1  **Genome evolution of a non-parasitic secondary heterotroph, the diatom *Nitzschia***

2  ***putrida***

3

4  Ryoma Kamikawa[1], Takako Mochizuki[2], Mika Sakamoto[2], Yasuhiro Tanizawa[2],

5  Takuro Nakayama[3], Ryo Onuma[4], Ugo Cenci[5], Daniel Moog[6,*], Samuel Speak[7],

6  Krisztina Sarkozi[7], Andrew Toseland[7], Cock van Oosterhout[7], Kaori Oyama[8], Misako

7  Kato[8], Keitaro Kume[9], Motoki Kayama[10], Tomonori Azuma[10], Ken-ichiro Ishii[10],

8  Hideaki Miyashita[10], Bernard Henrissat[11,12,13], Vincent Lombard[11,12], Joe Win[14],

9  Sophien Kamoun[14], Yuichiro Kashiyama[15], Shigeki Mayama[16], Shin-ya Miyagishima[4],

10  Goro Tanifuji[17], Thomas Mock[7], Yasukazu Nakamura[2]

11

12  [1]Graduate School of Agriculture, Kyoto University, Kyoto 606-8502, Japan;

13  [2]Department of Informatics, National Institute of Genetics, Research Organization of

14  Information and Systems, Shizuoka 411-8540, Japan; [3]Graduate School of Life

15  Sciences, Tohoku University, Sendai 980-8578, Japan; [4]Department of Gene Function

16  and Phenomics, National Institute of Genetics, Shizuoka 411-8540, Japan; [5]Univ. Lille,

17  CNRS, UMR 8576 – UGSF – Unité de Glycobiologie Structurale et Fonctionnelle, F-

18  59000 Lille, France; [6]Laboratory for Cell Biology, Philipps University Marburg, Karl-

19  von-Frisch-Str. 8, SYNMIKRO Research Center, Hans-Meerwein-Str. 6, 35032,

20  Marburg, Germany; [7]School of Environmental Sciences, University of East Anglia,

21  Norwich Research Park, Norwich, UK; [8]Graduate School of Humanities and Sciences,

22  Ochanomizu University, Tokyo, Japan; [9]Department of Clinical Medicine, Faculty of

23  Medicine, University of Tsukuba, Ibaraki 305-8572, Japan; [10]Graduate School of

24  Human and Environmental Studies, Kyoto University, Kyoto 606-8501, Japan;

1    [11]Architecture et Fonction des Macromolécules Biologiques (AFMB), CNRS,

2    Université Aix-Marseille, 163 Avenue de Luminy, 13288 Marseille, France; [12]INRA,

3    USC 1408 AFMB, 13288, Marseille, France; [13]Department of Biological Sciences,

4    King Abdulaziz University, Jeddah, 21589, Saudi Arabia; [14]The Sainsbury Laboratory,

5    University of East Anglia, Norwich Research Park, Norwich, UK; [15]Graduate School of

6    Engineering, Fukui University of Technology, Fukui, Japan; [16]Advanced Support

7    Center for Science Teachers, Tokyo Gakugei University, Koganei, Tokyo, Japan;

8    [17]Department of Zoology, National Museum of Nature and Science, Tsukuba 305-0005,

9    Japan

10

11   Corresponding author (RK): Graduate School of Agriculture, Kyoto University,

12   Kitashirakawa oiwake cho, Sakyo ku, Kyoto, Kyoto 606-8502, Japan

13   Email: kamikawa.ryoma.7v@kyoto-u.ac.jp

14   *current address: Max Planck Institute for Terrestrial Microbiology, Karl-von-Frisch-

15   Str. 10, 35043, Marburg, Germany

16

17   **Abstract**

18   Secondary loss of photosynthesis is observed across almost all plastid-bearing branches

19   of the eukaryotic tree of life. However, genome-based insights into the transition from a

20   phototroph into a secondary heterotroph have so far only been revealed for parasitic

21   species. Free-living organisms can yield unique insights into the evolutionary

22   consequence of the loss of photosynthesis, as the parasitic lifestyle requires specific

23   adaptations to host environments. Here we report on the diploid genome of the free-

24   living diatom *Nitzschia putrida* (35 Mbp), a non-photosynthetic osmotroph whose

2

1    photosynthetic relatives contribute ca. 40% of net oceanic primary production.

2    Comparative analyses with photosynthetic diatoms revealed that a combination of genes

3    loss, the horizontal acquisition of genes involved in organic carbon degradation, a

4    unique secretome and the rapid divergence of conserved gene families involved in cell

5    wall and extracellular metabolism appear to have facilitated the lifestyle of a non-

6    parasitic, free-living secondary heterotroph.

7

8    **Main**

9    The loss of photosynthesis in photoautotrophs seems to be accomplished if such loss is

10    compensated by a competitive advantage arising from the availability of an extracellular

11    energy source. Some secondary heterotrophs have evolved as parasites (Freese and Lane

12    2017; Hadariová et al. 2018; Janouškovec et al. 2019), relying on sufficient resources

13    provided by their hosts. Well studied examples are the Apicomplexa (e.g., Kissinger et

14    al. 2002), which have lost photosynthesis secondarily. However, examples of loss of

15    photosynthesis found in free-living secondary heterotrophs are as common as those of

16    parasites (Kamikawa et al. 2015; Hadariová et al. 2018; Dorrell et al. 2019; Kayama et

17    al. 2020a; Kayama et al. 2020b), and thereby such examples provide novel insights into

18    evolutionary processes required to thrive without photosynthesis and independently of a

19    resource-providing host. Given that a parasitic lifestyle accelerates the rate of evolution

20    (*cf.* Red Queens hypothesis, Van Valen 1974) and of loss of conserved orthologous

21    genes (e.g., Sun et al. 2018), the genome analysis of a non-parasitic secondary

22    heterotroph can provide insights uncompromised by parasite-specific adaptations. The

23    diatom *Nitzschia putrida* distributing in the mangrove estuaries is the ideal model to test

24    these hypotheses because it is a rare example of a free-living secondary heterotroph (Li

3

1    and Volcani 1987; Kamikawa et al. 2015a) within the diverse group of largely

2    photoautotrophic diatoms (Field et al. 1998; Mann 1999). As several genomes of the

3    latter have recently become available including close phylogenetic relatives (Armbrust

4    et al. 2004; Bowler et al. 2008; Mock et al. 2012), a genome-based comparative

5    metabolic reconstruction of *N. putrida* promises to reveal novel insights into what is

6    required to thrive as a free-living secondary heterotroph. Here, we have analysed the

7    draft genome sequence of a non-photosynthetic, obligately heterotrophic diatom species

8    (Bacillariophyceae). We provide insight into evolutionary processes underpinning

9    lifestyle shifts from photoautotrophy to free-living heterotrophy in the context of the

10    surface ocean ecosystem.

