

**A genomics approach reveals the global genetic polymorphism, structure and functional diversity of ten accessions of the marine model diatom *Phaeodactylum tricornutum***

Achal Rastogi<sup>1</sup>, Vieira FRJ<sup>1</sup>, Anne-Flore Deton-Cabanillas<sup>1</sup>, Alaguraj Veluchamy<sup>1,§</sup>, Catherine Cantrel<sup>1</sup>, Gaohong Wang<sup>2</sup>, Pieter Vanormelingen<sup>3</sup>, Chris Bowler<sup>1</sup>, Gwenael Piganeau<sup>4</sup>, Hanhua Hu<sup>2,\*</sup>, and Leila Tirichine<sup>1,€\*</sup>

<sup>1</sup>Institut de biologie de l'Ecole normale supérieure (IBENS), Ecole normale supérieure, CNRS, INSERM, PSL Université Paris 75005 Paris, France

<sup>2</sup>Key Laboratory of Algal Biology, Institute of Hydrobiology, Chinese Academy of Sciences, Wuhan, 430072 China

<sup>3</sup>Ghent University, Department of Biology, Research Group Protistology and Aquatic Ecology Krijgslaan 281/S8 9000 Gent, Belgium

<sup>4</sup> Sorbonne Universités, UPMC Univ Paris 06, CNRS, Biologie Intégrative des Organismes Marins (BIOM), Observatoire Océanologique, F-66650 Banyuls/Mer, France

<sup>§</sup>Current affiliation: Biological and Environmental Sciences and Engineering Division, Center for Desert Agriculture, King Abdullah University of Science and Technology, Thuwal 23955-6900, Saudi Arabia

<sup>€</sup>Current affiliation : Université de Nantes, CNRS, UFIP, UMR 6286, F-44000 Nantes, France

**\*Corresponding authors**

**Name:** Leila Tirichine

**Address:** Université de Nantes, CNRS, UFIP, UMR 6286, F-44000 Nantes, France

**Phone:** +33-276645058

**Email:** tirichine-l@univ-nantes.fr

**Name:** Hanhua Hu

**Address:** Key Laboratory of Algal Biology, Institute of Hydrobiology, Chinese Academy of Sciences, Wuhan, 430072 China

**Phone:** +86-27-68780078

**Email:** hanhuahu@ihb.ac.cn

**Running title:** *Phaeodactylum tricornutum* population genomics

## 41    **Abstract**

42

43    Diatoms emerged in the Mesozoic period and presently constitute one of the main  
44    primary producers in the world's ocean and are of a major economic importance. In  
45    the current study, using whole genome sequencing of ten accessions of the model  
46    diatom *Phaeodactylum tricornutum*, sampled at broad geospatial and temporal scales,  
47    we draw a comprehensive landscape of the genomic diversity within the species. We  
48    describe strong genetic subdivisions of the accessions into four genetic clades (A-D)  
49    with constituent populations of each clade possessing a conserved genetic and  
50    functional makeup, likely a consequence of the limited dispersal of *P. tricornutum* in  
51    the open ocean. We further suggest dominance of asexual reproduction across all the  
52    populations, as implied by high linkage disequilibrium. Finally, we show limited yet  
53    compelling signatures of genetic and functional convergence inducing changes in the  
54    selection pressure on many genes and metabolic pathways. We propose these findings  
55    to have significant implications for understanding the genetic structure of diatom  
56    populations in nature and provide a framework to assess the genomic underpinnings  
57    of their ecological success and impact on aquatic ecosystems where they play a major  
58    role. Our work provides valuable resources for functional genomics and for  
59    exploiting the biotechnological potential of this model diatom species.

60

## 61    **Introduction**

62    Diatoms are unicellular predominantly diploid and photosynthetic eukaryotes. They  
 63    belong to the large group known as Ochrophytes (plastid-bearing members of the  
 64    stramenopiles), constituents of CASH (Cryptomonads, Alveolates, Strameopiles, and  
 65    Haptophytes) lineages, and are believed to have evolved from serial endosymbiosis  
 66    involving green and red algal symbionts [1-3]. Ehrenberg first discovered diatoms in  
 67    the 19th century in dust samples collected by Charles Darwin in the Azores.  
 68    According to the earliest fossil records, they are believed to be in existence since at  
 69    least 190 million years [4] and their closest sister group are the Bolidomonads. In  
 70    nature, most diatoms likely live in obligate relationships with bacteria [5] but many,  
 71    like *Phaeodactylum tricornutum*, can be propagated in axenic conditions. In spite of  
 72    its low abundance in the open ocean [6], *P. tricornutum* is extensively used as a  
 73    model to study and characterize diatom metabolism, and to understand diatom  
 74    evolution [1, 7-11].

75    *P. tricornutum* is a coastal diatom found under highly unstable environments like  
 76    estuaries and rock-pools. Although it has never been reported to undergo sexual  
 77    reproduction, factors such as small cell size, discontinuous sexual phases, and the  
 78    sensitivity of sexual reproduction to many nonspecific abiotic components in diatoms  
 79    [12-14] limit our ability to constrain the sexual cycle of these organisms. Since the  
 80    discovery of *P. tricornutum* by Bohlin in 1897 and the characterization of different  
 81    morphologies, denoted fusiform, triradiate, oval, round and cruciform, 10 strains from  
 82    9 different geographic locations (sea shores, estuaries, rock pools, tidal creeks) around  
 83    the world, from sub-polar to tropical latitudes, have been accessioned (Fig S1) [well  
 84    described in [15]]. These accessions have been collected within the time frame of  
 85    approximately one century, from 1908 (Plymouth strain, Pt2/3) to 2000 (Dalian strain,

Pt10) (Fig S1) [15]. All the strains have been maintained either axenically or with native bacterial populations in different stock centers and have been cryopreserved after isolation. Previous studies have reported distinct functional behaviors of different accessions as adaptive responses to various environmental cues [16-19], but very little is known about their genetic diversity. However, based on sequence similarity of the ITS2 region within the 28S rDNA repeat sequence, the accessions can be divided into four genotypes (Genotype A: Pt1, Pt2, Pt3 and Pt9; Genotype B: Pt4; Genotype C: Pt5 and Pt10; Genotype D: Pt6, Pt7 and Pt8), with genotypes B and C being the most distant [15]. *P. tricornutum* is among the few diatom species with a whole genome sequence available to the community [20], and the only diatom for which extensive state-of-the-art functional and molecular tools have been developed over the past few decades [21-34]. These resources have advanced *P. tricornutum* as a model diatom species and provided a firm platform for future genome-wide structural and functional studies.

The accumulated effects of diverse evolutionary forces such as recombination, mutation, and selection have been found to dictate the structure and diversity of genomes in a wide range of species [35-38]. The existence of genomic diversity within a species reflects its potential to adapt to a changing environment. Exploring the genomic diversity within a species not only provides information about its evolution, it also offers opportunities to understand the role of various biotic and abiotic interactions in structuring a genome [39]. Such studies in diatoms are rare and estimates of genetic diversity within diatom populations are mostly inferred using microsatellite-based genotyping approaches [40-42]. Although these techniques have revealed a wealth of information about diatom evolution, their dispersal and reproductive physiology [39], additional insights can be obtained using state-of-the-

art whole genome comparative analysis techniques [42]. Deciphering the standing  
genomic variation of *P. tricornutum* across different accession populations, sampled  
at broad geospatial scale, is an important first step to assess the role of various  
evolutionary forces in regulating the adaptive capacities of diatoms in general  
(e.g.[43]). To understand the underlying genomic diversity within different accessions  
of *P. tricornutum* and to establish the functional implications of such diversity, we  
performed deep whole genome sequencing of the 10 most studied accessions, referred  
to as Pt1 to Pt10 [15, 18, 44]. We present a genome-wide diversity map of  
geographically distant *P. tricornutum* accessions, describing a stable genetic structure  
in the environment. This work further provides the community with whole genome  
sequences of the accessions, which will be a valuable genetic resource for functional  
studies of accession-specific ecological traits in the future.

123

## 124 **Results**

125

### 126 **Reference-assisted assembly reveals low nucleotide diversity across multiple** 127 **accessions of *P. tricornutum***

128 We sequenced the whole genomes of ten accessions of *P. tricornutum* using Illumina  
129 HiSeq 2000, and performed a reference-based assembly using the genome sequence  
130 of the reference strain Pt1 8.6 [1]. Across all accessions, the percentage of sequence  
131 reads mapped on the reference genome ranged between ~65% to ~80% (Table 1),  
132 with an alignment depth ranging between 26X and 162X, covering 92% to 98% of the  
133 reference genome (Table 1). Many regions on the reference genome that are observed  
134 as being unmapped by reads from individual ecotypes are annotated as being rich in

transposable elements (TEs) (Fig S2). At >90% identity the repeated proportion of unmapped reads varies between ~38% (Pt1) and 75% (Pt4).

Following the assembly, we performed variant calling using Genome Analysis Toolkit [45] and discovered 462,514 (depth  $\geq 4x$ ) single nucleotide polymorphisms (SNPs) including ~25% singleton sites, 573 insertions (of varying lengths from 1 bp to 312 bp) and 1,801 deletions (of lengths from 1 bp to 400 bp) (Fig 1A), across all the accessions. The spectrum of SNPs across all the accessions further reveals a higher rate of transitions (Ts) over transversions (Tv) ( $Ts/Tv = 1.6$ ). In total, compared to the reference alleles from Pt1.8.6, six possible types of single nucleotide changes could be distinguished, among which G:C  $\rightarrow$  A:T and A:T  $\rightarrow$  G:C accounted for more than ~60% of the observed mutations (Fig S3A). Further, most SNPs and INDELs (insertions and deletions) are shared between different accessions, except for Pt4, which possesses the highest proportion of specific SNPs (~35%) and INDELs (~75%) (Fig 1B). Interestingly, we found that most of the SNPs are heterozygous, and the proportion of heterozygous variants across all the accessions varies between ~45% (in Pt5 and Pt10) to ~98% (in Pt1, Pt2 and Pt3) (Fig 1C). Most of the variant alleles in the accessions with high proportions of heterozygous variants were further found to be significantly deviated from Hardy-Weinberg equilibrium (HWE) (chi-square test, P-value  $< 0.05$ ) (Fig 1C), possibly linked to prolonged asexual reproduction [46]. Surprisingly, despite significant differences in the proportion of heterozygote variant alleles between the accessions, which ranges between 45% to 98%, the average pairwise synonymous nucleotide diversity ( $\pi_s$ ) estimated from genes with callable sites across all the accessions is 0.007 per synonymous site. This indicates that any two homologous sequences taken at random across different populations will on average differ by only ~0.7% on synonymous positions. The non-synonymous

pairwise diversity ( $\pi_N$ ) over the same genes is 0.003, consistent with an excess of non-synonymous mutations being deleterious. Linkage disequilibrium (LD) analysis using only homozygous SNP sites revealed, on average, high LD ( $>0.7$ ) over pairs of variations, genome wide (Fig S3B). Further, based on the difference in the allelic frequencies of the SNPs, the pairwise *Fst* between the populations ranges from  $\sim 0.005$  (between Pt1 and Pt3) to  $\sim 0.4$  (between Pt4 and Pt10) (Fig 1D). Considering *Fst* as a measure of genetic differentiation or structuring between the populations, the ten *P. tricornutum* accessions can be clustered into 4 genetic groups/Clades with Pt1, Pt2, Pt3 and Pt9 in Clade A; Pt4 in Clade B; Pt5, Pt10 in Clade C; and Pt6, Pt7, Pt8 in Clade D, reflecting low intra-group *Fst* ( $\sim 0.02$ ) and high inter-group *Fst* ( $0.2 - 0.4$ ) (Fig 1D).

#### **Four genetic Clades of *P. tricornutum***

With the exception of Pt4, where we found the maximum number of variant alleles to be accession-specific, most of the variant alleles are shared between at least two accessions, indicating close genetic relatedness (Fig 1B). Therefore, in order to cluster the accessions based on the genome structure shared among them, we used Bayesian clustering approach by applying *Markov Chain Monte Carlo* (MCMC) estimations, programmed within the ADMIXTURE software [47]. As a result, population clustering of the ten accessions revealed four genetic clusters with Pt1, Pt2, Pt3, and Pt9 in one, Pt4 in a second, Pt5, Pt10 in a third, and Pt6, Pt7, Pt8 in a fourth cluster (Fig. 2A). Each cluster has a distinctive genetic makeup, which is also shared among different accessions in different proportions (represented with different colors in Fig. 2B). These four genetic clusters (Fig 2A) are also in broad agreement with *Fst*-based genetic Clades (Fig 1D), phylogenetic clusters inferred using ribosomal marker genes

185 (18S (Fig S4A), and ITS2 (Fig S4B)), as also reported previously [15], and at whole  
 186 genome scale (this study) as inferred by a phylogenetic tree generated using  
 187 maximum likelihood algorithm based on all (Fig 2C) and only homozygous  
 188 polymorphic sites (SNVs and INDELS) (Fig S4C). Based on the consensus from  
 189 different clustering algorithms, the four phylogenetic Clades were defined as Clade A  
 190 (Pt1, Pt2, Pt3 and Pt9), Clade B (Pt4), Clade C (Pt5 and Pt10), and Clade D (Pt6, Pt7,  
 191 and Pt8).

