

# Systematic and functional analysis of horizontal gene transfer events in diatoms.

Emmelen Vancaester<sup>1,2</sup>, Thomas Depuydt<sup>1,2</sup>, Cristina Maria Osuna-Cruz<sup>1,2</sup>, Klaas Vandepoele<sup>1,2,3,\*</sup>

<sup>1</sup>Department of Plant Biotechnology and Bioinformatics, Ghent University, Technologiepark 71, 9052 Ghent, Belgium

<sup>2</sup>VIB Center for Plant Systems Biology, Technologiepark 71, 9052 Ghent, Belgium

<sup>3</sup>Bioinformatics Institute Ghent, Ghent University, Technologiepark 71, 9052 Ghent, Belgium

\*Corresponding author

## 1           **1. Abstract**

2   Diatoms are a diverse group of mainly photosynthetic algae, responsible for 20% of worldwide oxygen  
3   production, which can rapidly respond to favourable conditions and often outcompete other  
4   phytoplankton. We investigated the contribution of horizontal gene transfer (HGT) to its ecological  
5   success. A systematic phylogeny-based bacterial HGT detection procedure across nine sequenced  
6   diatoms showed that 3-5% of their proteome has a horizontal origin and a large influx occurred at the  
7   ancestor of diatoms. More than 90% of HGT genes are expressed, and species-specific HGT genes in  
8   *Phaeodactylum tricornutum* undergo strong purifying selection. They are implicated in several  
9   processes including environmental sensing, and expand the metabolic toolbox. Cobalamin (vitamin  
10   B12) is an essential cofactor for roughly half of the diatoms and is only produced by bacteria. Genes  
11   involved in its final synthesis were detected as HGT, including five consecutive enzymes in  
12   *Fragilariopsis cylindrus*. This might give diatoms originating from the Southern Ocean, a region typically  
13   depleted in cobalamin, a competitive advantage. Overall, we show that HGT is a prevalent mechanism  
14   that is actively used in diatoms to expand its adaptive capabilities.

## 1           2. Introduction

2   Horizontal, also dubbed lateral, gene transfer (HGT) is the transfer of genetic information between  
3   reproductively isolated species by a route other than direct exchange from parent to progeny.  
4   Although HGT events are widespread and well documented among prokaryotes, they are much rarer  
5   in eukaryotes. Nevertheless, recently several examples of HGT from archaea or bacteria into  
6   eukaryotes have been reported. Functional HGT events have been described for almost all unicellular  
7   eukaryotic lineages, including fungi<sup>1,2</sup>, extremophilic red algae<sup>3</sup>, green algae<sup>4</sup>, rumen-associated  
8   ciliates<sup>5</sup>, oomycetes<sup>6</sup> and photosynthetic diatoms<sup>7,8</sup>. Next to events involving the maintenance of pre-  
9   existing functions, which occur mainly in endosymbiotic relationships, innovative events have been  
10  described which provide the recipient with new functions or an altered phenotype<sup>9</sup>. Although the  
11  uptake of genetic material happens by chance, fixation does not, making HGT predominantly  
12  important in the following processes: i) the alteration of iron uptake and metabolism<sup>8,10,11</sup>, ii)  
13  adaptation to an anaerobic lifestyle<sup>12,13</sup>, iii) nucleotide import and synthesis<sup>1,14</sup>, iv) novel defence  
14  mechanisms<sup>15,16</sup>, v) mechanisms to cope with stressors such as salt<sup>17,18</sup>, temperature<sup>4</sup> and heavy-metal  
15  concentrations<sup>3</sup> and vi) expansion of its metabolic capacities<sup>2,5,6</sup>.

16  Diatoms (Bacillariophyta) are one of the most abundant and species-rich groups of phytoplankton and  
17  release between 20-25% of the global amount of oxygen<sup>19</sup>. They can rapidly adapt to local conditions,  
18  outcompete other photosynthetic eukaryotes and dominate oceanic spring blooms, as long as silicon  
19  is not limited<sup>20</sup>. Moreover, they are found throughout every aquatic photic zone of this planet, such as  
20  oceans, intertidal zones, freshwater bodies, soil and even ice ecosystems<sup>21</sup>. Molecular clock evidence  
21  suggests that diatoms emerged between 225 and 200 million years ago<sup>22</sup> and their origin may be  
22  related to the end-Permian mass extinction which occurred around 250 million years ago. In the early  
23  Cretaceous, between 150 and 130 million years ago, diatoms split into the centric and pennate lineage.  
24  Several whole-genome sequences of representatives from polar centrics (*Thalassiosira pseudonana*<sup>23</sup>,  
25  *Thalassiosira oceanica*<sup>24</sup>, *Cyclotella cryptica*<sup>25</sup>), araphid pennates (*Synedra acus*<sup>26</sup>) and raphid pennates  
26  (*Phaeodactylum tricornutum*<sup>7,27</sup>, *Fistulifera solaris*<sup>28</sup>, *Fragilariopsis cylindrus*<sup>29</sup>, *Pseudo-nitzschia*  
27  *multistriata*<sup>30</sup>) have become available in recent years, which allows the analysis of the evolutionary  
28  history within diatoms. It is not fully understood how HGT has contributed to the ecological success of  
29  this environmentally important group of organisms. Moreover, diatoms evolved from several  
30  endosymbiotic events and their plastid is thought to have originated from a red alga, which has also  
31  contributed to their genetic set-up.

32  Although HGT detection has been previously performed in diatoms within the context of genome  
33  projects<sup>7,24,25,27,30</sup>, they were based on different methodologies and criteria and are therefore not  
34  directly comparable. While some studies used phylogenetics<sup>7,30</sup>, others relied purely on sequence  
35  homology searches<sup>24,25,27</sup>. In this study, we sought to systematically detect HGT events simultaneously  
36  across all sequenced diatoms. We delineated genes from horizontal descent using a high-throughput  
37  gene family phylogenetics-based approach, which allows dating transfer events. Here, we  
38  comprehensively explore the functional bias of HGT genes in diatoms and for the first time gain insight  
39  into their expression dynamics and patterns of selection.

### 3. Results

#### Detection and phylogenetic distribution of diatom HGT candidates

Twenty unicellular eukaryotic species (Table S1) were selected to deduce the contribution of bacterial-derived HGT. The identification of bacterial-to-eukaryotic HGT was achieved by building a phylogenetic tree per gene family. Therefore, all protein-coding genes from 20 eukaryotic species (Figure 1a) were clustered in 145,601 gene families, of which 32% are genes lacking similarity to any other protein in this dataset, followed by phylogenetic tree construction for 8,476 gene families having similarity to bacterial proteins. Also the species topology of these 20 unicellular eukaryotes was constructed, both based on single-locus trees and a concatenation-based approach of 156 near-single copy gene families (138,948 amino acids) (Figure 1a), with the haptophyte *Emiliana huxleyi* as an outgroup. Having the species tree available, allows for the dating of HGT candidates.

To avoid the misclassification of contaminating DNA present in the genome assembly as genomic regions originating by HGT, several quality analyses were performed. The guidelines proposed by Richards and Monier<sup>31</sup> were followed to exclude incorrect inference of HGT. Therefore, the gene origin was determined by phylogenetic tree construction followed by inspection of species-specific HGT genes. Also the percentage GC and the integration of HGT genes across chromosomes was assessed. First, the fraction of species-specific HGT was compared among all diatoms. More than 75% (2146/2844) of the predicted HGT genes in *S. acus* were only detected in this genome, while in all other species this fraction was drastically lower (11.58 +/- 9.25%) (Figure 1b). A donor analysis of these genes revealed that many were derived from *Sphingomonas sp.*, which has been described to be associated with *S. acus* in culture<sup>32</sup>. Contigs flagged to be contaminant based on a nucleotide sequence similarity search against all available Sphingomonadales genomes were clearly separable from *S. acus* based on their significantly lower percentage GC (Figure 2a) (42.1% vs 63.3%, p-value <2x10<sup>-16</sup>). Therefore, all 695 nuclear contigs having a GC content above 50% were removed, reducing the nuclear *S. acus* genome size by 4 Mb to 94.38 Mb and retaining 23,719 genes. Interestingly, the HGT detection procedure succeeded in both flagging the contaminant as detecting HGT events in the *S. acus* genome (Figure 2b). Despite the fact that in several other diatoms the GC content was significantly different between genes from horizontal and vertical descent, the mean difference never exceeds two percentage points (Figure S1).

Next, the enrichment of HGTs per contig or chromosome was evaluated to assess whether certain regions are derived from contamination, yielding no clear examples of clustering of HGT genes on specific genomic locations. The distribution of HGTs across the chromosome-level genomes of *P. tricornutum* and *T. pseudonana* is plotted in Figure S2 and shows an unbiased distribution of HGT genes. As it has been proposed that the transfer of transposable elements could be associated with facilitating gene transfer<sup>33</sup>, the distance between every gene and its closest transposable element was calculated in *P. tricornutum*. Species-specific HGT genes were significantly closer to transposable elements (TEs) (p-value 1.6x10<sup>-03</sup>), while the same was also true for vertically descended species-specific genes (p-value 2.7x10<sup>-14</sup>). This suggests that novel genes are more likely to integrate and become fixed close to repetitive regions.

Except for a fraction of genes in the *S. acus* data set, we could not identify genomic properties indicating that the identified HGT genes are caused by contamination. In total, 7,461 diatom genes were defined as having HGT origin, covering 1,979 gene families. This reflects 509 to 1,741 genes per species, making 3 to 5 percent of the diatom gene repertoire predicted to be HGT (Figure 1a). This is similar to previous phylogenetic-based estimations of HGT content in diatoms, which ranged from 587 genes (4.8%) in *P. tricornutum*<sup>7</sup> to 438 in *P. multistriata*<sup>30</sup> (3.6%) and is slightly higher than what was

46 reported in the anaerobic gut parasite *Blastocystis hominis*<sup>34</sup> (2.5%), where next to bacterial HGT also  
47 other transfers were described. The lower frequency in this stramenopile could be due to its  
48 constrained and reduced genome size as a result of its parasitic lifestyle. On average, a HGT gene family  
49 consisted out of 3.76 diatom genes and 2.55 diatom species. In total, only 106 HGT families were  
50 present in all nine diatoms. For 69 gene families the HGT copies were significantly expanded in at least  
51 one species, of which notably 26 and 21 gene families were expanded respectively in *S. robusta* and *S.*  
52 *acus* (Table S2). Indeed, gene family expansion by duplication has been observed before following HGT  
53 integration in eukaryotes<sup>1,35</sup> and this could be a strategy to diversify the original acquired function.

54 The age of all gene families of vertical descent was determined based on the lowest common ancestor  
55 of the observed species. Similarly, the most likely time point of integration for every HGT was  
56 determined using the species composition of the acceptor branch in the phylogenetic tree (Figure 1a).  
57 The large number of HGT gene families that can be attributed to the ancestor of diatoms is striking,  
58 ranging from 15% in *S. robusta* to 30% in *T. pseudonana* (Figure 1b). Another study<sup>27</sup> also detected a  
59 continuous flux of genes from prokaryotes during the evolutionary history of *P. tricornutum*. However,  
60 they claimed that most influx occurred at ancestor of the photosynthetic Stramenopiles (Ochrophytes),  
61 while our results indicate this happened more recently in the ancestor of the diatom clade.