11

12    **Results**

13    **Genome assembly**

14    K-mer-based GenomeScope analysis (Vurture et al. 2017) with 150 bp-long Illumina

15    short reads suggested the genome of *Nitzschia putrida* (Fig. 1A) to be diploid

16    (Supplementary Fig. 1A). To provide a high-quality genome with long-range contiguity,

17    PacBio sequencing (RSII platform) was performed resulting in $\geq$ 40-fold coverage. Due

18    to the confirmed diploid nature of the *N. putrida* genome, we have applied the Falcon

19    assembler and Falcon_unzip ver. 0.5 (Chin et al. 2016) to provide a first draft genome

20    of this species. Based on this assembly, we estimated a genome size of 35 Mbps,

21    including 87 scaffolds with an N50 of 860.9 kb. The longest scaffold was 3.8 Mbps. The

22    heterozygous regions of the genome (alternate contigs) estimated by the Falcon

23    assembler resulted in 12 Mbps, with an N50 of 121 kbps (Table S1). The Falcon

24    assembly was error corrected and polished by approximately 150-fold coverage of

4

1   Illumina short reads, which were subsequently used for generating the final assembly

2   with Pilon 1.2.2 (Walker et al. 2016) including manual curation.

3   According to the k-mer assessed diploid nature of the *N. putrida* genome, the

4   read coverage of the homozygous regions is approximately two-fold higher than the

5   read coverage for the heterozygous regions, suggesting the presence of diverged alleles

6   as previously identified in the genome of the photoautotroph diatom *Fragilariopsis*

7   *cylindrus* (Supplementary Figs. 1A & 1B). Thus, most of the diverged allelic variants

8   can be found in the heterozygous regions characterized by the presence of alternate

9   contigs (Supplementary Fig. 1B). On the basis of the analysis with Braker2 v.2.0.3

10  (Hoff et al. 2016), the *Nitzschia* genome comprises 15,003 and 5,767 inferred protein-

11  coding loci on the primary and alternate contigs, respectively (Table S1). Almost 40%

12  of loci in the genome of *N. putrida* appear to be characterised by diverged alleles. A

13  BUSCOv3 analysis (Waterhouse et al. 2018) revealed the genome to be complete at a

14  level of 90.1% based on the haploid set of genes.

15

16  **The loss of photosynthesis**

17  The haploid set of genes was used to reconstruct the nuclear-encoded plastid proteome

18  of *N. putrida* and therefore to reveal the extent of gene loss including key genes of

19  photosynthesis. A comparative analysis of the *N. putrida* plastome (Gruber et al. 2015)

20  with its photosynthetic counterparts revealed that more than 50% of nuclear encoded

21  plastid proteins have been lost (Fig. 1B). More than 500 orthogroups (Orthofinder,

22  Emms and Kelly 2015) of nuclear-encoded plastid proteins which are usually shared

23  between photosynthetic diatoms (Gruber et al. 2015) are missing in the predicted plastid

24  proteome of *N. putrida* (Fig. 1C). In the missing part of the plastid proteome were genes

1  encoding for proteins of light-harvesting antenna including fucoxanthin-chlorophyll *a*/*c*

2  protein (*fcp*), photosystem II and I (e.g. *psbA*, *psbC*, *psbO*, *psaA*, *psaB*, and *psaD*), the

3  cytochrome *b6*/*f* complex (e.g., *petA*), and carbon fixation (e.g. *rbcS*, *rbcL*) in addition

4  to genes of the Calvin cycle (e.g. phosphoribulokinase (*prk*)). Furthermore, a significant

5  number of key genes were missing for the biosynthesis of chlorophyll, carotenoids, and

6  plastoquinones (Fig. 1D).

7

8  Despite the loss of some of these key photosynthesis genes, unexpectedly, there is still a

9  significant number of genes left encoding common plastid metabolic pathways as

10  known from photosynthetic diatoms, including the generation of ATP by a nuclear-

11  encoded ATPase subunit (Kamikawa et al. 2015b). Almost all genes encoding for

12  plastid enzymes to synthesize essential amino acids are still encoded in the nuclear

13  genome of *N. putrida*. Furthermore, all genes of the heme pathway have been found,

14  and *N. putrida* appears to be able to synthesise riboflavin. The presence of plastid-

15  targeted transporters (Moog et al. 2020) enable the transport of phosphoenolpyruvate

16  (PEP), 3-phosphoglycerate (G3P), and dihydroxyacetone-phosphate (DHAP) across the

17  plastid membranes. Additionally, our genome-based reconstruction of plastid

18  metabolism detects the biosynthesis pathway for lipids and the ornithine cycle in *N.*

19  *putrida* (Fig. 1E; Supplementary Fig. 2). The latter has neither been reported in previous

20  transcriptome-based studies with this species (Kamikawa et al. 2017) nor in any other

21  secondary heterotrophs (e.g., Dorrell et al. 2019; Kayama et al. 2020).

22

23  **Communication between organelles and light-dependent gene expression**

1    The lack of $CO_2$ fixation in plastids of *N. putrida*, which reduces the amount of amino

2    acids, lipids and other metabolites to be synthesised, appears to be partially

3    compensated by the remodelling of metabolic interactions with mitochondria and

4    peroxisomes as well as by the keeping nitrogen recycling (Fig. 1; Fig. 2A;

5    Supplementary Fig. S3&S4). It appears that the non-photosynthetic plastid of *N. putrida*

6    still exchanges glutamine and ornithine, both of which are important intermediates of

7    the ornithin cycle. Indeed, all genes for the ornithine-urea cycle have been retained in

8    the *N. putrida* genome. The ornithine-urea cycle is indispensable for nitrogen recycling

9    in photosynthetic diatoms (Allen et al. 2011; Smith et al. 2019), and even after the loss

10   of photosynthesis, nitrogen recycling appears to be essential in *N. putrida* (Fig. 2A) due

11   to its osmotrophic lifestyle. Usually, the ornithine-urea cycle is tightly linked with TCA

12   cycle and/or photorespiration in photosynthetic diatoms (Allen et al. 2011; Smith et al.