192 Further sequential assessment of the 18S and ITS2 rDNA gene sequences across  
 193 different clades indicated the presence of multiple variations, including both  
 194 heterozygous and homozygous variant alleles (Fig S4D and S4E). Because the  
 195 ribosomal DNA region including 18S and ITS2 is highly repetitive, which is on  
 196 average ~4 times more than non-ribosomal genes (Fig 3A), these differences can be  
 197 understood as intra-genomic variations within the genome. However, taxonomists and  
 198 ecologists use differences within 18S gene sequences as a measure of species  
 199 assignation and to estimate species delineation [6]. This latter practice has been  
 200 previously shown to be very conservative as no differences in the 18S gene were  
 201 found between reproductively isolated species [48]. Alternatively, the possibility of  
 202 sub-populations or cryptic populations cannot be ignored, as previously reported in  
 203 planktonic foraminifers [49] and coccolithophores [50].

204 We examined the possible presence of sub-populations on 18S gene heterozygosity in  
 205 the Pt8 accession. In particular, we confirmed the expression of all the heterozygous  
 206 alleles within the 18S rDNA gene using whole genome and total-RNA sequencing of  
 207 a monoclonal culture (propagated from a single cell) from Pt8 population (constituent  
 208 of Clade D), referred to as Pt8Tc (Fig S4D), indicating that the cultures are a single  
 209 population.



210 Next, concerning the observed polymorphisms within the 18S ribosomal marker gene,  
 211 we investigated whether the four clades can be considered as different species. We  
 212 looked for the existence of compensatory base changes (CBCs) within secondary  
 213 structures of the ITS2 gene between all pairs of accessions. The presence of CBCs  
 214 within ITS2 has been recently suggested to account for reproductive isolation in  
 215 multiple plant species [51] and between diatom species [52, 53]. By comparing the  
 216 ITS2 secondary structure from all the accessions, we did not find any CBCs between  
 217 any given pair of accessions (Fig S5). As a control, we compared the ITS2 secondary  
 218 structure of all the *P. tricornutum* accessions with the ITS2 sequences of other diatom  
 219 species (*Cyclotella meneghiniana*, *Pseudo-nitzschia delicatissima*, *Pseudo-nitzschia*  
 220 *multiseries*, *Fragilariopsis cylindrus*) that have significant degrees of evolutionary  
 221 divergence as depicted previously using multiple molecular marker genes [20, 54],  
 222 and found multiple CBCs in them (Fig S5).

223

## 224 **Close genetic relatedness depicted by large structural genomic variations among** 225 **accessions**

226 Next, using a normalized measure of read depth (see Materials and Methods), we  
 227 found that 259 and 590 genes, representing ~2% and ~5% of the total gene content,  
 228 respectively, have been lost or exhibit copy number variation (CNV), across the ten  
 229 accessions (Fig 3A, 3B) (File S1). Multiple randomly chosen loci were also validated  
 230 by PCR for their loss from certain accessions compared to the reference strain Pt1 8.6  
 231 (Fig S6). Compared to the reference, approximately 70% of the genes that are either  
 232 lost or show CNV are shared among multiple accessions with an exception of Pt10,  
 233 which displays the maximum number of lost genes and accession-specific genes  
 234 exhibiting CNV (Fig 3B). In addition, we detected 207 TEs (~6% of the total

235 annotated TEs) (File S2) showing CNVs across one or more accessions (Fig 3C, 3D),  
 236 80% of which are shared among two or more accessions, with Pt10 again possessing  
 237 the maximum number of accession-specific TEs exhibiting CNVs (Fig 3C). Not  
 238 surprisingly, across all the accessions, class I-type TEs, which undergo transposition  
 239 via a copy-and-paste mechanism, show more variation in the estimated number of  
 240 copies than class II-type TEs (Fig 3D, S7) that are transposed by a cut-and-paste  
 241 mechanism. Euclidean distance estimated between accessions, based on the variation  
 242 in the number of copies of different genes and TEs displaying CNVs, followed by  
 243 hierarchical clustering, depicted three genetic clusters: Pt1, Pt2, Pt3, Pt9 in cluster1;  
 244 Pt5, Pt10 in cluster 2, and Pt4, Pt6, Pt7, Pt8 in cluster 3 (Fig 3A, 3D). These clusters  
 245 are in broad agreement with the ones described by *Fst*, and indicate the closer genetic  
 246 makeup between accessions within a cluster than between the clusters. Further,  
 247 biological processes can only be traced for ~40% of the genes exhibiting accession-  
 248 specific CNVs. Among all the enriched biological processes (chi-square test,  $P < 0.01$ )  
 249 (File S1), a gene associated to nitrate assimilation (Phatr3\_EG02286) is observed to  
 250 have higher copy number specifically in Pt4. Likewise, each accession can be  
 251 characterized by specific genetic features, represented by ~0.3% to ~28% accession-  
 252 specific CNVs (Fig 3B), possibly linked to the explicit functional behavior of some  
 253 accessions in response to various environmental cues, as reported previously [16-18].

254

## 255 **Functional characterization of the genetic diversity within *P. tricornutum* clades**

256 Species are under continuous pressure to adapt to a changing environment over time.  
 257 We therefore wanted to understand the functional consequences of the genetic  
 258 diversity between the accessions. Localization of the polymorphic sites over genomic  
 259 features (genes, TEs, and intergenic regions) revealed highest rate of variation within

genes (Fig 2C), specifically on exons, and was consistent across all the studied accessions. An average non-synonymous to synonymous variant ratio ( $\pi_N / \pi_S$ ) was estimated to be  $\sim 0.43$ , which is higher than in the *Ostreococcus tauri*,  $\pi_N / \pi_S = 0.2$ [55]. We further identified genes within different phylogenetic clades experiencing different selection pressure based on lowest and highest  $\pi_N / \pi_S$  ratios. Across all the accessions, 241 genes displaying  $\pi_N / \pi_S > 1$  and a higher frequency of non-synonymous as compared to synonymous polymorphism, as expected under balancing selection [56](File S3). Furthermore, many genes (902) were found to have loss-of-function (LoF hereafter) variant alleles (Fig 4A), including frame-shift mutations and mutations leading to theoretical start/stop codon loss and/or gain. Based on the presence of functional domains (Pfam domains), all *P. tricornutum* annotated genes [57] were grouped into 3,020 gene families. These families can be as large as the reverse transcriptase gene family, which is highly abundant in marine plankton [58], representing 149 candidate genes having reverse transcriptase domains, or as small as families that constitute single gene candidates. Across all the accessions, we observed that the majority of genes experiencing LoF mutations belong to large gene families (Fig 4B). This is consistent with a previous observation of the existence of functional redundancy in gene families as a balancing mechanism for null mutations in yeast [59]. Therefore, to estimate an unbiased effect of any evolutionary pressure (LoF allele mutations) on different gene families, we calculated a ratio, termed the effect ratio (EfR, see Materials and Methods), which normalizes that if any gene family has enough candidates to buffer the effect on some genes influencing evolutionary pressure, it will be considered as being less affected compared to those for which all or most of the constituents are under selection pressure. From this analysis, each genetic clade displayed a specific set of gene

285 families as being under selection (Fig 5). Functional characterization of constrained  
 286 genes revealed genes enriched in two families, (1) AAA family proteins that often  
 287 perform chaperone like functions that assist in the assembly or disassembly of  
 288 proteins complexes, protein transport and degradation as well as other functions such  
 289 as replication, recombination, repair and transcription [60], (2) tetratricopeptide-like  
 290 repeats known for their role in a variety of biological processes, such as cell cycle  
 291 regulation, organelle targeting and protein import, vesicle fusion and  
 292 biomineralization [61]. A redox class of enzymes are common to both groups of  
 293 genes and a significant proportion of unknown function proteins is found in the group  
 294 of genes under balancing selection (File S3).

295 Consistent with the population structure, accessions within individual clades are more  
 296 closely related than the accessions belonging to other clades (Fig S8A and S8B),  
 297 suggesting variation in functional relatedness between different proposed  
 298 phylogenetic clades.

### 299 **Selection of *MetH* facilitated methionine biosynthesis over *MetE***

300 Apart from the genetic clade-specific families that are under selection pressure, across  
 301 all the accessions a group of gene families associated with methionine biosynthesis  
 302 (*MetH*, Phatr3\_J23399) was also observed as experiencing higher  $\pi_N / \pi_S$  ratio (Fig 5).  
 303 In *P. tricornutum*, *MetE* (cobalamin-independent methionine synthase) and *MetH*  
 304 (cobalamin-dependent methionine synthase) are known to catalyze conversion of  
 305 homocysteine to methionine in the presence of symbiotic bacteria and vitamin B12 in  
 306 the growth media, respectively. Previous reports have suggested that growing axenic  
 307 cultures in conditions of high cobalamin (vitamin B12) availability results in  
 308 repression of *MetE*, leading to its loss of function and high expression of the *MetH*  
 309 gene in *P. tricornutum* and *C. reinhardtii* [62-64]. In accordance with these results,

we observed a high expression of *MetH* in axenically grown laboratory cultures (Fig 6A) compared to its expression in cells cultured with their natural co-habitant bacteria. However, we were not able to trace any significant signature for the loss of *MetE* gene although its expression is significantly lower in axenic cobalamin-containing cultures (Fig 6B). Similar observations were obtained for *CBA1* and *SHMT* genes (Fig 6C and 6D), which under cobalamin scarcity enhance cobalamin acquisition and manage reduced methionine synthase activity, respectively [63].

## Discussion

Using whole genome sequence analysis of accessions sampled across multiple geographic locations around the world (Fig S1), the aim of this study was to describe the global genetic and functional diversity of the model diatom *Phaeodactylum tricornutum*. By defining a comprehensive landscape of natural variations across multiple accessions, we could investigate genetic structure between *P. tricornutum* populations, and a summary of our results is presented in Fig 7. In order to do so, we first performed reference-based assembly and found consistently high genome coverage (>90%) mapped by sequencing reads from respective accessions, where some accessions have more coverage (>98%, Pt1, Pt2, Pt3, and Pt9) than others (Table 1). This difference is independent of the size of the sequencing library, as it does not correlate with the genome coverage (Table 1), and a portion of unmapped reads is likely a consequence of the incomplete reference genome, which contains several gaps [1]. Additionally, given the redundant nature of unmapped reads together with the fact that the unmapped reference genome is annotated as being rich in TEs (Fig S2), a major portion of unmapped reads likely account for large structural variability within the genomes of individual accessions. This explanation is most clear

335 in Pt10, which is shown to have the largest number of gene losses (Fig 3B) and the  
336 highest number of accession-specific TEs with high copy numbers (Fig 3C), and  
337 covers the least (92%) of the reference genome (Table 1). This suggests the role of  
338 TEs in creating substantial genetic diversity as also shown in many species of plants  
339 and animals [65, 66].

340 Next, based on patterns of variations discovered over the whole genomes and on the  
341 molecular marker genes (18S and ITS2) of all the accessions, and by using various  
342 clustering algorithms (see Results), the ten accessions could be grouped into four  
343 genetic clades. Clade A clusters Pt1, Pt2, Pt3, Pt9; Clade B includes Pt4; Clade C  
344 clusters Pt5, Pt10; and Clade D clusters Pt6, Pt7, Pt8. Most of the structural variants  
345 discovered, both small (SNPs and INDELS) and large (CNV and Gene Loss), are  
346 shared among populations within a clade rather than between clades. This suggests  
347 high intra-clade relatedness over a variety of structural, functional and possibly  
348 ecological traits.

349 *P. tricornutum* is a coastal species with limited dispersal potential, which is consistent  
350 with the reports of its absence in the open ocean from *Tara* Oceans data [6].  
351 Consequently, the Fixation index (*Fst*) between different genetic clades is very high  
352 (0.2 – 0.4), (Fig 1D) confirming the existence of strong population subdivisions into  
353 four genetic clades. As expected for an organism with limited dispersal potential,  
354 most of the populations are structured geographically (Fig 2A), with the exception of  
355 Pt5 and Pt9. In addition, the fact that the subdivisions do not correlate with the  
356 sampling time (Fig. 7, Fig S1), which spans approximately a century, suggests long  
357 and stable genetic populations, which is in line with reports from other diatom species  
358 [40, 41]. Although there exist strong genetic structuring within the accessions, the  
359 average nucleotide diversity ( $\pi$ ), estimated across all the accessions, is remarkably

low (0.2%) compared to the diversity estimates in other unicellular eukaryotes [35, 55, 67-69] but in line with previous estimations in marine phytoplanktonic eukaryotes [70]. The high linkage disequilibrium ( $>0.7$ ) observed across all the accessions (Fig S3B) can be explained by prolonged asexual reproduction [71], a common behavior among diatoms [72].