62 Finally, several structural gene features were evaluated according to their mode of inheritance. The  
63 coding gene length of vertically descended species-specific genes in all diatom species was significantly  
64 shorter compared to all other genes (p-value  $<2 \times 10^{-16}$ ) and significantly shorter to the species-specific  
65 HGT genes in all diatoms, except for *T. pseudonana* (Figure S3). In yeast, it has also been observed that  
66 *de novo* genes were on average shorter than conserved and horizontally transferred genes<sup>36</sup>. Species-  
67 specific HGT genes on the other hand, were significantly shorter to all other genes in *C. cryptica* (p-  
68 value  $2.1 \times 10^{-02}$ ) and *S. robusta* (p-value  $1.4 \times 10^{-05}$ ). Given that introns are a typical eukaryotic gene  
69 feature, HGT genes are expected to have a shorter total intron length, especially for recent acquisitions  
70 as HGT genes adapt to their recipient genome. The intron length of HGT genes was significantly shorter  
71 in several pennate diatoms (*F. solaris*:  $1.8 \times 10^{-03}$ , *P. tricornutum*:  $3.4 \times 10^{-02}$ , *S. robusta*:  $1.1 \times 10^{-07}$  and *S.*  
72 *acus*:  $4.6 \times 10^{-03}$ ) and for several diatoms the young species-specific HGT genes had shorter introns than  
73 the rest of the gene repertoire (*F. cylindrus*:  $1.1 \times 10^{-03}$ , *F. solaris*:  $9.8 \times 10^{-03}$ , *P. tricornutum*:  $1.1 \times 10^{-02}$ , *S.*  
74 *robusta*:  $2 \times 10^{-09}$  and *C. cryptica*:  $4.9 \times 10^{-02}$ ). These results indicate that introns become an emerging  
75 property of HGT genes after integration.

## 76 **The functional landscape of diatom HGT genes**

77 To gain insight in the functional repertoire of HGT genes, a gene ontology (GO) and functional domain  
78 (Interpro) enrichment was performed. Out of the 7,461 diatom HGT genes, 6,024 (81%) were  
79 annotated with an Interpro domain and 3,893 (52%) with a GO term. The only GO term which was  
80 enriched for HGT genes in all nine diatom species is pseudouridine synthesis (GO:0001522), while  
81 enriched protein domains covered pseudouridine synthase (IPR006145), S-adenosyl-L-methionine-  
82 dependent methyltransferase (IPR029063) and nitroreductase (IPR029479). An overview of the  
83 enriched functional categories across different ages can be found in (Figure S4, Figure S5). A more in-  
84 depth exploration of several functional categories is given in Supplementary Note 1, while an overview  
85 of all discussed functions and their corresponding gene families can be found in Table S3.

## 86 **Cobalamin uptake**

87 Cobalamin (vitamin B12) is a complex molecule composed out of a central cobalt-containing corrin  
88 ring, a lower ligand of 5,6-dimethylbenzimidazole (DMB) and an upper axial ligand that can either be  
89 an hydroxy-, cyano-, methyl or adenosyl group. Vitamin B12 acts as a coenzyme in three enzymes in  
90 eukaryotes: methylmalonyl-CoA mutase (MCM), type II ribonucleotide reductase (RNRII) and

91 methionine synthase (METH). Despite that more than half of the algal species surveyed (171/326)<sup>37</sup>  
92 are auxotrophic for vitamin B12, including 37 out of 58 diatoms, *de novo* synthesis has only been  
93 described to occur in prokaryotes. Therefore cobalamin availability alters the composition of marine  
94 phytoplankton communities<sup>38</sup>. The exchange of cobalamin in return for organic compounds is believed  
95 to underpin the close mutualistic interactions between heterotrophic bacteria and auxotrophic algae<sup>39</sup>  
96 A correlation was detected between the scattered phylogenetic pattern of absence of a cobalamin-  
97 independent methionine synthase (METE) and auxotrophy for this vitamin<sup>40</sup>. It has been suggested  
98 that this loss has a biogeographical basis as there is a tendency for diatoms occurring in the Southern  
99 Ocean to retain METE more often<sup>41</sup>. Moreover, it has been recently proven that cyanobacteria produce  
100 the chemical variant pseudo-cobalamin, where adenine substitutes DMB as the lower ligand, which is  
101 less bioavailable to eukaryotic algae<sup>42</sup>. However, some species, including *P. tricornutum* and *E. huxleyi*<sup>43</sup>  
102 can remodel this to cobalamin using CobT, CobS and CobC via the nucleotide loop assembly<sup>39,42</sup>. Here  
103 BluB, necessary for DMB production<sup>44</sup>, was detected to have originated by HGT from  
104 alphaproteobacteria in *F. cylindrus*, *P. multistriata* and *P. multiseriis* (Figure 3a). More than 90 percent  
105 of the cobalamin-producing alpha- and gammaproteobacteria encode BluB<sup>39</sup>. Moreover, five HGT  
106 genes were detected in the final synthesis of the cobalamin biosynthesis pathway, which can also  
107 function as scavenging and repair genes: CobN, CobA/CobO, CobQ/CbiP, CobD/CbiB and CobU/CobP  
108 (Figure 3a). These genes were previously also detected in diatom metatranscriptomes and *P. granii*<sup>45,46</sup>,  
109 where CobN, CobS and CobU were more highly expressed under iron replete conditions. Interestingly,  
110 for all diatoms, except for *C. cryptica*, the CbiB gene also contains the CbiZ domain, which is involved  
111 in the removal of the lower ligand<sup>43</sup>. Only *F. cylindrus* and *P. multiseriis* contain the full suite of these  
112 detected HGT genes in their cobalamin pathway, while *P. tricornutum*, *F. solaris* and *S. acus* possess  
113 none (Figure 3a). Interestingly, while *P. multiseriis* is auxotrophic for cobalamin, *F. cylindrus* is not.  
114 Thus, despite the presence of the cobalamin independent methionine synthase METE, *F. cylindrus*  
115 expanded its repertoire of cobalamin synthesis genes and prefers to maximally optimize its uptake to  
116 perform methionine synthesis by the more efficient METH. By querying the metatranscriptomic TARA  
117 Oceans data it was clear that CobU, CbiZ+CbiB and CobQ are significantly correlated with nitrate  
118 concentration and day length (Figure 3b) (Figure S6, Figure S7), while CobU and CbiZ+CbiB are anti-  
119 correlated with temperature (Figure S8) and CobU and CobQ are anti-correlated with iron (Figure S9).  
120 HGT genes in the cobalamin pathway are particularly abundant in the Southern and Pacific Ocean  
121 (Figure 3b). The lower production rate of bacteria in low temperature and the photodegradation of  
122 cobalamins, which could be of particular importance during arctic summers, might explain the  
123 cobalamin limitation and the specific expression of vitamin B12-related genes in these regions of the  
124 ocean.

## 125 Environmental adaptation to light sensing and cold protection

126 Diatoms employ photosensory proteins to gain information about their environment and respond to  
127 changing light conditions. Proteorhodopsins (PR) perceive light to drive ATP generation and are  
128 especially important when photosynthesis is comprised during iron-limiting conditions. This study  
129 confirms the bacterial origin of the PR-genes in *F. cylindrus* and *Pseudo-nitzschia granii*<sup>47</sup>, next to brown  
130 algae, dinophytes and haptophytes. Furthermore, the red/far-red light sensing phytochrome DPH1<sup>48</sup>  
131 was detected as HGT in *P. tricornutum* (1 copy), *S. robusta* (4), *S. acus* (7), *C. cryptica* (1) and *T.*  
132 *pseudonana* (1), and formed together with brown algae an independent branch from green algal and  
133 fungal DPHs, similar as in previous reports<sup>48,49</sup> and was predicted to have originated from HGT.

134 Arctic diatoms such as *F. cylindrus* undergo periods of prolonged darkness, low temperature and high  
135 salinity. Their ability to thrive in these conditions could be partially attributed to cryoprotectants that  
136 interfere with the growth of ice<sup>50</sup>. Ice-binding proteins were found to be laterally transferred from a



137 basidiomycete lineage to *Fragilariopsis curta* and *F. cylindrus*<sup>51</sup>. Also the phylogenetic tree inferred in  
138 this study detected relatedness between fungal and diatom antifreeze proteins, but wasn't classified  
139 as HGT as this pipeline only detects bacterial-to-eukaryotic HGT. However, a second gene family of *F.*  
140 *cylindrus* proteins containing the ice-binding protein domain (IPR021884), was found to be transferred  
141 from *Cryobacterium*.

#### 142 **Carbon and nitrogen metabolism**

143 Diatoms can rapidly recover from prolonged nitrogen limitation due to presence of the urea cycle that  
144 allows for carbon fixation into nitrogenous compounds<sup>52</sup>. Two genes in the metabolic branches derived  
145 from this pathway, carbamate kinase and ornithine cyclodeaminase were found to be laterally  
146 transferred, both here as in previous studies<sup>7,52</sup>. The latter enzyme is responsible for the conversion of  
147 ornithine to proline, which is the main osmolyte during salt stress in diatoms. Another way of a  
148 nitrogen storage and translocation is the catabolism of purines to urate that can be further degraded  
149 to allantoin. It was found that plants and diatoms independently evolved a fusion protein (Urah-Urad  
150 domain; allantoin synthase) to perform the second and third step in this urate degradation pathway<sup>53</sup>.  
151 Exactly as in<sup>53</sup>, this gene was detected to be laterally transferred from alphaproteobacteria, where this  
152 fusion event occurred, to the ancestor of haptophytes and stramenopiles.

153 Moreover, several genes in carbohydrate metabolism were found to be laterally transferred. The  
154 acetyl-CoA conversion to acetate occurs in a two-step process where phosphate acetyltransferase  
155 (PTA) adds a phosphate group to form acetylphosphate, that is in turn is catalyzed to acetate by acetate  
156 kinase (ACK)<sup>54</sup>. The PTA gene family was found to have bacterial origins and emerged in the ancestor  
157 of haptophytes and stramenopiles. In all diatoms, except for *F. cylindrus* and *P. multistriata*, multiple  
158 copies were found of this gene. Also acetate kinase was detected as a HGT gene in the pennate diatoms  
159 *P. tricornutum*, *S. robusta* and *S. acus*. Furthermore, this enzyme was predicted to be involved in the  
160 bifid shunt<sup>54</sup>. Here, the key enzyme XPK cleaves xylulose-5-phosphate to acetyl-phosphate and  
161 glyceraldehyde-3-phosphate, followed by conversion of acetyl-phosphate to acetate by ACK. Also XPK  
162 was laterally transferred in the pennate diatoms, single-copy in *P. tricornutum* and significantly  
163 expanded to five copies in *S. acus*. XPK and ACK are syntenic in *P. tricornutum*, what was already  
164 suggested to point to a bacterial origin as this spatial organization is also detected in Proteobacteria  
165 and Cyanobacteria<sup>54</sup>. Interestingly, *S. acus* has also conserved the physical association of XPK and ACK  
166 and maintained a bidirectional promoter, although an inversion of the gene order occurred (Figure  
167 S10). Furthermore, phosphofructokinase and the cytosolic fructose-bisphosphate aldolase Fba4 in the  
168 glycolysis<sup>55</sup>, phosphopentose epimerase<sup>56</sup> in the pentose phosphate pathway and a putative D-lactate  
169 dehydrogenase are enzymes that were predicted to be transferred from bacteria present in diatoms.  
170 Finally, also bacterial xylanases, glucanases and glucosidases expanded the carbohydrate metabolic  
171 repertoire in diatoms.