13   2019). However, *N. putrida* is not likely to perform photorespiration (Fig. 2A). The

14   metabolic exchange with the peroxisome through glycolate likely has ceased as

15   phosphoglycolate phosphatase and peroxisomal glycolate oxidase are missing. Thus,

16   photorespiration is unlikely to take place in non-photosynthetic plastids of *N. putrid* due

17   to the lack of Ribulose 1, 5-bisphosphate carboxylase/oxygenase (RuBisCO) and other

18   key enzymes of the Calvin cycle (Fig. 1). Nevertheless, peroxisomes still appear to play

19   a role in *N. putrida* for the production of malate or glyoxylate, which feed into

20   respiratory pathways of the mitochondria to support ATP and NADPH production (Fig.

21   2A).

22

23   Light in photosynthetic organisms does not only play a significant role for

24   photosynthesis generating ATP and NADPH, it also regulates cell division, diel cycles

7

1    and different signalling processes unlike in many heterotrophic organisms (Ashworth et

2    al. 2013; Chauton et al. 2013; Smith et al. 2017). Thus, we identified remaining

3    photoreceptors and cell-cycle regulators and their effect on regulating diel cycles and

4    genome-wide light-dependent gene expression. Although we found that all the cyclin

5    and cyclin-dependent kinases (Huysman et al. 2010) were still encoded and expressed in

6    the genome of *N. putrida* (Fig. 2B; Supplementary Fig. S5A&S5B), we were unable to

7    identify a diel cycle in cell division over a growth period between 12 hours light and 12

8    hours dark condition and 48 hours darkness with cells previously acclimatised to a

9    periodic change of light and dark for 12 hours each (Fig. 2C). This suggests these cell-

10   cycle regulators potentially have neo/subfunctionalized and therefore have a different

11   redulatory role in *N. putrida* unrelated to the diel cycle. The loss of the transcription

12   factor *bHLH-1a* (RITMO1), which has been identified as a master regulator of diel

13   periodicity (Annunziata et al. 2019), corroborates our finding that *N. putrida* has lost the

14   ability to perform diel cycles. In addition, most of the other photoreceptors known from

15   photosynthetic diatoms have also been lost (Fig. 2B) such as the blue-light sensing

16   Aureochromes 1a/b, both of which are transcription factors responsible for

17   photoacclimation (Kroth et al. 2017). Despite the lack of light-dependent cell-cycle

18   regulation, a few remaining photoreceptors were identified including bHLH1b_PAS,

19   Aureochrome 1c, and Cryptochrome-DASH/CPF2 (Fig. 2B) (Coesel et al. 2009;

20   Ashworth et al. 2013). Basic ZIP transcription factors possessing potentially light-

21   sensitive PAS domains (bZIP-PAS) (Fortunato et al. 2016), were also identified in the

22   *N. putrida* genome such as homologues to bZIP6 and bZIP7 of *Phaeodactylum*

23   *tricornutum* (Rayko et al. 2010). The latter has been duplicated and diversified in *N.*

24   *putrida* (Supplementary Fig. S5C). The presence of bZIP-PAS protein in a heterotrophic

8

1    eukaryote is not unprecedented as some oomycetes, non-photosynthetic parasites, have

2    been reported to encode them in their genomes (e.g., Kong et al. 2020). Although their

3    role in regulating gene expression remains to be investigated in *N. putrida*, light still

4    appears to influences the expression of some genes in this heterotrophic species.

5    Comparative transcriptome analyses every four hours during a shift from a light phase

6    to darkness (Fig. 2D) revealed eight clusters characterised by different expression

7    patterns, and there is no cluster explicitly representing the light-dependent gene

8    expression patterns as seen in photosynthetic algae (e.g., Ashworth et al. 2013; Fujiwara

9    et al. 2020). However, one of the clusters suggested some genes were expressed only in

10   the mid light phase: cluster 7 containing 90 genes (0.6% total). EuKaryotic Orthologous

11   Groups (KOGs) analysis with these genes detected 44 genes with known functional

12   domains, and of them, 21 are responsible for substrate import and carbon metabolisms

13   (Supplementary Fig. S5D). However, the photoreceptor homologues above,

14   bHLH1b_PAS, Aureochrome 1c and Cryptochrome-DASH/CPF2, were not part of this

15   cluster, and there was no explicit trend in their gene expression patterns with respect to

16   changes between light and dark conditions.

17

18   **The acquisition of genes through horizontal gene transfer**

19   Recent studies estimated that 3-5% of genes in diatom genomes were acquired through

20   species-specific horizontal gene transfer (HGT) (Vancaester et al. 2020). For *N. putrida*,

21   we identified 73 genes potentially acquired via HGT based on phylogenetic tree

22   reconstruction. Three of these genes were only shared with photosynthetic and non-

23   photosynthetic species of the genus *Nitzschia*. Hence, they were likely acquired by the

24   last common ancestor of the *Nitzschia* genus (Supplementary Fig. S6&S7). These three

9

1    genes are one sulfate transporter and two serine hydrolases, which contribute to the

2    sulfate assimilation and degradation of various compounds, respectively. Of the

3    remaining 70 HGT genes, 25 were of bacterial origin and 34 appear to have been

4    originated in other microbial eukaryotes such as fungi. The remaining 11 genes were

5    either of viral or of ambiguous origin (Fig. 3A). As the majority of the HGT genes in *N.*

6    *putrida* were not identified in any other diatom species, it suggests that they play a role

7    in the heterotrophic lifestyle and therefore might have been essential for the transition

8    from photoautotrophy to secondary heterotrophy. This hypothesis is corroborated by the

9    fact that many of these genes (ca. 50%) are involved in catabolic processes such as

10   carbohydrate metabolism (such as glucose 1P dehydrogenase, GH27 α-galactosidase)

11   and degradation of amino acids (such as arginase). One HGT gene (trichothecene 3-O-

12   acetyltransferase) even appears to convey tolerance against mycotoxin (Khatibi et al.

13   2011) (Figs. 3B-3E). Thus, the acquisition of these genes may have provided the

14   ancestor of *N. putrida* with the potential to utilize new substrates and to compete fungi,

15   both of which can be considered functional traits essential for a successful transition to

16   becoming a secondary heterotroph. Only five HGT genes encode secreted proteins

17   according to our *in-silico* analysis (see below), which suggests that most of the others

18   contribute to intracellular metabolic processes.