Given the observation that there exist a large proportion of heterozygous variant alleles (Fig 1C), the high *Fst* between the clades, and the low nucleotide diversity across the accessions, we propose that allele frequency plays a significant role in the genetic differentiation of the clades. The difference in allele frequencies is possibly linked to adaptive selection. This phenomenon has recently been studied in diatoms where allele-specific expression of numerous loci has been demonstrated to be a significant source of adaptive evolution in the cold-adapted diatom species *Fragilariopsis cylindrus* [73]. Furthermore, high proportions of heterozygous variant alleles in some Clades (Clade A, 98%, Fig 7) compared to others (Clade B, 45%, Fig 7) suggests a high selection pressure in the Clade B accession Pt4, which is further supported by the strong signals of balancing selection in Pt4 (Fig 4A, 7). Conversely, the large number of alleles that are deviated from HWE within Clade A member populations (Fig 7) is most likely due to prolonged asexual reproduction, which is also associated with generating high linkage disequilibrium across all the isolates [71]. Asexual reproduction results in higher proportions of divergent alleles within loci with less genetic variation among individuals and a significant deviation from HWE [71]. Therefore, it is possible that Pt4 undergoes sexual reproduction, reasonably often, as it possesses the least number of heterozygous variant alleles, most of which follow HWE. Besides, recent reports suggest induction of sex in diatoms under low nutrient [13] and low light [14] conditions, which resembles the

385 natural niche of Pt4 [17]. Therefore, it would be interesting to explore Pt4 as a model  
386 for investigating sexual phases in *P. tricornutum*.  
387 Interestingly, despite high variability in the levels of heterozygosity between different  
388 accessions (Fig 7), the mutational spectrum, compared to the reference, and across all  
389 the accessions consisted of high G:C → A:T and A:T → G:C transitions (Fig S3A).  
390 Deamination of cytosines dominantly dictates C to T transitions in both plants and  
391 animals [74, 75], and CpG methylation potential of the genome is greatly influenced  
392 by heterozygous SNPs in CpG dinucleotides [76]. Previous studies have demonstrated  
393 low DNA methylation in *P. tricornutum*, using Pt1 8.6, a monoclonal strain  
394 accessioned from a Pt1 single cell as a reference [31, 77]. Because there exist  
395 significant differences in the proportion of heterozygote variant alleles between the  
396 accessions (45% - 98%), testing for DNA methylation patterns across different  
397 accessions may provide an interesting opportunity to dissect cross-talk between loss  
398 of heterozygosity and DNA methylation in the selection of certain traits [78].  
399 High *F<sub>st</sub>*, and yet low nucleotide diversity across all the accessions, suggests some  
400 degree of genetic and functional convergence among the accessions. This can be  
401 explained as a consequence of maintaining all the accessions in lab cultures. The  
402 hypothesis is supported by our observation that the *MetH* gene is under balancing  
403 selection in all the accessions (Fig. 5), due to the functional selection of *MetH*-  
404 dependent biosynthesis of methionine over *MetE* in the presence of high Vitamin B12  
405 in lab-grown cultures [62, 63]. However, such genetic and functional convergence is  
406 limited to certain gene families and pathways, as each clade possesses numerous  
407 clade-specific gene-families that are under balancing selection (Fig 7), possibly linked  
408 to local adaptive traits (Fig 5).



409 It is also worth considering that genetic homogenization, i.e., low nucleotide  
 410 diversity, or high linkage-disequilibrium, across the meta-population, can also be a  
 411 consequence of continuous gene flow between the accessions. However, in the case of  
 412 *P. tricornutum*, gene flow seems limited as highly differentiated populations structure  
 413 geographically, except Pt5 of Clade C and Pt9 of Clade A (Fig. 2A). In addition, *P.*  
 414 *tricornutum* is not known to reproduce sexually, although various components (genes)  
 415 of the meiosis pathway are conserved in *P. tricornutum* as well as in other diatom  
 416 species known to undergo sexual reproduction [79]. Furthermore, the absence of  
 417 contemporary base changes (CBC) within ITS2 secondary structure between all the  
 418 accessions compared to the presence of many CBCs between *P. tricornutum*  
 419 accessions and other diatom species suggests that the accessions may be able to  
 420 exchange genetic material sexually. However, because *P. tricornutum* is a coastal  
 421 diatom with only limited dispersal capacity, which is further supported by its apparent  
 422 absence in the open ocean [6], the possibility of gene flow within different  
 423 populations is likely to be limited at best.

424 The four genetic clades are further supported by functional specialization of grouped  
 425 populations, nicely illustrated with Pt4 in Clade B. Pt4 shows a low non-  
 426 photochemical quenching capacity (NPQ) [17], which has been proposed to be an  
 427 adaptive trait to low light conditions. Specifically, this accession has been proposed to  
 428 establish an up-regulation of a peculiar light harvesting protein LHCX4 in extended  
 429 dark conditions [17, 19]. In line with these observations, a gene involved in nitrate  
 430 assimilation (Phatr3\_EG02286) in Pt4 shows high copy numbers, suggesting an  
 431 altered mode of nutrient acquisition. Nitrate assimilation was shown to be regulated  
 432 extensively under low light or dark conditions to overcome nitrate limitation of  
 433 growth in *Thalassiosira weissflogii* [80]. Pt4 is likely adapted to the low light and

434 highly seasonal environment that characterizes the high latitudes where it was found,  
 435 which may well affect its nitrate assimilation capacity [81, 82]. Additional functions  
 436 emerging from Clade C (Pt5 and Pt10) include vacuolar sorting and vesicle-mediated  
 437 transport gene-families to be under balancing selection, which could be an indication  
 438 of altered intracellular trafficking [83].

439 In conclusion, the study presents pan-genomic diversity of the model diatom *P.*  
 440 *tricornutum*. This is the first study within diatoms that provides a comprehensive  
 441 landscape of diversity at whole genome sequence level and brings new insights to our  
 442 understanding of diatom functional ecology and evolution. Given our observation that  
 443 *P. tricornutum* accessions possess high numbers of heterozygous alleles, it would be  
 444 interesting to think of possible selective functional preferences of one allele over the  
 445 other under different environmental conditions or during the life/cell cycle. In the  
 446 future, such studies could be crucial for deciphering the mechanisms underpinning  
 447 allele divergence and selection within diatoms. Likewise, more than answers, our  
 448 study delivers more questions, which should help address the genetic basis of diatom  
 449 success in diverse ocean ecosystems. Finally, this study provides the community with  
 450 genomic sequences of *P. tricornutum* accessions that can be useful for functional  
 451 studies.

452

## 453 **Experimental procedures**

454

### 455 **Sample preparation, sequencing and mapping**

456 Ten different accessions of *P. tricornutum* were obtained from the culture collections  
 457 of the Provasoli-Guillard National Center for Culture of Marine Phytoplankton  
 458 (CCMP, Pt1=CCMP632, Pt5=CCMP630, Pt6=CCMP631, Pt7=CCMP1327,

Pt9=CCMP633), the Culture Collection of Algae and Protozoa (CCAP, Pt2=CCAP 1052/1A, Pt3= CCAP 1052/1B, Pt4= CCAP 1052/6), the Canadian Center for the Culture of Microorganisms (CCCM, Pt8=NEPCC 640), and the Microalgae Culture Collection of Qingdao University (MACC, Pt10=MACC B228). All of the accessions were grown axenically in batch cultures with a photon fluency rate of 75  $\mu\text{mol photons m}^{-2} \text{ s}^{-1}$  provided by cool-white fluorescent tubes in a 12:12 light: dark (L:D) photoperiod at 20 °C. Exponentially growing cells were harvested and total DNA was extracted with the cetyltrimethylammonium bromide (CTAB) method. At least 6  $\mu\text{g}$  of genomic DNA from each accession was used to construct a sequencing library following the manufacturer's instructions (Illumina Inc.). Paired-end sequencing libraries with a read size of 100 bp and an insert size of approximately 400 bp were sequenced on an Illumina HiSeq 2000 sequencer at Berry Genomics Company (China). The corresponding data can be accessed using bioSample accessions: SAMN08369620, SAMN08369621, SAMN08369622, SAMN08369623, SAMN08369624, SAMN08369625, SAMN08369626, SAMN08369627, SAMN08369628, SAMN08369629. Low quality read-pairs were discarded using FASTQC with a read quality (Phred score) cutoff of 30. Using the genome assembly published in 2008 as reference [1], we performed reference-assisted assembly of all the accessions. We used BOWTIE (-n 2 -X 400) for mapping the high quality NGS reads to the reference genome followed by the processing and filtering of the alignments using SAMTOOLS and BEDTOOLS. Detailed methods are provided in File S5.

## Discovery of small polymorphisms and large structural variants

483 GATK [45], configured for diploid genomes, was used for variant calling, which  
 484 included single nucleotide polymorphisms (SNPs), small insertions and deletions  
 485 ranging between 1 and 300 base pairs (bp). The genotyping mode was kept default  
 486 (genotyping mode = DISCOVERY), Emission confidence threshold (-  
 487 stand\_emit\_conf) was kept 10 and calling confidence threshold (-stand\_call\_conf)  
 488 was kept at 30. The minimum number of reads per base, to be called as a high quality  
 489 SNV, was kept to 4 (read-depth  $\geq 4x$ ). Following this filtration step, the number of  
 490 sites in the protein coding genes covered for all 10 accessions, and therefore callable  
 491 to estimate the genome wide synonymous and non-synonymous polymorphism, added  
 492 up 11.0 Mbp. The average pairwise synonymous and non-synonymous diversity  $\pi_s$   
 493 and  $\pi_N$  [84] were estimated for all genes using in R house scripts.

494 Next, considering Z-score as a normalized measure of read-depth, gene and TE  
 495 candidates showing multiple copies (representing CNV) or apparently being lost  
 496 (representing gene loss) were determined. For TE CNV analysis, TEs that are more  
 497 than 100 bp lengths were considered. We measured the fold-change ( $F_c$ ) by dividing  
 498 normalized read depth per genomic feature (Z-score per gene or TE) by average of  
 499 normalized read depth of all the genes/TEs (average Z-score), per sample. Genes or  
 500 TEs with log2 scaled fold change  $\geq 2$  were reported and considered to exist in more  
 501 than one copy in the genome. Genes where the reads from individual accession  
 502 sequencing library failed to map on the reference genome were considered as  
 503 potentially lost within that accession and reported. Detailed method is provided in File  
 504 S5. Later, some randomly chosen loci were picked and validated for the loss in the  
 505 accessions compared to the reference genome by PCR analysis.

506

# 507 **Validation of gene loss and quantitative PCR analysis**

In order to validate gene loss, DNA was extracted from all the accessions as described previously [21] and PCR was performed with the primers listed in Table S1. PCR products were loaded in 1% agarose gel and after migration gels were exposed to UV light and photographs were taken using a gel documentation apparatus to visualize the presence and absence of amplified fragment. To assess gene expression, RNA was extracted as described in [22] from accessions grown axenically in Artificial Sea Water (ASW) [85] supplemented with vitamins as well as in the presence of their endemic bacteria in ASW without vitamins. qPCR was performed as described previously [22].

517

#### 518 ***P. tricornutum* population structure**

519 **Haplotype analysis:** First, to cluster the accessions as haplogroups, ITS2 gene (chr13: 42150-43145) and 18S gene (chr13: 43553-45338) were used. Polymorphic sites across all the accessions within ITS2 and 18S genes were called and used to generate their corresponding accession specific sequences, which were then aligned using CLUSTALW. The same approach was employed to perform haplotype analysis at the whole genome scale. Later, a maximum likelihood algorithm was used to generate the 18S, ITS2 and, whole genome tree with bootstrap values of 1,000. We used MEGA7 [86] to align and deduce the phylogenetic trees.

527 **CBC analysis:** CBC analysis was done by generating the secondary structure of ITS2 sequences, using RNAfold [87], across all *P. tricornutum* accessions and other diatom species. The other species include one centric diatom species *Cyclotella meneghiniana* (AY906805.1), and three pennate diatoms *Pseudo-nitzschia delicatissima* (EU478789.1), *Pseudo-nitzschia multiseriata* (DQ062664.1), *Fragilariopsis cylindrus* (EF660056.1). The centroid secondary structures of ITS2

gene with lowest minimum free energy were used for CBC analysis. We used 4SALE [88] for estimating the presence of CBCs between the secondary structure of ITS2 gene across all the species.