172 The biosynthetic aspartate-derived pathway to synthesize the four amino acids, lysine, threonine,  
173 methionine and isoleucine was completed due to HGT<sup>57</sup>. Aspartate semialdehyde dehydrogenase (*asd*)  
174 performs the second step in this pathway and is derived from Proteobacteria. The end product L-  
175 aspartate 4-semialdehyde can either be used by dihydrodipicolinate synthase (*dapA*) towards lysine  
176 biosynthesis, or by homoserine dehydrogenase (*thrA*) towards threonine and methionine. Both genes  
177 were laterally transferred from bacteria. The metabolic pathways of other amino acids were also  
178 affected, the last step in tryptophan synthesis is achieved by tryptophan synthase. While in diatoms  
179 the alpha and beta subunit of this enzyme are merged, in *P. tricornutum* an extra copy of the beta  
180 subunit is present<sup>58</sup> (Phatr3\_J52286) that was deemed bacterial. Also alanine racemase, arginine  
181 biosynthesis ArgJ, leucyl-tRNA synthetase *leuRS2*, glycyl-tRNA synthetase *glyRS2* and tyrosine-tRNA  
182 ligase *tyrRS2* were laterally transferred.

### 183 Selection pattern for HGT genes in *P. tricornutum*

184 Genomic sequence information from ten *Phaeodactylum* accessions, belonging to four clades sampled  
185 across the world<sup>59</sup>, was used to determine the maintenance and selection pattern across the  
186 detected HGT genes. The retention of species-specific HGT genes across different strains confirmed  
187 their horizontally derived origin and did not point to contamination (for more details, see  
188 Supplementary Note 2). Moreover, analyzing gene selection patterns gives an indication on the  
189 strength of functional conservation. Variant calling resulted in a data set of 585,715 high-confidence  
190 bi-allelic SNPs. The total number of SNPs per strain across the genome was low and ranged from 0.96%  
191 to 1.37% (Table S4). To detect selective pressure,  $\pi N/\pi S$ , was calculated. This metric compares the  
192 fraction of synonymous and nonsynonymous mutations within a coding open-reading frame across  
193 strains. A gene experiencing neutral evolution has a  $\pi N/\pi S$  value of 1, whilst a value smaller than 1  
194 signifies negative purifying selection. The smaller the ratio of non-synonymous and synonymous  
195 nucleotide diversity, the stronger is the level of purifying selection acting on the gene. The average  
196 synonymous nucleotide diversity ( $\pi S$ ) across all accessions is 0.009, while the non-synonymous  
197 nucleotide diversity ( $\pi N$ ) is 0.003, thus the genome-wide average  $\pi N/\pi S$  ratio is 0.3. This value is similar  
198 to what was described in Rastogi et al.<sup>59</sup> and means most genes undergo strong purifying selection.  
199 The average  $\pi N/\pi S$  for genes of vertical descent is 0.302, while for HGT genes it is significantly lower  
200 at 0.268 (p-value  $5.9 \times 10^{-4}$ ). When comparing the  $\pi N/\pi S$  ratios for HGT and vertical genes across age  
201 classes (Figure 4), it is apparent that this difference is due to the youngest gene categories, being the  
202 *Phaeodactylum*-, raphid pennate diatoms and pennate diatoms specific genes, where vertical genes  
203 are less constrained than HGT genes in those age categories. To the best of our knowledge, this is the  
204 first time the selection pressure of bacterial HGT genes is assessed in unicellular eukaryotes and  
205 compared with vertically descended genes while taking age into account. Although it has already been  
206 observed that *de novo* genes display patterns of rapid evolution and the strength of purifying selection  
207 increases with age<sup>36</sup>, it is remarkable to observe that HGT genes deviate from this pattern. Unlike  
208 recent innovations from vertical descent, young HGT genes are quickly integrated in the biological  
209 network exemplified by their high levels of purifying selection.

### 210 Expression and co-expression network analysis of HGT genes

211 The availability of RNA sequencing experiments in several diatoms, allows for the construction of  
212 genome-wide gene expression atlases quantifying gene expression levels across a wide range of  
213 conditions. These compendia consisted out of 13 to 76 conditions per species, all having biological  
214 replicates per condition (Table S5). The vast majority of HGT genes are expressed: in *P. tricornutum* all  
215 HGT genes are expressed, in *T. pseudonana* 558 out of 580 (96%), in *S. robusta* 1597 out of 1741 (92%)  
216 and in *F. cylindrus* 741 out of 762 (97%). Given that most HGT genes are kept under purifying selection  
217 in *P. tricornutum* and are transcribed in diatoms, this is indicative that they are functional and can play  
218 a vital role in expanding the functional repertoire. Indeed, 64% of the predicted HGT genes in *P.*  
219 *tricornutum* were translated into proteins in a proteogenomic analysis<sup>60</sup>. This is similar to 63% of all  
220 proteins in the genome that were detected to be translated.

221 Next, the expression specificity was calculated per gene, where a low value signifies broad expression  
222 in many (or even all) conditions and a value close to one indicates expression in one or a few. Species-  
223 specific genes have a higher mean condition-specific expression, both for vertical as horizontal derived  
224 genes in *P. tricornutum* (p-value  $< 2 \times 10^{-16}$ ,  $1.1 \times 10^{-06}$ ), *S. robusta* ( $< 2 \times 10^{-16}$ ,  $1.6 \times 10^{-13}$ ), *F. cylindrus* ( $< 2 \times 10^{-16}$ ,  
225  $2.2 \times 10^{-03}$ ) and *T. pseudonana* ( $< 2 \times 10^{-16}$ ,  $3.7 \times 10^{-02}$ ). A declining trend of condition-specificity was  
226 observed over time. Whenever there was a significant difference in condition specificity between HGT  
227 and vertical genes within the same age category, HGT genes consistently displayed on average a more  
228 specific expression pattern (Figure 5). A high tissue-specificity for species-specific genes which

229 decreases over time has also been observed in mouse<sup>61,62</sup>. Interestingly, the selection pressure in *P.*  
230 *tricornutum* across all age classes and per mode of inheritance is not correlated with expression  
231 specificity (Figure S11), showing that genes having a highly specific expression are not necessarily  
232 under less purifying selection, defying the trend that was previously observed in mammals<sup>63</sup>.

233 Based on a global co-expression *P. tricornutum* network constructed using an expression atlas  
234 comprising 211 samples, for every gene in *P. tricornutum* the co-expression neighbourhood was  
235 defined as a module. Subsequently, these modules and the known gene function for genes part of this  
236 module were used, through guilt-by-association analysis, to gain functional insights in the detected  
237 HGT genes. For 19 HGT genes the co-expression modules confirmed enrichment for at least one known  
238 function. Fructose-bisphosphate aldolase Fba4 is enriched in its co-expression module for genes  
239 involved in a carbohydrate metabolic process and aspartate semialdehyde dehydrogenase (*asd*) has  
240 enrichment for amino acid biosynthesis. New functions were attributed to 320 out of 509 HGT genes  
241 based on significant GO enrichment of the co-expression modules. For example, two HGT proteins  
242 involved in amino acid synthesis - tryptophan synthase  $\beta$  chain (p-value  $2.2 \times 10^{-16}$ ) and ArgJ (p-value  
243  $1 \times 10^{-15}$ ) - are predicted to be coregulated with photosynthetic genes, although both proteins do not  
244 contain the chloroplast targeting peptide. The co-expression neighbourhood of phosphofructokinase  
245 is significantly enriched (p-value  $3.7 \times 10^{-06}$ ) to be involved in the Krebs cycle, while this protein is part  
246 of the glycolysis and directly upstream of the citric acid cycle. Also Phatr3\_J40382, which contains a  
247 pyruvate kinase-like domain, is enriched for this GO term, and this corroborates its metabolic function.  
248 Methionine sulfoxide reductase MsrB (Phatr3\_J13757) was predicted to have a similar expression  
249 pattern to genes partaking in iron-sulphur cluster assembly (p-value  $1.1 \times 10^{-03}$ ) and metal ion transport  
250 (p-value  $1.4 \times 10^{-02}$ ). In yeast these proteins were already shown to have a protective role for FeS  
251 clusters during oxidative stress<sup>64</sup>. The far-red light phytochrome DPH1 is enriched for genes involved  
252 in transcriptional regulation (p-value  $5.5 \times 10^{-03}$ ), which could point to its primary role in the light sensing  
253 cascade. These results demonstrate that co-expression network analysis offers a pragmatic means to  
254 predict the biological processes HGT genes are involved in.

#### 255 4. Discussion

256 Through the application of phylogeny-based HGT detection, we identified 1,979 gene families with a  
257 horizontal origin in diatoms. Although HGT detection has been previously performed in diatoms, this  
258 is the first large-scale and systematic approach of HGT detection across all available sequenced  
259 diatoms. While some previous studies were based on phylogenetics<sup>7,30</sup>, most relied purely on sequence  
260 homology searches<sup>24,25,27</sup>, while it has been shown that the degree of gene similarity does not  
261 necessary necessarily reflect phylogenetic relationships<sup>65,66</sup>. Although HGT had been previously  
262 predicted in *P. multistriata*, *C. cryptica*, *T. oceanica* and twice in *P. tricornutum*, only a fraction of HGT  
263 genes were confirmed across these studies, going from 6 to 45% (Figure S12). This could be due to the  
264 usage of different methods, criteria and underlying databases. For example, horizontal genes were  
265 defined in *T. oceanica* if they didn't show high similarity with any other stramenopile and thus contain  
266 *T. oceanica*-specific genes from both bacterial- and eukaryotic-to eukaryotic origin. Nonetheless, the  
267 overlap between all methods is still significantly higher than expected by chance for all species.

268 Assessing the strength of purifying selection using genome-wide nucleotide diversity information  
269 showed that lateral gene transfer is quickly followed by fixation compared to the retention of new  
270 genes originating from non-coding regions, so-called *de novo* genes. Moreover, most species-specific  
271 HGT genes are present across all strains in *P. tricornutum*, suggesting they were acquired prior the  
272 divergence of these strains and are actively maintained in the population. This underlines the  
273 importance of HGT in diatoms and indicates that fixation of a laterally transferred gene takes place  
274 quickly after the initial uptake. Indeed, in grasses it was recently shown that several plant-to-plant LGT



275 fragments were rapidly integrated and spread across the population, after which erosion occurred on  
276 neutrally selected genes within those fragments<sup>67</sup>.

277 Among the HGT events, we detected the transfer of five concurrent genes of the vitamin B12  
278 biosynthetic pathway. A cobalamin addition experiment in a high-nutrient low chlorophyll (HNLC)  
279 region in the Gulf of Alaska significantly altered the species composition, going from diatom-  
280 dominated plankton to an increased fraction of ciliates and dinoflagellates<sup>68</sup>. This could be explained  
281 by the presence of these HGT genes and the corresponding enhanced uptake mechanism of vitamin  
282 B12 and its analogues, which give diatoms a competitive advantage during limiting conditions. In  
283 conclusion, our results support a high genetic plasticity and ability for local adaptation in diatoms due  
284 to HGT.