19

20   **The genetic toolkit for the evolution of non-parasitic secondary heterotrophy**

21   Despite the loss of many nuclear genes and their families, the genome size of *N. putrida*

22   is not significantly different to photosynthetic relatives such as *Fragilariopsis cylindrus*,

23   *Phaedoctylum tricornutum* and the more distantly related diatom *Thalassiosira*

24   *pseudonana* (Table S1). By comparing KOGs of paralog proteins, there was no

1    significant difference in the number of unique KOG IDs between these four diatom

2    species (Fig. 3F&3G). However, when we compared the number of paralog proteins

3    assigned to each KOG ID, there were several KOG categories for which *N. putrida* had

4    a higher number of paralogous proteins compared to the other diatom species:

5    nucleotide transport (F), transcription (K), signal transduction (T), intracellular

6    trafficking, secretion, vesicular transport (U), and cytoskeleton (Z) (Fig. 3H). Even after

7    normalization by total gene numbers, nucleotide transport (F), signal transduction (T),

8    and cytoskeleton (Z) genes were more abundant in the *N. putrida* genome

9    (Supplementary Fig. S8).  This observation was corroborated by *N. putrida*-specific

10    enrichment of important Pfam domains in these functions such as Adenylate/Guanylate

11    cyclase and Cyclic nucleotide esterase, Leucine rich repeat (LRR), and

12    glycosyl/galactosyl transferase domains (Supplementary Fig. S8).

13

14    A microbial heterotroph either acquires nutrients by phagotrophy, the preferred

15    nutrition of many parasites, or by osmotrophy, which is the uptake of dissolved organic

16    compounds by osmosis as realised by bacteria and fungi, for instance (Richards and

17    Talbot 2013; Richards and Talbot 2018). As *N. putrida* grows well under axenic

18    conditions (Kamikawa et al. 2015; Ishii & Kamikawa 2017), it is likely an osmotroph,

19    dependent on the uptake of dissolved organic compounds across the silicified cell wall

20    and the plasma membrane. As realised by osmotrophic fungi, *N. putrida* may even be

21    able to degrade higher molecular weight compounds extracellularly to be subsequently

22    taken up as individual molecules by specific transporters or even osmosis (Richards and

23    Talbot 2013; Richards and Talbot 2018). Thus, it is likely that cell wall, membrane, and

24    secreted proteins diversified in *N. putirda* compared to photosynthetic diatoms to

11

1     facilitate osmotrophy. We analysed the enrichment of paralog proteins and differences

2     in nutrient transporters involved in the uptake of dissolved organic compounds such as

3     solute carriers.

4        Although the *N. putrida* genome does not differ in the abundance of nutrient

5     transporters relative to photosynthetic diatoms (Fig. 4A), we did find a significant

6     difference in the composition of these genes. For instance, the number of genes

7     encoding silicon transporters, solute symporters, and the resistance-nodulation-cell

8     division superfamily were more than twice as large in *N. putrida* compared to

9     photosynthetic diatom species (Fig. 4B; Supplementary Fig. S9).

10       Expansion of those gene families may, at least partly, have been achieved by

11     recent tandem duplications (Fig. 4C). To gain insight into when the expansion had

12     occurred, we performed a coalescence analysis, which revealed that silicon transporters

13     (SITs) in *N. putrida* began to expand around 3.3 Mya, while divergence from another

14     non-photosynthetic diatom *N. alba* is estimated to have occurred around 6.67 Mya

15     (Supplementary Fig. S10). Thus, their recent expansion suggests

16     neo/subfunctionalisation of the gene family in response to the change in lifestyle. The

17     divergence rate of SIT genes was much larger than that of control genes (e.g., myosin),

18     indicating that SIT diversification might have contributed to the adaptation of the

19     heterotrophic lifestyle. In support of this hypothesis, we detected several sites under

20     positive selection in different members of the SIT family (Table S2), which implies that

21     the evolution of those genes may have been driven by diversifying selection.

22

23       The solute sodium symporters are estimated to have diverged around 7.3 Mya,

24     markedly earlier than the SIT gene family. Although the divergence rate is also larger

1    than that of control genes (Supplementary Fig. S10), we did not find evidence of

2    diversifying selection in this gene family. The differences between these two families of

3    transporters suggest that their expansion might have occurred in a stepwise manner and

4    driven by different evolutionary forces.

5         Furthermore, although the overall carbohydrate-active enzymes (CAZyme)

6    family composition of *Nitzschia* was not remarkably different from that of

7    photosynthetic diatoms (Supplementary Fig. S11), families encoding β-glycoside

8    hydrolase (GH8), laminarinase (GH16_3), pectinase (GH28), β-glucanase (GH72), α-

9    mannan hydrolyzing enzymes (GH99), and β-1,2-glucan hydrolytic enzymes (GH114)

10   were also enriched in *N. putrida* compared to photosynthetic species (Fig. 4D).

11   Expansion of these families might, at least partly, have been achieved by recent tandem

12   duplications (Fig. 4E). Notably, more than one third of proteins assigned to the above

13   six families are predicted to be secreted in *N. putrida*.

14

15   **The predicted secretome of the non-parasitic, free-living secondary heterotroph *N.***

16   ***putrida***

17   Given that the secretome plays an important role for substrate degradation and

18   subsequent uptake of low-molecular weight compounds in osmotrophs (Richards and

19   Talbot 2013), we conducted a comparative analysis to predict secreted proteins of *N.*

20   *putrida in-silico* by idintifying proteins with N-terminal signal peptides and a lack of

21   transmembrane domains. The resulting proteins were clustered using TribeMCL

22   (Enright et al. 2002), and plastid- and lysosome-localised proteins were subsequently

23   removed using ASAFind according to their characteristic targeting motifs (Gruber et al.

24   2005) and Pfam domains. The number of putatively secreted proteins is 978, 998, 596,

13

1 and 718, in *N. putrida*, *F. cylindrus*, *P. tricornutum*, and *T. pseudonana*, respectively,

2 which corresponds to between 5 and 7% of the total number of genes in their genomes

3 (Supplementary Fig. S12A). Nevertheless, there were significant differences when we

4 compared the diversity of proteins for each between the four diatom species (Figs. 5A &

5 5B); *N. putrida* on average had a significantly higher number of proteins per tribe than

6 any of the other diatom species (Two-sided Wilcoxon signed rank test; $p < 0.01$; Fig.

7 5C). Especially, proteins involved in heterotrophy such as organic matter

8 degradation/modification including CAZymes and peptidases were more abundant in *N.*

9 *putrida* than in the photosynthetic diatom genomes (188 in *N. putrida*, 142 in *F.*

10 *cylindrus*, 118 in *P. tricornutum*, and 101 in *T. pseudonana*; Supplementary Fig. S12A).