**Population genetics:** Further, we measured various population genetic functions to estimate the effect of evolutionary pressure in shaping the diversity and resemblance between different accession populations. Within individual accessions, by using approximate allelic depths of reference/alternate alleles, we calculated the alleles that are deviated from Hardy Weinberg equilibrium (HWE). We used chi-square estimation to evaluate alleles observed to deviate significantly ( $P$ -value  $< 0.05$ ) from the expected proportion as per  $p^2$  (homozygous) +  $2pq$  (heterozygous) +  $q^2$  (homozygous) = 1) and should be 0.25% + 0.50% + 0.25%. Alleles were considered heterozygous if the proportion of ref/alt allele is between 20-80%. The proportion of ref/alt allele was calculated by dividing the number of reads supporting ref/alt base change by total number of reads mapped at the position. We evaluated average  $R^2$  as a function to measure the linkage disequilibrium with increasing distance (1 kb, 5 kb, 10 kb, 20 kb, 30 kb, 40 kb and 50 kb) between any given pair of mutant alleles across all the accessions using expectation-maximization (EM) algorithm deployed in the VCFtools. Although no recombination was observed within the accessions, attempts were made to look for recombination signals using LDhat [89] and RAT [90]. Genetic differentiation or variability between the accessions was further assessed using the mathematical function of Fixation index ( $F_{ST}$ ), as described by Weir and Cockerham 1984 [91].

**Genetic clustering:** Genetic clustering of the accessions was done using Bayesian clustering approach by applying *Markov Chain Monte Carlo* (MCMC) estimation programmed within ADMIXTURE (version linux-1.3.0) [47]. Accessory tools like

558 PLINK (version 1.07-x86\_64) [92] and VCFtools (version 0.1.13) [93] were used to  
 559 format the VCF files to ADMIXTURE accepted formats. In the absence of data from  
 560 individuals of each accession/sample, we assumed the behavior of each individual in a  
 561 sample to be coherent. Conclusively, instead of estimating the genetic structure within  
 562 an accession, we compared it across all the accessions. We first estimated the possible  
 563 clusters of genomes, (K), across all the accessions, by using cross-validation error  
 564 (CV error) function of ADMIXTURE [94]. Finally, we used ADMIXTURE with 200  
 565 bootstraps, to estimate the genome clusters within individual accessions by  
 566 considering the possible number of genomes derived via CV-error function.

567

# **568 Functional characterization of polymorphisms**

569 snpEff [95] and KaKs [96] calculator were used to annotate the functional nature of  
 570 the polymorphisms. Along with the non-synonymous, synonymous, loss-of-function  
 571 (LOF) alleles, transition to transversion ratio and mutational spectrum of the single  
 572 nucleotide polymorphisms were also measured.  $\pi_N/\pi_S$  ratios were calculated for  
 573 5232 protein coding genes containing more than 10 SNP. 10% of genes with lower  
 574  $\pi_N/\pi_S$  were considered as under strong purifying selection on amino-acid  
 575 composition (File S3). Genes with  $\pi_N/\pi_S$  higher than 1 and average frequency on non-  
 576 synonymous polymorphism higher than the average frequency of synonymous  
 577 polymorphism were considered as candidate genes under balancing selection on  
 578 amino-acid composition (File S3). Various in-house scripts were also used at different  
 579 levels for analysis and for plotting graphs. Data visualization and graphical analysis  
 580 were performed principally using ClicO [97], CYTOSCAPE [98], IGV [99] and R  
 581 (<https://www.r-project.org/about.html>). Based on the presence of functional domains  
 582 all the Phatr3 genes [57] were grouped into 3,020 gene families. Subsequently, the

583 constituents of each gene family were checked for being either affected by loss-of-  
 584 function mutations or experiencing balancing selection. To estimate an unbiased  
 585 effect of any evolutionary pressure (LoF allele or balancing selection mutations) on  
 586 different gene families, induced because of high functional redundancies in the gene  
 587 families, a normalized ratio named as effect ratio (EfR), was calculated. Precisely, the  
 588 EfR normalizes the fact that if any gene family have enough candidates to buffer the  
 589 effect on some genes influencing evolutionary pressures, it will be considered as less  
 590 affected compared to the situation where all or most of the constituents are under  
 591 selection pressure. The ratio was estimated as shown below and gene families with  
 592 EfR larger than 1 were considered as being significantly affected.

*Effect Ratio (EfR)*

$$= \frac{\frac{\text{Number of genes affected within the given gene family}}{\text{Total number of genes in the given gene family}}}{\frac{\text{Total number of genes affected in all the gene families}}{\text{Total number of genes in all the gene families}}}$$

593  
 594 Additionally, significantly enriched (chi-square test, P-value < 0.05) biological  
 595 processes associated within genes experiencing LoF mutations, purifying selection,  
 596 balancing selection (BS), or showing CNV, or being lost (GnL), were estimated by  
 597 calculating observed to expected ratio of their percent occurrence within the given  
 598 functional set (BS, LoF, CNV) and their occurrence in the complete annotated Phatr3  
 599 ([http://protists.ensembl.org/Phaeodactylum\\_tricornutum/Info/Index](http://protists.ensembl.org/Phaeodactylum_tricornutum/Info/Index)) biological  
 600 process catalog. Later, considering gene family EfR as a function to measure the  
 601 association rate, we deduced Pearson pairwise correlations between different  
 602 accessions. The correlation matrix describes that if many equally affected gene  
 603 families are shared between any given pair of accessions, they will have higher



604 correlation compared to others. Finally, hierarchical clustering using Pearson pairwise  
605 correlation matrix assessed the association between the accessions.

606

## 607 **Acknowledgements**

608 HH acknowledges support from National Natural Science Foundation of China (grant  
609 No. 91751117). GW acknowledges the Strategic Priority Research Program of the  
610 Chinese Academy of Sciences (grant No. XDA17010502). CB acknowledges funding  
611 from the ERC Advanced Award ‘Diatomite’, the LouisD Foundation of the Institut de  
612 France, the Gordon and Betty Moore Foundation, and the French Government  
613 ‘Investissements d’Avenir’ programmes MEMO LIFE (ANR-10-LABX-54), PSL\*  
614 Research University (ANR-1253 11-IDEX-0001-02), and OCEANOMICS (ANR-11-  
615 BTBR-0008). CB also thanks the Radcliffe Institute of Advanced Study at Harvard  
616 University for a scholar’s fellowship during the 2016-2017 academic year. LT  
617 acknowledges funds from the CNRS and MEMO LIFE (ANR-10-LABX-54). AR was  
618 supported by an International PhD fellowship from MEMO LIFE (ANR-10-LABX-  
619 54).

620

## 621 **Conflict of interest**

622 The authors declare no conflicts of interest.

623

## 624 **References**

625

- 626 1. Bowler C, Allen AE, Badger JH, Grimwood J, Jabbari K, Kuo A, et al: **The**  
627 **Phaeodactylum genome reveals the evolutionary history of diatom**  
628 **genomes. *Nature* 2008, **456**:239-244.**
- 629 2. Moustafa A, Beszteri B, Maier UG, Bowler C, Valentin K, Bhattacharya D:  
630 **Genomic footprints of a cryptic plastid endosymbiosis in diatoms. *Science***  
631 **2009, **324**:1724-1726.**

- 632 3. Dorrell RG, Gile G, McCallum G, Meheust R, Baptiste EP, Klinger CM, et al:  
633 **Chimeric origins of ochrophytes and haptophytes revealed through an**  
634 **ancient plastid proteome.** *Elife* 2017, **6**.
- 635 4. Armbrust EV: **The life of diatoms in the world's oceans.** *Nature* 2009,  
636 **459**:185-192.
- 637 5. Amin SA, Parker MS, Armbrust EV: **Interactions between diatoms and**  
638 **bacteria.** *Microbiol Mol Biol Rev* 2012, **76**:667-684.
- 639 6. Malviya S, Scalco E, Audic S, Vincent F, Veluchamy A, Poulain J, Wincker P,  
640 Iudicone D, de Vargas C, Bittner L, et al: **Insights into global diatom**  
641 **distribution and diversity in the world's ocean.** *Proc Natl Acad Sci U S A*  
642 2016.
- 643 7. Allen AE, Dupont CL, Obornik M, Horak A, Nunes-Nesi A, McCrow JP, et al:  
644 **Evolution and metabolic significance of the urea cycle in photosynthetic**  
645 **diatoms.** *Nature* 2011, **473**:203-207.
- 646 8. Huysman MJ, Fortunato AE, Matthijs M, Costa BS, Vanderhaeghen R, Van den  
647 Daele H, et al: **AUREOCHROME1a-mediated induction of the diatom-specific**  
648 **cyclin dsCYC2 controls the onset of cell division in diatoms (Phaeodactylum**  
649 **tricornutum).** *Plant Cell* 2013, **25**:215-228.
- 650 9. Morrissey J, Sutak R, Paz-Yepes J, Tanaka A, Moustafa A, Veluchamy A, et al:  
651 **A novel protein, ubiquitous in marine phytoplankton, concentrates iron at**  
652 **the cell surface and facilitates uptake.** *Curr Biol* 2015, **25**:364-371.
- 653 10. Tanaka T, Maeda Y, Veluchamy A, Tanaka M, Abida H, Marechal E, et al: **Oil**  
654 **accumulation by the oleaginous diatom *Fistulifera solaris* as revealed by the**  
655 **genome and transcriptome.** *Plant Cell* 2015, **27**:162-176.
- 656 11. Fortunato AE, Jaubert M, Enomoto G, Bouly JP, Raniello R, Thaler et al:  
657 **Diatom phytochromes reveal the existence of far-red-light-based sensing in**  
658 **the ocean.** *Plant Cell* 2016, **28**:616-628.
- 659 12. Godhe A, Kremp A, Montresor M: **Genetic and microscopic evidence for**  
660 **sexual reproduction in the centric diatom *Skeletonema marinoi*.** *Protist*  
661 2014, **165**:401-416.
- 662 13. Moore ER, Bullington BS, Weisberg AJ, Jiang Y, Chang J, Halsey KH:  
663 **Morphological and transcriptomic evidence for ammonium induction of**  
664 **sexual reproduction in *Thalassiosira pseudonana* and other centric diatoms.**  
665 *PLoS One* 2017, **12**:e0181098.
- 666 14. Mouget JL, Gastineau R, Davidovich O, Gaudin P, Davidovich NA: **Light is a**  
667 **key factor in triggering sexual reproduction in the pennate diatom *Haslea***  
668 ***ostrearia*.** *FEMS Microbiol Ecol* 2009, **69**:194-201.
- 669 15. De Martino AM, A. Juan Shi, K.P. Bowler, C.: **Genetic and phenotypic**  
670 **characterization of *Phaeodactylum tricornutum* (Bacillariophyceae)**  
671 **accessions.** *J Phycol* 2007, **43**:992–1009.
- 672 16. Stanley MS, Callow JA: **Whole cell adhesion strength of morphotypes and**  
673 **isolates of *Phaeodactylum tricornutum* (Bacillariophyceae).** *European*  
674 *Journal of Phycology* 2007, **42**:191-197.
- 675 17. Bailleul B, Rogato A, de Martino A, Coesel S, Cardol P, Bowler C, et al: **An**  
676 **atypical member of the light-harvesting complex stress-related protein**  
677 **family modulates diatom responses to light.** *Proc Natl Acad Sci U S A* 2010,  
678 **107**:18214-18219.

