

1 Genome wide assessment of genetic diversity and transcript variations in 17 accessions of the
2 model diatom *Phaeodactylum tricornutum*

3

4 Chaumier Timothée^{1¥}, Feng Yang^{1¥}, Eric Manirakiza¹, Ouardia Ait-Mohamed², Yue Wu¹,
5 Uditia Chandola¹, Bruno Jesus³, Gwenael Piganeau⁴, Agnès Groisillier¹, and Leila Tirichine^{1*}

6

7 ¹Nantes Université, CNRS, US2B, UMR 6286, F-44000 Nantes, France

8

9 ² Immunity and Cancer Department, Institut Curie, PSL Research University, INSERM U932,
10 75005 Paris, France

11

12 ³ Nantes Université, Institut des Substances et Organismes de la Mer, ISOMer, UR 2160, F-
13 44000 Nantes, France

14

15 ⁴Integrative Biology of Marine Organisms (BIOM), Sorbonne University, CNRS,
16 Oceanological Observatory of Banyuls, Banyuls-sur-Mer, France

17

18

19 [¥]Equal contribution

20

21

22 *Correspondence: tirichine-l@univ-nantes.fr; Tel.: +33-276645058

23

24 Abstract

25 Diatoms, a prominent group of phytoplankton, have a significant impact on both the oceanic
26 food chain and carbon sequestration, thereby playing a crucial role in regulating the climate.
27 These highly diverse organisms show a wide geographic distribution across various latitudes.
28 In addition to their ecological significance, diatoms represent a vital source of bioactive
29 compounds that are widely used in biotechnology applications. In the present study, we
30 investigated the genetic and transcriptomic diversity of 17 accessions of the model diatom
31 *Phaeodactylum tricornutum* including those sampled a century ago as well as more recently
32 collected accessions. The analysis of the data reveals a higher genetic diversity and the
33 emergence of novel clades, indicating an increasing diversity within the *P. tricornutum*
34 population structure, compared to the previous study and a persistent long-term balancing
35 selection of genes in old and newly sampled accessions. However, the study did not establish a
36 clear link between the year of sampling and genetic diversity, thereby, rejecting the hypothesis
37 of loss of heterozygosity in cultured strains. Transcript analysis identified novel transcript
38 including non-coding RNA and other categories of small RNA such as PiwiRNAs.
39 Additionally, transcripts analysis using differential expression as well as Weighted Gene
40 Correlation Network Analysis has provided evidence that the suppression or downregulation of
41 genes cannot be solely attributed to loss of function mutations. This implies that other
42 contributing factors, such as epigenetic modifications, may play a crucial role in regulating gene
43 expression. Our study provides novel genetic resources, which are now accessible through the
44 platform PhaeoEpiView (<https://PhaeoEpiView.univ-nantes.fr>), that offer both ease of use and
45 advanced tools to further investigate microalgae biology and ecology, consequently enriching
46 our current understanding of these organisms.

47

48

49 **Introduction**

50 Photosynthetic microalgae are important components of life in the oceans providing
51 organic biomass and fueling a range of key biogeochemical processes. Diatoms in particular
52 are widely recognized as one of the most significant phylum of phytoplankton, owing to their
53 substantial contribution to primary productivity, carbon fixation, and biogeochemical cycling
54 of essential nutrients such as nitrogen and silicon [1, 2]. In addition to their ecological
55 importance, diatoms are a rich source of bioactive compounds with diverse applications in
56 various industries, including nutraceuticals, nanotechnology, pharmaceuticals, and food and
57 feed industries [3, 4]. In recent years, using model species in diatoms has dramatically increased
58 our knowledge about the biology and ecology of these important organisms [5]. Particularly,
59 one species, the diatom *Phaeodactylum tricornutum* has proven to be a robust model for
60 research, yielding a wealth of knowledge and advancing our understanding in this field of
61 research.

62 *P. tricornutum*, a marine pennate diatom is commonly found in coastal waters, including
63 tidal areas, estuaries, rock pools and shallow waters exposing the species to important
64 fluctuations in light intensity and salinity. The diatom is a well-established model with a
65 genome that has been fully assembled and well annotated, along with an expanding molecular
66 toolbox using the reference strain Pt1 8.6 [6-10]. Genome wide sequencing of ten accessions of
67 *P. tricornutum* (Pt1 to Pt10) using Illumina, identified throughout the genome diverse
68 variations, including single nucleotide and insertion deletion polymorphisms (SNPs, INDEls)
69 and copy number variations [11]. This study provided important insights into the genetic
70 diversity of the isolates clustering them into four distinct clades with a conserved genetic and
71 functional makeup. Previous studies have revealed distinguishing features among different
72 accessions. Pt4 displayed a low non-photochemical quenching (NPQ), Pt5 demonstrated higher
73 adhesion, Pt6 exhibited substantial lipid accumulation, Pt8, Pt3 and Pt9 have different cell
74 morphologies and Pt3 demonstrated increased tolerance to variations in salinity, among other
75 traits [12-14].

76 The 10 sequenced accessions of *P. tricornutum* were mostly collected more than a century
77 ago and have been preserved as either lab cultures, or frozen stocks in culture collections for
78 extended periods of time. Therefore, their genetic composition may have been affected,
79 potentially favoring genes that are adapted to lab conditions [15-20]. Seven more recent isolates
80 were collected from the environment and sequenced, Pt11 (HongKong, China), Pt12, Pt13 (both

81 from Bourneuf Bay, West Atlantic, France), Pt14 (Gulf of Salerno, Italy), Pt15 (East China
82 sea), Pt16 (Helgoland, Atlantic Ocean, North Sea, Germany) and Pt17 (Banyuls Bay, Gulf of
83 Lion, France) (Figure 1, Table S1). Assessment of genetic diversity within natural accessions
84 of a model diatom, such as *P. tricornutum* is critical to our understanding of fundamental
85 questions relevant to diatom's biology and ecology. A high genetic diversity within the *P.*
86 *tricornutum* species presents substantial implications, including valuable insights into their
87 adaptive strategies across diverse ecological niches, alongside the identification of pivotal
88 genetic determinants governing responses to environmental factors.

89 DNA sequence polymorphism can lead to phenotypic variations but it remains only the
90 first step in understanding how these polymorphisms can affect the phenotype. Variations in
91 transcript levels are another proxy to better understand the contribution of genes to phenotypic
92 variations. Diatoms have developed sophisticated sensory and gene regulatory mechanisms to
93 detect and respond to environmental cues. They employ transcription factors and regulatory
94 elements to control the initiation and rate of transcription. Furthermore, post-transcriptional
95 mechanisms, such as RNA splicing, RNA transport, stability, and translation, play essential
96 roles in determining the abundance and activity of specific gene products. These integrated
97 processes collectively enable diatoms to fine-tune gene expression in response to changing
98 environmental conditions [7]. An illustrative example is the diel and circadian rhythms in gene
99 expression, which are synchronized with light and dark cycles. These rhythmic gene regulation
100 mechanisms enable diatoms to optimize their metabolic processes and growth in a time-
101 dependent manner, playing a significant role in their overall physiology and ecology [21].

102 Differences in gene expression can be attributed to different DNA sequence
103 polymorphisms including SNPs and INDELS that can nullify gene function or induce variations
104 in splicing. It is important to ask whether the genes that show differences in expression are
105 under selective pressure and whether there is a link between transcript level variations and DNA
106 sequence polymorphisms. DNA sequence may not explain differences in gene expression, cases
107 rather known to be the consequences of epigenetic factors including DNA methylation, post-
108 translational modifications of histones and small and long non-coding RNAs [22, 23].

109 In the present study, we analysed the genetic diversity of 17 accessions of the model
110 diatom *P. tricornutum* by examining both DNA sequences and transcript levels. These
111 accessions were collected from various coastal regions of world seas and oceans including
112 recently sampled accessions that were not included in the previous study that had sequenced
113 accessions collected over a span of 100 years. Our findings indicate a higher genetic diversity

114 that defined more distinct clades and long-term balancing selection of genes in old and newly
115 sampled accessions. However, our analysis failed to establish a clear link between the temporal
116 factor of sample collection year and the extent of genetic diversity. Consequently, the
117 hypothesis proposing a decline in genetic diversity, specifically the loss of heterozygosity in
118 cultured strains over time, could not be supported [24]. Furthermore, our study identified novel
119 transcripts among which various non-coding RNA species and provides insights into the
120 regulation of genes mediated by genetic and transcript diversity. Our study offers easy and
121 valuable access to these novel genetic resources, particularly focusing on a model species,
122 through the PhaeoEpiView platform (<https://PhaeoEpiView.univ-nantes.fr>) [10]. This
123 unprecedented accessibility provides a multitude of opportunities for exploring diverse
124 ecological functions by leveraging the genetic diversity of this model organism, thereby
125 expanding our understanding of the biology and ecology of microalgae.

126

127 **Material and methods**

128 **Material used and growth conditions**

129 Eighteen different accessions of *P. tricornutum* were acquired from the Provasoli-Guillard
130 National Center for Culture of Marine Phytoplankton, Roscoff and Nantes culture collections
131 (Table S1). All of the accessions were grown axenically using Enhanced Artificial Sea Water
132 (EASW) [25] in batch cultures at 19°C, under 12/12 light dark period with a light intensity of
133 70 $\mu\text{mol photons m}^{-2} \text{ s}^{-1}$.

134 **Growth curves**

135 The cultures were grown in 30 ml of EASW with an initial concentration of 10^5 cells / ml. Cell
136 counts were measured using flow cytometry (CytoFLEX, Becman, USA), every 2 days for 20
137 days, with 1 ml of sample taken from each accession culture each time. After 7 days of culture,
138 1 ml from each of Pt11 to Pt17 were used for light microscopy to describe variant shapes and
139 cell size using Axio inverted microscope (ZEISS, Germany). Photos were further analyzed
140 using Zeiss software (ZEN 2.6).

141 **Pulse amplitude modulated variable *chlorophyll a* measurements**

142 Variable fluorescence measurements were carried out using an Imaging PAM fluorometer
143 (Walz) with a blue measuring light (450 nm), controlled by the software ImagingWin v2.46i
144 (Heinz Walz GmbH, Effeltrich, Germany). The actinic and saturating light were also blue and

145 provided by fluorometer LED panel. The saturation pulse intensity was $6000 \mu\text{mol photons m}^{-2}$
146 s^{-1} for 0.8 s. Samples were dark-adapted for one hour before carrying out any measurements.
147 For the construction of RLC (Rapid Light-response Curves) [26], the samples were exposed to
148 nine incremental intensities of actinic light with an irradiance step duration of 30 s. The PAR
149 (photosynthetically active radiation; 400 – 700 nm) steps used were: 0, 5, 19, 31, 37, 42, 47,
150 56, 75, 143, 280 and $519 \mu\text{mol photons m}^{-2} \text{s}^{-1}$. The first point of the RLC corresponds to the
151 dark-adapted state, yielding the minimum fluorescence yield (F_o) and the maximum
152 fluorescence yield (F_m), allowing the calculation of the maximum PSII quantum efficiency
153 (F_v/F_m) as $F_v/F_m = (F_m - F_o)/F_m$. The remaining light steps measured the fluorescence yield (F'),
154 the maximum fluorescence yield (F_m') in the light-exposed state and the effective PSII quantum
155 yield at each experimental light level (E) as $F_q'/F_m' = (F_m' - F')/F_m'$. Relative PSII electron
156 transport rates were calculated as $\text{rETR}(E) = F_q'/F_m'(E) \times E$ and non-photochemical quenching
157 as $\text{NPQ} = (F_m - F_m')/F_m'$. Maximum relative electron transport rates (rETR_{max}) were estimated
158 by fitting the RLCs with the photosynthesis-light response model of [27] and maximum NPQ
159 (NPQ_{max}) by fitting the NPQ-light response model of [28].

160 **DNA extraction and PCR protocol**

161 After 7 days of culture, cells were centrifuged at $4000 \times g$ for 20 min and washed twice with
162 1XPBS. DNA was isolated using a CTAB protocol as described previously [29, 30]. A volume
163 of 1.5 mL of CTAB buffer (450 μL 10% CTAB, 420 μL of 5 M NaCl, 60 μL of 0.5 M EDTA,
164 150 μL of 1 M Tris HCL) preheated to 65°C , was placed into a 2 ml plastic tube together with
165 diatom pellet and incubated for 1 hour at 65°C , then DNA was isolated using Chloroform
166 Isoamyl (24/1) after centrifugation for 10 minutes ($12000 \times g$). The upper phase was removed
167 and incubated with 3.2 μg of RNase A for 1 h at 37°C . DNA was isolated again using
168 Chloroform Isoamyl (24/1) after centrifugation ($12000 \times g$) to remove protein and RNA. Same
169 volume of isopropanol and 8% volume of ice-cold ammonium acetate 7.5M, were used for
170 precipitation at -20°C overnight. Nucleic acids were recovered after centrifugation ($12000 \times g$)
171 at 4°C and purified by absolute ethanol, then washed with 70% ethanol. DNA concentration
172 was measured using NanoDrop ND-1000 spectrophotometer (Thermo Scientific, Wilmington,
173 DE, USA). PCR amplification was carried out on Mastercycler[®]nexus $\times 2$ (eppendorf, Germany)
174 in 20 μL total volume including 9.6 μL Go Taq, 50 ng DNA and 10 μM forward and reverse
175 primers. The PCR program consisted of 95°C for 5 min, then 35 cycles of 95°C denaturation
176 for 30s, annealing at appropriate annealing temperature (56°C - 62°C) for 30 s, 72°C extension
177 for 30 s, and a final extension step at 72°C for 10 min. PCR products were electrophoresed on

178 1% agarose gel, and the gel images were acquired using EBOX CX5 System (VILBER BIO
179 IMAGING, France). Primer sequences are listed in Table S2.

180 **RNA extraction**

181 A total of 300 ml exponentially growing cells were centrifuged at 4000 x g, 4°C for 20 min and
182 immediately re-suspended in 500 µL TRIzol® Reagent (Invitrogen, Thermo Fisher Scientific,
183 USA) and vortexed vigorously before being stored at - 80 °C. RNA was isolated using Trizol
184 reagent as described previously [31]. Purity and quantity of RNA were assessed using
185 NanoDrop ND-1000 spectrophotometer (Thermo Scientific, USA). To remove genomic DNA,
186 RNA samples were treated with Ambion™ DNase I (Invitrogen, Thermo Fisher Scientific,
187 USA), according to manufacturer's instructions. RNA was quantified using Qubit™ RNA BR
188 Assay Kit, 500 assays (Invitrogen, Thermo Fisher Scientific, USA).