11

12 The most common secreted proteins in *N. putrida* are leucine rich repeat-containing

13 (LRR) proteins (Supplementary Fig. 12B), many of which contain additional domains

14 such as tegument and glycoprotein domains, suggesting an increased functional

15 diversity (Fig. 5D). Interestingly, only very few LRR-containing proteins were

16 identified in the predicted secretomes of the photosynthetic diatoms, indicating that

17 signal-peptide-dependent secretion of abundant and diverse LRR-containing proteins

18 maybe an essential requirement in this secondary heterotroph. In addition to LRR-

19 containing proteins, the top ten most enriched proteins in *N. putrida* were Von

20 Willebrand factor type D involved in adhesion or clotting, two types of endopeptidases,

21 trypsin and leishmanolysin (cell surface peptidase of the human parasite *Leishmania*),

22 intradiol ring-cleavage dioxygenase protein for degradation of aromatic compounds,

23 methyltransferase, and four proteins with unknown function (Supplementary Fig. 12B).

24

14

1    Transcriptional dynamics of the predicted secretome over a diel cycle (Fig. 2) revealed

2    the presence of four different clusters. Genes in cluser 1 were transcribed at the

3    beginning of the first light phase and genes in cluster 2 at the end of the dark phase and

4    into the second light phase (Fig. 5E). Genes of clusters 3 were most strongly expressed

5    in the middle and end of the first light phase, whereas genes in cluster 4 were relatively

6    weakly expressed throughout day and night. These results suggest that some stimuli

7    including light conditions or nutrient conditions play a role in the regulation of these

8    genes, which might either be a relict from the photosynthetic ancestor or a response to

9    diel cycles of organic substances in the aquatic system occupied by *N. putrida*. As only

10   five of the secreted proteins potentially have been acquired via horizontal gene transfer

11   (Supplementary Fig. 7), the origin of the secretome in *N. putrida* most likely is derived

12   from a photosynthetic ancestor.

13

14   **Discussion**

15   *N. putrida* experienced a series of genetic adaptations towards a heterotrophic lifestyle.

16   *N. putrida* took a step backwards in one of the major evolutionary transitions, from

17   photoautotrophs to heterotrophs, potentially relaxing selection on some of the now

18   redundant gene networks and functions. As expected, more than 50% of nuclear

19   encoded plastid proteins have been lost in the *N. putrida* plastid proteome in

20   comparison to its photosynthetic counterparts (Gruber et al. 2015). However, the total

21   number of genes (~15,000) fell within the range of photosynthetic microalgae, and we

22   found no evidence of pseudogene formation, genome streamlining (e.g., Wolf and

23   Koonin 2013), gene family contraction (cf. birth-and-death hypothesis, Nei and Kumar

24   2000), or reductive genome evolution (Black Queen Hypothesis, Morris et al. 2012).

15

1    The relatively large genome size is not unexpected given that *N. putrida* is a free-living

2    osmotroph. This free-living lifestyle in a complex and highly variable coastal marine

3    environment likely is the reason why a significant number of genes including some

4    photoreceptors, cell-cycle regulators, and common plastid metabolic pathways present

5    in photosynthetic diatoms have remained. Even though some of the latter genes were

6    still expressed, *N. putrida* appears to lack a diel growth cycle, which suggests that these

7    cell-cycle regulators have neo/subfunctionalized. However, as a cetain number of genes

8    still appears to be regulated by light, osmotrophy potentially benefits from diel

9    fluctuations of resouces such as dissolved organic carbon in aquatic environments

10   (Ottesen et al. 2014; Aylward et al. 2015; Frischkorn et al. 2018). For photoautotrophs,

11   it is important to regulate the cell cycle in accordance with diel cycles for optimising

12   photosynthesis and therefore cell proliferation (Ashworth et al. 2013; Chauton et al.

13   2013; Ottesen et al. 2013; Smith et al. 2017; Hernández Limón et al. 2020). Without

14   being reliant on light as its primary energy source, the osmotroph *N. putrida* no longer

15   requires coordinating its cell cycle with diel cycles. Thus, after the loss of

16   photosynthesis, the strict light-dependent regulation of gene expression might have

17   become less important and gene expression therefore may have become predominantly

18   regulated by other stimuli. Indeed, many photoreceptors are missing but duplication of

19   genes for bZIP transcription factors with PAS domains and genes for signal

20   transduction and cellular regulatory roles such as adenyl/guanyl cyclase and cyclic

21   nucleotide esterase domains were enriched in the *N. putrida* genome. Furthermore, the

22   peroxisome-plastid interaction is no longer requried after the loss of photosynthesis

23   giving rise to loss of carbon fixation in the context of glycolate recycling. In contrast,

24   the ornithin-urea cycle likely remains to be functional to faciliate nitrogen recycling.

1

2  Gene family expansions and neo/subfunctionalizations appear to have played a

3  prominent role in the adaptation to its new lifestyle given that many proteins predicted

4  to be secreted have diversified in *N. putirda*, possibly to facilitate osmotrophy.

5  Altogether, the drastic change of lifestyle associated with the "devolution" in light of

6  autotrophic lifestyles did not result in the reductive genome evolution such as loss in

7  gene number, predicted from non-photosynthetic plastid-bearing parasites.

8

9  **Methods**

10  **Cultivation, DNA and RNA extraction, and sequencing**

11  *Nitzschia putrida* NIES-4239 was cultivated in the daigo' IMK medium (Wako)

12  including 1% Luria–Bertani medium based on the artificial seawater made with

13  MARINE ART SF-1 (Osaka Yakken Co.) at 20°C under the 12 hours light and 12 hours

14  dark conditions: 50 µmol photons/m$^2$/s with Plant cultivation LED light (BC-BML3,

15  Biomedical Science). DNA was extracted with Extrap Soil DNA Kit Plus ver. 2

16  (Nippon Steel). Total DNA was subjected to library construction with TruSeq DNA

17  PCR Free (350; Illumina) and to 151 bp paired-end sequencing by HiseqX, resulting in

18  660 million paired-end reads, and to PacBio RSII, with SMRT cell 8Pac V3 and DNA

19  Polymerase Binding Kit P6 v2, in Macrogen, resulting in 1.3 Gb subreads. Total RNA

20  was extracted with Trizol (Sigma) according to the manufacturer's instruction and was

21  subjected to library construction with TruSeq RNA Sample Prep Kit v2 (Illumina) and

22  101 bp paired-end sequencing by Hiseq2500, resulting in 107.5 million paired-end

23  reads.