1 **5. Figure legends**

2

3 **Figure 1: Overview of HGT events across diatoms.** **a**, Species phylogeny determined by IQ-Tree with  
4 values at the internal nodes that denote bootstrap, gene- and site-concordance factors respectively.  
5 Branches are coloured according to their phylogenetic classification. Bold and framed values at internal  
6 nodes reflect the number of predicted HGT events. For diatoms, the number of species-specific GFs of  
7 HGT origin, the total number of HGT genes and its fraction of the proteome to be originating from HGT  
8 is tabularised. For *S. acus* the number of HGT genes both prior and after removal of contamination is  
9 mentioned. **b**, Distribution of the age classes of HGTs across the nine investigated diatom species.

10 **Figure 2: Contamination of Sphingomonadales in genome of *S. acus*.** **a**, Percentage GC of *S. acus*  
11 contigs versus nucleotide sequence similarity to Sphingomonadales based on at least 70% identity and  
12 25% alignment coverage. **b**, Percentage GC of *S. acus* contigs versus presence of a HGT candidate within  
13 each contig.

14 **Figure 3: Cobalamin pathway in diatoms.** **a**, Overview of the cobalamin biosynthesis pathway and its  
15 presence in diatoms. Genes in bold are of horizontal descent and the presence of a gene is displayed  
16 by a filled circle, diamond or square depending on its position in the pathway. **b**, Expressed number of  
17 diatom sequences of several HGT genes involved in the cobalamin pathway across stations sampled  
18 worldwide during the TARA Oceans project in the surface layer, coloured according to their nitrate  
19 concentration.

20 **Figure 4: Selective pressure over time in *P. tricornutum*.** Distribution of selective pressure, measured  
21 by  $\pi_N/\pi_S$ , across age classes sorted from young to old and per origin in *P. tricornutum*. Number of  
22 genes is indicated in between brackets. The asterisks denote a statistical difference per type within the  
23 same age category and have the following confidence range for p-values; \* : 0.05, \*\* : 0.01, \*\*\* : 0.001,  
24 \*\*\*\*: 0.0001.

25 **Figure 5: Expression specificity over time in diatoms.** Distribution of expression specificity across age  
26 classes sorted from young to old and per origin in four diatoms. The number of genes is indicated in  
27 between brackets. The asterisks denote a statistical difference per type within the same age category  
28 and have the following confidence range for p-values; \* : 0.05, \*\* : 0.01, \*\*\* : 0.001, \*\*\*\*: 0.0001.

29

## 30 **6. Material and methods**

### 31 **Gene family construction**

32 The publicly available genomes of 17 stramenopiles, 1 alveolate, 1 rhizarian and 1 haptophyte (listed  
33 in Table S1) were downloaded and their nuclear proteomes, totalling to 398,001 protein-coding genes,  
34 were searched for similarity in an all-against-all fashion with BLASTp (version 2.6+) using an e-value  
35 cut-off of  $10^{-5}$  and retaining maximum 4,000 hits. Next, clustering of these protein-coding genes was  
36 performed using OrthoFinder (version 2.1.2)<sup>69</sup>.

### 37 **Species tree phylogeny**

38 To delineate the species phylogeny for all SAR members, using *Emiliana huxleyi* as an outgroup,  
39 OrthoFinder gene families where all species have a copy number of either 1 or 2 genes were selected,  
40 and one gene sequence was randomly picked in case of duplication. MUSCLE<sup>70</sup> was used to build a  
41 concatenated sequence alignment. Afterwards IQ-Tree (version 1.7.0b7)<sup>71</sup> was used to build a  
42 concatenated tree using 1,000 bootstraps, estimate the single-locus trees and finally to calculate the  
43 gene- and site-concordance factors of the inferred species tree<sup>72</sup>.

### 44 **HGT detection**

45 The NCBI non-redundant protein database (download date 08/06/2018) was complemented with the  
46 proteomes of 20 species (Table S1). Diamond (version 0.9.18.119)<sup>73</sup> searches were performed in  
47 sensitive mode against this database for all proteins of these 20 species, retaining maximum 1,000 hits  
48 per query. Hits were reduced to maximum five sequences for each order and 15 sequences per phylum.  
49 Genes families with at least one copy in a diatom and at least one third of the diatom members having  
50 a bacterial hit were analysed. The hits of all diatom members were combined and clustered using CD-  
51 HIT (version 4.6.1)<sup>74</sup> based on a 95% identity cut-off. Next, the sequences were aligned with MAFFT  
52 (version 7.187)<sup>75</sup> in automatic mode. Maximum likelihood trees are produced using IQTree (version  
53 1.6.5)<sup>71</sup> including a test for the best fitting protein model (-mset JTT, LG, WAG, Blosum62, VT,  
54 Dayhoff)<sup>76</sup>. The FreeRate model was used to account for rate heterogeneity across sites (-mrate R),  
55 empirical base frequencies were calculated (-mfreq F) and 1,000 rounds of ultra-parametric  
56 bootstrapping (-bb 1000) (UFBoot2)<sup>77</sup> were run.

57 Phylogenetic trees were reordered based on midpoint rooting, unless the whole eukaryotic fraction  
58 formed a cluster and then this cluster was used as a subtree for rooting. For every node having a  
59 bootstrap support  $\geq 90$ , and consisting out of a bacterial and eukaryotic subtree the last common  
60 ancestor (LCA) was defined. When the eukaryotic subtree consisted out of more than 20 sequences,  
61 they could only compose at most 85% of the total number of sequences belonging to that node. When  
62 several nodes complied to these rules having the same eukaryotic fraction, the last common ancestor  
63 of the bacterial subset was considered the donor of this event. To avoid classifying endosymbiotic gene  
64 transfer (EGT) incorrectly as HGT, only bacterial-to-eukaryotic events were analysed and events older  
65 than SAR + Haptophytes, also dubbed SAR+CCTH, were discarded.

### 66 **Gene family expansion**

67 Expanded families were delineated by calculating the Z-score profile of the gene copy number per HGT  
68 family across all diatoms excluding the allodiploid *F. solaris*. Families where the variance is larger than  
69 two and the Z-score for a particular species is larger than three, were deemed expanded in that species.

### 70 **Contamination detection**

71 Structural genomic annotation features for sequenced diatoms were retrieved and GC content, coding  
72 sequence length, number of introns per gene and intron length were compared between horizontally  
73 and vertically transferred genes among several age categories. Also the distance for every gene in *P.*  
74 *tricornutum* to the closest transposable element as defined by<sup>27</sup>, centromeric and telomeric regions  
75 elucidated by<sup>78</sup> was calculated and compared among the different origins. Statistical significance was  
76 calculated by the Wilcoxon rank sum test. For the diatoms *T. pseudonana* and *P. tricornutum*, whose  
77 genomes are resolved on chromosome-scale level, the distribution of HGT genes was plotted using R.

78 To assess the degree of contamination from *Sphingomonas sp.* in *S. acus*, 914 genomes of the order  
79 *Sphingomonadales* were retrieved from NCBI and a nucleotide blast against the *S. acus* genome was  
80 performed. Contigs having at least 70% identity and 25% alignment coverage were deemed to have  
81 *Sphingomonadales* origin.

## 82 **Functional interpretation of HGT genes**

83 The proteomes of all species were functionally annotated using Interproscan (version 60)<sup>79</sup> in order to  
84 obtain functional domain annotations and Gene Ontology (GO) terms. KEGG orthology identifiers<sup>80</sup>  
85 were attained using EggNOG-mapper<sup>81</sup>. For all diatoms, the chloroplast targeting signal was predicted  
86 using ASAFind (version 1.1.7)<sup>82</sup>. Only GO terms within the subtree 'biological process' were taken into  
87 account. These terms were expanded to also contain all ancestral functional information. GO and  
88 Interpro domain enrichment was performed on the HGT genes per species using hypergeometric  
89 testing, and multiple hypothesis testing was constrained using Benjamini–Hochberg correction (q  
90 value < 0.05). Functional enrichments found in at least two species were visualized using the  
91 ComplexHeatmap package<sup>83</sup> (R version 3.4) and clustered using the complete linkage method.

92 Tandem duplicates were defined as genes belonging to the same gene family and located within 15  
93 genes of each other and identified using i-ADHoRe v3.0<sup>84</sup> (alignment method: gg2, gap size 15, tandem  
94 gap 15, cluster gap 15, q-value 0.85, probability cut-off 0.01, anchor\_points 3, level\_2\_only FALSE, FDR  
95 as method for multiple hypothesis correction).

## 96 **Metatranscriptome analysis**

97 For several selected HGT gene families involved in cobalamin synthesis a HMM profile was created  
98 using hmmer3 v3.1b2<sup>85</sup> and these were uploaded to the Ocean Gene Atlas webserver<sup>86</sup> to query the  
99 eukaryotic MATOU gene dataset<sup>87</sup> (blastp, evalue cut-off 10<sup>-10</sup>) linked the metatranscriptomic TARA  
100 Oceans data. Only sequences taxonomically assigned as diatoms were further analysed. The  
101 abundance was estimated as the number of sequence per station and depth.

## 102 **Population genetics**

103 Data from ten resequencing strains was downloaded from the public repository SRA  
104 (<https://www.ncbi.nlm.nih.gov/sra>) (SRR6476693-SRR6476702) and these reads were mapped to the  
105 *Phaeodactylum tricornutum* genome using BWA-mem (version 0.7.17)<sup>88</sup>. The read alignments per  
106 strain were filtered to only include unique mappings without chimeric alignments using samtools  
107 (version 1.6)<sup>89</sup>. SGSGeneLoss (version 0.1)<sup>90</sup> was run in a relaxed mode to determine to  
108 presence/absence pattern of all genes across the ten strains (minCov=1, lostCutoff=0.05, thus  
109 requiring only one read and 5% gene coverage to be perceived as present). The resulting phylogenetic  
110 pattern of HGT genes was visualized using the ComplexHeatmap package<sup>83</sup> (R version 3.4) and  
111 clustered using the complete linkage method.

112 SNP calling was performed per strain using HaplotypeCaller, after which SNPs were integrated using  
113 GenotypeGVCFs. Both methods are available within the GATK framework (version 3.7)<sup>91</sup>. SNPs were

114 filtering using the GATK recommended hard filters (QD<2.0; FS>60.0; MQ<40.0; MQRankSum ≤ 12.5;  
115 ReadPosRankSum ≤ 8.0)<sup>92</sup> and only bi-allelic SNPs were retained.

116 To estimate the degree of negative purifying pressure across the proteome, only coding positions  
117 having a read depth ≥ 10 across all strains were considered, calculated using SAMtools mpileup<sup>93</sup>. In  
118 total, 89% of all genic positions could be analysed and 272,235 SNPs were observed in these regions.  
119 We used SnpEff (version 4.3t)<sup>94</sup> to predict the individual effect per SNP and  $\pi N/\pi S$  was calculated taking  
120 only the callable positions for complete codons into account and correcting for the allele frequency of  
121 the mutation in the population. Statistical significance of difference in selective pressure across mode  
122 of inheritance and age classes was calculated by the Wilcoxon rank sum test.