189 **DNA and RNA sequencing**

190 Extracted DNA for Pt1 8.6 and Pt11 to Pt17 was sequenced on Illumina Novaseq 6000 platform,
191 using 250 bp paired-end reads. Yields for Pt1 8.6, Pt11, Pt12, Pt13, Pt14, Pt15, Pt16 and Pt17
192 were 5.9, 13.2, 6.3, 5.7, 8.0, 6.6, 5.9 and 6.6 million read pairs, respectively. Messenger RNAs
193 for Pt1 8.6 and Pt1 to Pt17 were sequenced in duplicates on Illumina Novaseq 6000 platform,
194 using 150bp paired-end reads. The RNA libraries were enriched for matured RNAs and
195 sequenced in stranded mode, yielding between 15.1 and 25.7 million read pairs.

196 **Bioinformatics analysis**

197 **Variant calling analysis**

198 Paired-end Illumina libraries from each ecotype (Pt1 to Pt17) were first trimmed using
199 Trimmomatic [32] with “ILLUMINACLIP:adapters.fa:2:30:10:2:keepBothReads
200 LEADING:20 TRAILING:20 SLIDINGWINDOW:4:20 MINLEN:40”. Reads were then
201 mapped with BWA-mem2 2.2.8 [33] on *P. tricornutum* Phatr2 assembly (accession
202 GCA_000150955.2). Mapping rates ranged from 94.70% for Pt8 to 99.36% for Pt13. Variant
203 calling and filtering were performed with GATK package version 4.2.2.0, following GATK
204 best practices [34]. In short, HaplotypeCaller module was run with a call confidence of 30, a
205 sample ploidy of 2 and double precision was activated for pair-HMM algorithm. Variants that
206 were called were functionally annotated by snpEff [35] with *P. tricornutum* database v5.0 and
207 transposable elements annotation from [7], an upstream/downstream region size set at 2kb and
208 gene putative loss of function (LOF) annotation activated. Only SNPs and insertions/deletions

209 (INDELs) were retained from annotated VCF files. Finally, GATK's VariantFiltration module
210 was used on SNPs with the following filters: "QD<2.0; QUAL<30; SOR>3.0; FS>60.0;
211 MQ<40" and on INDELs with: "QD<2.0; QUAL<30; FS>200.0". INDELs of size above 50bp
212 were extracted with SelectVariant module and some of the longest were validated by PCR.

213 **Fixation index computation**

214 Fixation index (Fst) was computed with ANGSD 0.939 [36] between all 17 accessions possible
215 pairs. First, in order to compare only regions where data were present for all samples, the
216 callable genome size was defined where read coverage on reference was no less than 10X across
217 all accessions. Allele frequencies were computed for all ecotypes using ANGSD "-doSaf 1 -GL
218 2 -minMapQ 1 -minQ 20", then folded site frequency spectrum (2DSFS) was determined with
219 "realSFS -maxIter 100 -fold 1" for all ecotypes combinations. Finally, we computed all pairwise
220 Fst on the resulting indexes with "realSFS fst index -fold 1". All Fst values were gathered in a
221 matrix and displayed as a heatmap using the R [37] package Pheatmap 1.0.12
222 (https://scicrunch.org/resolver/RRID:SCR_016418).

223 **Population clustering**

224 Callable SNPs and INDELs were analyzed with ADMIXTURE 1.3.0 [38] with a random seed
225 of 12345679, cross-validation (CV) activated and a bootstrap value of 200. Numbers of
226 ancestral populations (K-value) were tested from 1 to 17 and cross-validation error was plotted
227 in order to select K leading to the lowest CV error. A PCA was then performed on Q-estimates
228 for K=15 and estimated ancestral fractions were plotted with R.

229 **CNV and gene loss analysis**

230 For each ecotype (Pt1 to Pt17), raw number of mapped fragments from the Variant Calling
231 Analysis BAM files were counted on each Phatr3 gene [7] using featureCounts [39] in
232 unstranded paired-end mode and reads were assigned to all their overlapping features ("-O"
233 option). Genes with no counts were deemed as possibly lost. Raw counts were then normalized
234 for each gene following FPKM formula: $FPKM_normalized_count = (gene_raw_count \times 10^9) / (gene_length \times total_sample_counts)$. Similar to previous work [11], binary logarithm Fold
235 Change (log2FC) was calculated as the log2 ratio of normalized count for each gene to the
236 average (mean) normalized count of all the genes per accession. Genes with a log2FC ≥ 2
237 were considered as showing putative Copy Number Variation (CNV) compared to the reference
238 strain. Finally, lost genes and genes exhibiting CNV in only one accession were marked as
239

240 ecotype-specific. Heatmap plots were made in R using Pheatmap 1.0.12
241 (https://scicrunch.org/resolver/RRID:SCR_016418) and UpsetR 1.4.0 [40] packages.

242 **Genes with loss of function**

243 After variant annotation by snpEff, we used an in-house script to find the total and specific
244 number of genes affected by loss-of-function (LoF) mutations for each ecotype (Pt1-Pt17). First,
245 we selected genes with LoF variants alleles retained in the VCF annotation file whether they
246 are homozygous or not. Then, we searched for the genes that are specific to each accession and
247 considered the non-accession specific genes common if shared by two or more accessions.

248 **Site Frequency Spectrum (SFS) analysis**

249 A matrix of allele counts per ecotype was created from the variants called previously on the
250 callable genome. Briefly, for each biallelic SNP, a value of 0, 1 or 2 was determined for each
251 ecotype, depending on its ploidy (homozygous on reference allele, heterozygous
252 reference/alternate alleles or homozygous on alternate allele, respectively). Moreover,
253 functional annotation as described previously and the affected gene (if applicable) were added
254 for each SNP. Folded SFS was then calculated for each functional category of SNPs (nonsense,
255 non-synonymous, synonymous, intergenic) and the resulting data was plotted with R.

256

257 **Searching for signatures of Balancing Selection (BS) on non-synonymous SNPs**

258

259 One of the signature of balancing selection is the excess of nonsynonymous polymorphisms
260 segregating at intermediate frequencies [41]. Genes with less than 10 synonymous (S) + non-
261 synonymous (NS) SNPs were filtered out from the allele counts matrix (see “*Site Frequency*
262 *Spectrum (SFS) analysis*”). The ratio non-synonymous versus synonymous diversity was
263 estimated by Watterson’s theta θ assuming twice as many non-synonymous than synonymous
264 sites ($\theta_{wNS} / \theta_{wS}$), defined as “(Number of NS/2) / (Number of S)” was calculated for each of
265 the remaining 9267 genes. The 91 genes with an excess of non-synonymous SNPs (showing a
266 θ ratio over 3), were extracted and further investigated. Finally, the same process was performed
267 ecotype-wise, with a number of genes with a θ ratio > 3 ranging from 34 in Pt14 to 62 in Pt4.

268

269 **Phylogeny of the 17 accessions**

270 A matrix of genome-wide biallelic SNPs and INDELS allele counts per ecotype (0, 1 or 2
271 depending on the called ploidy of the variant, see “*Site Frequency Spectrum (SFS) analysis*”),
272 was computed for 640,454 variants found in the population Pt1 to Pt17. Canberra distance and
273 average linkage functions were identified to produce the tree that represented best the matrix

274 by the “find_dend()” method from R library “dendextend” 1.16.0 [42]. Then, an unrooted
275 neighbor-joining tree was built using “phangorn” R package 2.10.0 [43] on the Canberra
276 distance matrix and colored according to the ecotypes clades.

277

278 **Expression analysis**

279 After filtering raw data with the removal of adapters and low quality reads, clean reads were
280 aligned against the reference genome using HISAT2 2.0.5 [44]. Reads were assigned to each
281 transcript using the FPKM metric which normalizes for differences in library size and gene
282 length. In order to compare gene expression levels in different accessions, the graphical
283 representation of the distribution of gene expression and FPKM levels in different samples has
284 been performed using the ggplot2 R package (v3.4.0)[45]. In order to differentiate between
285 coding and non-coding transcripts, the Coding Potential Assessment Tool (CPAT) (DOI:
286 10.1093/nar/gkt006) which is a convenient and rapid method to categorize transcripts based on
287 their coding scores, was used. This algorithm relies on specific models to assign coding
288 potential scores to individual transcripts. In our research, we employed models from human,
289 mouse, and zebrafish. Consequently, a table was generated, presenting the coding potential
290 outcomes for each input transcript.

291

292

293 **Principal Component Analysis**

294 To elucidate the relationships between distinct accessions, we conducted Principal Component
295 Analysis (PCA) on the gene expression values (FPKM) of all the samples. Specifically, we first
296 computed the average FPKM value for two replicates of each sample, followed by a logarithmic
297 transformation of this average value ($\log_2(\text{FPKM}+1)$). Ultimately, we represented the samples
298 according to their expression levels. In our analysis, we evaluated how well the samples were
299 represented in PCA using the cos2 (square cosine) metric as a measure of quality. A cos2 value
300 closer to one indicates a stronger representation of the variable by the two displayed
301 components.

302 **Repeats detection in novel transcripts**

303 Reads from both replicates of Pt1 to Pt17 RNAseq libraries were mapped with HISAT2 2.0.5
304 [44] using the default parameters and all the mapping information were combined. The reads
305 were then assembled with StringTie 1.3.3b [46] and novel transcripts were kept. We screened

306 these 656 novel transcripts for repeats with RepeatMasker (*RepeatMasker Open-4.0*. 2013-
307 2015 <http://www.repeatmasker.org>) Galaxy Tool wrapper version 4.0.9 (slow settings with
308 matrices for 43% GC content), using a manually curated database of 71 reference transposable
309 elements (TE) of *P. tricornutum*. RepeatMasker was run on the public server at
310 <https://usegalaxy.org> [47].

311 **WGCNA network analysis**

312 The network was obtained using Weighted Gene Correlation Network Analysis (WGCNA) R
313 package (version 1.17) [48]. Before constructing the co-expression network, we filtered out
314 genes having a row median less than 10 reads. The expression matrix was transformed with the
315 vst (Variance Stabilizing Transformation) function from DESeq2 R package (version 1.32.0)
316 [49]. The sequencing steps for the network construction for *P. tricornutum* accessions have
317 been described previously [50].

318 For Network construction, the WGCNA R package [48] was used to identify network modules
319 from 36 RNA-Seq datasets representing expression data from 18 *P. tricornutum* accessions
320 (two replicates per accession). First, the quality of the raw counting matrix was checked. A
321 hierarchical clustering analysis based on the "average" method allowed us to identify the Pt1R2
322 as an outlier, so this sample was filtered out from further analysis.

323 **Results**

324 **Phenotypic traits characterization**

325 To assess phenotypic differences among the 17 accessions, we monitored their growth,
326 cell morphology and photosynthetic features. Significant differences in growth rate were
327 recorded at day 4 of the exponential phase (Fig.2a). In most of the cultures, cell growth rate
328 dropped after day 11 and entered the stationary phase. Pt8, Pt3 and Pt10 showed faster growth
329 rate and higher final concentrations than other accessions, with 1165×10^4 , 1032×10^4 and
330 960×10^4 cells mL^{-1} respectively in the end of exponential phase ($p < 0.05$). On the other hand,
331 Pt4 and Pt9 showed slower growth rate than the others and the lowest final concentrations, with
332 419×10^4 and 559×10^4 cells mL^{-1} respectively.

333 Cell dimensions were measured for only 7 accessions (Pt11 to Pt17) and compared to the
334 previously published Pt1 to Pt10 accessions [12]. Among all accessions, the previously
335 measured Pt5 had the longest length (25-30 μm) and Pt14 cells were the shortest (10-15 μm).

336 Pt12 were the thinnest cells with the lowest length/width ratio (9.9 ± 1.1) and Pt17 were the
337 largest (4.6 ± 0.8) (Table S1, Fig. S1). Different morphotypes were observed in Pt16 (a mix of
338 75% fusiform and 25% of triradiate). As reported previously [11], we found few oval cells
339 mixed with fusiform in Pt3 and Pt9 (7.3% and 7.6% respectively). Triradiate were reported in
340 Pt8 [12] but we did not observe triradiate cells in our conditions. Triradiate morphotype was
341 reported to be instable in this accession [51].

342 **Measurements of photosynthetic parameters**

343 To assess photosynthetic capacities of *P. tricornutum* accessions, we measured maximal
344 PSII quantum yield (F_v/F_m) and maximal relative electron transport rates (rETRmax). Among
345 all the accessions, Pt12 showed the highest F_v/F_m followed by Pt13 ($p = 0.027$) and Pt11 ($p =$
346 0.0047) (Fig. 2b). These three accessions also showed the best photosynthetic performances
347 based on rETRmax (Fig. 2c), while Pt4, Pt5, and Pt6 demonstrated the lowest values in PSII
348 quantum efficiency (Fig. 2c). The English Channel accessions, Pt1, Pt2 and Pt3 showed similar
349 results.

350 To evaluate the response to excess light energy, we measured non-photochemical
351 quenching (NPQ). Pt9 and Pt17 showed the highest NPQ capacity, c. 8.2 and c. 9.2 respectively
352 while Pt4 showed a lower NPQ capacity of around 2.5 (Fig. 2d). These observations suggest
353 that Pt9 and Pt17 can tolerate environments with higher light intensity compared to other
354 accessions. This is consistent with their geographical distribution in latitudes that receive
355 greater amounts of solar radiation. Similarly, Pt4 NPQ reflects an adaptation to its sampling
356 location in higher latitudes, specifically the Norwegian Fjords.

357 **Variant calling analysis**

358 We performed variant calling analysis on previously published sequences of Pt1 to Pt10
359 [11] and the newly sequenced Pt11 to Pt17 accessions using the reference strain Pt1 8.6 genome
360 sequences. All the accessions had a good sequence coverage allowing a confident variant
361 calling. We identified 731,357 single nucleotide polymorphisms (SNPs), 44,470 insertions
362 (from 1 to 422 bp length) and 52,867 deletions (varying from 1 to 274 bp) (Fig. 3 a,b). Site
363 frequency spectrum (SFS) which reflects the numbers of variants segregating at different
364 frequencies showed the expected excess of low frequency alleles (Fig. 3c). The highest increase
365 in low frequency SNPs was observed in non-sense polymorphisms. Non-synonymous
366 polymorphisms showed the second highest increase, when compared to intergenic and

367 synonymous polymorphisms. Most of the SNPs (59-63%) were found in genes while INDELS
368 were mostly detected in intergenic regions (Fig.3 d). Our analysis identified 22,4253 additional
369 SNPs and 74,918 additional INDELS in Pt1 to Pt10 compared to our previous study [11].

370 Despite the higher number of discovered SNPs and INDELS, the overall trend of their
371 distribution among the 10 previously analyzed accessions remains the same. Across all the
372 accessions most of the SNPs were heterozygous (HetSNPs) with >95% in Pt1, Pt2, Pt3, Pt9,
373 Pt15 and Pt17 and 65% to 68% in Pt6, Pt7, Pt8 and Pt16 and the lowest proportion of HetSNPs
374 were found in Pt4, Pt5 and Pt10 to Pt14 (<49%). Across all the accessions, the numbers of
375 INDELS was similar except for Pt4, Pt6, Pt7 and Pt8, which showed the highest number of
376 INDELS (Fig. 3b, Table S3). Most INDELS were shared among the accessions. A total of 14
377 INDELS were validated by PCR for randomly chosen loci (Fig. S2). To further assess genetic
378 diversity, we investigated copy number variations (CNVs) and gene loss (GL). A total of 284
379 and 180 genes show CNV or GL respectively. Most of CNVs are shared and eleven accessions
380 out of 17 show specific CNVs with 40 genes in Pt10 followed by Pt4, Pt6, Pt14 and Pt16 with
381 15, 15, 8 and 7 genes respectively (Fig. 3e, Table S4, Table S5). Randomly chosen loci were
382 validated by PCR for gene loss (Fig. S3).