24

1    **Genome assembly and construction of gene models**

2    PacBio reads were assembled into contigs using Falcon (ver. 0.7.0) with a length cutoff

3    of 7,000 bp for seed reads and an estimated genome size of 33 Mbp. Genome size

4    estimation was performed on the GenomeScope web server

5    (http://qb.cshl.edu/genomescope/) based on the K-mer frequency distribution of

6    Illumina reads calculated by JellyFish ver. 2.2.6 with a K-mer size of 21. The resultant

7    primary and associate contigs were then subjected to Falcon_unzip (ver. 0.5.0),

8    generating partially haplotype-phased contigs (primary contigs) and fully phased

9    contigs (haplotigs). The assembly was polished using PacBio reads and Quiver

10   program, followed by SNP and short indel error correction using Pilon (ver. 1.2.2) with

11   Illumina reads mapped by BWA (ver. 0.7.15). Indel errors in the vicinity of hetero

12   SNPs were further fixed manually, as they were difficult to be automatically corrected.

13   Contigs derived from plastid and mitochondrial genomes were identified using

14   BLASTN and separated from contigs derived from the nuclear genome.

15           RNA-seq reads were trimmed under the parameters of

16   ILLUMINACLIP:TruSeq3-PE.fa:2:30:10, LEADING:20 TRAILING:20,

17   SLIDINGWINDOW:4:15 and MINLEN:75 using Trimmomatic (ver. 0.36) (Bolger et

18   al. 2014). The trimmed reads aligned to the assembled contigs using HISAT2 (ver.

19   2.0.4) (Kim et al. 2015). They were provided to BRAKER2 gene annotation pipeline

20   (ver 2.0.3) as training data to be used for *ab initio* prediction of protein-coding genes.

21   Supplementarily, PASA (ver. 2.3.3) (Haas et al., 2003) was used to generate transcript-

22   based gene models by integrating *de novo* transcriptome assembly and genome-guided

23   assembly using Trinity (ver. 2.5.0) (Grabherr et al. 2011). The genome-guided assembly

24   used the mapping result from HISAT2 with --dta option. TransDecoder (ver. 5.0.2)

18

1    (Haas et al. 2013) was employed to extract protein coding regions from PASA result

2    with the alignment files from BlastP (ver 2.7.1) (Camacho et al. 2009) against UniRef90

3    with -evalue 1e-5 option and hmmscan (http://hmmer.org/, ver 3.1b2) against pfam (El-

4    Gebali et al. 2019) database.

5

6    The gene models that overlapped with the results from Braker were removed using

7    BlastP with evalue 1e-5 option, and remaining gene models were merged with the

8    Braker gene models to generate the final gene annotation. Transposable elements in the

9    NIES-4239 genome were searched by RepeatMasker (ver. 4.9.0) with using Dfam3.1

10    and RepBase-20170127 as reference repeat libraries.

11        The integrity of gene annotation was assessed by BUSCO (ver. 3.0.2) (Simão

12    et al. 2015) and the Eukaryota odb9 (ver. 2) dataset. The manipulation of sam/bam file

13    was used by samtools (ver. 1.9). The sequence files of gene region from gff file were

14    used by gffread. (ver. 0.9.11) (Pertea and Pertea 2020).

15        Organellar genome annotation was performed as described in previous studies

16    (Kamikawa et al. 2015b; Kamikawa et al. 2018).

17        Assembled genomes were deposited to DNA Data Bank of Japan under the

18    accession numbers BLYE01000001-BLYE01000234 for the nuclear genome,

19    LC600866 for the mitochondrial genome, and LC600867 for the plastid genome.

20

21    **Functional annotation**

22    The predicted protein coding genes were annotated using InterProScan, and RPS-

23    BLAST search was performed against KOG (EuKaryotic Orthologous Groups)

24    database. KO identifiers for Kyoto Encyclopedia of Genes and Genomes (KEGG)

19

1   metabolic pathways were assigned using KEGG Automatic Annotation Server (KAAS,

2   Moriya et al., 2004). Transporter proteins were annotated with TransportTP (Li et al.

3   2009). Reference proteome datasets for three photosynthetic diatom species were

4   obtained from the JGI Genome Portal: *Phaeodactylum tricornutum* CCAP 1055/1 v2.0

5   (Phatr2_bd_unmapped_GeneModels_FilteredModels1_aa.fasta,

6   Phatr2_chromosomes_geneModels_FilteredModels2_aa.fasta, 10,402 protein sequences

7   in total), *Thalassiosira pseudonana* CCMP 1335

8   (Thaps3_bd_unmapped_GeneModels_FilteredModels1_aa.fasta,

9   Thaps3_chromosomes_geneModels_FilteredModels2_aa.fasta, 11,776 sequences),

10  *Fragilariopsis cylindrus* CCMP 1102 (Fracy1_GeneModels_FilteredModels3_aa.fasta,

11  21,066 sequences). KEGG and KOG annotation was performed with them in the same

12  manner as NIES-4239. For plastid and mitochondrial genomes assembled in the above

13  procedures, annotation was performed by MFANNOT (insert ref).

14

15  **CAZyme annotation**

16  We performed a manual annotation of CAZymes (Lombard et al. 2014) using a

17  combination of BLAST (Altschul et al. 1997) and HMM searches (Mistry et al. 2013),

18  similar to that done previously (Curtis et al. 2012; Cenci et al. 2018). To assess the

19  similarity between the CAZyme family profiles of the two species, we generated heat

20  maps derived from an average linkage hierarchical clustering based on Bray-Curtis

21  dissimilarity matrix distances and Ward's method (Bray and Curtis 1957; Ward 1963;

22  Cenci et al. 2018). The heat maps were computed with Rstudio software

23  (https://www.rstudio.com/) using vegan in the R package

1   (http://cc.oulu.fi/~jarioksa/softhelp/vegan/html/vegdist.html) (Oksanen 2014) with

2   vegdist and hclust commands.

3

4   **Annotation of Cyclins, Cyclin-dependent kinases, bZIP transcription factors, and**

5   **photoreceptor proteins**

6   Cyclins, cyclin-dependent kinases, transcription factors were retrieved from the results

7   of Pfam annotation (above). In addition, those proteins were surveyed by homology-

8   based search with homologues of *P. tricornutum* and *T. pseudonana* (Huysman et al.

9   2010; Rayko et al. 2010) as queries. We also specifically surveyed photoreceptor

10   proteins with diel cycle-based expression, according to the list of Annunziata et al.

11   (2019). Detected sequences were again confirmed as proteins of our interest by

12   reciprocal blastP search against non-redundant database of GenBank. For Cyclins,

13   CDKs, and bZIP transcription factors in *N. putrida*, *P. tricornutum*, and *T. pseudonana*

14   were subjected to MAFFT (version 7.394; Kato and Standley 2013), followed by

15   removal of ambiguously aligned sites by Bioedit (Hall 1999). The resultant datasets

16   were subjected to the maximum likelihood analysis with IQ-TREE (Nguyen et al. 2015)

17   under the LG+$\Gamma$+F model with 100 non-parametric bootstrap analyses.