### 123 **Expression and co-expression analysis**

124 An expression atlas for every diatom species which has RNA-Seq expression data available was  
125 generated. First, relevant experiments were searched using Curse<sup>95</sup>, which also allows the user to  
126 identify and curate replicates. The experiments listed in (Table S5) were used to generate the  
127 expression compendia. Next, the atlas was generated using Prose<sup>95</sup>, which uses the SRA toolkit to  
128 download the raw data locally, FastQC to perform quality control and adapter detection, Trimmomatic  
129 for automatic read trimming and finally kallisto for expression quantification in transcripts per million  
130 (TPM). Genes were deemed expressed when having a TPM value of at least 3. The condition-specificity,  
131 also known as tau<sup>96</sup>, of every gene was calculated as follows, where x is the TPM value per condition,  
132 max is the maximal expression of a gene and n is the number of conditions in the expression  
133 compendium:

$$134 \quad \tau = \frac{\sum(1 - (x/\max))}{n}$$

135 Condition-specific genes were defined as having a tau value of bigger or equal than 0.9.

136 The generated expression atlas for *P. tricornutum* was also used to define co-expression clusters. The  
137 Pearson correlation was calculated in a pairwise manner between all genes and the highest reciprocal  
138 rank (HRR)<sup>97,98</sup> was determined at 23, by maximizing the recovery of known GO annotations, while  
139 restraining the number of novel predictions. A cluster was defined for every gene based on this cut-off  
140 and GO enrichment using hypergeometric testing was run per cluster. Multiple hypothesis testing was  
141 constrained using Benjamini–Hochberg correction (q value < 0.05).

### 142 **Data availability**

143 All gene families, phylogenetic trees of horizontal descent and the dating of the HGT events within  
144 these gene families are available on Zenodo (<https://zenodo.org/record/3555201>).

### 145 **7. Author contributions**

146 E.V. wrote the manuscript, performed species topology delineation, HGT detection analysis, functional  
147 interpretation, metagenomic and population genomic analysis. T.D. aided in HGT delineation and  
148 performed co-expression analysis. C.M.O.-C. aided in population genomic analysis and performed  
149 expression analysis generation of *S. robusta*. K.V. supervised the project. All authors read, edited and  
150 approved the manuscript.

### 151 **8. Acknowledgments**

152 E.V. wants to acknowledge the funding obtained by the BOF project GOA01G01715.



## 9. References

1. Dean, P. *et al.* Transporter gene acquisition and innovation in the evolution of Microsporidia intracellular parasites. *Nat Commun* **9**, 1709 (2018).
2. Gonçalves, C. *et al.* Evidence for loss and reacquisition of alcoholic fermentation in a fructophilic yeast lineage. *eLife* **7**, (2018).
3. Schönknecht, G. *et al.* Gene Transfer from Bacteria and Archaea Facilitated Evolution of an Extremophilic Eukaryote. *Science* **339**, 1207–1210 (2013).
4. Krasovec, M. *et al.* Genome Analyses of the Microalga Picochlorum Provide Insights into the Evolution of Thermotolerance in the Green Lineage. *Genome Biol Evol* **10**, 2347–2365 (2018).
5. Ricard, G. *et al.* Horizontal gene transfer from Bacteria to rumen Ciliates indicates adaptation to their anaerobic, carbohydrates-rich environment. *BMC Genomics* **7**, 22 (2006).
6. Savory, F., Leonard, G. & Richards, T. The Role of Horizontal Gene Transfer in the Evolution of the Oomycetes. *Plos Pathog* **11**, e1004805 (2015).
7. Bowler, C. *et al.* The Phaeodactylum genome reveals the evolutionary history of diatom genomes. *Nature* **456**, 239 (2008).
8. Marchetti, A. *et al.* Ferritin is used for iron storage in bloom-forming marine pennate diatoms. *Nature* **457**, 467 (2008).
9. Husnik, F. & McCutcheon, J. Functional horizontal gene transfer from bacteria to eukaryotes. *Nat Rev Microbiol* **16**, 67 (2017).
10. Tsaousis, A. D. *et al.* Evolution of Fe/S cluster biogenesis in the anaerobic parasite Blastocystis. *Proceedings of the National Academy of Sciences* **109**, 10426–10431 (2012).
11. Kominek, J. *et al.* Eukaryotic Acquisition of a Bacterial Operon. **176**, 1356-1366.e10 (2019).
12. Stairs, C. W., Roger, A. J. & Hampl, V. Eukaryotic Pyruvate Formate Lyase and Its Activating Enzyme Were Acquired Laterally from a Firmicute. *Mol Biol Evol* **28**, 2087–2099 (2011).
13. Stairs, C. *et al.* Microbial eukaryotes have adapted to hypoxia by horizontal acquisitions of a gene involved in ridoquinone biosynthesis. *Elife* **7**, (2018).
14. Alexander, W. G., Wisecaver, J. H., Rokas, A. & Hittinger, C. T. Horizontally acquired genes in early-diverging pathogenic fungi enable the use of host nucleosides and nucleotides. *Proc Natl Acad Sci USA* **113**, 4116–4121 (2016).
15. Strese, Å., Backlund, A. & Alsmark, C. A recently transferred cluster of bacterial genes in *Trichomonas vaginalis* - lateral gene transfer and the fate of acquired genes. *BMC Evol Biol* **14**, 119 (2014).
16. Chou, S. *et al.* Transferred interbacterial antagonism genes augment eukaryotic innate immune function. *Nature* **518**, 98 (2015).
17. Harding, T., Roger, A. J. & Simpson, A. G. Adaptations to High Salt in a Halophilic Protist: Differential Expression and Gene Acquisitions through Duplications and Gene Transfers. *Frontiers in Microbiology* **8**, 944 (2017).
18. Foflonker, F., Mollegard, D., Ong, M., Yoon, H. S. & Bhattacharya, D. Genomic Analysis of Picochlorum Species Reveals How Microalgae May Adapt to Variable Environments. *Mol Biol Evol* **35**, 2702–2711 (2018).
19. Field, C. B., Behrenfeld, M. J., Randerson, J. T. & Falkowski, P. G. Primary production of the biosphere: integrating terrestrial and oceanic components. *Science* **281**, 237–240 (1998).
20. Winder, M. & Cloern, J. E. The annual cycles of phytoplankton biomass. *Phil. Trans. R. Soc. B* **365**, 3215–3226 (2010).
21. Janech, M. G., Krell, A., Mock, T., Kang, J.-S. & Raymond, J. A. Ice-binding proteins from sea ice diatoms (Bacillariophyceae). *J Phycol* **42**, 410–416 (2006).
22. Nakov, T., Beaulieu, J. & Alverson, A. Accelerated diversification is related to life history and locomotion in a hyperdiverse lineage of microbial eukaryotes (Diatoms, Bacillariophyta). *New Phytol* **219**, 462–473 (2018).
23. Armbrust, E. *et al.* The Genome of the Diatom *Thalassiosira Pseudonana*: Ecology, Evolution, and Metabolism. *Science* **306**, 79–86 (2004).