383 To understand the functional impact of genetic variations among *P. tricornutum*
384 accessions, we examined the loss of function mutations (LoF) such as premature stop codons,
385 frameshifts and start loss. A total of 31,536 LoF was found, among which 588 were shared
386 across the accessions (Fig. 3f). Accession specific LoFs were mostly found in Pt4 (61), Pt14
387 (20) and Pt16 (17) (Table S6). LoF mutations were enriched in gene ontology (GO) categories
388 of molecular function type (Table S6) and in genes that belong to large gene families as
389 previously shown [11].

390 **Population structure and phylogeny of *P. tricornutum* accessions**

391 To examine the global population structure of *P. tricornutum* accessions, we used
392 pairwise fixation index (Fst), a measure of genetic differentiation revealing groups with low Fst
393 index (< 5%) (Figure 4a). To further estimate genetic relatedness among *P. tricornutum*
394 accessions, we used admixture proportion inference which allows the assignation of individual
395 genetic variations into clusters based on shared allele frequency patterns [38]. We ran
396 ADMIXTURE with various plausible values of *K* which is the number of source populations
397 and found a stable admixture pattern proportions with *K*=15, reflecting the number of ancestral

398 populations (Fig. S4, Table S7). Based on individual ancestry with similarity of clusters
399 between accessions, we distinguished 6 clades: Pt1, 2, 3, 9, 15 and 17 in clade 1 with most of
400 the clusters (up to 11) reflecting a clade with multiple ancestral populations, Pt16 in clade 2,
401 Pt4 in clade 3, Pt5, 10 and 11 in clade 4 with only 3 clusters, Pt6, 7 and 8 in clade 5 and Pt12,
402 Pt13 and Pt14 in clade 6 (Figure 4b).

403 To confirm admixture analysis clades, we performed a Principle Component Analysis
404 (PCA) which revealed similar results reflecting common ancestry except for Pt14 which is far
405 from its admixture defined clade 6 (Fig. 4c). Of note, cluster composition proportions of Pt14
406 are different from Pt12 and Pt13. To further assess the segmentation among *P. tricornutum*
407 accessions, we used CNVs to run a hierarchical clustering and found that the 17 accessions fall
408 into 6 clusters supporting further both PCA and admixture analyses: Pt1, Pt2, Pt3, Pt9, Pt15 and
409 Pt17 in cluster 1, Pt16 in cluster 2, Pt4 in cluster 3, Pt6, Pt7 and Pt8 in cluster 4, Pt5, Pt10 and
410 Pt11 in cluster 5 and Pt12, Pt13 and Pt14 in cluster 6 (Fig. 4d). Phylogenetic analysis at whole
411 genome scale using genetic variations (SNPs and INDELS) of the 17 accessions supported
412 further the clustering into six clades observed with Fst and PCA analyses (Fig. 4e).

413 **Balancing and relaxed selection in Pt clades**

414 We calculated the ratio of nonsynonymous to synonymous nucleotide site diversity using
415 Watterson's estimate of theta (θ_{wN}/θ_{wS}) [52] as a measure of the efficiency of natural
416 selection. Ratios > 3 suggest that selective pressure maintains non-synonymous
417 polymorphisms, a signature for balancing selection (BS) which refers to selective processes by
418 which alleles are maintained in a population at frequencies larger than expected from genetic
419 drift alone [53], while θ_{wN}/θ_{wS} ratios < 1 refer to non-synonymous polymorphisms being
420 counter-selected pointing to a purifying selection. We identified 91 common genes under BS
421 and 2422 under purifying selection (Table S8). Most of the genes under BS are of unknown
422 function. However, those with known function were enriched in genes coding for functions
423 such as cell proliferation and growth, perception, transmembrane transport activity, stress
424 responses and adaptation to the environment.

425 **Identification of transcript level variations and co-expression network modules in *P.*** 426 ***tricornutum* accessions**

427 Differences in gene expression is known to control inter and intra specific phenotypic
428 variations providing living organisms with abilities to colonize different ecological niches. To

429 identify differences in transcriptomes, RNAs from *P. tricornutum* accessions (Pt1 to Pt17) were
430 sequenced and mapped to the reference genome. The mapping of mRNA-Seq reads was above
431 85% for all replicates except for one, Pt16R2, which had a mapping rate of 41.27%. Pearson
432 correlation coefficients between each of the two replicates was mostly around 0.98 (Fig. S5).
433 Interestingly, a total of 656 genes including 25 from the chloroplast were novel. These genes
434 are widely distributed over the genome among which some were specific to each of the ecotypes
435 and 385 genes were found to be common to all of them (Table S9). They showed an average
436 length of 709.42 bp, with 334 genes categorized as sense and 322 genes as antisense. Except
437 from few genes that were annotated, most of these novel genes were of unknown function
438 (Table S9). The majority of novel genes were downregulated compared to the average gene
439 expression but their expression remains significant and cannot be considered as part of a
440 spurious phenomenon of background low-level transcription.

441 The analysis of novel transcripts using RepeatMasker revealed that 12.26% of them
442 contained repetitive elements, primarily Copia LTR_retrotransposon of class I and rare MuDR2
443 Terminal Inverted Repeats of class II (Table S10). Additionally, approximately 0.48% of the
444 transcripts contained simple repeats. Using a coding potential assessment tool, we identified
445 several non-coding RNA with sizes varying between 202 bp and 8379 bp (Table S11).
446 Furthermore, we detected several other RNA types, including sRNA, miRNA, tRNA, snoRNA,
447 antisense and piwiRNAs (Table S11).

448 Then, we examined differentially expressed genes (DEGs) under our standard growth
449 conditions by analyzing RNA-Seq data across the 17 accessions and comparing them with the
450 reference strain Pt18.6. This strain was derived from Pt1, which displayed the lowest number
451 of DEG (1308 genes) as expected (Figure 5 a). In contrast, Pt7 showed the highest number of
452 DEG (5086 genes). The remaining ecotypes showed variable numbers of DEG, with Pt5
453 exhibiting the lowest number at 1890 genes. On average, most ecotypes had approximately
454 4000 DEG. With the exception of Pt1, which exhibited a bias towards upregulation (twice as
455 many upregulated genes as downregulated genes), the other ecotypes displayed a balanced ratio
456 of upregulated and downregulated genes. The majority of upregulated genes, their
457 $\log_2(\text{FoldChange})$ values varies between 1 and 3, while the majority of downregulated genes, the
458 value of $\log_2(\text{FoldChange})$ varies between -3 and -1. A substantial fraction of DEGs, regardless
459 of their upregulation or downregulation status are found to be specific to one or multiple
460 ecotypes suggesting an ecotype-dependent gene expression pattern that likely underlies ecotype

461 specific trait (Figure 5b, c, Table S12, S13). On the other hand, only a minor subset of DEGs
462 displaying upregulation or downregulation were shared across all the ecotypes. We conducted
463 a principal component analysis using the average replicates expression to evaluate whether the
464 ecotypes exhibited comparable expression profiles. Our analysis distinguished five clusters that
465 partially aligned with the clades defined in this study, implying a correlation to some extent
466 between genetic diversity and expression patterns (Figure 5d).

467 To explore the relationship between phenotypic traits, specifically photosynthesis and
468 transcripts, we closely examined the NPQ response and the genes related to its regulation in
469 response to light. Indeed, the NPQ capacity depends on transthylakoidal proton gradient, but
470 also on antenna proteins called Light-Harvesting Complex Protein X (LHCX) and on the
471 Diatoxanthin Xanthophyll cycle [54], equivalent to the Zeaxanthin cycle in land plants [55].
472 Among the LHCX genes, only LCHX1 shows strong expression in all ecotypes, confirming its
473 constitutive role in low light, while LHCX2 and 3 are involved in high light and LHCX4
474 expression increases in the dark [56] (Fig. S6). Similarly, all genes involved in the Diatoxanthin
475 Xanthophyll cycle show negligible expression in our low light culture conditions.

476 To understand the underlying nature of the conserved transcriptomic responses, we
477 analyzed the enrichment of GO terms for both upregulated and downregulated DEGs (Figure
478 S7). Additionally, we performed GO enrichment analysis on genes that were specifically
479 upregulated or downregulated in a single accession as well as per clade, where applicable. Only
480 few GO terms emerged from the analysis of accession specific DEGs, namely photosynthesis
481 GO related terms (light harvesting, protein chromophore linkage) in Pt8, lipid metabolic
482 processes and translation in Pt3 for downregulated genes, while upregulated genes were
483 enriched in ribosome biogenesis and rRNA processing in Pt7, protein transport in Pt16 and
484 glucose metabolism in Pt11 (Figure S7, Table S12, S13).

485 The WGCNA package was used to construct gene co-expression network of transcripts
486 from an expression matrix of ~432,000 transcripts derived from 36 RNA-seq samples, with 2
487 replicates collected from the 18 accessions including the reference strain Pt1 8.6. This approach
488 yielded in 33 distinct co-expressed modules (labeled by different colors) with dark slate blue
489 and plum2 modules containing each the smallest number of genes (106) and green yellow with
490 the largest number of genes (1599) (Figure 6c, Table S14). These modules were constituted by
491 genes demonstrating analogous expression profiles, which may or may not be consistent among
492 different clades suggesting the existence of additional factors besides genetic polymorphisms

493 that could modulate transcript levels (Figure S8). The modules were further categorized into
494 six distinct clusters, each characterized by a group of genes exhibiting comparable expression
495 patterns thus implying their involvement in shared pathways (Fig. 6c, Table S15). Based on the
496 GO annotation analysis, we identified significant functional enrichments in different groups.
497 Group I displayed a substantial increase in oxidoreductase activity, while Group II showed an
498 enrichment in calcium ion binding activity. Group III was characterized by an enrichment in
499 chlorophyll binding and light harvesting, whereas Group IV was associated with RNA
500 processing and translation. Group V showed a significant enrichment in photosynthesis and cell
501 redox homeostasis, while Group VI exhibited an enrichment in cell division and DNA binding
502 activity (Table S15).

503 **Discussion**

504
505 The present study was designed to comprehensively characterize the phenotypic, genetic,
506 and transcriptomic diversity among seventeen distinct *P. tricornutum* accessions, which were
507 collected from various locations across the world's oceans and included more recently sampled
508 accessions compared to those examined in prior studies [11, 12]. Growth dynamics, cell
509 morphology, and photosynthetic traits were monitored, and significant inter-accession
510 differences were recorded. Notably, the Pt4 and Pt9 strains exhibited a distinct growth pattern,
511 implying a probable correlation to their respective sampling locations. More specifically, Pt4,
512 sampled from the Norwegian fjords, seems to be adapted to a lower light intensity regime than
513 that employed in our study, while Pt9, a tropical strain, showed a slower growth at 19°C
514 compared to the presumed higher temperature in its geographical location. Moreover, the
515 evaluation of photosynthetic abilities across various accessions corroborated the association
516 with the sampling sites. Pt4 displayed the smallest $rETR_{max}$ and PSII maximum quantum
517 efficiency, while Pt12, Pt13, and Pt14 collected from the Mediterranean Sea and Atlantic side
518 demonstrated higher photosynthetic performance, as evidenced by their important F_v/F_m and
519 $rETR_{max}$ values. Strains within the same clade show similar growth and photosynthetic
520 performances, but there is no apparent correlation pattern with the year of sampling. Each clade
521 includes accessions from both older and more recently collected samples. Measuring cell
522 dimensions in the recently acquired accessions and their comparison with previously
523 characterized ones revealed notable variations. Specifically, Pt12, Pt14, and Pt17 exhibited
524 significant deviations, with Pt12 having the shortest cell length, Pt14 displaying the smallest
525 length-to-width ratio, and Pt17 showing the largest cells when compared to the remaining
526 accessions. Additionally, our investigation identified Pt16 as a new accession with a

527 combination of triradiate and fusiform morphotypes. These disparities in cell sizes may
528 potentially confer an advantageous trait in terms of enhanced gliding capabilities and improved
529 photosynthetic efficiency [57]. The observed differences in cell dimensions and, at times,
530 morphology are not surprising, given that *P. tricornutum* does not rely on silica for growth.
531 This lack of dependence on silica may confer flexibility in morphogenesis, a trait not typically
532 observed in silicified diatom species. The variations in cell sizes signify an adaptation to local
533 environments, highlighting an environmental-induced control of morphogenesis rather than a
534 genetic one, as demonstrated previously [51]. Drill-core records from Lake Titicaca in Peru
535 revealed a strong correlation between size trends in the diatom *Cyclostephanos andinus* and
536 environmental variables. This suggests that diatom size responds to regional environmental
537 changes driven by global processes that affect lake level and thermal stratification [58]. This,
538 in turn, implies that environmentally mediated epigenetic changes modulate phenotypic traits
539 within the same species

540 Variants calling analysis showed a larger number of SNPs and INDELS than previously
541 reported in Pt1 to Pt10. This is due to the gapped alignment mode used for SNPs and INDELS
542 calling which performs better than ungapped mapping [59]. Improvements made to
543 HaplotypeCaller's algorithm since 2018 were also likely playing a role in the gain of sensitivity
544 we noticed. Most of the SNPs were located in coding regions, while INDELS were mostly
545 found in intergenic regions as a consequence of their highly deleterious effects within coding
546 regions. Most of the SNPs were found to be heterozygous, indicating that *P. tricornutum* has a
547 high level of heterozygosity. Our previous study has demonstrated a substantial level of
548 heterozygosity in *P. tricornutum* populations. The recent sampling of genetically distinct
549 accessions has reaffirmed the persistence of this trait, despite not being related to the previously
550 identified highly heterozygous accessions Pt1, Pt2, and Pt3, thereby supporting the reliability
551 of the heterozygosity measure. This high level of heterozygosity is intriguing in *P. tricornutum*
552 considering that the species is not known to reproduce sexually suggesting an advantage of
553 heterozygous alleles or the detrimental homozygous alleles that get selected against, as reported
554 in inbred population of clonal honey bees, *Apis mellifera capensis* which retained after 20 years
555 of inbreeding high heterozygosity throughout its genome due to selection against homozygotes
556 [60]. Similar examples of heterozygosity advantage through its maintenance at high proportions
557 of the genome were reported in other species, isolated wolf populations and a hermaphrodite
558 worm after several generations of selfing [61, 62]. An alternative explanation for the observed
559 high heterozygosity could be due to the mutations that occurred in the ancestral lineage of Pt1,
560 Pt2, Pt3, Pt9, Pt15, which was revealed through admixture analysis, indicating that these

561 accessions share a common ancestry and are closely related to Pt17, which also exhibits high
562 heterozygosity. In contrast, accessions with lower heterozygosity display a distinct ancestry
563 pattern.