- 679 18. Abida H, Dolch LJ, Mei C, Villanova V, Conte M, Block MA, et al: **Membrane**  
680 **glycerolipid remodeling triggered by nitrogen and phosphorus starvation in**  
681 **Phaeodactylum tricornutum**. *Plant Physiol* 2015, **167**:118-136.
- 682 19. Taddei L, Stella GR, Rogato A, Bailleul B, Fortunato AE, Annunziata R, et al:  
683 **Multisignal control of expression of the LHCX protein family in the marine**  
684 **diatom Phaeodactylum tricornutum**. *J Exp Bot* 2016, **67**:3939-3951.
- 685 20. Tirichine L, Rastogi A, Bowler C: **Recent progress in diatom genomics and**  
686 **epigenomics**. *Curr Opin Plant Biol* 2017, **36**:46-55.
- 687 21. Falciatore A, Casotti R, Leblanc C, Abrescia C, Bowler C: **Transformation of**  
688 **nonselectable reporter genes in marine diatoms**. *Mar Biotechnol (NY)* 1999,  
689 **1**:239-251.
- 690 22. Siaux M, Heijde M, Mangogna M, Montsant A, Coesel S, Allen A, et al:  
691 **Molecular toolbox for studying diatom biology in Phaeodactylum**  
692 **tricornutum**. *Gene* 2007, **406**:23-35.
- 693 23. De Riso V, Raniello R, Maumus F, Rogato A, Bowler C, Falciatore A: **Gene**  
694 **silencing in the marine diatom Phaeodactylum tricornutum**. *Nucleic Acids*  
695 *Res* 2009, **37**:e96.
- 696 24. Huysman MJ, Martens C, Vandepoele K, Gillard J, Rayko E, Heijde M, et al:  
697 **Genome-wide analysis of the diatom cell cycle unveils a novel type of**  
698 **cyclins involved in environmental signaling**. *Genome Biol* 2010, **11**:R17.
- 699 25. Maheswari U, Jabbari K, Petit JL, Porcel BM, Allen AE, Cadoret JP, et al: **Digital**  
700 **expression profiling of novel diatom transcripts provides insight into their**  
701 **biological functions**. *Genome Biol* 2010, **11**:R85.
- 702 26. Maheswari U, Mock T, Armbrust EV, Bowler C: **Update of the Diatom EST**  
703 **Database: a new tool for digital transcriptomics**. *Nucleic Acids Res* 2009,  
704 **37**:D1001-1005.
- 705 27. Kaur S, Spillane C: **Reduction in carotenoid levels in the marine diatom**  
706 **Phaeodactylum tricornutum by artificial microRNAs targeted against the**  
707 **endogenous phytoene synthase gene**. *Mar Biotechnol (NY)* 2015, **17**:1-7.
- 708 28. Diner RE, Bielinski VA, Dupont CL, Allen AE, Weyman PD: **Refinement of the**  
709 **Diatom Episome Maintenance Sequence and Improvement of Conjugation-**  
710 **Based DNA Delivery Methods**. *Front Bioeng Biotechnol* 2016, **4**:65.
- 711 29. Nymark M, Sharma AK, Sparstad T, Bones AM, Winge P: **A CRISPR/Cas9**  
712 **system adapted for gene editing in marine algae**. *Sci Rep* 2016, **6**:24951.
- 713 30. Rastogi A, Murik O, Bowler C, Tirichine L: **PhytoCRISP-Ex: a web-based and**  
714 **stand-alone application to find specific target sequences for CRISPR/CAS**  
715 **editing**. *BMC Bioinformatics* 2016, **17**:261.
- 716 31. Veluchamy A, Lin X, Maumus F, Rivarola M, Bhavsar J, Creasy T, et al: **Insights**  
717 **into the role of DNA methylation in diatoms by genome-wide profiling in**  
718 **Phaeodactylum tricornutum**. *Nat Commun* 2013, **4**.
- 719 32. Veluchamy A, Rastogi A, Lin X, Lombard B, Murik O, Thomas Y, et al: **An**  
720 **integrative analysis of post-translational histone modifications in the**  
721 **marine diatom Phaeodactylum tricornutum**. *Genome Biol* 2015, **16**:102.
- 722 33. Daboussi F, Leduc S, Marechal A, Dubois G, Guyot V, Perez-Michaut C, et al:  
723 **Genome engineering empowers the diatom Phaeodactylum tricornutum for**  
724 **biotechnology**. *Nat Commun* 2014, **5**:3831.

- 725 34. Serif M, Dubois G, Finoux AL, Teste MA, Jallet D, Daboussi F: **One-step**  
726 **generation of multiple gene knock-outs in the diatom *Phaeodactylum***  
727 ***tricornutum* by DNA-free genome editing.** *Nat Commun* 2018, **9**:3924.
- 728 35. Flowers JM, Hazzouri KM, Pham GM, Rosas U, Bahmani T, Khraiweh B, et al:  
729 **Whole-Genome Resequencing Reveals Extensive Natural Variation in the**  
730 **Model Green Alga *Chlamydomonas reinhardtii*.** *Plant Cell* 2015, **27**:2353-  
731 2369.
- 732 36. Cao J, Schneeberger K, Ossowski S, Gunther T, Bender S, Fitz J, et al: **Whole-**  
733 **genome sequencing of multiple *Arabidopsis thaliana* populations.** *Nat*  
734 *Genet* 2011, **43**:956-963.
- 735 37. Liti G, Carter DM, Moses AM, Warringer J, Parts L, James SA, et al: **Population**  
736 **genomics of domestic and wild yeasts.** *Nature* 2009, **458**:337-341.
- 737 38. Lachance J, Tishkoff SA: **Population Genomics of Human Adaptation.** *Annu*  
738 *Rev Ecol Syst* 2013, **44**:123-143.
- 739 39. Godhe A, Ryneerson T: **The role of intraspecific variation in the ecological**  
740 **and evolutionary success of diatoms in changing environments.** *Philos Trans*  
741 *R Soc Lond B Biol Sci* 2017, **372**.
- 742 40. Harnstrom K, Ellegaard M, Andersen TJ, Godhe A: **Hundred years of genetic**  
743 **structure in a sediment revived diatom population.** *Proc Natl Acad Sci U S A*  
744 2011, **108**:4252-4257.
- 745 41. Whittaker KA, Ryneerson TA: **Evidence for environmental and ecological**  
746 **selection in a microbe with no geographic limits to gene flow.** *Proc Natl*  
747 *Acad Sci U S A* 2017, **114**:2651-2656.
- 748 42. Rengefors K, Kremp A, Thorsten BH, Reusch A, Wood M: **Genetic diversity**  
749 **and evolution in eukaryotic phytoplankton: revelations from population**  
750 **genetic studies.** *Journal of plankton research* 2017, **39**:165-179.
- 751 43. Matuszewski S, Hermisson J, Kopp M: **Catch Me if You Can: Adaptation from**  
752 **Standing Genetic Variation to a Moving Phenotypic Optimum.** *Genetics*  
753 2015, **200**:1255-1274.
- 754 44. Bailleul B, Berne N, Murik O, Petroutsos D, Prihoda J, Tanaka A, et al:  
755 **Energetic coupling between plastids and mitochondria drives CO2**  
756 **assimilation in diatoms.** *Nature* 2015, **524**:366-369.
- 757 45. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al:  
758 **The Genome Analysis Toolkit: a MapReduce framework for analyzing next-**  
759 **generation DNA sequencing data.** *Genome Res* 2010, **20**:1297-1303.
- 760 46. Chen G, Ryneerson TA: **Genetically distinct populations of a diatom co-exist**  
761 **during the North Atlantic spring bloom.** *LIMNOLOGY and OCEANOGRAPHY*  
762 2016:2165-2179.
- 763 47. Alexander DH, Novembre J, Lange K: **Fast model-based estimation of**  
764 **ancestry in unrelated individuals.** *Genome Res* 2009, **19**:1655-1664.
- 765 48. Piganeau G, Eyre-Walker A, Jancek S, Grimsley N, Moreau H: **How and why**  
766 **DNA barcodes underestimate the diversity of microbial eukaryotes.** *PLoS*  
767 *One* 2011, **6**:e16342.
- 768 49. de Vargas C, Norris R, Zaninetti L, Gibb SW, Pawlowski J: **Molecular evidence**  
769 **of cryptic speciation in planktonic foraminifers and their relation to oceanic**  
770 **provinces.** *Proc Natl Acad Sci U S A* 1999, **96**:2864-2868.

- 771 50. Saez AG, Probert I, Geisen M, Quinn P, Young JR, Medlin LK: **Pseudo-cryptic**  
772 **speciation in coccolithophores.** *Proc Natl Acad Sci U S A* 2003, **100**:7163-  
773 7168.
- 774 51. Wolf M, Chen S, Song J, Ankenbrand M, Muller T: **Compensatory base**  
775 **changes in ITS2 secondary structures correlate with the biological species**  
776 **concept despite intragenomic variability in ITS2 sequences--a proof of**  
777 **concept.** *PLoS One* 2013, **8**:e66726.
- 778 52. Kaczmarek I, Mather L, Luddington I, Muise F, Ehrman J: **Cryptic diversity in**  
779 **a cosmopolitan diatom known as *Asterionellopsis glacialis* (Fragilariaceae):**  
780 **Implications for ecology, biogeography, and taxonomy.** *American Journal of*  
781 *Botany* 2014.
- 782 53. Amato A, Kooistra WH, Ghiron JH, Mann DG, Proschold T, Montresor M:  
783 **Reproductive isolation among sympatric cryptic species in marine diatoms.**  
784 *Protist* 2007, **158**:193-207.
- 785 54. Medlin LK: **A timescale for diatom evolution based on four molecular**  
786 **markers: reassessment of ghost lineages and major steps defining diatom**  
787 **evolution.** *Vie Milieu / Life & Environment* 2015.
- 788 55. Blanc-Mathieu R, Krasovec M, Hebrard M, Yau S, Desgranges E, Martin J, et  
789 al: **Population genomics of picophytoplankton unveils novel chromosome**  
790 **hypervariability.** *Sci Adv* 2017, **3**:e1700239.
- 791 56. Bitarello BD, de Filippo C, Teixeira JC, Schmidt JM, Kleinert P, Meyer D, et al:  
792 **Signatures of Long-Term Balancing Selection in Human Genomes.** *Genome*  
793 *Biol Evol* 2018, **10**:939-955.
- 794 57. Rastogi A, Maheswari U, Dorrell RG, Vieira FRJ, Maumus F, Kustka A, et al:  
795 **Integrative analysis of large scale transcriptome data draws a**  
796 **comprehensive landscape of *Phaeodactylum tricornutum* genome and**  
797 **evolutionary origin of diatoms.** *Sci Rep* 2018, **8**:4834.
- 798 58. Lescot M, Hingamp P, Kojima KK, Villar E, Romac S, Veluchamy A, et al:  
799 **Reverse transcriptase genes are highly abundant and transcriptionally**  
800 **active in marine plankton assemblages.** *ISME J* 2016, **10**:1134-1146.
- 801 59. Gu Z, Steinmetz LM, Gu X, Scharfe C, Davis RW, Li WH: **Role of duplicate**  
802 **genes in genetic robustness against null mutations.** *Nature* 2003, **421**:63-66.
- 803 60. Ogura T, Wilkinson AJ: **AAA+ superfamily ATPases: common structure--**  
804 **diverse function.** *Genes Cells* 2001, **6**:575-597.
- 805 61. Zeytuni N, Zarivach R: **Structural and functional discussion of the tetra-trico-**  
806 **peptide repeat, a protein interaction module.** *Structure* 2012, **20**:397-405.
- 807 62. Helliwell KE, Wheeler GL, Leptos KC, Goldstein RE, Smith AG: **Insights into the**  
808 **evolution of vitamin B12 auxotrophy from sequenced algal genomes.** *Mol*  
809 *Biol Evol* 2011, **28**:2921-2933.
- 810 63. Bertrand EM, Allen AE, Dupont CL, Norden-Krichmar TM, Bai J, Valas RE, et al:  
811 **Influence of cobalamin scarcity on diatom molecular physiology and**  
812 **identification of a cobalamin acquisition protein.** *Proc Natl Acad Sci U S A*  
813 2012, **109**:E1762-1771.
- 814 64. Helliwell KE, Collins S, Kazamia E, Purton S, Wheeler GL, Smith AG:  
815 **Fundamental shift in vitamin B12 eco-physiology of a model alga**  
816 **demonstrated by experimental evolution.** *ISME J* 2015, **9**:1446-1455.