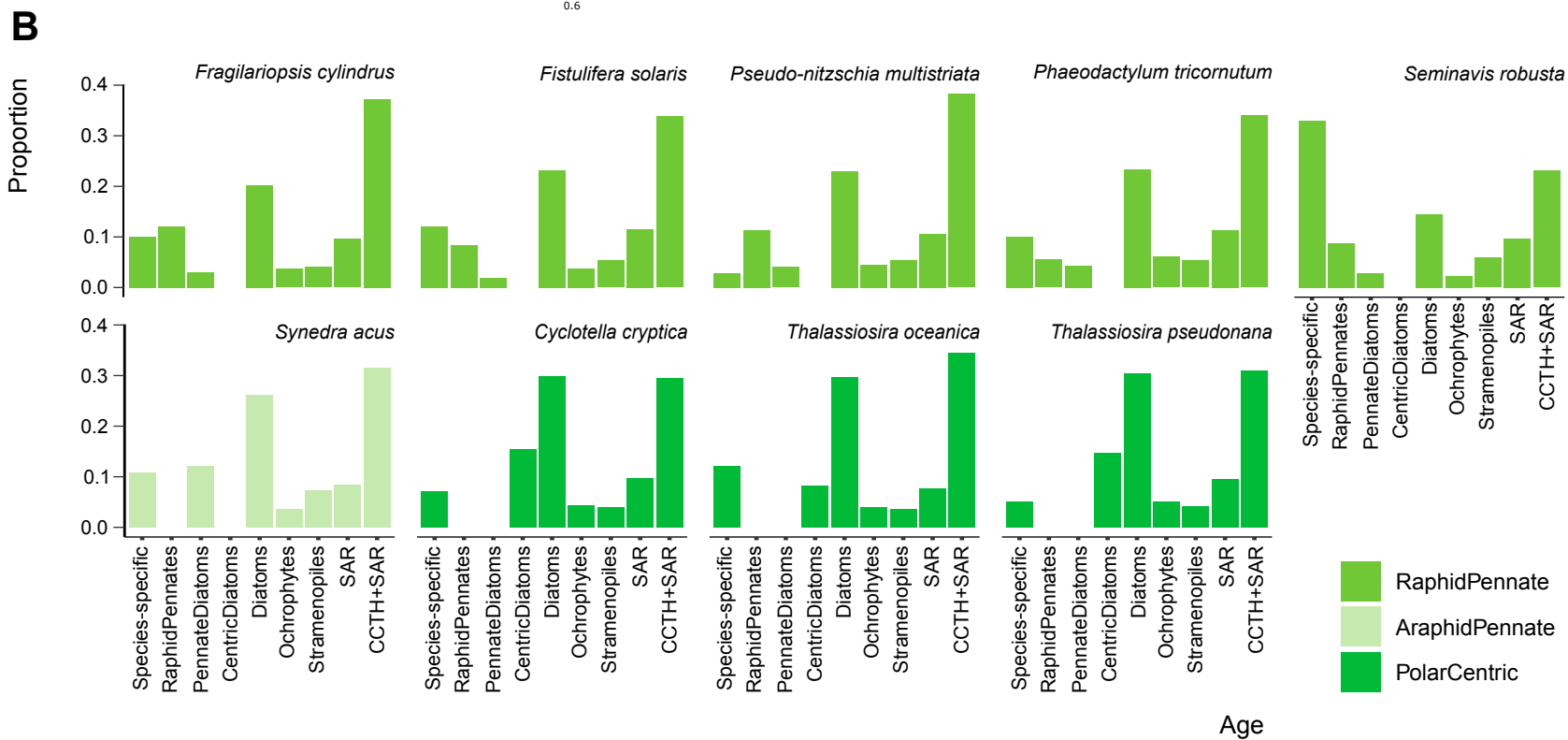
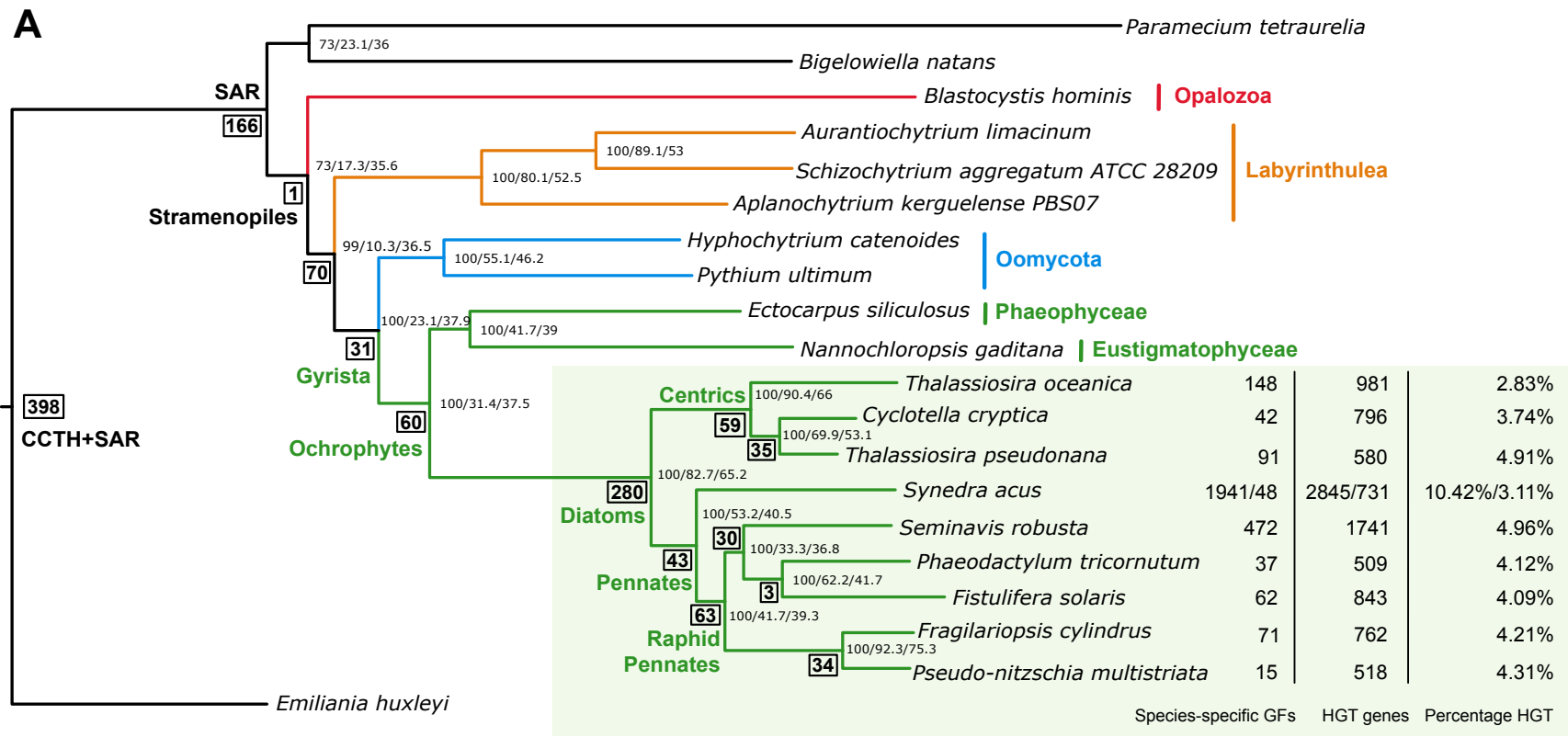
24. Lommer, M. *et al.* Genome and low-iron response of an oceanic diatom adapted to chronic iron limitation. *Genome Biol.* **13**, R66 (2012).
25. Traller, J. C. *et al.* Genome and methylome of the oleaginous diatom *Cyclotella cryptica* reveal genetic flexibility toward a high lipid phenotype. *Biotechnol Biofuels* **9**, (2016).
26. Galachyants, Y. *et al.* Sequencing of the complete genome of an araphid pennate diatom *Synedra acus* subsp. *radians* from Lake Baikal. *Dokl Biochem Biophys* **461**, 84–88 (2015).
27. Rastogi, A. *et al.* Integrative analysis of large scale transcriptome data draws a comprehensive landscape of *Phaeodactylum tricornutum* genome and evolutionary origin of diatoms. *Sci Rep* **8**, 4834 (2018).
28. Tanaka, T. *et al.* Oil Accumulation by the Oleaginous Diatom *Fistulifera solaris* as Revealed by the Genome and Transcriptome. *Plant Cell* **27**, 162–176 (2015).
29. Mock, T. *et al.* Evolutionary genomics of the cold-adapted diatom *Fragilariopsis cylindrus*. *Nature* **541**, 536 (2017).
30. Basu, S. *et al.* Finding a partner in the ocean: molecular and evolutionary bases of the response to sexual cues in a planktonic diatom. *New Phytol.* **215**, 140–156 (2017).
31. Richards, T. & Monier, A. A tale of two tardigrades. *Proc National Acad Sci* **113**, 4892–4894 (2016).
32. Zakharova, Yu., Adel'shin, R., Parfenova, V., Bedoshvili, Ye. & Likhoshway, Ye. Taxonomic characterization of the microorganisms associated with the cultivable diatom *Synedra acus* from Lake Baikal. *Microbiology+* **79**, 679–687 (2010).
33. Keeling, P. J. & Palmer, J. D. Horizontal gene transfer in eukaryotic evolution. *Nat Rev Genet* **9**, 605–618 (2008).
34. Eme, L., Gentekaki, E., Curtis, B., Archibald, J. & Roger, A. Lateral Gene Transfer in the Adaptation of the Anaerobic Parasite *Blastocystis* to the Gut. *Curr Biol* **27**, 807–820 (2017).
35. Murphy, C. L. *et al.* Horizontal Gene Transfer as an Indispensable Driver for Evolution of Neocallimastigomycota into a Distinct Gut-Dwelling Fungal Lineage. *Appl Environ Microbiol* **85**, e00988-19 (2019).
36. Vakirlis, N. *et al.* A Molecular Portrait of De Novo Genes in Yeasts. *Mol. Biol. Evol.* **35**, 631–645 (2018).
37. Croft, M. T., Lawrence, A. D., Raux-Deery, E., Warren, M. J. & Smith, A. G. Algae acquire vitamin B12 through a symbiotic relationship with bacteria. *Nature* **438**, 90–93 (2005).
38. Bertrand, E. M. *et al.* Phytoplankton–bacterial interactions mediate micronutrient colimitation at the coastal Antarctic sea ice edge. *Proc Natl Acad Sci USA* **112**, 9938–9943 (2015).
39. Heal, K. R. *et al.* Two distinct pools of B<sub>12</sub> analogs reveal community interdependencies in the ocean. *Proc Natl Acad Sci USA* **114**, 364–369 (2017).
40. Helliwell, K. E., Wheeler, G. L., Leptos, K. C., Goldstein, R. E. & Smith, A. G. Insights into the Evolution of Vitamin B12 Auxotrophy from Sequenced Algal Genomes. *Mol Biol Evol* **28**, 2921–2933 (2011).
41. Ellis, K. A., Cohen, N. R., Moreno, C. & Marchetti, A. Cobalamin-independent Methionine Synthase Distribution and Influence on Vitamin B12 Growth Requirements in Marine Diatoms. *Protist* **168**, 32–47 (2017).
42. Helliwell, K. E. *et al.* Cyanobacteria and Eukaryotic Algae Use Different Chemical Variants of Vitamin B12. *Curr. Biol.* **26**, 999–1008 (2016).
43. Helliwell, K. E. The roles of B vitamins in phytoplankton nutrition: new perspectives and prospects. *New Phytol* **216**, 62–68 (2017).
44. Campbell, G. R. O. *et al.* *Sinorhizobium meliloti* bluB is necessary for production of 5,6-dimethylbenzimidazole, the lower ligand of B12. *Proceedings of the National Academy of Sciences* **103**, 4634–4639 (2006).
45. Cohen, N. R. *et al.* Iron and vitamin interactions in marine diatom isolates and natural assemblages of the Northeast Pacific Ocean: Iron and vitamin interactions in diatoms. *Limnol. Oceanogr.* **62**, 2076–2096 (2017).

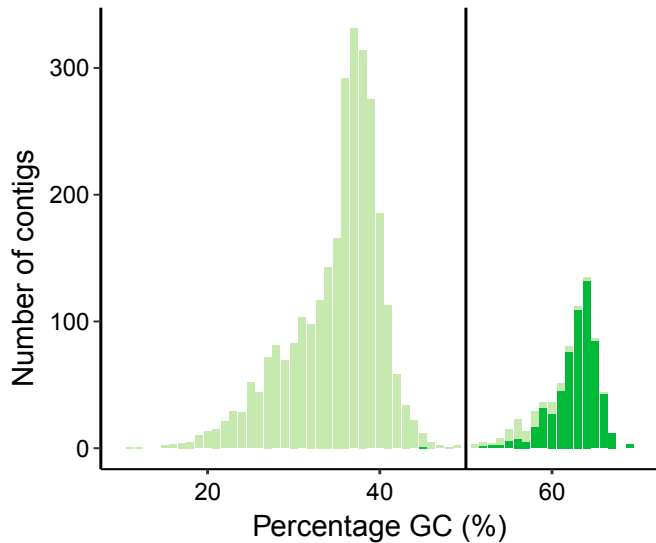
46. Cohen, N. R. *et al.* Iron storage capacities and associated ferritin gene expression among marine diatoms: Iron storage and ferritin expression in diatoms. *Limnol. Oceanogr.* **63**, 1677–1691 (2018).
47. Marchetti, A., Catlett, D., Hopkinson, B., Ellis, K. & Cassar, N. Marine diatom proteorhodopsins and their potential role in coping with low iron availability. *ISME J* **9**, 2745 (2015).
48. Fortunato, A. E. *et al.* Diatom Phytochromes Reveal the Existence of Far-Red-Light-Based Sensing in the Ocean. *Plant Cell* **28**, 616–628 (2016).
49. Montsant, A. *et al.* Identification and comparative genomic analysis of signaling and regulatory components in the diatom *Thalassiosira pseudonana*. *Journal of Phycology* **43**, 585–604 (2007).
50. Bayer-Giraldi, M. *et al.* Growth suppression of ice crystal basal face in the presence of a moderate ice-binding protein does not confer hyperactivity. *Proc Natl Acad Sci USA* **115**, 7479–7484 (2018).
51. Sorhannus, U. Evolution of Antifreeze Protein Genes in the Diatom Genus *Fragilariopsis*: Evidence for Horizontal Gene Transfer, Gene Duplication and Episodic Diversifying Selection. *Evol Bioinform Online* **7**, EBO.S8321 (2011).
52. Allen, A. *et al.* Evolution and metabolic significance of the urea cycle in photosynthetic diatoms. *Nature* **473**, 203 (2011).
53. Oh, J. *et al.* Diatom Allantoin Synthase Provides Structural Insights into Natural Fusion Protein Therapeutics. *ACS Chem. Biol.* **13**, 2237–2246 (2018).
54. Fabris, M. *et al.* The metabolic blueprint of *Phaeodactylum tricornutum* reveals a eukaryotic Entner–Doudoroff glycolytic pathway. *Plant J* **70**, 1004–1014 (2012).
55. Allen, A. E. *et al.* Evolution and Functional Diversification of Fructose Bisphosphate Aldolase Genes in Photosynthetic Marine Diatoms. *Mol Biol Evol* **29**, 367–379 (2012).
56. Whitaker, J. W., McConkey, G. A. & Westhead, D. R. The transferome of metabolic genes explored: analysis of the horizontal transfer of enzyme encoding genes in unicellular eukaryotes. *Genome Biol* **10**, R36 (2009).
57. Sun, G. & Huang, J. Horizontally acquired DAP pathway as a unit of self-regulation: Gene transfer and metabolic network. *Journal of Evolutionary Biology* **24**, 587–595 (2011).
58. Jiroutová, K., Horák, A., Bowler, C. & Oborník, M. Tryptophan Biosynthesis in Stramenopiles: Eukaryotic Winners in the Diatom Complex Chloroplast. *J Mol Evol* **65**, 496–511 (2007).
59. Rastogi, A. *et al.* A genomics approach reveals the global genetic polymorphism, structure, and functional diversity of ten accessions of the marine model diatom *Phaeodactylum tricornutum*. *ISME J* (2019) doi:10.1038/s41396-019-0528-3.
60. Yang, M., Lin, X., Liu, X., Zhang, J. & Ge, F. Genome annotation of a model diatom *Phaeodactylum tricornutum* using an integrated proteogenomic pipeline. *Mol Plant* **11**, 1292–1307 (2018).
61. Lehner, B. & Fraser, A. G. Protein domains enriched in mammalian tissue-specific or widely expressed genes. *Trends in Genetics* **20**, 468–472 (2004).
62. Freilich, S. *et al.* Relationship between the tissue-specificity of mouse gene expression and the evolutionary origin and function of the proteins. *Genome Biol.* **6**, R56 (2005).
63. Zhang, L. & Li, W.-H. Mammalian Housekeeping Genes Evolve More Slowly than Tissue-Specific Genes. *Mol Biol Evol* **21**, 236–239 (2004).
64. Sideri, T. C., Willetts, S. A. & Avery, S. V. Methionine sulphoxide reductases protect iron-sulphur clusters from oxidative inactivation in yeast. *Microbiology* **155**, 612–623 (2009).
65. Koski, L. B. & Golding, G. B. The Closest BLAST Hit Is Often Not the Nearest Neighbor. *J Mol Evol* **52**, 540–542 (2001).
66. Philippe, H. *et al.* Resolving Difficult Phylogenetic Questions: Why More Sequences Are Not Enough. *PLoS Biol* **9**, e1000602 (2011).
67. Olofsson, J. K. *et al.* Population-Specific Selection on Standing Variation Generated by Lateral Gene Transfers in a Grass. *Current Biology* **29**, 3921–3927.e5 (2019).