564 The SFS analysis provided compelling evidence consistent with the predictions of the
565 nearly neutral theory of evolution, revealing an excess of low frequency alleles. This prevalence
566 of lower frequency alleles in non-sense and non-synonymous polymorphisms can be attributed
567 to the effects of purifying selection, which acts against deleterious mutations. Consequently,
568 our investigation aimed to examine whether there was an elevated occurrence of non-
569 synonymous mutations in the Pt1 8.6 reference strain, in contrast to both the original Pt1
570 accession and the closely related Pt2 strain within the same clade. This analysis sought to
571 determine if the Pt1 8.6 strain had undergone a process of "domestication." However,
572 unexpectedly, we did not observe differences in non-synonymous mutations. Instead, we
573 observed a remarkable predominance of LOF mutations in the reference strain Pt1 8.6,
574 suggesting an adaptation to laboratory culture conditions facilitated through LOF mediated
575 mechanisms. The majority of these LOF mutations resulted in the repression or reduced
576 expression of targeted genes. However, a substantial number of genes showed moderate to high
577 expression levels, implying that these LOF mutations may function as an evolutionary
578 mechanism for generating new functional genes, serving as an adaptive response to culture
579 conditions [63-65]. An illustrative example is the domestication of maize where most of the
580 mutations are loss of function and the selection for a variety of traits has led to fixation of loss
581 of function alleles in today's crops [66, 67]. Another example is the human loss of function
582 mutations in the promoter of a red blood cell chemoreceptor, DARC that resulted in the
583 protection of human against malaria caused by *Plasmodium vivax* [68]. It is important to note
584 that not all LOF mutations lead to complete functional knockout. For instance, LOF mutations
585 at the 5' or 3' regions of genes may not entirely abolish their functions, and truncated proteins
586 resulting from such mutations could act as dominant-negative factors [69, 70]. Interestingly,
587 about 10.75% (331 out of 3078 genes) LOF showed moderate to high expression. Some specific
588 examples of these genes with known functions include: (i) a mitochondrial enzyme known as
589 glutamate dehydrogenase (Phatr3 J13951), which has been reported to play a crucial role in
590 carbon and nitrogen metabolism as well as energy supply under abiotic stresses in *Arabidopsis*
591 and the red alga *Pyropia haitanensis* [71-74] ; (ii) an LCH15 protein (Phatr3 J48882), which
592 functions as a chlorophyll binding protein and potentially contributes to the adaptation to
593 different light conditions in laboratory cultures and (iii) a heat shock transcription factor
594 (Phatr3 J49594) which similarly may contribute to the adaptation to lab culture temperatures.

595 Admixture analysis, PCA, and hierarchical clustering all identified six clusters, which
596 suggests that the samples had shared ancestry and similar geographical origins. However, not
597 all accessions within the same cluster had shared geographical origins. The English Channel
598 and East China Sea populations clustered together, indicating that *P. tricornutum* accessions
599 may have been dispersed by various means, or that similar ecological niches across the
600 sampling sites led to convergent evolution. Genome-wide phylogeny analysis confirmed six
601 clades, with some accessions falling into previously identified clades and others forming two
602 new clades, which suggests genetic divergence.

603 Several loci that are believed to be under balancing selection were found to have a high
604 level of genetic diversity. Notably, genes coding for stress-inducible proteins
605 (Phatr3_EG00471, Phatr3_J54019), outer membrane receptors (Phatr3_EG01193), and cell
606 cycle genes (Phatr3_J34920, Phatr3_EG00817) displayed an excess of polymorphism,
607 reflecting their role in protecting cells from stresses such as high temperatures, starvation, and
608 infection, as well as in cell division and growth. Consistent with this, transmembrane protein
609 have been previously observed to evolve faster than protein without a transmembrane domain
610 in unicellular eukaryotes such as yeast [75] and *Ostreococcus* [76]. Interestingly, the genes
611 identified under balancing selection were found in multiple clades and were specific to one or
612 more accessions, suggesting that spatially varying selection forces may be related to local
613 niches. It is expected that these same selection forces will apply to accessions from similar
614 ecological niches or with a shared origin and sampling locations.

615 Comparison of genes under BS between accessions sampled at divergent time points,
616 namely 1910, 1930, 1956, 1989 and 2016 revealed several identical genes suggesting the
617 persistence of long term balancing selection acting on these genomic regions (Figure 6a, Table
618 S8). Notably, the majority of these genes were functionally associated with stress resistance
619 and fundamental cellular processes, highlighting their potential significance in adaptation to
620 various habitats. Long term balancing selection was found at genes involved in diverse
621 processes such as disease resistance, self-incompatibility, and heat stress providing advantages
622 and enhancing fitness in natural populations [77-80].

623
624 Profiling transcript levels in the 17 accessions identified novel genes in the assayed
625 growth conditions suggesting condition and accession specific genes that were not identified in
626 the numerous growth conditions previously reported [7]. Interestingly, our analysis revealed
627 that genes carrying LOF mutations displayed a significant decrease in expression levels when
628 compared to their non-LOF counterparts (Figure 6b), implying the role of DNA sequence

629 variations in shaping gene expression patterns. Nonetheless, we also noticed a considerable
630 number of LOF mutations that did not result in downregulation. The observed LOF mutations
631 with no effect on gene expression is likely due to the robustness of the genome through gene
632 duplication and compensatory mechanisms allowing for the tolerance of many LOF variants,
633 resulting in the majority of these mutations being silent and having little to no impact. Multiple
634 LOF mutations were observed across clades as well as within them, particularly among
635 accessions that exhibited extreme phenotypes (short cell size versus long ones, low versus high
636 photosynthetic performance). This suggests that there may be variations in genetic backgrounds
637 and/or epigenetic factors among these accessions. For instance, Pt12 and Pt14, despite having
638 vastly different cell morphologies (very long versus short cells), share 1273 LOF mutations.
639 Similarly, Pt2 and Pt9, as well as Pt6 and Pt12, share 1479 and 1120 LOF mutations
640 respectively, but exhibit distinct photosynthetic performances. In general, no clear association
641 can be established between genetic diversity including SNP, INDELS and LOF mutations and
642 the phenotypic traits investigated in this study.

643 Interestingly, we observed LOF mutations that led to upregulation of genes instead,
644 suggesting the potential existence of Gain of Function (GOF) mutations. These GOF mutations
645 known to occur mostly in unstructured regions may be attributed to the emergence of
646 transcription factor binding sites, miRNA binding sites, an RNA binding protein or new
647 functional domains [81, 82]. Since genes and their products do not operate in isolation but rather
648 in biological networks, these newly acquired domains may acquire functionality through their
649 interactions with other proteins. Additionally, compensatory mechanisms for LOF or GOF may
650 involve epigenetic processes that serve as a platform for interacting with specific proteins or
651 complexes. WGCNA analysis revealed several network modules that were further merged into
652 six clusters with similar expression patterns indicating co-regulated genes and pathways. Both
653 differential expression and WGCNA analysis corroborated the presence of differences in
654 transcript levels, which cannot be solely attributed to genetic variations. This observation
655 implies the involvement of other regulatory mechanisms, such as epigenetics, that are known
656 to modulate gene transcription [6, 83, 84].

657

658 **Conclusions**

659 Our study provides a comprehensive assessment of the genetic and transcriptional
660 diversity among 17 natural accessions of the model diatom *P. tricornutum*. Our investigation
661 has uncovered novel clades, which are likely indicative of previously unexplored ecological
662 niches. Moreover, we have identified new genes that expand the existing transcriptome

663 repertoire of this species. By incorporating recently sampled accessions, we have further
664 confirmed a persistent long-term balancing selection and the high level of heterozygosity in *P.*
665 *tricornutum* through population genetic analyses, suggesting that this characteristic arises from
666 a heterozygous advantage. Our findings establish a crucial groundwork for future research that
667 utilizes sequencing data from various *P. tricornutum* accessions which we made available via
668 PhaeoEpiView platform (<https://PhaeoEpiView.univ-nantes.fr>) [10] for easy and
669 comprehensive use. This will enhance our understanding of diatom biology, foster
670 advancements in biotechnology applications, and optimize trait selection.

671 **Acknowledgements**

672 We thank Carine Pruvost for media preparation. We are grateful to Achal Rastogi for his helpful
673 discussions on the bioinformatics methodology related to phylogeny. LT acknowledges support
674 from the region of Pays de la Loire and Nantes métropole (ConnecTalent EPIALG project),
675 Epicycle ANR project (ANR-19-CE20- 0028-02) and Région Pays de la Loire ImpulseAlgae
676 projects. FY was supported by CSC Grant 201906310152. We are grateful to the bioinformatics
677 core facility of Nantes University (BiRD Biogenouest) for its technical support.

678

679 **Authors contribution**

680 LT conceived and designed the study. TC performed and coordinated the bioinformatics anal-
681 ysis. FY conducted most of the experiments. AG extracted RNA and performed the PAM study.
682 EM contributed to the bioinformatics analysis. BJ supervised the PAM study and assisted with
683 data analysis. GP and LT supervised the genetic population study. OAM performed the
684 WGCNA analysis. YW and UC assisted with DNA extraction. TC, FY, AG, GP and LT ana-
685 lysed and interpreted the data. LT supervised and coordinated the study. LT wrote the manu-
686 script with input from TC, FY, AG, EM and OAM. All authors read and edited the manuscript.

687

688 **Data availability**

689 The data that support the findings of this study are openly available in BioProjects
690 PRJNA430316 and PRJNA971163

691

692 **Competing interests**

693 None of the authors have any competing interests

694

695 **References**

- 696 1. Falkowski PG: **Evolution of the nitrogen cycle and its influence on the biological**
697 **sequestration of CO₂ in the ocean.** *Nature* 1997, **387**:272-275.
- 698 2. Treguer PJ, De La Rocha CL: **The world ocean silica cycle.** *Ann Rev Mar Sci* 2013, **5**:477-501.
- 699 3. Lauritano C AJ, Hansen E, Albrigtsen M, Escalera L, Esposito F, Helland K, Hanssen KØ,
700 Romano G and Ianora A: **Bioactivity Screening of Microalgae for Antioxidant, Anti-**
701 **Inflammatory, Anticancer, Anti-Diabetes, and Antibacterial Activities.** . *Front Mar Sci* 2016,
702 **3**:68.
- 703 4. Rabiee N, Khatami, M., Jamalipour Soufi, G., Fatahi, Y., Iravani, S., & Varma, R. S. : **Diatoms**
704 **with invaluable applications in nanotechnology, biotechnology, and biomedicine: recent**
705 **advances.** . *ACS Biomaterials Science and Engineering* 2021, **7**:3053-3068.

- 706 5. Falciatore A, Jaubert M, Bouly JP, Bailleul B, Mock T: **Diatom Molecular Research Comes of**
707 **Age: Model Species for Studying Phytoplankton Biology and Diversity.** *Plant Cell* 2020,
708 **32:547-572.**
- 709 6. Veluchamy A, Rastogi A, Lin X, Lombard B, Murik O, Thomas Y, Dingli F, Rivarola M, Ott S, Liu
710 X, et al: **An integrative analysis of post-translational histone modifications in the marine**
711 **diatom *Phaeodactylum tricornutum*.** *Genome Biol* 2015, **16:102.**
- 712 7. Rastogi A, Maheswari U, Dorrell RG, Vieira FRJ, Maumus F, Kustka A, McCarthy J, Allen AE,
713 Kersey P, Bowler C, Tirichine L: **Integrative analysis of large scale transcriptome data draws**
714 **a comprehensive landscape of *Phaeodactylum tricornutum* genome and evolutionary**
715 **origin of diatoms.** *Sci Rep* 2018, **8:4834.**
- 716 8. Rastogi A, Murik O, Bowler C, Tirichine L: **PhytoCRISP-Ex: a web-based and stand-alone**
717 **application to find specific target sequences for CRISPR/CAS editing.** *BMC Bioinformatics*
718 **2016, 17:261.**
- 719 9. Nymark M, Sharma AK, Sparstad T, Bones AM, Winge P: **A CRISPR/Cas9 system adapted for**
720 **gene editing in marine algae.** *Sci Rep* 2016, **6:24951.**
- 721 10. Wu Y, Chaumier T, Manirakiza E, Veluchamy A, Tirichine L: **PhaeoEpiView: an epigenome**
722 **browser of the newly assembled genome of the model diatom *Phaeodactylum***
723 ***tricornutum*.** *Sci Rep* 2023, **13:8320.**
- 724 11. Rastogi A, Vieira FRJ, Deton-Cabanillas AF, Veluchamy A, Cantrel C, Wang G, Vanormelingen
725 P, Bowler C, Piganeau G, Hu H, Tirichine L: **A genomics approach reveals the global genetic**
726 **polymorphism, structure, and functional diversity of ten accessions of the marine model**
727 **diatom *Phaeodactylum tricornutum*.** *ISME J* 2020, **14:347-363.**
- 728 12. De Martino AM, A. Juan Shi, K.P. Bowler, C.: **Genetic and phenotypic characterization of**
729 ***Phaeodactylum tricornutum* (Bacillariophyceae) accessions.** *J Phycol* 2007, **43:992-1009.**
- 730 13. Bailleul B, Rogato A, de Martino A, Coesel S, Cardol P, Bowler C, Falciatore A, Finazzi G: **An**
731 **atypical member of the light-harvesting complex stress-related protein family modulates**
732 **diatom responses to light.** *Proc Natl Acad Sci U S A* 2010, **107:18214-18219.**
- 733 14. Abida H, Dolch LJ, Mei C, Villanova V, Conte M, Block MA, Finazzi G, Bastien O, Tirichine L,
734 Bowler C, et al: **Membrane glycerolipid remodeling triggered by nitrogen and phosphorus**
735 **starvation in *Phaeodactylum tricornutum*.** *Plant Physiol* 2015, **167:118-136.**
- 736 15. Sprouffske K, Aguilar-Rodriguez J, Wagner A: **How Archiving by Freezing Affects the**
737 **Genome-Scale Diversity of *Escherichia coli* Populations.** *Genome Biol Evol* 2016, **8:1290-**
738 **1298.**
- 739 16. Riesco MF, Robles V: **Cryopreservation Causes Genetic and Epigenetic Changes in Zebrafish**
740 **Genital Ridges.** *PLoS One* 2013, **8:e67614.**
- 741 17. Wing KM, Phillips MA, Baker AR, Burke MK: **Consequences of Cryopreservation in Diverse**
742 **Natural Isolates of *Saccharomyces cerevisiae*.** *Genome Biol Evol* 2020, **12:1302-1312.**
- 743 18. Kram KE, Geiger C, Ismail WM, Lee H, Tang H, Foster PL, Finkel SE: **Adaptation of *Escherichia***
744 ***coli* to Long-Term Serial Passage in Complex Medium: Evidence of Parallel Evolution.**
745 *mSystems* 2017, **2.**
- 746 19. Russo MT, Aiese Cigliano R, Sanseverino W, Ferrante MI: **Assessment of genomic changes in**
747 **a CRISPR/Cas9 *Phaeodactylum tricornutum* mutant through whole genome resequencing.**
748 *PeerJ* 2018, **6:e5507.**
- 749 20. Bulankova P, Sekulic M, Jallet D, Nef C, van Oosterhout C, Delmont TO, Vercauteren I, Osuna-
750 Cruz CM, Vancaester E, Mock T, et al: **Mitotic recombination between homologous**
751 **chromosomes drives genomic diversity in diatoms.** *Curr Biol* 2021, **31:3221-3232 e3229.**
- 752 21. Hafker NS, Andreatta G, Manzotti A, Falciatore A, Raible F, Tessmar-Raible K: **Rhythms and**
753 **Clocks in Marine Organisms.** *Ann Rev Mar Sci* 2023, **15:509-538.**
- 754 22. Tirichine L, Bowler C: **Decoding algal genomes: tracing back the history of photosynthetic**
755 **life on Earth.** *Plant J* 2011, **66:45-57.**
- 756 23. Cruz de Carvalho MH, Sun HX, Bowler C, Chua NH: **Noncoding and coding transcriptome**
757 **responses of a marine diatom to phosphate fluctuations.** *New Phytol* 2016, **210:497-510.**