18

19   **Annotation of peroxisomal proteins**

20   To identify potential peroxisomal proteins, the protein models of the diatom were

21   screened via KEGG annotation and BLAST analyses. First, annotation of all the protein

22   sequences was performed via GhostKOALA (KEGG Automatic Annotation Server;

23   https://www.kegg.jp/ghostkoala/). Peroxisomal protein candidates were subsequently

24   identified using the "KEGG Mapper – Reconstruct Pathway" tool

1    (https://www.genome.jp/kegg/tool/map_pathway.html). In addition, all the *Nitzschia*

2    proteins were screened for the presence of a peroxisomal targeting signal of type 1

3    (PTS1, a C-terminal tri-peptide) using a local command line script identifying those

4    entries, which contain the amino acids [SAC][KRHS][LM] within the last three

5    positions of a protein sequence. Detected proteins were then functionally annotated with

6    WebMGA (http://weizhongli-lab.org/metagenomic-analysis/server/kog/) as well as

7    BlastKOALA (https://www.kegg.jp/blastkoala/) and their sequences were further

8    investigated for the presence of additional targeting signals using SignalP 3.0

9    (http://www.cbs.dtu.dk/services/SignalP-3.0/) and 4.1

10    (http://www.cbs.dtu.dk/services/SignalP-4.1/), TargetP 1.1

11    (http://www.cbs.dtu.dk/services/TargetP-1.1/index.php), PredSL

12    (http://aias.biol.uoa.gr/PredSL/), Predotar (https://urgi.versailles.inra.fr/predotar/) and

13    TMHMM 2.0 (http://www.cbs.dtu.dk/services/TMHMM/). To identify factors for

14    peroxisomal biogenesis/maintenance (peroxins) as well as for photorespiration and the

15    glyoxylate cycle, manual BLAST analyses against the *Nitzschia* proteins were

16    conducted using protein queries from the diatoms *Phaeodactylum tricornutum* and

17    *Thalassiosira pseudonana,* the cryptophyte *Guillardia theta* (Davis et al. 2017,

18    Gonzalez et al. 2011, Mix et al. 2018) as well as yeast (peroxin identification) and an e-

19    value cut-off of e-4. For verification, identified candidates were analyzed via reciprocal

20    BlastP against the NCBI nr database (https://blast.ncbi.nlm.nih.gov) and using NCBI

21    Conserved Domain Search (https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi) for

22    identification of conserved domains within the protein sequences. Missing proteins

23    present in peroxisomes of many organisms were especially surveyed in the

24    transcriptome data (see above).

1

2 **Annotation of mitochondrial proteins**

3 To identify mitochondrial proteins, we retrieved mitochondrial reference pathways and

4 proteins of *Homo sapiens*: from a stand-alone reactome server (Sidiropoulos et al. 2017;

5 Fabregat et al. 2018a; Fabregat et al. 2018b; Jassal et al. 2020) on our systems, Uniprot

6 (UniProt Consortium 2019) ids of proteins which are annotated that its subcellar

7 locations are 'mitochondrial matrix', 'mitochondrial inner membrane', 'mitochondrial

8 outer membrane' or 'mitochondrial intermembrane space' were extracted using the

9 cypher query language. To search and annotate for mitochondrial proteins, we

10 performed PSI-BLAST (Altschul et al. 1990; Altschul et al. 1997; Camacho et al. 2009)

11 search with proteins which have above uniprot ids as queries for the assemblies of

12 genome and transcriptome. As validation of those annotations, we performed three

13 analyses: 1) psi-blast search with above blast-hit sequences for swissprot (UniProt

14 Consortium 2019), 2) KEGG orthology (KO) number assignment by kofamscan

15 (Aramaki et al. 2020), 3) mitochondrial import signal analysis by Mitofates (Fukasawa

16 et al. 2015) and Nommpred (Kume et al. 2018). We manually refined our annotations

17 based on these results. Missing proteins present in mitochondria of many organisms

18 and/or in each metabolic pathway were especially surveyed in the transcriptome data

19 (see above).

20

21 **Annotation of plastid proteins**

22 To identify plastid proteins, protein sequences were subjected to SignalP4.0 (Petersen et

23 al. 2011) followed by ASAFind (Gruber et al. 2015). Functions and related metabolisms

24 for proteins predicted to be localized in the plastid were estimated by their KEGG ID

1  (see above) and KEGG mapper (Kanehisa et al. 2012). Missing proteins in each

2  metablic pathway were especially surveyed in the transcriptome data (see above).

3  Plastid proteins in the photosynthetic diatoms *Phaeodactylum tricornutum* and

4  *Thalassiosira pseudonana* have been already published in Gruber et al. (2005). We

5  retrieved them and clustered with the plastid proteins of *N. putrida* as orthogroups by

6  Orthofinder with default settings (Emms and Kelly 2015).

7

8  **Secretome analysis**

9  We detected protein sequences with N-terminal signal peptides and no internal

10  transmembrane domain by evaluation with signalP4.0 (Petersen et al. 2011) and

11  HMTMM (Sonnhammer et al., 1998; Krogh et al., 2001), respectively. Those sequences

12  highly likely comprise secretome and plastid sequences. The sequences were clustered

13  by TribeMCL (Enlight et al. 2002) into Tribes, with homology under e-30 criterion. If a

14  tribe comprises sequences predicted to be localized in plastids with the "high

15  confidence" category by ASAFind (Gruber et al. 2015), more than or equal to those

16  predicted to be non-plastidal or localized in plastids with the "low confidence" category,

17  we removed the tribe. Functional categorization was performed with domain annotation

18  by Pfam, and we also removed tribes if the included sequences were predicted to have a

19  domain for apparent organellar proteins such as plastid protein translocons, plastid

20  heme biosynthesis, components of photosystems, organellar transporters, and lysosomal

21  proteins by KO definition and/or Pfam. To evaluate whether the above procedure to

22  detect secreted proteins is appropriate for diatoms, we made a benchmark dataset of

23  diatoms including 21 secreted proteins experimentally confirmed (Bruckner et al. 2011;

24  Buhmann et al. 2016; Lachnit et al. 2019; Dell'Aquila et al. 2020) and 62 of their

1   homologues, 182 plastid proteins, 29 mitochondrial proteins, and 29 proteins localized

2   in other compartments such as cytosol and nucleus. From the benchmark set, our

3   method identified 77 secreted proteins as secretome proteins, indicating 92.8% recovery

4   rate. But no non-secreted proteins was not identified as secretome proteins in our

5   benchmark set, suggesting that our secretome dataset for whole protein data in the four

6   diatoms less likely comprise high propotion of non-secreted proteins. In addition, the

7   benchmark dataset included 17 secreted protein sequences, experimentally verified, of

8   *Phaeodactyum*, and of them, 13 protein sequences were retrieved by our genome-wide

9   secretome analysis. Four proteins not identified were appeared to lack signal peptides

10  detactable either SignalP4.0 and SignalP3.0, indicating those sequences might be

11  secreted by signal peptide-independnet ways.