68. Koch, F. *et al.* The effect of vitamin B<sub>12</sub> on phytoplankton growth and community structure in the Gulf of Alaska. *Limnol. Oceanogr.* **56**, 1023–1034 (2011).
69. Emms, D. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biol* **16**, 157 (2015).
70. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
71. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**, 268–274 (2015).
72. Minh, B. Q., Hahn, M. & Lanfear, R. *New methods to calculate concordance factors for phylogenomic datasets.* <http://biorxiv.org/lookup/doi/10.1101/487801> (2018).
73. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat Methods* **12**, 59–60 (2015).
74. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
75. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002).
76. Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A. & Jermini, L. S. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods* **14**, 587–589 (2017).
77. Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Mol Biol Evol* **35**, 518–522 (2018).
78. Diner, R. E. *et al.* Diatom centromeres suggest a mechanism for nuclear DNA acquisition. *Proc Natl Acad Sci USA* **114**, E6015–E6024 (2017).
79. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification. *Bioinformatics* **30**, 1236–1240 (2014).
80. Ogata, H. *et al.* KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Res.* **27**, 29–34 (1999).
81. Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper. *Mol. Biol. Evol.* **34**, 2115–2122 (2017).
82. Gruber, A., Rocap, G., Kroth, P., Armbrust, E. & Mock, T. Plastid proteome prediction for diatoms and other algae with secondary plastids of the red lineage. *Plant J* **81**, 519–528 (2015).
83. Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in multidimensional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).
84. Proost, S. *et al.* i-ADHoRe 3.0—fast and sensitive detection of genomic homology in extremely large data sets. *Nucleic Acids Res.* **40**, e11–e11 (2012).
85. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
86. Villar, E. *et al.* The Ocean Gene Atlas: exploring the biogeography of plankton genes online. *Nucleic Acids Res.* **46**, W289–W295 (2018).
87. Carradec, Q. *et al.* A global ocean atlas of eukaryotic genes. *Nat Commun* **9**, 373 (2018).
88. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
89. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
90. Golicz, A. A. *et al.* Gene loss in the fungal canola pathogen *Leptosphaeria maculans*. *Funct Integr Genomics* **15**, 189–196 (2015).
91. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research* **20**, 1297–1303 (2010).
92. Van der Auwera, G. A. *et al.* From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best Practices Pipeline. *Current Protocols in Bioinformatics* **43**, (2013).

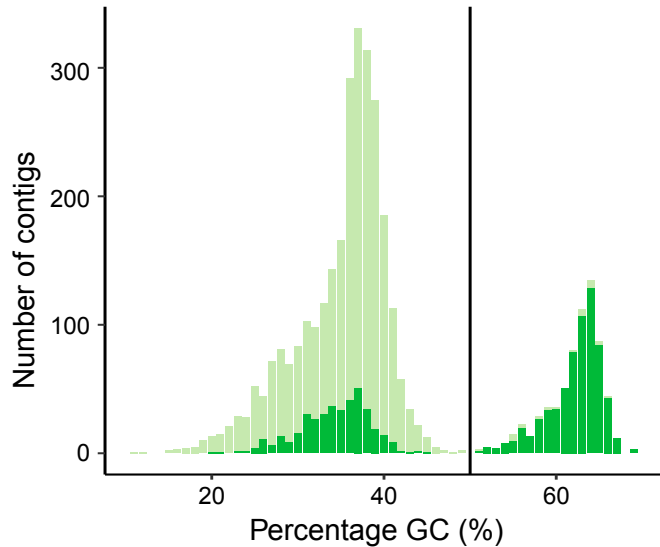
93. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
94. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w<sup>1118</sup>; iso-2; iso-3. *Fly* **6**, 80–92 (2012).
95. Vanechoutte, D. & Vandepoele, K. Curse: building expression atlases and co-expression networks from public RNA-Seq data. *Bioinformatics* **35**, 2880–2881 (2019).
96. Kryuchkova-Mostacci, N. & Robinson-Rechavi, M. A benchmark of gene expression tissue-specificity metrics. *Brief Bioinform* **18**, 205–214 (2016).
97. Liesecke, F. *et al.* Ranking genome-wide correlation measurements improves microarray and RNA-seq based global and targeted co-expression networks. *Sci Rep* **8**, 10885 (2018).
98. Mutwil, M. *et al.* Assembly of an Interactive Correlation Network for the Arabidopsis Genome Using a Novel Heuristic Clustering Algorithm. *Plant Physiol.* **152**, 29–43 (2010).

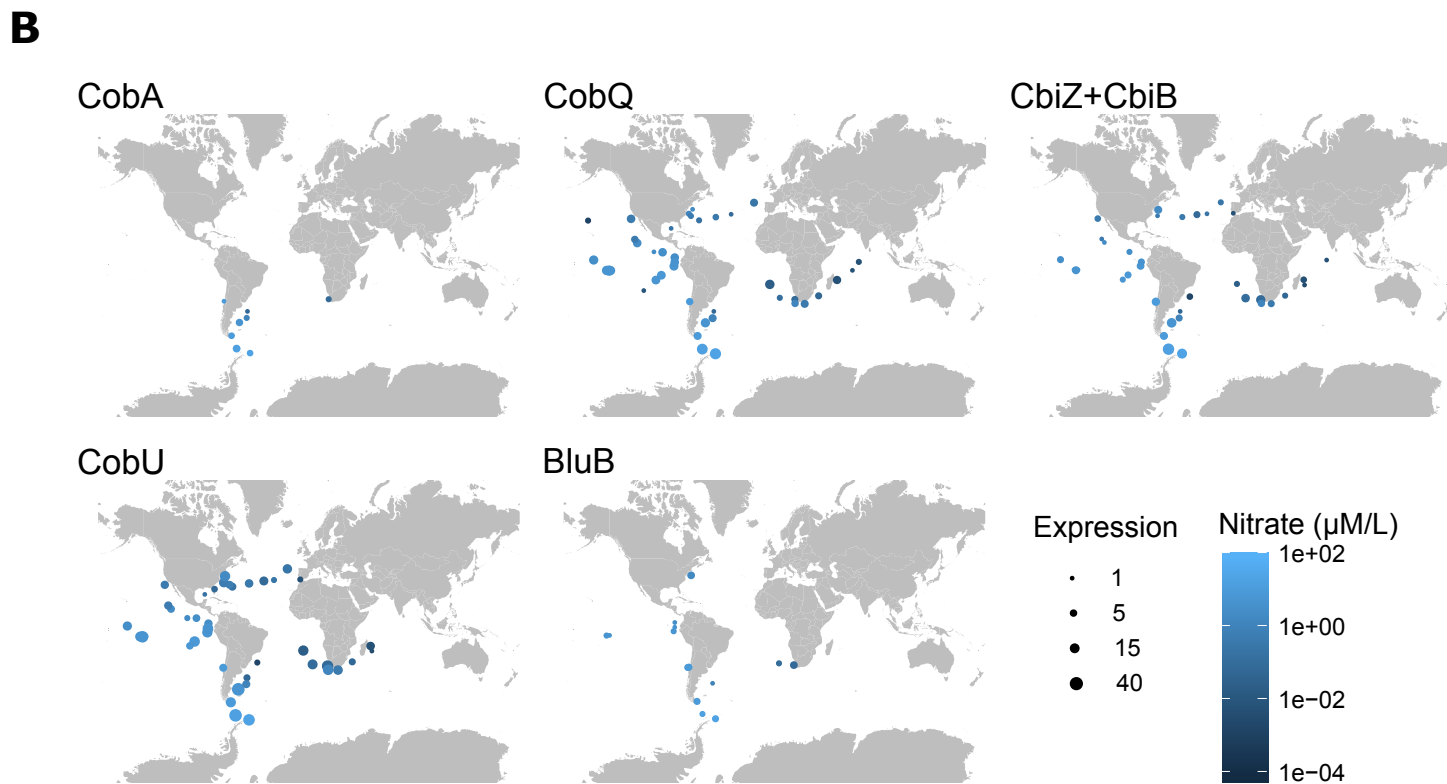
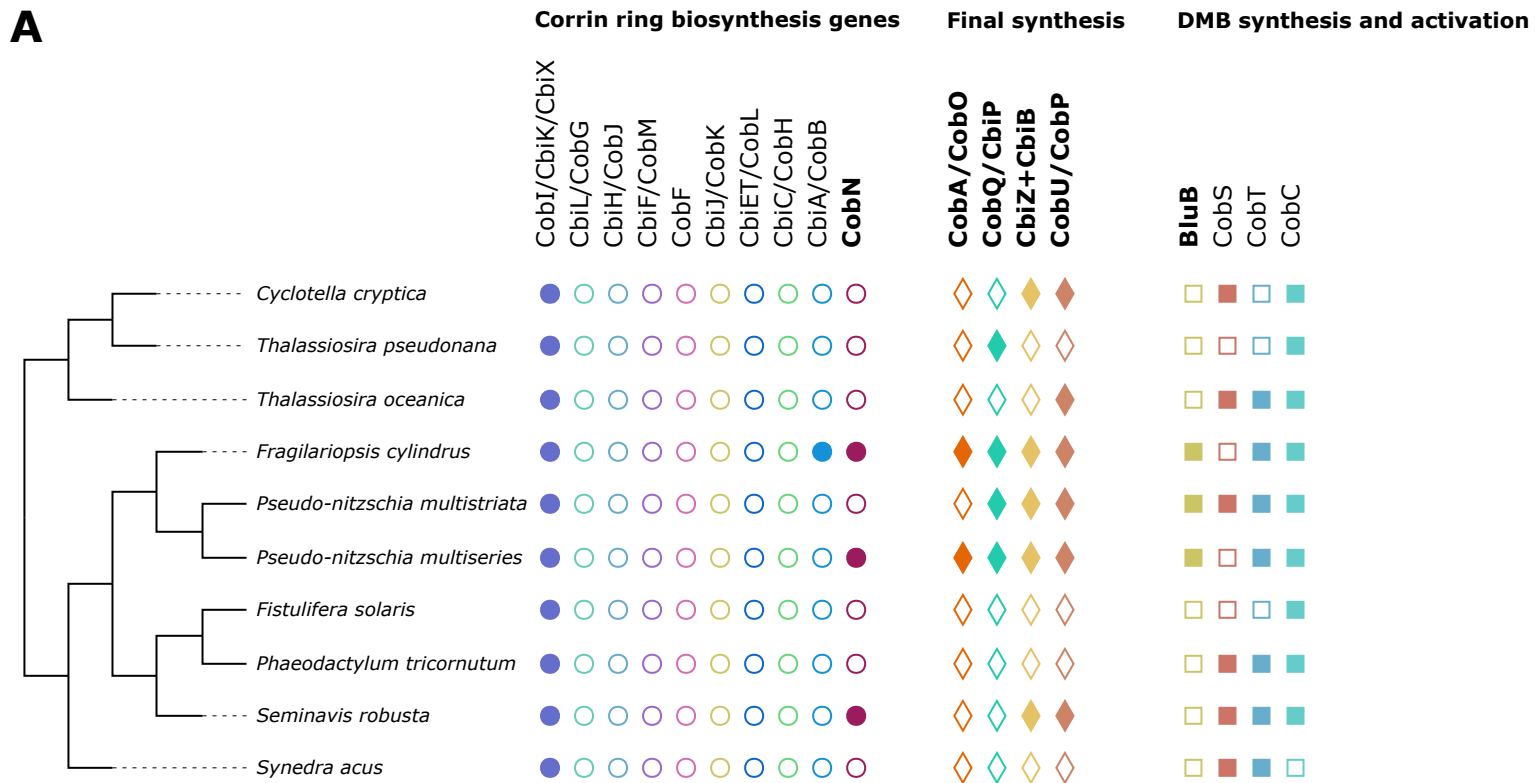