- 758 24. Nguyen DT, Wu B, Long H, Zhang N, Patterson C, Simpson S, Morris K, Thomas WK, Lynch M,
759 Hao W: **Variable Spontaneous Mutation and Loss of Heterozygosity among Heterozygous**
760 **Genomes in Yeast.** *Mol Biol Evol* 2020, **37**:3118-3130.
- 761 25. Vartanian M, Descles J, Quinet M, Douady S, Lopez PJ: **Plasticity and robustness of pattern**
762 **formation in the model diatom *Phaeodactylum tricornutum*.** *The New phytologist* 2009,
763 **182**:429-442.
- 764 26. Ralph P.J. GR: **Rapid light curves: A powerful tool to assess photosynthetic activity.** *Aquatic*
765 *Botany* 2005, **82**:222-237.
- 766 27. Platt T, Gallegos, C.L., Harrison, W.G. : **Photoinhibition of photosynthesis in natural**
767 **assemblages of marine phytoplankton.** *J Mar Res* 1980, **38**:687-701.
- 768 28. Serodio J, Lavaud J: **A model for describing the light response of the nonphotochemical**
769 **quenching of chlorophyll fluorescence.** *Photosynth Res* 2011, **108**:61-76.
- 770 29. Richards E RM, Rogers S. : **Preparation of genomic DNA from plant tissue.** *Current protocols*
771 *in molecular biology* 1994, **27**:1-23.
- 772 30. Nguyen TN, Berzano M, Gualerzi CO, Spurio R: **Development of molecular tools for the**
773 **detection of freshwater diatoms.** *J Microbiol Methods* 2011, **84**:33-40.
- 774 31. Siaut M, Heijde M, Mangogna M, Montsant A, Coesel S, Allen A, Manfredonia A, Falciatore A,
775 Bowler C: **Molecular toolbox for studying diatom biology in *Phaeodactylum tricornutum*.**
776 *Gene* 2007, **406**:23-35.
- 777 32. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina sequence data.**
778 *Bioinformatics* 2014, **30**:2114-2120.
- 779 33. Vasimuddin M, Misra, S., Li, H. and Aluru, S.: **Efficient Architecture-Aware Acceleration of**
780 **BWA-MEM for Multicore Systems.** In *IEEE International Parallel and Distributed Processing*
781 *Symposium (IPDPS)*. pp. 314-324. Rio de Janeiro, Brazil; 2019:314-324.
- 782 34. van der Auwera G, O'Connor, B.D.: *Genomics in the Cloud: Using Docker, GATK, and WDL in*
783 *Terra*. 2020.
- 784 35. Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM: **A**
785 **program for annotating and predicting the effects of single nucleotide polymorphisms,**
786 **SnEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3.** *Fly*
787 *(Austin)* 2012, **6**:80-92.
- 788 36. Korneliussen TS, Albrechtsen A, Nielsen R: **ANGSD: Analysis of Next Generation Sequencing**
789 **Data.** *BMC Bioinformatics* 2014, **15**:356.
- 790 37. **R: A language and environment for statistical computing** [[https://www.R-project.org/.](https://www.R-project.org/)]
- 791 38. Alexander DH, Novembre J, Lange K: **Fast model-based estimation of ancestry in unrelated**
792 **individuals.** *Genome Res* 2009, **19**:1655-1664.
- 793 39. Liao Y, Smyth GK, Shi W: **featureCounts: an efficient general purpose program for assigning**
794 **sequence reads to genomic features.** *Bioinformatics* 2014, **30**:923-930.
- 795 40. Conway JR, Lex A, Gehlenborg N: **UpSetR: an R package for the visualization of intersecting**
796 **sets and their properties.** *Bioinformatics* 2017, **33**:2938-2940.
- 797 41. Fijarczyk A, Babik W: **Detecting balancing selection in genomes: limits and prospects.** *Mol*
798 *Ecol* 2015, **24**:3529-3545.
- 799 42. Galili T: **dendextend: an R package for visualizing, adjusting and comparing trees of**
800 **hierarchical clustering.** *Bioinformatics* 2015, **31**:3718-3720.
- 801 43. Schliep KP: **phangorn: phylogenetic analysis in R.** *Bioinformatics* 2011, **27**:592-593.
- 802 44. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL: **Graph-based genome alignment and**
803 **genotyping with HISAT2 and HISAT-genotype.** *Nat Biotechnol* 2019, **37**:907-915.
- 804 45. Ito K, Murphy D: **Application of ggplot2 to Pharmacometric Graphics.** *CPT Pharmacometrics*
805 *Syst Pharmacol* 2013, **2**:e79.
- 806 46. Perteau M, Perteau GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL: **StringTie enables**
807 **improved reconstruction of a transcriptome from RNA-seq reads.** *Nat Biotechnol* 2015,
808 **33**:290-295.

- 809 47. Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Cech M, Chilton J, Clements D, Coraor
810 N, Gruning BA, et al: **The Galaxy platform for accessible, reproducible and collaborative**
811 **biomedical analyses: 2018 update.** *Nucleic Acids Res* 2018, **46**:W537-W544.
- 812 48. Langfelder P, Horvath S: **WGCNA: an R package for weighted correlation network analysis.**
813 *BMC Bioinformatics* 2008, **9**:559.
- 814 49. Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for RNA-**
815 **seq data with DESeq2.** *Genome Biol* 2014, **15**:550.
- 816 50. Ait-Mohamed O, Novak Vanclova AMG, Joli N, Liang Y, Zhao X, Genovesio A, Tirichine L,
817 Bowler C, Dorrell RG: **PhaeoNet: A Holistic RNAseq-Based Portrait of Transcriptional**
818 **Coordination in the Model Diatom Phaeodactylum tricornutum.** *Front Plant Sci* 2020,
819 **11**:590949.
- 820 51. Zhao X, Rastogi A, Deton Cabanillas AF, Ait Mohamed O, Cantrel C, Lombard B, Murik O,
821 Genovesio A, Bowler C, Bouyer D, et al: **Genome wide natural variation of H3K27me3**
822 **selectively marks genes predicted to be important for cell differentiation in Phaeodactylum**
823 **tricornutum.** *New Phytol* 2021, **229**:3208-3220.
- 824 52. Watterson GA: **On the number of segregating sites in genetical models without**
825 **recombination.** *Theor Popul Biol* 1975, **7**:256-276.
- 826 53. Siewert KM, Voight BF: **Detecting Long-Term Balancing Selection Using Allele Frequency**
827 **Correlation.** *Mol Biol Evol* 2017, **34**:2996-3005.
- 828 54. Lepetit B, Sturm S, Rogato A, Gruber A, Sachse M, Falciatore A, Kroth PG, Lavaud J: **High light**
829 **acclimation in the secondary plastids containing diatom Phaeodactylum tricornutum is**
830 **triggered by the redox state of the plastoquinone pool.** *Plant Physiol* 2013, **161**:853-865.
- 831 55. Jahns P, Holzwarth AR: **The role of the xanthophyll cycle and of lutein in photoprotection of**
832 **photosystem II.** *Biochim Biophys Acta* 2012, **1817**:182-193.
- 833 56. Taddei L, Stella GR, Rogato A, Bailleul B, Fortunato AE, Annunziata R, Sanges R, Thaler M,
834 Lepetit B, Lavaud J, et al: **Multisignal control of expression of the LHCX protein family in the**
835 **marine diatom Phaeodactylum tricornutum.** *J Exp Bot* 2016, **67**:3939-3951.
- 836 57. Malerba ME, Palacios MM, Palacios Delgado YM, Beardall J, Marshall DJ: **Cell size,**
837 **photosynthesis and the package effect: an artificial selection approach.** *New Phytol* 2018,
838 **219**:449-461.
- 839 58. Spanbauer TL, Fritz, S.C. and Baker, P.A.: **Punctuated changes in the morphology of an**
840 **endemic diatom from Lake Titicaca.** *Paleobiology* 2018, **44**:89 - 100.
- 841 59. Pfeifer SP: **From next-generation resequencing reads to a high-quality variant data set.**
842 *Heredity (Edinb)* 2017, **118**:111-124.
- 843 60. Smith NMA, Wade C, Allsopp MH, Harpur BA, Zayed A, Rose SA, Engelstadter J, Chapman NC,
844 Yagound B, Oldroyd BP: **Strikingly high levels of heterozygosity despite 20 years of**
845 **inbreeding in a clonal honey bee.** *J Evol Biol* 2019, **32**:144-152.
- 846 61. Kardos M, Akesson M, Fountain T, Flagstad O, Liberg O, Olason P, Sand H, Wabakken P,
847 Wikenros C, Ellegren H: **Genomic consequences of intensive inbreeding in an isolated wolf**
848 **population.** *Nat Ecol Evol* 2018, **2**:124-131.
- 849 62. Guo L, Zhang S, Rubinstein B, Ross E, Alvarado AS: **Widespread maintenance of genome**
850 **heterozygosity in Schmidtea mediterranea.** *Nat Ecol Evol* 2016, **1**:19.
- 851 63. Behe MJ: **Experimental evolution, loss-of-function mutations, and "the first rule of adaptive**
852 **evolution".** *Q Rev Biol* 2010, **85**:419-445.
- 853 64. Caseys C: **Loss of Function, a Strategy for Adaptation in Arabidopsis.** *Plant Cell* 2019, **31**:935.
- 854 65. Hottes AK, Freddolino PL, Khare A, Donnell ZN, Liu JC, Tavazoie S: **Bacterial adaptation**
855 **through loss of function.** *PLoS Genet* 2013, **9**:e1003617.
- 856 66. Meyer RS, Purugganan MD: **Evolution of crop species: genetics of domestication and**
857 **diversification.** *Nat Rev Genet* 2013, **14**:840-852.
- 858 67. Murray AW: **Can gene-inactivating mutations lead to evolutionary novelty?** *Curr Biol* 2020,
859 **30**:R465-R471.

- 860 68. Tournamille C, Colin Y, Cartron JP, Le Van Kim C: **Disruption of a GATA motif in the Duffy**
861 **gene promoter abolishes erythroid gene expression in Duffy-negative individuals.** *Nat*
862 *Genet* 1995, **10**:224-228.
- 863 69. de Valles-Ibanez G, Hernandez-Rodriguez J, Prado-Martinez J, Luisi P, Marques-Bonet T,
864 Casals F: **Genetic Load of Loss-of-Function Polymorphic Variants in Great Apes.** *Genome Biol*
865 *Evol* 2016, **8**:871-877.
- 866 70. MacArthur DG, Balasubramanian S, Frankish A, Huang N, Morris J, Walter K, Jostins L,
867 Habegger L, Pickrell JK, Montgomery SB, et al: **A systematic survey of loss-of-function**
868 **variants in human protein-coding genes.** *Science* 2012, **335**:823-828.
- 869 71. Labboun S, Terce-Laforgue T, Roscher A, Bedu M, Restivo FM, Velanis CN, Skopelitis DS,
870 Moschou PN, Roubelakis-Angelakis KA, Suzuki A, Hirel B: **Resolving the role of plant**
871 **glutamate dehydrogenase. I. In vivo real time nuclear magnetic resonance spectroscopy**
872 **experiments.** *Plant Cell Physiol* 2009, **50**:1761-1773.
- 873 72. Grzechowiak M, Sliwiak J, Jaskolski M, Ruskowski M: **Structural Studies of Glutamate**
874 **Dehydrogenase (Isoform 1) From Arabidopsis thaliana, an Important Enzyme at the Branch-**
875 **Point Between Carbon and Nitrogen Metabolism.** *Front Plant Sci* 2020, **11**:754.
- 876 73. Terce-Laforgue T, Clement G, Marchi L, Restivo FM, Lea PJ, Hirel B: **Resolving the Role of**
877 **Plant NAD-Glutamate Dehydrogenase: III. Overexpressing Individually or Simultaneously**
878 **the Two Enzyme Subunits Under Salt Stress Induces Changes in the Leaf Metabolic Profile**
879 **and Increases Plant Biomass Production.** *Plant Cell Physiol* 2015, **56**:1918-1929.
- 880 74. Li S, Shao Z, Lu C, Yao J, Zhou Y, Duan D: **Glutamate Dehydrogenase Functions in Glutamic**
881 **Acid Metabolism and Stress Resistance in Pyropia haitanensis.** *Molecules* 2021, **26**.
- 882 75. Julenius K, Pedersen AG: **Protein evolution is faster outside the cell.** *Mol Biol Evol* 2006,
883 **23**:2039-2048.
- 884 76. Jancek S, Gourbiere S, Moreau H, Piganeau G: **Clues about the genetic basis of adaptation**
885 **emerge from comparing the proteomes of two Ostreococcus ecotypes (Chlorophyta,**
886 **Prasinophyceae).** *Mol Biol Evol* 2008, **25**:2293-2300.
- 887 77. Malaria Genomic Epidemiology N, Band G, Rockett KA, Spencer CC, Kwiatkowski DP: **A novel**
888 **locus of resistance to severe malaria in a region of ancient balancing selection.** *Nature*
889 2015, **526**:253-257.
- 890 78. Key FM, Teixeira JC, de Filippo C, Andres AM: **Advantageous diversity maintained by**
891 **balancing selection in humans.** *Curr Opin Genet Dev* 2014, **29**:45-51.
- 892 79. Segurel L, Thompson EE, Flutre T, Lovstad J, Venkat A, Margulis SW, Moyses J, Ross S, Gamble
893 K, Sella G, et al: **The ABO blood group is a trans-species polymorphism in primates.** *Proc*
894 *Natl Acad Sci U S A* 2012, **109**:18493-18498.
- 895 80. Junprung W. SP, Tassanakajon A., Van Stappen G., Peter Bossier: **Balancing selection at the**
896 **ATP binding site of heat shock cognate 70 (HSC70) contributes to increased**
897 **thermotolerance in Artemia franciscana.** *Acquaculture* 2021, **531**:735988.
- 898 81. Li XH, Babu MM: **Human Diseases from Gain-of-Function Mutations in Disordered Protein**
899 **Regions.** *Cell* 2018, **175**:40-42.
- 900 82. Meyer Kea: **Mutations in disordered regions can cause disease by creating dileucine motifs.**
901 *Cell* 2018, **175**:239–253.
- 902 83. Jaenisch R, Bird, A.: **Epigenetic regulation of gene expression: how the genome integrates**
903 **intrinsic and environmental signals.** *Nat Genet* 2003, **33**:245-254.
- 904 84. Veluchamy A, Lin X, Maumus F, Rivarola M, Bhavsar J, Creasy T, O'Brien K, Sengamalay NA,
905 Tallon LJ, Smith AD, et al: **Insights into the role of DNA methylation in diatoms by genome-**
906 **wide profiling in Phaeodactylum tricornutum.** *Nat Commun* 2013, **4**.