- 817 65. Quadrana L, Bortolini Silveira A, Mayhew GF, LeBlanc C, Martienssen RA,  
818 Jeddelloh JA, et al: **The Arabidopsis thaliana mobilome and its impact at the**  
819 **species level.** *Elife* 2016, **5**.
- 820 66. Bonchev G, Parisod C: **Transposable elements and microevolutionary**  
821 **changes in natural populations.** *Mol Ecol Resour* 2013, **13**:765-775.
- 822 67. Liti G: **The fascinating and secret wild life of the budding yeast *S. cerevisiae*.**  
823 *Elife* 2015, **4**.
- 824 68. Blanc-Mathieu R, Verhelst B, Derelle E, Rombauts S, Bouget FY, Carre I, et al:  
825 **An improved genome of the model marine alga *Ostreococcus tauri* unfolds**  
826 **by assessing Illumina de novo assemblies.** *BMC Genomics* 2014, **15**:1103.
- 827 69. Hirakawa MP, Martinez DA, Sakthikumar S, Anderson MZ, Berlin A, Gujja S, et  
828 al: **Genetic and phenotypic intra-species variation in *Candida albicans*.**  
829 *Genome Res* 2015, **25**:413-425.
- 830 70. Filatov DA: **Extreme Lewontin's Paradox in Ubiquitous Marine**  
831 **Phytoplankton Species.** *Mol Biol Evol* 2019, **36**:4-14.
- 832 71. Allen DE, Lynch M: **The effect of variable frequency of sexual reproduction**  
833 **on the genetic structure of natural populations of a cyclical parthenogen.**  
834 *Evolution* 2012, **66**:919-926.
- 835 72. Koester JA, Berthiaume CT, Hiranuma N, Parker MS, Iverson V, Morales R, et  
836 al: **Sexual ancestors generated an obligate asexual and globally dispersed**  
837 **clone within the model diatom species *Thalassiosira pseudonana*.** *Sci Rep*  
838 2018, **8**:10492.
- 839 73. Mock T, Otilar RP, Strauss J, McMullan M, Paajanen P, Schmutz J, et al:  
840 **Evolutionary genomics of the cold-adapted diatom *Fragilariopsis cylindrus*.**  
841 *Nature* 2017, **541**:536-540.
- 842 74. Becker C, Hagmann J, Muller J, Koenig D, Stegle O, Borgwardt K, et al:  
843 **Spontaneous epigenetic variation in the *Arabidopsis thaliana* methylome.**  
844 *Nature* 2011, **480**:245-249.
- 845 75. Rahbari R, Wuster A, Lindsay SJ, Hardwick RJ, Alexandrov LB, Al Turki S, et al:  
846 **Timing, rates and spectra of human germline mutation.** *Nat Genet* 2016,  
847 **48**:126-133.
- 848 76. Shoemaker R, Deng J, Wang W, Zhang K: **Allele-specific methylation is**  
849 **prevalent and is contributed by CpG-SNPs in the human genome.** *Genome*  
850 *Res* 2010, **20**:883-889.
- 851 77. Huff JT, Zilberman D: **Dnmt1-independent CG methylation contributes to**  
852 **nucleosome positioning in diverse eukaryotes.** *Cell* 2014, **156**:1286-1297.
- 853 78. Kanai Y, Ushijima S, Tsuda H, Sakamoto M, Hirohashi S: **Aberrant DNA**  
854 **methylation precedes loss of heterozygosity on chromosome 16 in chronic**  
855 **hepatitis and liver cirrhosis.** *Cancer Lett* 2000, **148**:73-80.
- 856 79. Patil S, Moeys S, von Dassow P, Huysman MJ, Mapleson D, De Veylder L, et al:  
857 **Identification of the meiotic toolkit in diatoms and exploration of meiosis-**  
858 **specific SPO11 and RAD51 homologs in the sexual species *Pseudo-nitzschia***  
859 **multistriata and *Seminavis robusta*.** *BMC Genomics* 2015, **16**:930.
- 860 80. Clark DP, Flynn KJ, Ownes NJ: **The large capacity for dark nitrate-assimilation**  
861 **in diatoms may overcome nitrate limitation of growth.** *New Phytologist*  
862 2002.

- 863 81. Ivanikova NV, McKay R, Bullerjahn GS: **Construction and characterization of a**  
864 **cyanobacterial bioreporter capable of assessing nitrate assimilatory capacity**  
865 **in freshwaters.** *Limnology and Oceanography* 2005, **3**:86-93.
- 866 82. Weiguo Li JW: **Influence of light and nitrate assimilation on the growth**  
867 **strategy in clonal weed *Eichhornia crassipes*.** *Aquatic Ecology* 2011.
- 868 83. Pickett-Heaps JD, Forer A: **Pac-Man does not resolve the enduring problem**  
869 **of anaphase chromosome movement.** *Protoplasma* 2001, **215**:16-20.
- 870 84. Nei M, Li WH: **Mathematical model for studying genetic variation in terms**  
871 **of restriction endonucleases.** *Proc Natl Acad Sci U S A* 1979, **76**:5269-5273.
- 872 85. Vartanian M, Descles J, Quinet M, Douady S, Lopez PJ: **Plasticity and**  
873 **robustness of pattern formation in the model diatom *Phaeodactylum***  
874 ***tricornutum*.** *The New phytologist* 2009, **182**:429-442.
- 875 86. Kumar S, Stecher G, Tamura K: **MEGA7: Molecular Evolutionary Genetics**  
876 **Analysis Version 7.0 for Bigger Datasets.** *Mol Biol Evol* 2016, **33**:1870-1874.
- 877 87. Lorenz R, Bernhart SH, Honer Zu, Siederdisen C, Tafer H, Flamm C, et al:  
878 **ViennaRNA Package 2.0.** *Algorithms Mol Biol* 2011, **6**:26.
- 879 88. Seibel PN, Muller T, Dandekar T, Schultz J, Wolf M: **4SALE--a tool for**  
880 **synchronous RNA sequence and secondary structure alignment and editing.**  
881 *BMC Bioinformatics* 2006, **7**:498.
- 882 89. Auton A, McVean G: **Recombination rate estimation in the presence of**  
883 **hotspots.** *Genome Res* 2007, **17**:1219-1227.
- 884 90. Etherington GJ, Dicks J, Roberts IN: **Recombination Analysis Tool (RAT): a**  
885 **program for the high-throughput detection of recombination.**  
886 *Bioinformatics* 2005, **21**:278-281.
- 887 91. Weir BS, Cockerham CC: **Estimating F-Statistics for the Analysis of**  
888 **Population Structure.** *Evolution* 1984, **38**:1358-1370.
- 889 92. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al:  
890 **PLINK: a tool set for whole-genome association and population-based**  
891 **linkage analyses.** *Am J Hum Genet* 2007, **81**:559-575.
- 892 93. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al: **The**  
893 **variant call format and VCFtools.** *Bioinformatics* 2011, **27**:2156-2158.
- 894 94. Alexander DH, Lange K: **Enhancements to the ADMIXTURE algorithm for**  
895 **individual ancestry estimation.** *BMC Bioinformatics* 2011, **12**:246.
- 896 95. Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, et al: **A program**  
897 **for annotating and predicting the effects of single nucleotide**  
898 **polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster***  
899 **strain w1118; iso-2; iso-3.** *Fly (Austin)* 2012, **6**:80-92.
- 900 96. Zhang Z, Li J, Zhao XQ, Wang J, Wong GK, Yu J: **KaKs\_Calculator: calculating**  
901 **Ka and Ks through model selection and model averaging.** *Genomics*  
902 *Proteomics Bioinformatics* 2006, **4**:259-263.
- 903 97. Cheong WH, Tan YC, Yap SJ, Ng KP: **ClicO FS: an interactive web-based**  
904 **service of Circos.** *Bioinformatics* 2015, **31**:3685-3687.
- 905 98. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al:  
906 **Cytoscape: a software environment for integrated models of biomolecular**  
907 **interaction networks.** *Genome Res* 2003, **13**:2498-2504.
- 908 99. Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G,  
909 et al: **Integrative genomics viewer.** *Nat Biotechnol* 2011, **29**:24-26.

910  
911  
912



## 913 Table Legends

914 **Table 1. Reference-assisted mapping statistics.** The table summarizes the origin  
915 and year of sampling of each accession of *P. tricornutum* along with the number of  
916 total reads mapped on the reference. Average depth (X=average number of reads  
917 aligned on each base covered across the entire genome) was estimated using the  
918 number of mapped read pairs and the horizontal coverage (aka. coverage breadth)  
919 across the whole genome.

920

921

922

## 923 Figure Legends

924

925 **Figure 1. Genetic diversity between *P. tricornutum* accessions.** (A) The bar plot  
926 represents total number of discovered SNPs, with the proportion of heterozygous  
927 SNPs (dark blue) and homozygous SNPs (light blue), INSERTIONS (orange) and  
928 DELETIONS (yellow) in each accession compared to the reference genome. (B) The  
929 stack bar plot represents the proportion of total vs specific polymorphic variant sites,  
930 including SNPs, insertions and deletions (from left to right, respectively) across all  
931 the accessions. (C) The world map indicates proportion of heterozygous (dark violet)  
932 and homozygous SNPs (violet) in each accession, represented as pie charts. The outer  
933 ring represents the proportion of variant alleles being significantly deviated from  
934 HWE (deep red). (D) The heat-map shows the genetic differentiation or association  
935 between all possible pairs of accessions. The colors indicate  $F_{ST}$  values, which range  
936 from 0.02 to 0.4, with a color gradient from yellow to green, respectively. Values  
937 closer to 0 signify close genetic makeup and values closer to 1 indicate strong genetic  
938 structuring between the populations.

939

940 **Figure 2: Clustering of *P. tricornutum* accessions.** (A) Principal component  
941 analysis (PCA) showing the distribution of the ten accessions based on their shared  
942 genome structure, revealing four genetic clusters referred to as Clades A, B, C, and D.  
943 (B) Pie charts showing the genetic make-up of the genomes of each accession. Each  
944 accession have a distinct genetic makeup, which is shared at both inter and intra Clade  
945 level, and represented with different colors. (C) Phylogenetic association of the  
946 accessions based on 468,188 genome-wide polymorphic sites (including SNP and  
947 INDELS) using a maximum likelihood approach. The numbers on the branches  
948 indicate the bootstrap values. Pie charts adjacent to each node of the whole genome  
949 tree correspond to the proportion of SNPs and INDELS over all functional features of  
950 the genome; GENES (blue), TEs (yellow), IGRs (Intergenic Regions, represented in  
951 grey).

952

953 **Figure 3. Large structural variations within accessions.** (A) The heat-map displays  
954 the fold-change (FC) of read depth between each reference gene and median of read  
955 depth of all the reference genes, within each accession. Using Z-score as a measure of  
956 normalized read depth,  $\log_2$  fold change (FC) is calculated as a ratio of Z-score per  
957 gene to the average normalized read depth of all the genes per accession. A blue to  
958 red color gradient in the heat-map represents low to high  $\log_2$ FC. From all the  
959 accessions only those genes are plotted where  $\log_2$ FC is more than 2 in at least one of  
960 the accessions and are considered to exhibit copy number variation (CNV). (B) The  
961 bar plots represent the total and specific numbers of genes, denoted on Y-axis, that

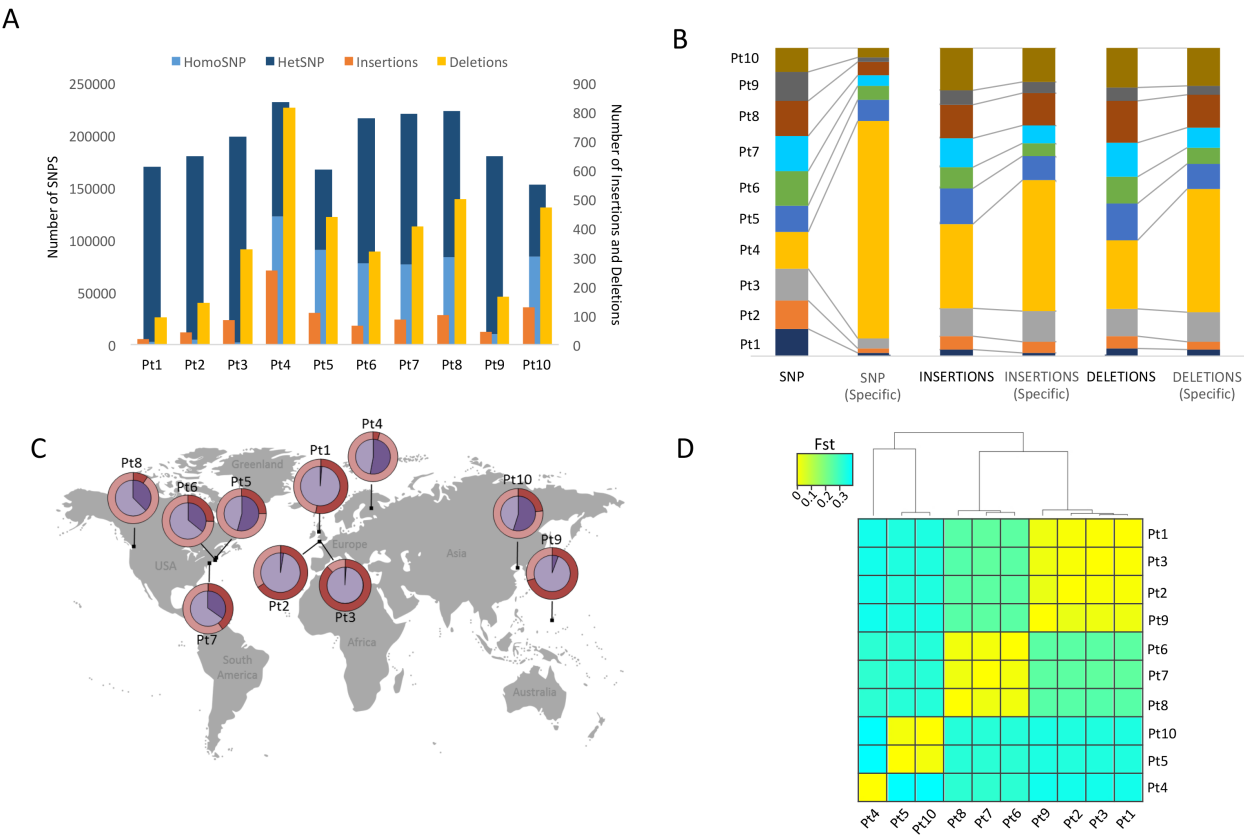
exhibit a loss or multiple copies (CNV) within one or more accessions. (C) The bar plots represent the number of total- and accession-specific TEs exhibiting CNVs across one or more accessions. (D) With similar principle aesthetics as panel (A) of this figure, the heat-map shows the patterns of log2FC only across all the accessions of those TEs exhibiting CNV in at least one of the ten accessions studied.