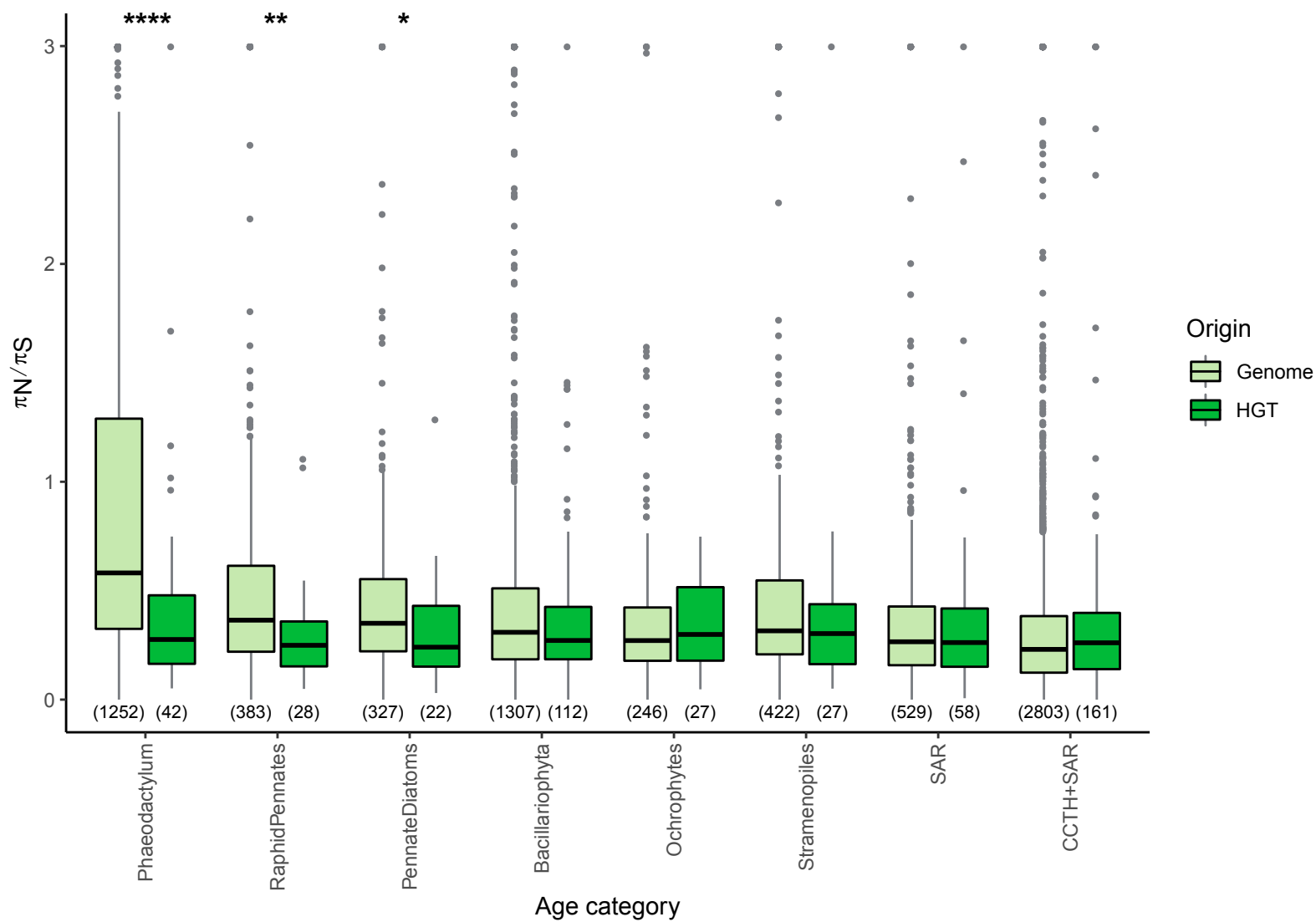


**A** Similarity to *Sphingomonadales*

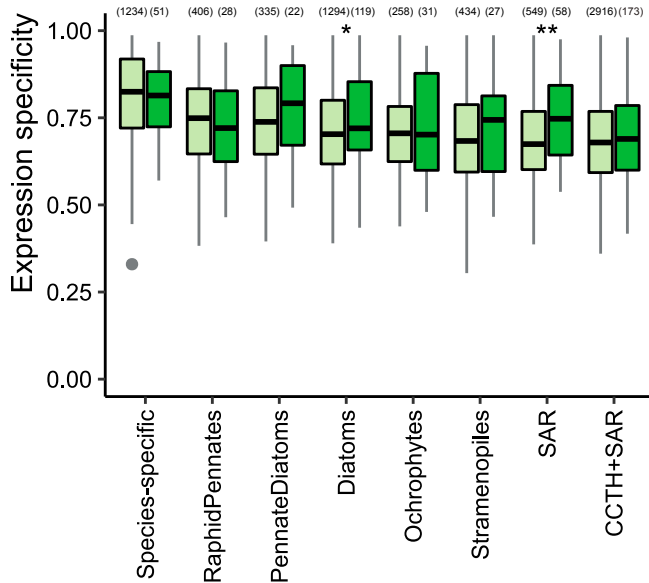
False True

**B** Presence of a predicted HGT gene

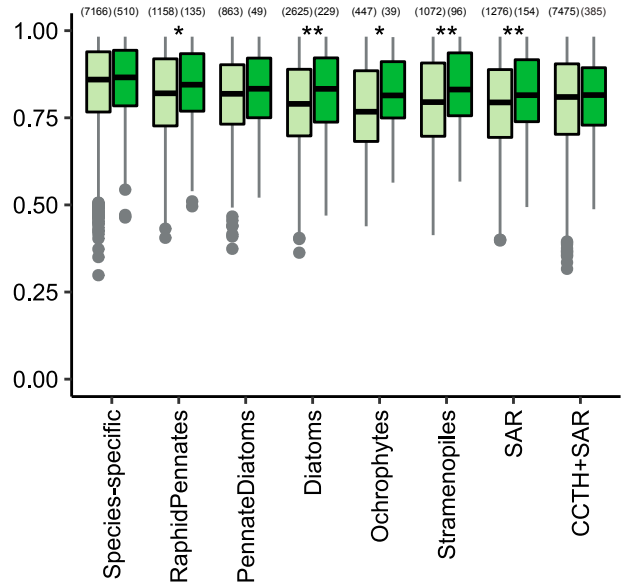




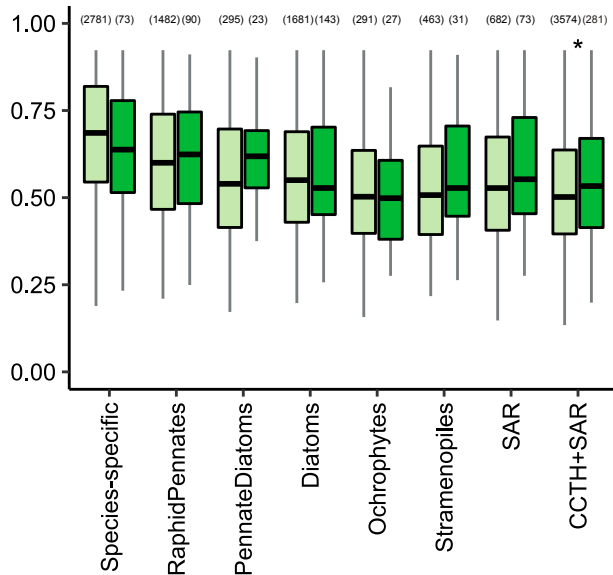
### *Phaeodactylum tricornutum*



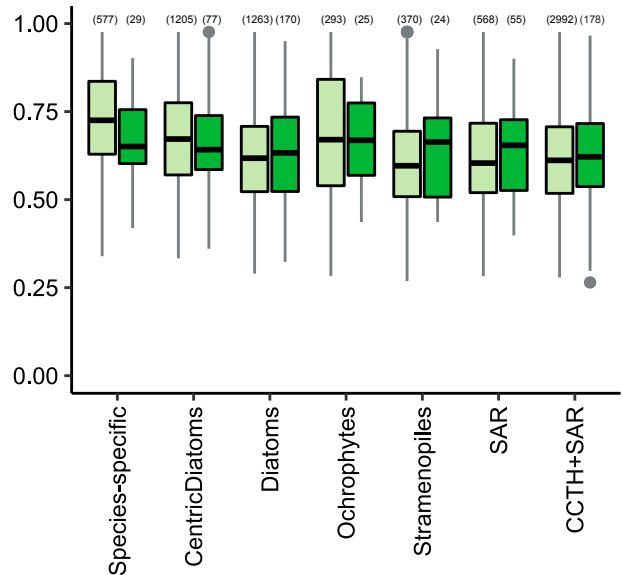
### *Seminaavis robusta*



### *Fragliariopsis cylindrus*



### *Thalassiosira pseudonana*



Origin Genome HGT

Age category