907

908

909 **Legend**

910 **Figure 1.** World map illustrating the sampling sites for the 17 *P. tricornutum* accessions
911 analysed in the study with the year of sampling indicated in red.

912

913 **Figure 2.** Growth and photosynthetic features for the 17 accessions. **a** Growth rates in the 17
914 accessions, with error bars indicating standard deviations based on triplicate cultures. **b.** Mean
915 maximum relative electron transport rate (rETR_{max}) for the 17 accessions. **c** The maximum
916 PSII photochemical efficiency (F_v/F_m). **d** Non-photochemical quenching (NPQ) for the same
917 17 accessions.

918

919 **Figure 3.** Genetic diversity across *P. tricornutum* accessions. **a** Composition of heterozygous
920 and homozygous SNPs, the stacked bar plot represents the number of SNPs discovered in Pt1
921 to Pt17 accessions, showing a contrasted proportion of homozygous and heterozygous SNPs.
922 Homozygous SNPs are displayed in light blue, heterozygous SNPs are shown in dark blue.
923 **b** The bar plot represents the number of insertions (red) and deletions (yellow) called in Pt1 to
924 Pt17 accessions. **c** Folded site frequency spectrum SFS of 1898, 166877, 206536 and 168969
925 non-sense, non-synonymous, intergenic and synonymous SNPs, respectively. In order to obtain
926 an unbiased SFS, only one accession was chosen to represent each group of genetically close
927 accessions. These variants were thus called on Pt1-4-7-9-11-12-14-16-17. **d** Pie charts represent
928 different proportion of SNPs and INDELS over all functional features of the genome; GENES
929 (blue), TEs (gray), IGRs (Intergenic Regions, represented in yellow). **e** The bar plots represent
930 the total number of genes considered to exhibit CNV per accession (left, colored by clade) and
931 those that are accession-specific (top, in grey). Only accessions having specific genes with CNV
932 are shown in the matrix (center, in black). **f** The bar chart represents total and specific numbers
933 of genes that are affected by loss-of-function (LoF) mutations for all ecotypes (Pt1-Pt17). The
934 total number of genes (blue color) is the number of all non-duplicated genes on which a single
935 variant (INDEL or SNP) was taken into account. The grey bars represent the number of genes
936 unique to a specific accession and not present in the others.

937 **Figure 4.** Clustering of *P. tricornutum* accessions into clades. **a** The heat-map shows the genetic
938 differentiation or association between all possible pairs of accessions. The colors indicate F_{st}
939 values, which range from 0.02 to 0.4, with a color gradient from yellow to green, respectively.

940 Values closer to 0 signify close genetic makeup and values closer to one indicate strong genetic
941 structuring between the populations. **b** ADMIXTURE plot representing the ancestry genome
942 fractions of Pt1 to Pt17 (in color) for K=15. **c** Principal component analysis (PCA) showing the
943 distance between the seventeen accessions based on their shared genome structure. **d** The
944 heatmap shows the log₂ fold change (log₂FC) of normalized reads counts between each
945 reference gene and average of the read counts of all the reference genes per accession. FPKMs
946 are used as normalized values, the log₂ ratio of each gene FPKM over the mean FPKM of each
947 accession being calculated. A blue to red color gradient in the heatmap represents low to high
948 log₂FC. The previously described clades A, B, C, D and the newly defined clades E and F are
949 shown as colored annotations on the top. Genes having a null FPKM in a given accession
950 (considered lost) are displayed in black. Only genes having a log₂FC over 2 in at least one
951 ecotype are plotted in this figure (222 genes) and are considered to exhibit Copy Number
952 Variation (CNV). **e** Phylogenetic association of the 17 accessions based on genome-wide
953 biallelic SNPs and Indels (640454 variants), built from a hierarchical clustering (canberra
954 distance and average linkage functions).

955

956 **Figure 5.** Transcription levels variations in the 17 accessions. **a** The bar plots represent the
957 number of differentially expressed genes (DEG) in the 17 accessions denoted on the Y axis. All
958 DEGs are displayed in blue. Downregulated genes are shown in grey and upregulated in green.
959 **b** The bar plots represent total number of upregulated genes in each accession (left, colored by
960 clade) and those that are specific to a given group of accessions (top, in grey). **c** The bar plots
961 represent the total number of downregulated genes in each accession (left, colored by clade)
962 and those that are specific to a given group of accessions (top, in grey). **d** Principal component
963 analysis (PCA) showing the distance between the seventeen accessions based on gene
964 expression values.

965

966 **Figure 6.** Balancing selection, LOF and WGCNA analyses. **a** The plot represents the total
967 number of genes showing a signature for balancing selection (BS) per accessions (left, colored
968 by clade) and those that are found only in a given group of accessions (top, in grey). Genes
969 under BS in all accessions are shown in blue. **b** Distribution of gene expression levels in genes
970 affected by loss-of-function mutations (blue) and unaffected genes (green) for each accession
971 using FPKM values. **c** Eigengene adjacency heatmap of the 33 merged modules of *P.*

972 *tricornutum* accessions network. Each row and column in the heatmap corresponds to one
973 module (represented by their color). The scale bar on the right represents the correlation
974 strength ranging from 0 (blue) to 1 (red).

975

976 **Supplementary figures**

977 **Figure S1.** Light microscopy images of *P. tricornutum* accessions depicting cell morphology
978 and size. The red-framed cells (Pt3 in **a**, and Pt9 in **b**) represent the proportions of oval
979 morphotype. **c** Cell size and morphology of Pt11 (**c**), Pt12 (**d**), Pt13 (**e**), Pt14 (**f**) Pt15 (**g**). **h** Cell
980 size, morphology and proportion of tri-radiate morphotype in Pt16. **i** Cell size and morphology
981 in Pt17. Red lines indicate the 30 cells for which width and length were measured.

982

983 **Figure S2.** Gel images of PCR product illustrating the validation of INDELs observed in Pt11
984 to Pt17. **a** Insertion validation. **b** Deletion validation. The molecular weight marker is 1 kb plus
985 DNA ladder (M), and the negative control is represented by (N).

986 **Figure S3.** Gel images of PCR product illustrating the validation of gene loss observed in Pt11
987 to Pt17. The molecular weight marker is 1 kb plus DNA ladder (M), and the negative control is
988 represented by (N).

989

990 **Figure S4.** The plot shows the error of ADMIXTURE cross-validation process for K ranging
991 from 1 to 17, from all accessions callable SNPs and INDELs. The lowest value (15) gives an
992 indication of the ancestral populations number.

993

994 **Figure S5.** Inter-samples correlation heat map. The correlation coefficient is represented by the
995 square of the Pearson correlation coefficient (R). The greater the value of R, the higher the
996 degree of similarity between samples.

997

998 **Figure S6.** Comparative transcriptional analysis of genes involved in non-photochemical
999 quenching (NPQ) capacity across accessions using Pt1 8.6 as reference strain. **a** Light-
1000 harvesting complex X1 (Lhcx1). **b** Light-harvesting complex X2 (Lhcx2). **c** Light-harvesting

1001 complex X3 (Lhcx3). **d** Light-harvesting complex X3 ((Lhcx4). **e** Zeaxanthin epoxidase 1
1002 (ZEP1). **f** Zeaxanthin epoxidase 2 (ZEP2). **g** Zeaxanthin epoxidase 3 (ZEP3).

1003 **Figure S7.** GO enrichment analysis. **a** GO terms enrichment in each of Pt1, Pt7, Pt9, Pt11, Pt12,
1004 Pt16 and Pt17. In other accessions, no significant enrichment of GO terms was found. **b** The
1005 enrichment of GO terms was performed on the set of specific downregulated DEGs for each
1006 ecotype. Bar charts show which GO terms are represented and in which ecotype (Pt3, Pt8, Pt10,
1007 Pt11, Pt14, Pt16, Pt17). In other ecotypes, no significant GO enrichment terms were found.

1008 **Figure S8.** Heatmap visualization of gene expression and corresponding eigengenes across
1009 accessions in the 33 modules. The x-axis of both the heatmap and barplot displays the *P.*
1010 *tricornutum* accessions in the following order: Pt1R1, Pt2R1, Pt2R2, Pt3R1, Pt3R2 until
1011 Pt17R2. Each row of the heatmap corresponds to genes belonging to the module. Red color
1012 denotes overexpression and green underexpression.

1013

1014 **Table S1.** Sampling location, date and morphological features of the 17 accessions of *P. trio-*
1015 *rutum*

1016 **Table S2.** List of primer sequences used in the study

1017 **Table S3.** INDEL loci identified in Pt11 to Pt17 accessions

1018 **Table S4.** Copy number variation (CNV) identified in Pt1 to Pt17 accessions

1019 **Table S5.** Gene loss in Pt1 to Pt17 accessions

1020 **Table S6.** Loss of function (LoF) loci in Pt1 to Pt17 accessions

1021 **Table S7.** Table of the Fst distance between the different populations

1022 **Table S8.** Genes under balancing selection identified in Pt1 to Pt17 accessions

1023 **Table S9.** List of the novel genes identified in Pt1 to Pt17 accessions

1024 **Table S10.** Repeats identified in novel transcripts

1025 **Table S11.** Non coding RNA and other categories of RNA identified in novel transcripts. The
1026 table is composed of Sequence Name :the name of the sequence, RNA size: the length of the
1027 original transcript, ORF size:the size of the potential ORF within the sequence, Fickett Score:the
1028 Fickett score is a linguistic feature that distinguishes protein-coding RNA and ncRNA accord-
1029 ing to the combinational effect of nucleotide composition and codon usage bias, Hexamer
1030 Score:the hexamer score is calculated using a log-likelihood ratio to measure differential hex-

1031 amer usage between coding and non-coding sequences, Coding Probability:the coding proba-
1032 bility assigned to each transcript(Human:coding probability(CP) \geq 0.364 indicates coding se-
1033 quence,Mouse:coding probability(CP) \geq 0.44 indicates coding sequence and Zebrafish:coding
1034 probability(CP) \geq 0.38 indicates coding sequence), Coding Label:marking for each sequence
1035 whether it is a coding, non-coding, or unknown coding potential transcript.

1036 **Table S12.** List of accession specific genes up or down regulated and their GO

1037 **Table S13.** List of genes up or down regulated and their GO in all accessions or per clade

1038 **Table S14.** *P. tricorutum* accessions modules after merging and their annotation

1039 **Table S15.** Module composition of identified clusters and their corresponding GO