**Figure 4. Evolutionary and functional consequences of polymorphisms.** (A) The bar plot represents total and specific numbers of genes that are subject to balancing selection, or experiencing loss-of-function (LoF) mutations. For each category, the accessions are plotted as stack plots with total and specific numbers of genes. Numbers of genes in each category are indicated. (B) The box plot represents the number of gene families affected by loss-of-function (LOF) mutations and suggests a bias of such mutations on the genes belonging to large gene families. Y-axis represents, as log scale, the number of genes in the gene families vs those that are not affected by LoF mutations.

**Figure 5. Balancing selection within each genetic Clade.** Based on the EfR metric, the network displays highly affected gene families experiencing balancing selection. Gene families associated with *MetH* genes under balancing selection in all the accessions are indicated within the blue circles. The red circles group individual accessions as Clades.

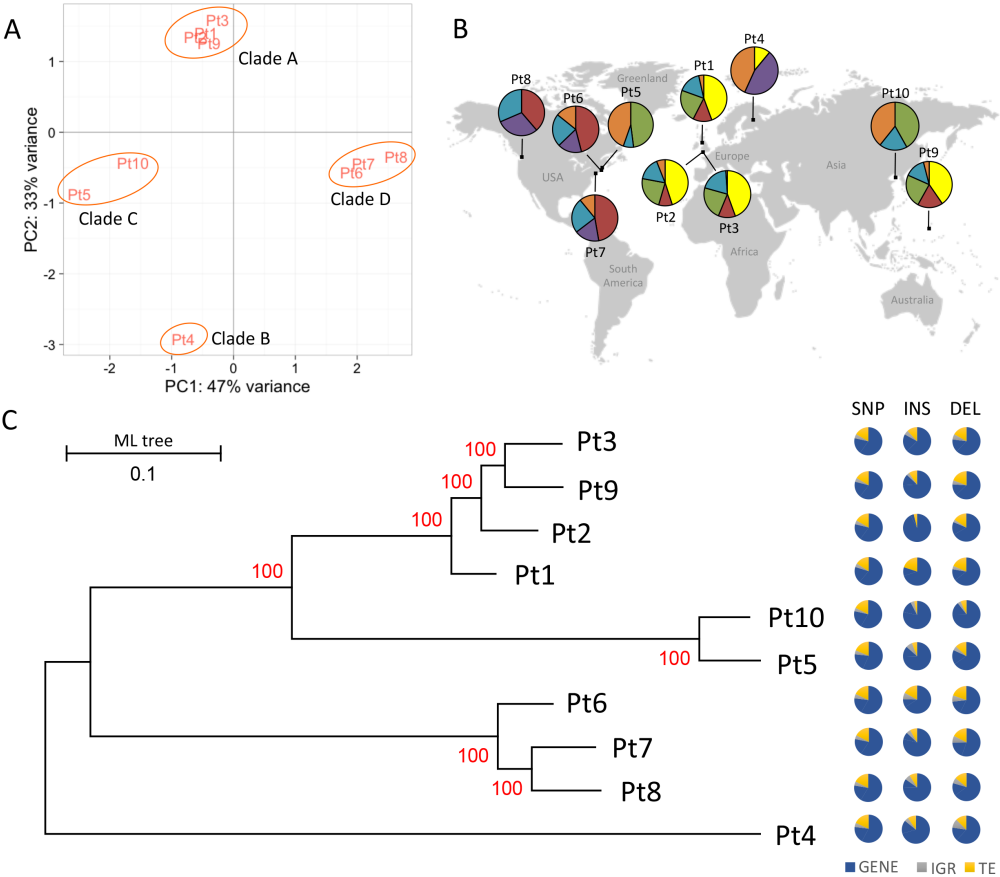
**Figure 6. Selection of MetH-facilitated methionine biosynthesis over MetE.** The bar plots represent relative expression of (A) *MetH*, (B) *MetE*, (C) *CBA1* and (D) *SHMT* genes in four (Pt2, Pt3, Pt4 and Pt8) of the ten accessions with the presence of vitamin B12 in axenic cultures (light gray), and with natural bacteria and no vitamin B12 in the growing media (black).

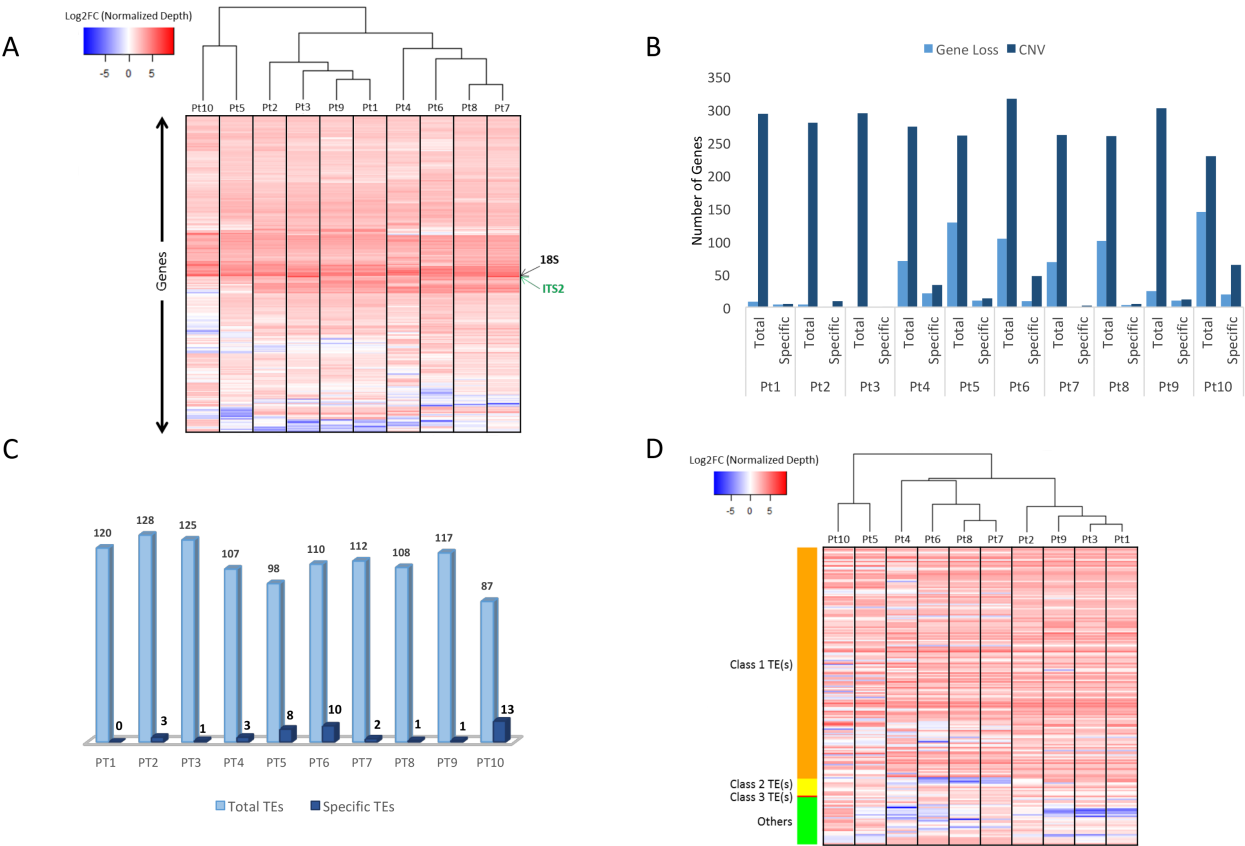
**Figure 7. Population structure of the ecotypes.** The color gradient from yellow to blue indicates low to high numerical values across each ecotype (indicated on top X-axis of panel A) within different functional categories indicated on Y-axis. These functional categories includes (from top to bottom), Year of sampling = Year in which the respective ecotype was sampled, Total SNP = Absolute number of SNPs found in each ecotype, Specific/Total SNP (%) = percentage of ecotype specific SNPs, Heterozygosity = Number of heterozygous SNPs from a set of total SNPs within each ecotype, Homozygosity = Number of homozygous SNPs from a set of total SNPs within respective ecotype, Sites deviated from HWE = Number of SNP sites predicted to be deviated from Hardy-Weinberg equilibrium (HWE), BS = Number of genes under Balancing Selection, LoF = Number of genes localizing Loss of Function variant sites.



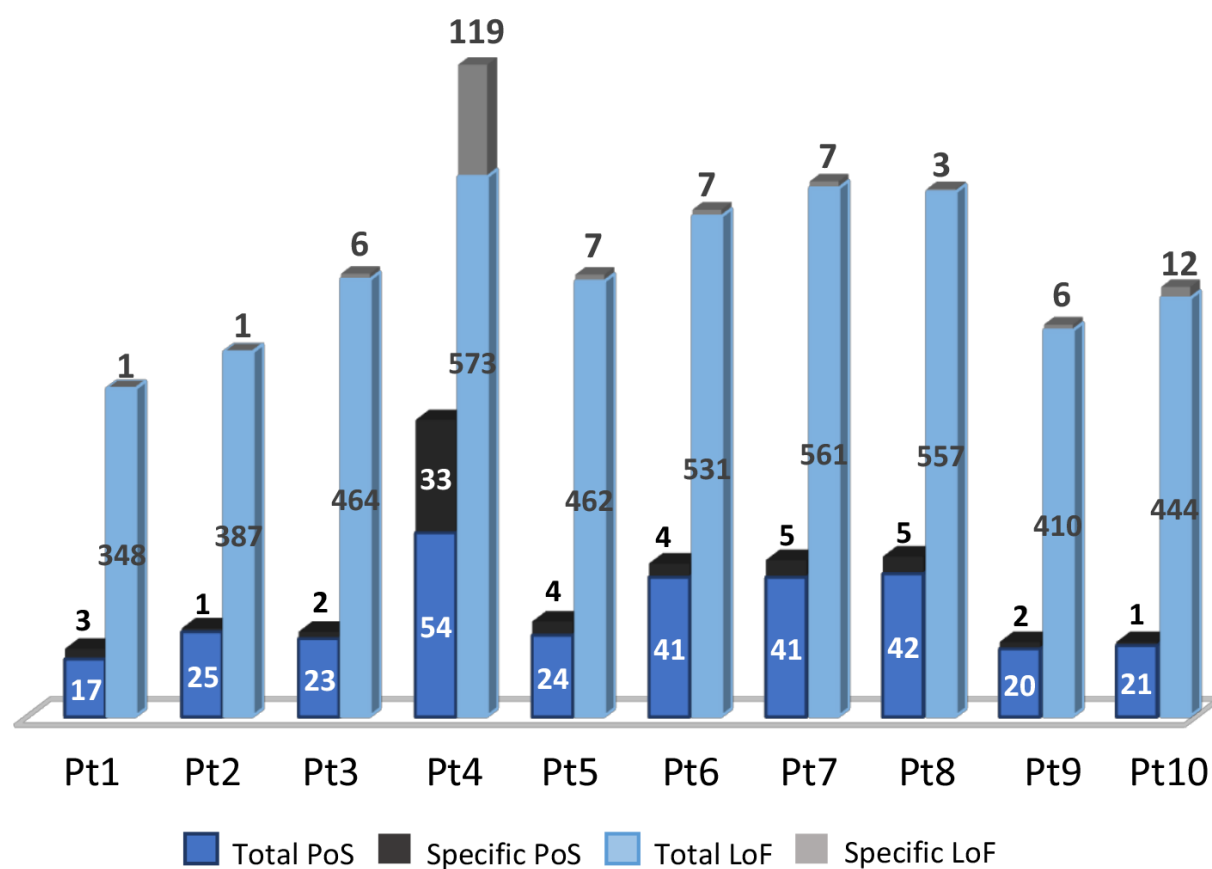
Library Name	Origin	Year of Isolation	Mapped Read-Pairs	% Mapped Read-Pairs	Alignment Depth (X)	Genome Coverage (%)
Pt1	Blackpool, UK	1956	3, 642,044	79.41	26.5	98.0
Pt2	Plymouth, UK	Prior to 1910	6, 016,241	78.23	43.8	98.0
Pt3	Plymouth, UK	1930s	6, 373,591	65.62	46.4	98.3
Pt4	Finland	1951	15, 583,665	67.31	113.5	94.0
Pt5	West Dennis, MA, USA	1972	5, 346,009	75.50	38.9	93.2
Pt6	MA, USA	1956	3, 922,830	64.50	28.5	94.1
Pt7	Long Island, NY, USA	1952	4, 937,516	67.30	35.9	94.9
Pt8	Vancouver, Canada	1987	22, 235,170	78.36	162.1	94.4
Pt9	Guam, Micronesia	1981	7, 551,099	74.68	55.2	97.5
Pt10	Dalian, China	2000	5, 436,057	72.59	39.6	92.1