1040

1041

1042

1043

1044

1045

1046

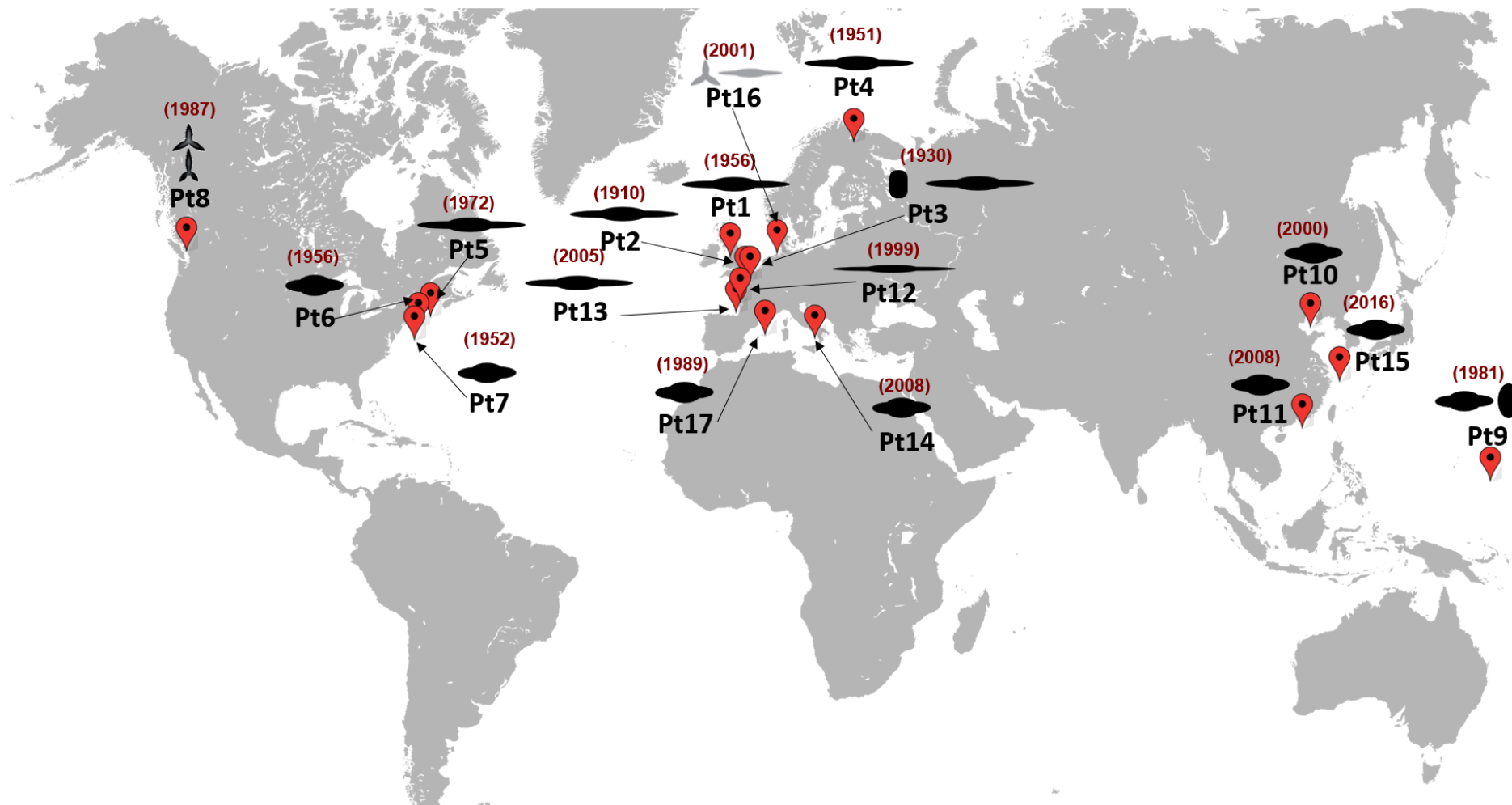
1047

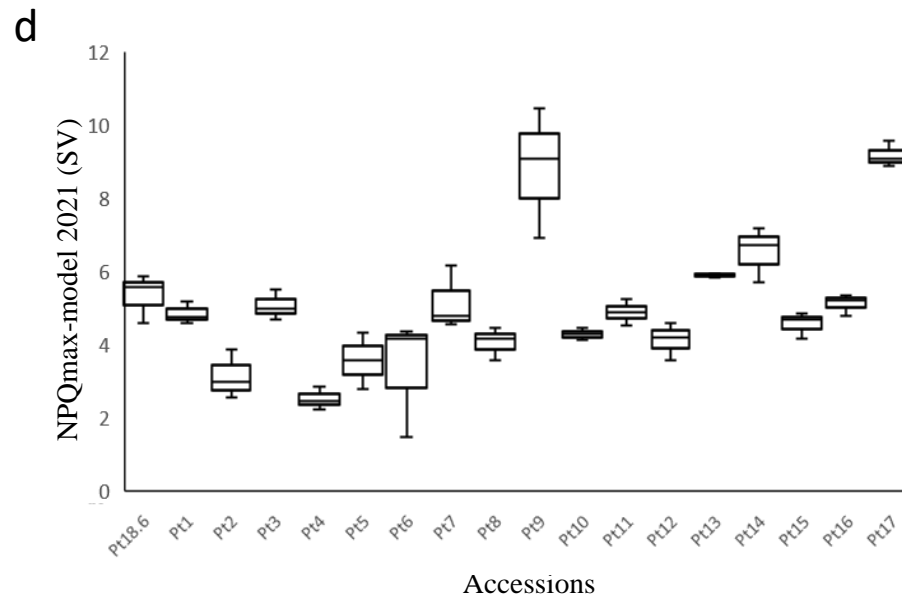
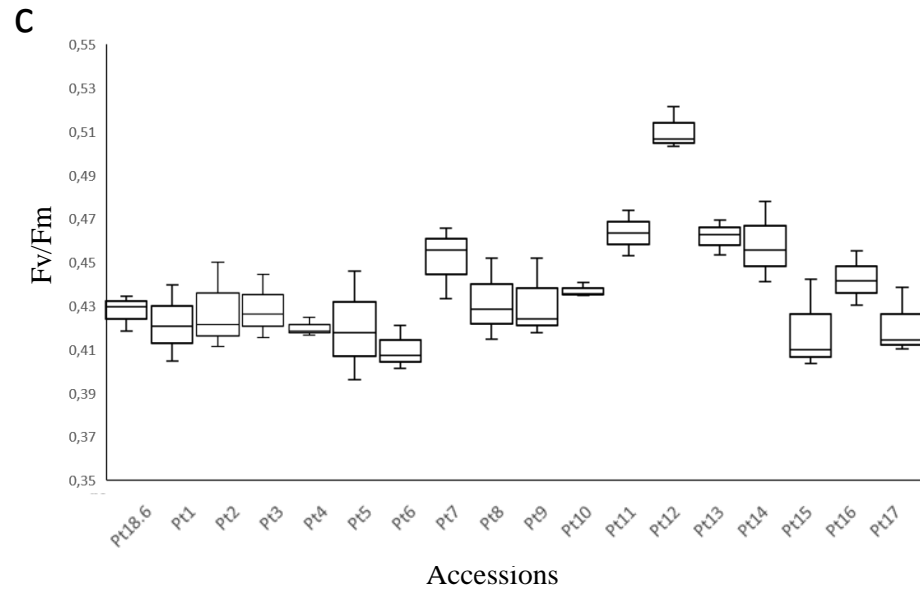
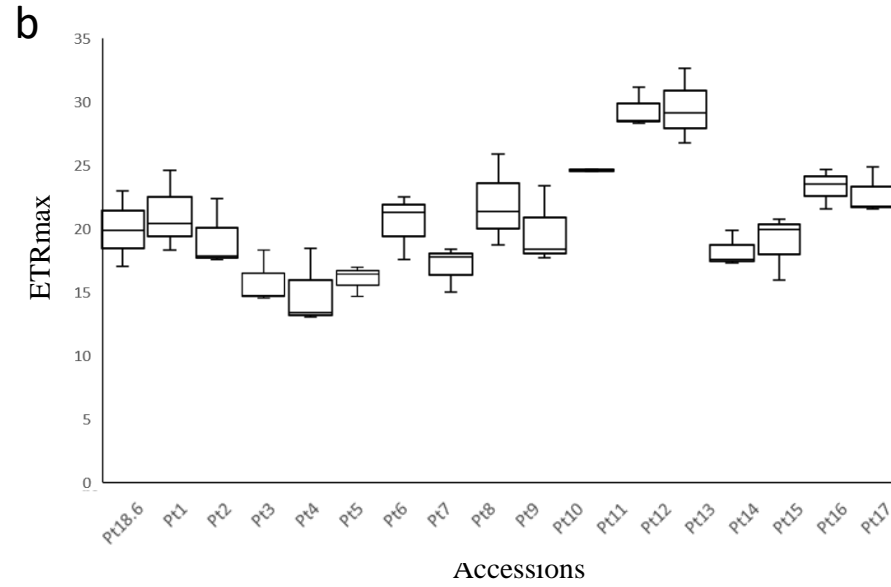
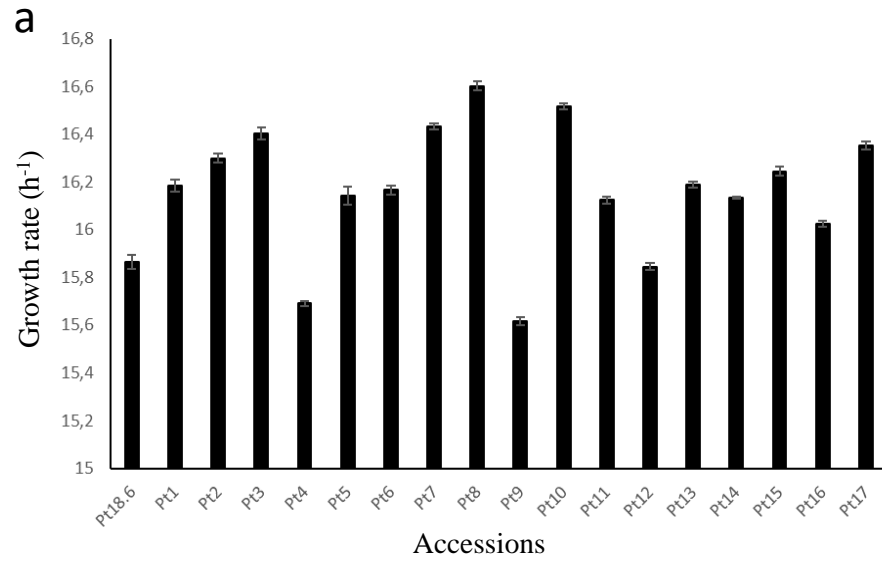
1048

1049

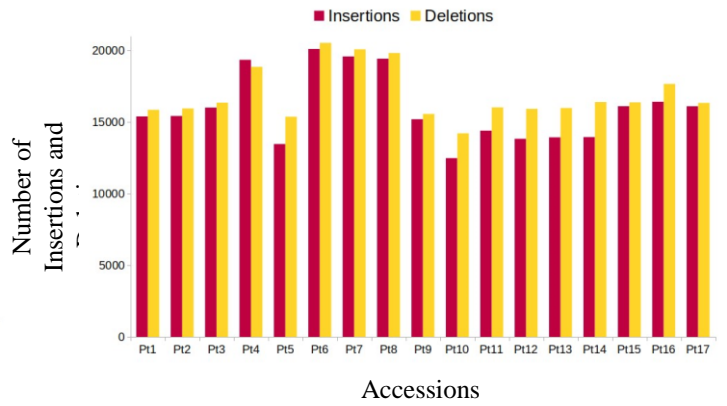
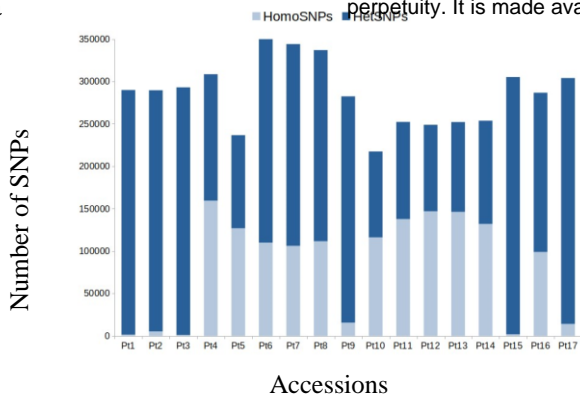
1050

1051

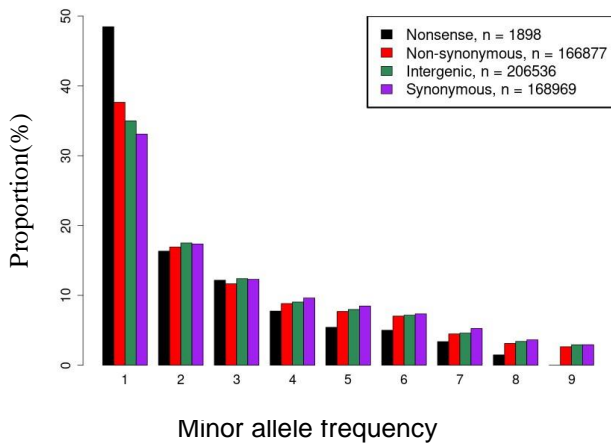




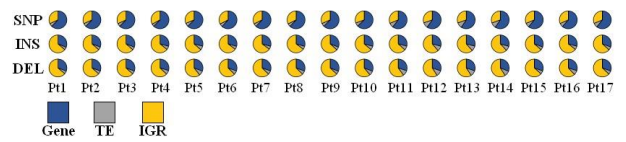
a



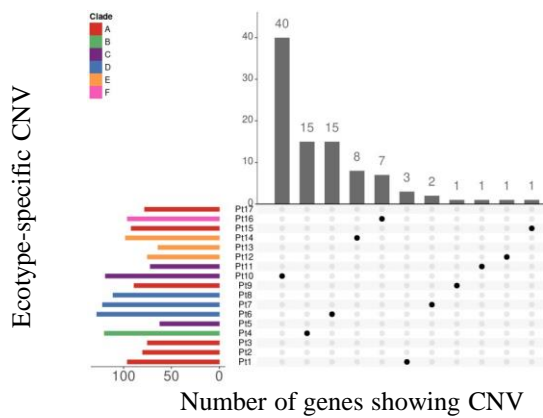
c



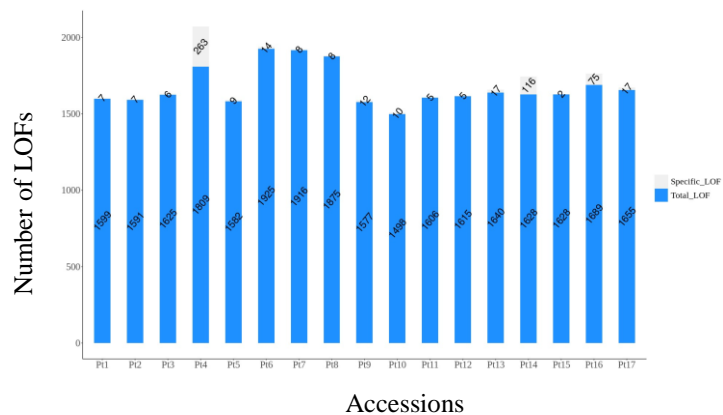
d

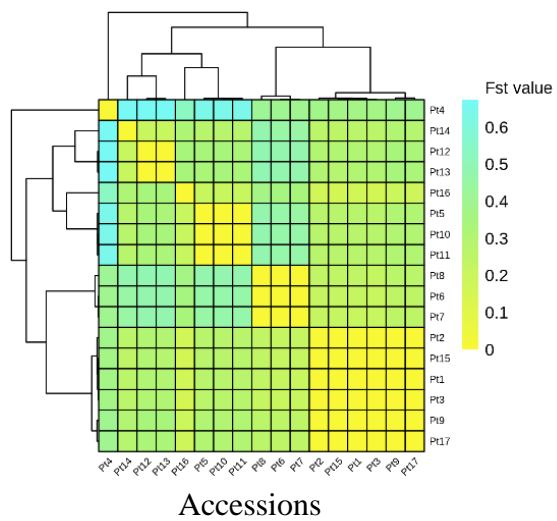
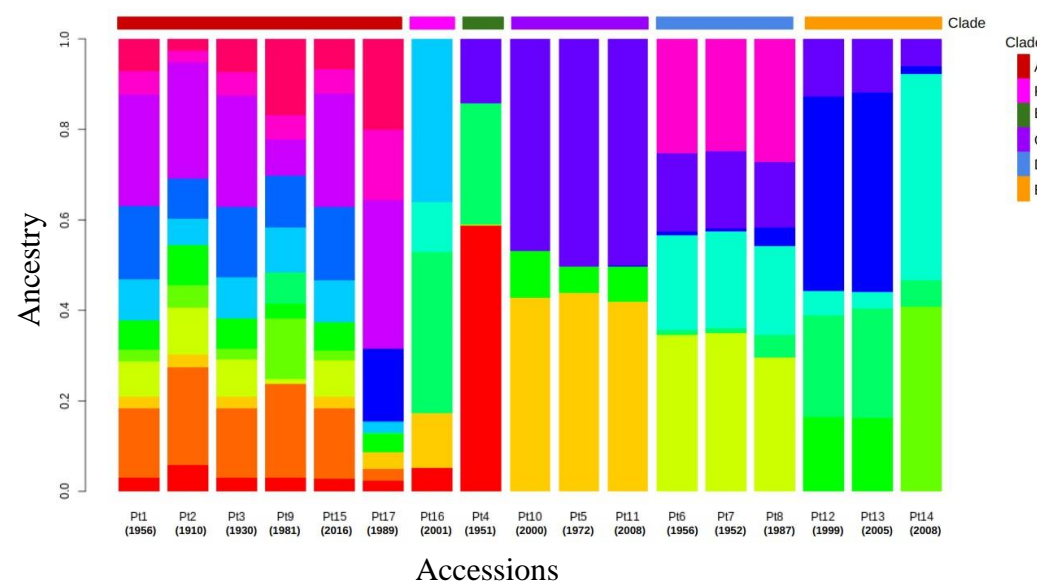
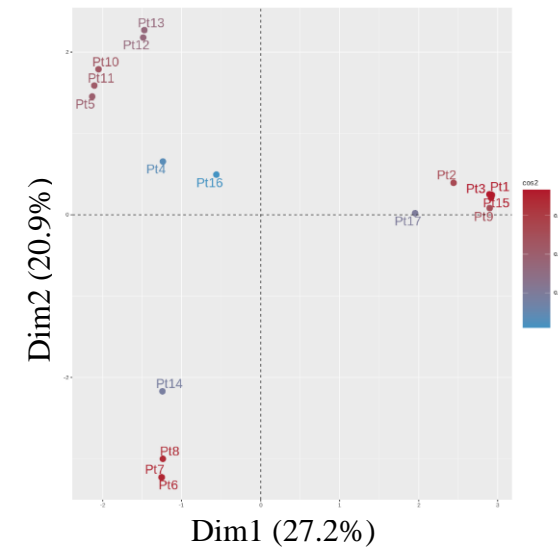
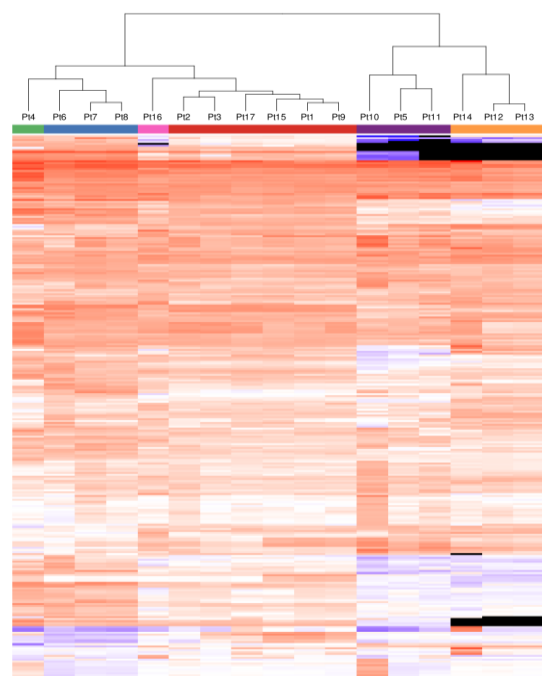
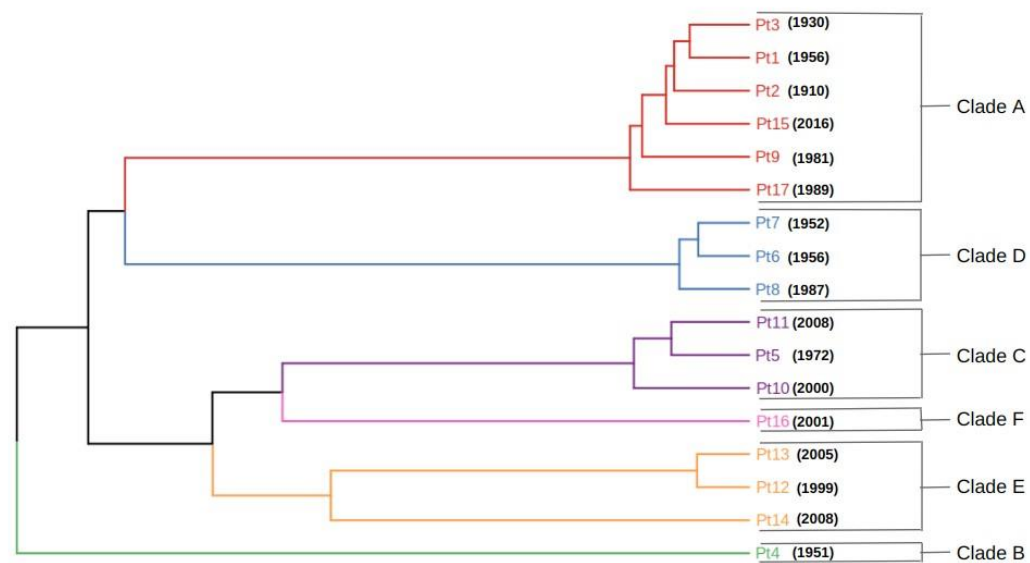


e

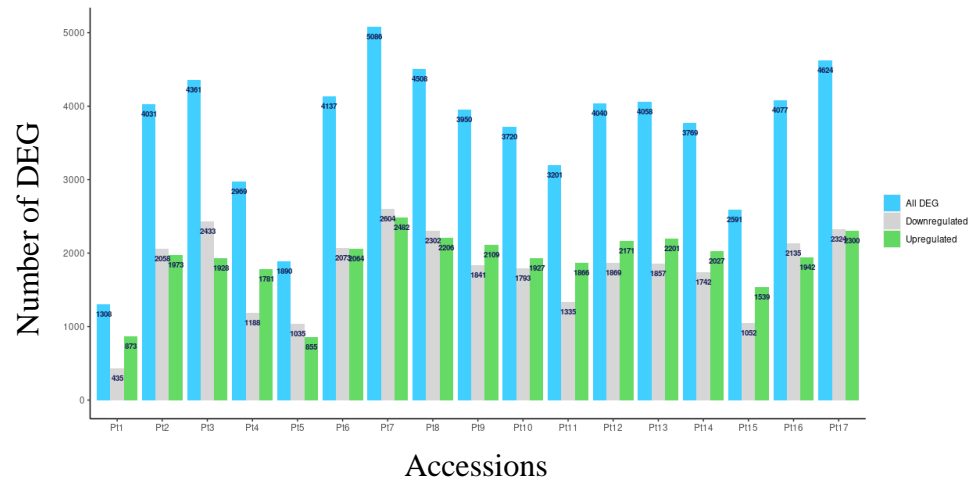


f

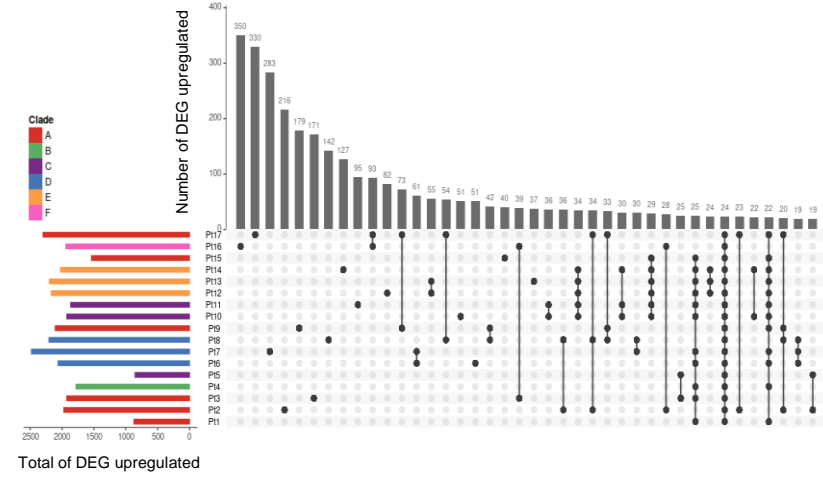


a**b****c****d****e**

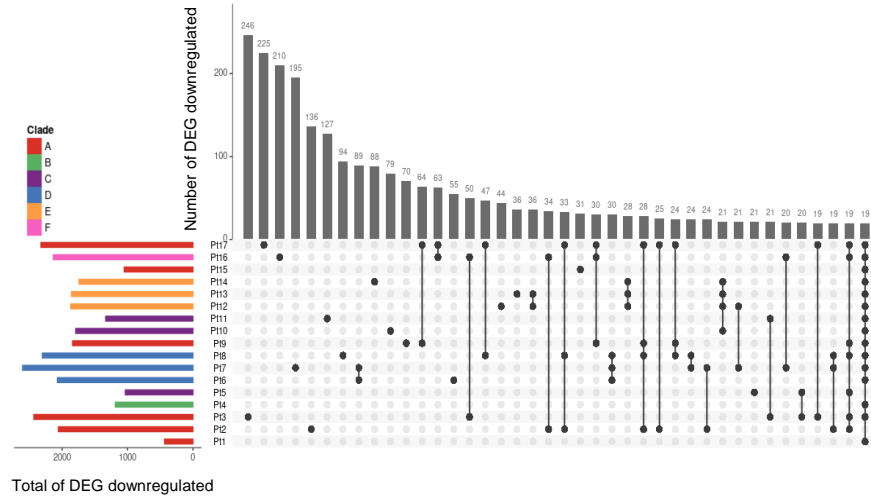
a



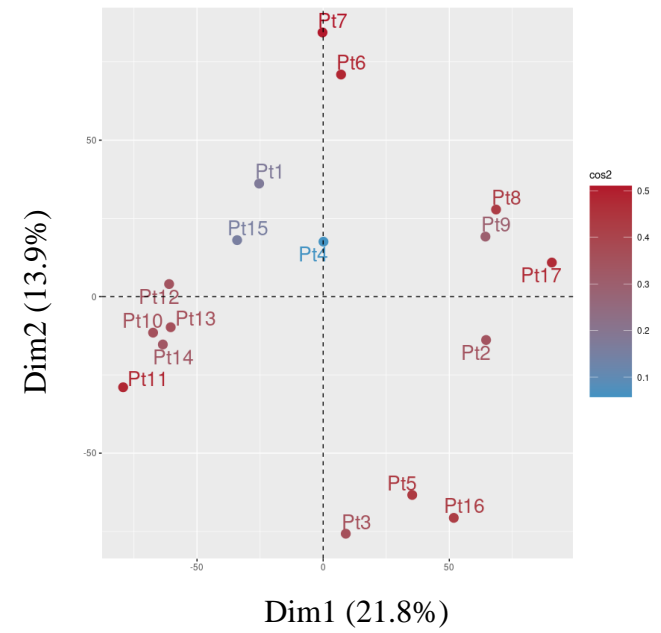
b



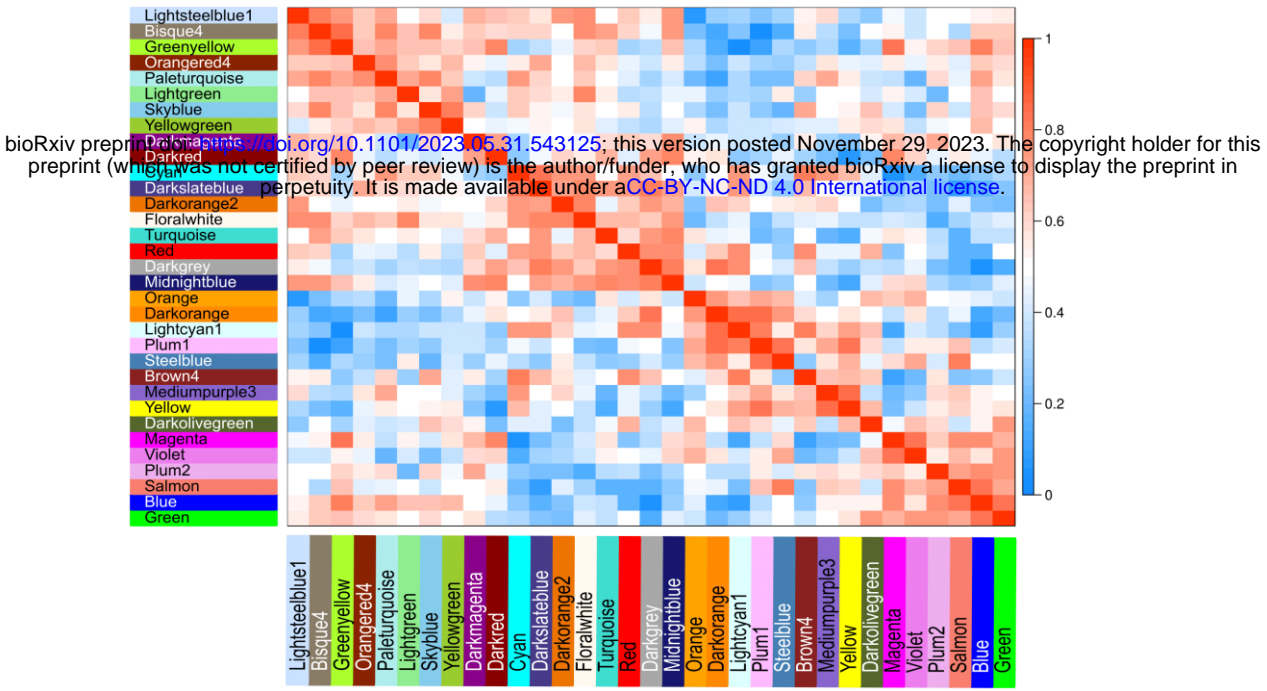
c



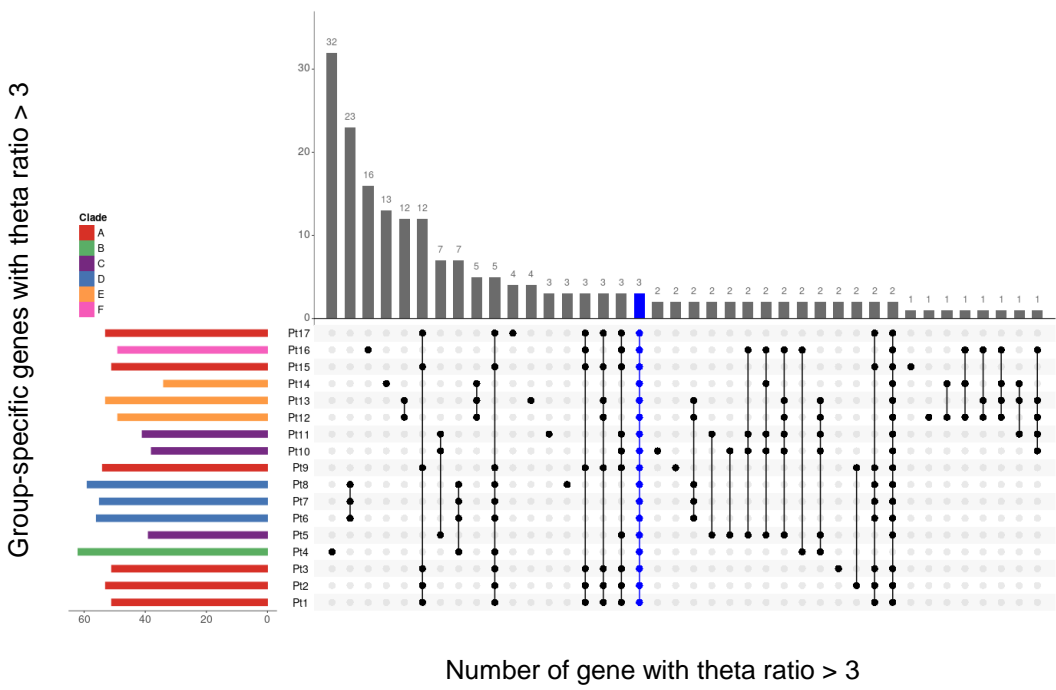
d



a



b



c