Table 1

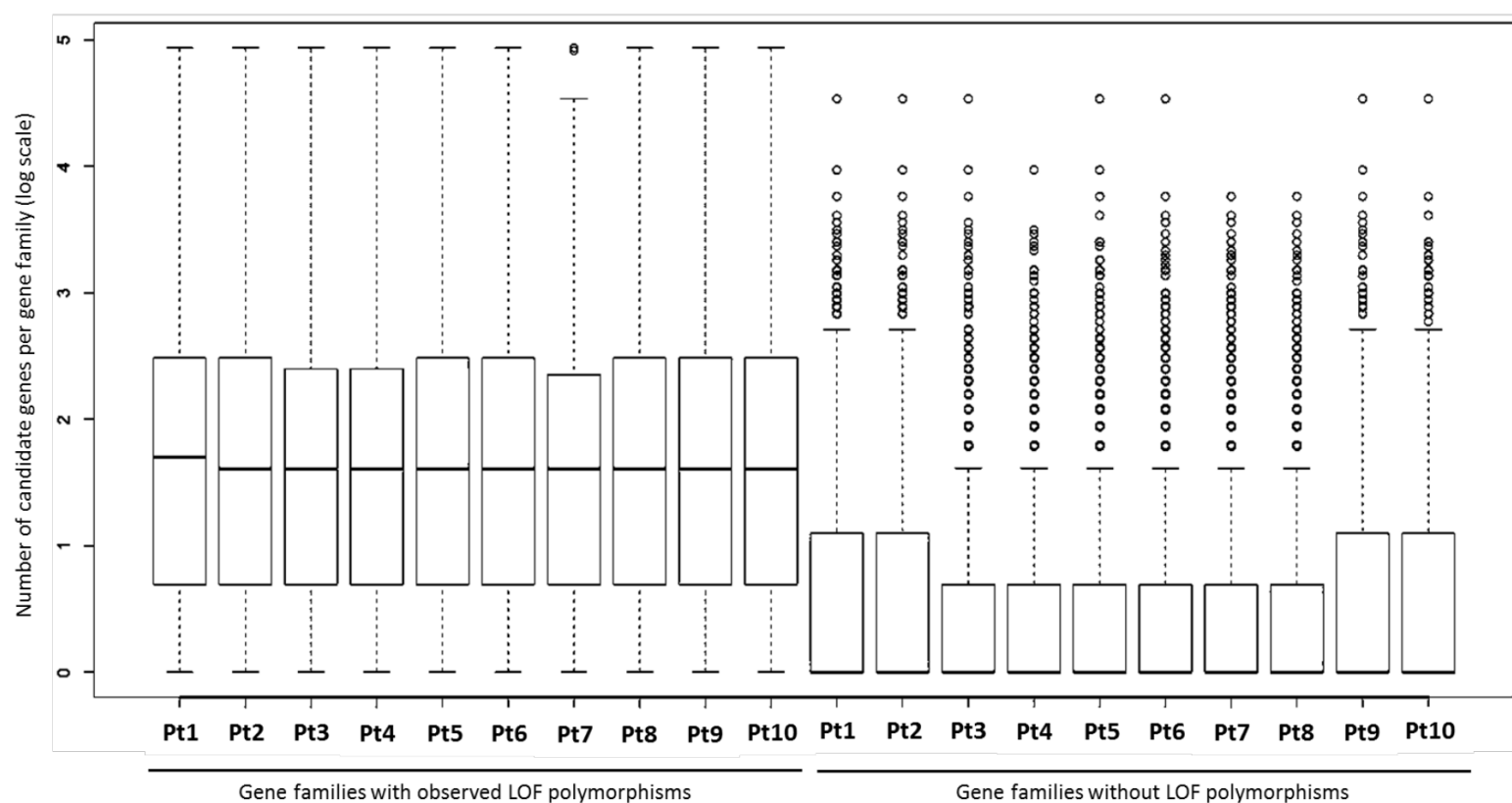


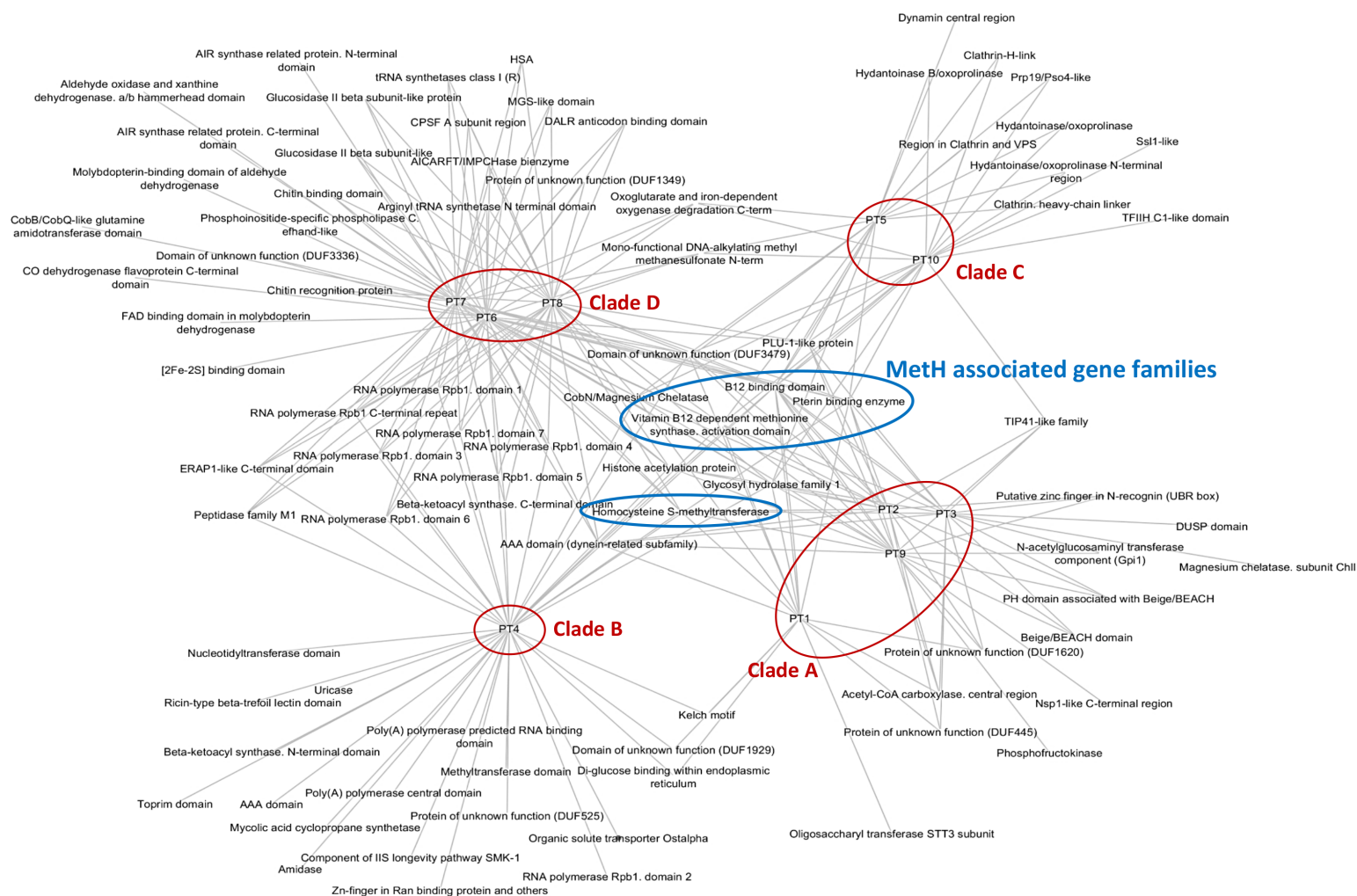


A



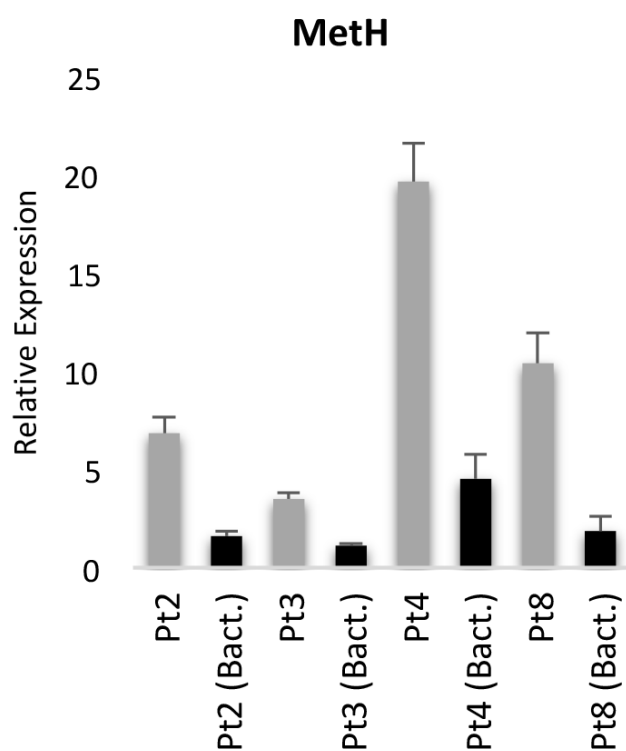
B



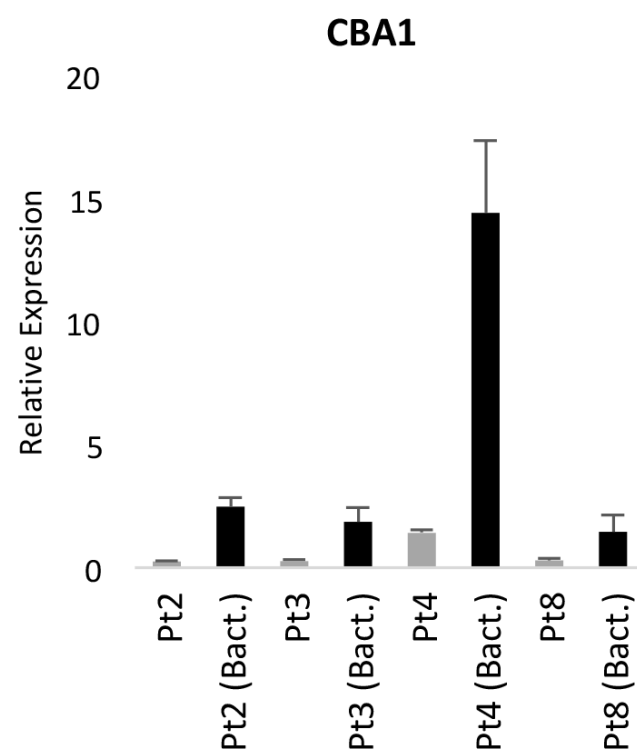




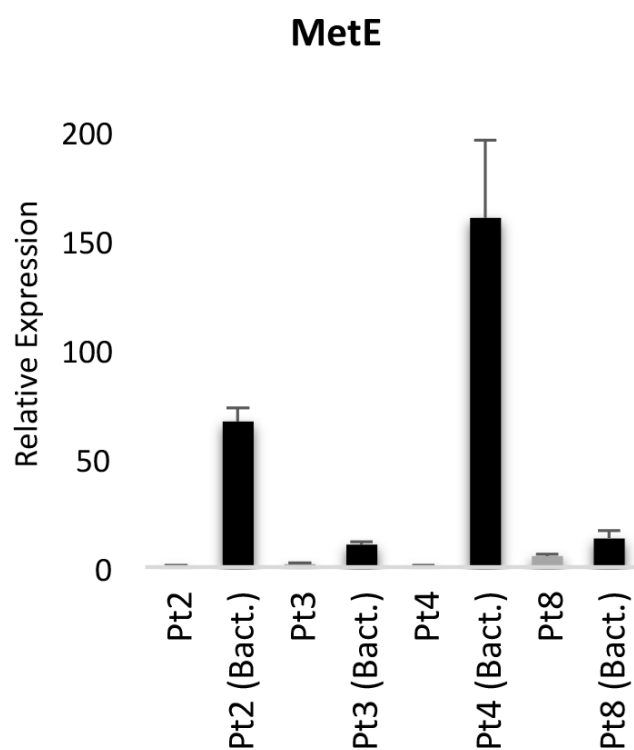
A



C



B



D