12        Difference in the distributions of the number of protein sequences in each

13  secretome tribe among species was tested by the Wilcoxon signed rank test. *p*-values for

14  three pairs of comparisons that include *N. putrida* as a counterpart were adjusted by the

15  Benjamini–Hochberg procedure for multiple-testing correction. The Wilcoxon signed

16  rank test and the Benjamini–Hochberg procedure were conducted by using 'wilcoxon'

17  function implemented in SciPy library (version 1.4.1) and 'multipletests' function

18  implemented in Statsmodels library (version 0.11.1) for Python, respectively.

19        For LRR, additional domains were also searched by NCBI Conserved Domain

20  Search. All the possible secreted LRR proteins in *N. putrida*, *P. tricornutum*, and *T.*

21  *pseudonana* were subjected to MAFFT (version 7.394; Kato and Standley 2013) and

22  ambiguously aligned sites were removed by Bioedit (Hall 1999). The resultant dataset

23  was subjected to IQ-TREE (Nguyen et al. 2015) under the LG+$\Gamma$+F model with 100

24  non-parametric bootstrap analyses.

1

## Horizontal gene transfers

3  Total *N. putrida* protein sequences were subjected to Orthofinder with default settings

4  (Emms and Kelly 2015) together with whole protein sets of the photosynthetic diatom

5  genomes *P. tricornutum*, *F. cylindrus*, and *T. pseudonana*, resulting in 4,920 proteins

6  unique to *N. putrida*. The 4,920 proteins were subjected to similarity search using

7  BLASTP (BLAST+ 2.6.0) against GenBank non-redundant protein sequence database

8  and protein sequence database provided by the Marine Microbial Eukaryote

9  Transcriptome Sequencing Project (https://github.com/dib-lab/dib-MMETSP; Johnson

10  et al. 2019) with an e-value threshold of 1.0E-15. BLASTP hits from two databases

11  were combined and sorted by their bit-scores. Top ten hits with higher bit-scores for

12  each query sequence were then inspected their taxonomic composition with two criteria:

13  1) if those hits do not include any diatom sequences, namely, sequences from diatom

14  species or dinoflagellate species bearing diatom-derived plastids, 2) if the ten hits do not

15  include the diatom sequences except sequences from diatoms belonging to genus

16  *Nitzschia*. Query sequences, whose top ten BLASTP hits match the two criteria, were

17  considered as initial candidates for 1) genes derived from *N. putrida* specific horizontal

18  gene transfer (HGT) event and 2) *Nitzchia*-lineage-specific HGT, respectively. Each of

19  candidates with HGT origin was aligned together with all of the hit sequences in

20  aforementioned BLASTP analyses by MAFFT (version 7.394; Katoh and Standley

21  2013). After ambiguously aligned positions were removed by BMGE (version 1.12;

22  Criscuolo and Gribaldo 2010), each single-gene alignment was subjected to a

23  maximum-likelihood (ML) phylogenetic analysis by IQ-TREE (Nguyen et al. 2015)

24  under LG+Γ+F model. The statistical supports for the bipartitions in ML trees were

1  assessed by a nonparametric ML bootstrap analysis (100 replicates) under the same

2  substitution model used for the ML trees. All topologies of ML trees for both of *N.*

3  *putrida*-specific HGT and *Nitzchia*-lineage-specific HGT candidates were manually

4  assessed.

5

6  **Transcriptome analyses**

7  The cells of *N. putrida* NIES-4239 cultivated under the 12 hours light and 12 hours dark

8  conditions under the same condition as described above were then cultivated under the

9  dark conditions for 48 hours. The cells possibly acclimated to the dark condition were

10  transferred to the fresh liquid medium and incuvated under the 12 hours light and 12

11  hours dark condition for 25 hours in which the last one hour was of the second light

12  term. RNA was extracted every 4 hours at the one hour, 5 hours, and 9 hours after

13  initiation of the first light condition, one hour, 5 hours, and 9 hours after initiation of the

14  dark condition, and one hour after initiation of the second light term. The cultivation

15  and RNA extraction were performed twice, samples which were called Replicates A and

16  B. The extracted total RNAs were subjected to 151 bp paired-end sequencing by

17  NextSeq 500 according to the manufacturer's instruction in Bioengineering Lab, Japan,

18  resulting in 9.8-13 million paired-end reads for each sample. Adapter trimming and

19  quality filtering were performed with fastX toolkit

20  (http://hannonlab.cshl.edu/fastx_toolkit/). In quality filtering, reads with quality scores

21  > 20 for at least 75% of their length were retained, resulting in 7.3-11.1 million paired-

22  end reads. To obtain gene expression scores, one side of the paired-end reads was

23  mapped to the reference by Bowtie2 ver. 2.3.4.1 (Langmead and Salzberg 2012).

24  SAMtools ver. 1.8 (Li et al. 2009), BEDtools ver. 2.19.1 (Quinlan and Hall 2010), and

27

1    R ver. 3.5.3 (Ihaka and Gentleman 1996) were used to calculate the reads per kilobase

2    of exon per million mapped reads (RPKM). To extract the reproducible change in

3    transcriptome of Replicate A and B, the genes with low value of Pearson's correlation

4    coefficient between two biological replicates ($\leqq$ 0.9) were omitted, resulting in 1,971

5    genes. To investigate expression patterns of 1,971 genes, the values of RPKM + 1 of

6    each gene were transformed to log2 and centered by median values, and $k$-means

7    clustering ($k$ = 8) was performed by using Cluster 3.0 (de Hoon et al. 2004) based on

8    Spearman correlation and complete linkage. The clustered expression patterns were

9    visualized by Java TreeView (Saldanha 2004) and R.

10         The transcriptome data were deposited in DNA Data Bank of Japan under the

11    accession number PRJDB11016.

12

13    **Divergence time estimation**

14    All sequences of interest in diatoms were first clustered using cd-hit/4.6.8, removing

15    sequences of 100% alignment and clustering at low identity. The largest clusters for

16    each gene family was then aligned using PRANK 170427 (Löytynoja, 2014), poorly

17    aligned sequences were removed with TrimAl 1.2 (Capella-Gutiérrez, Silla-Martínez

18    and Gabaldón, 2009) and put into gene blocks to remove gapped columns with Gblocks

19    v0.91b (Castresana 2000; Talavera and Castresana 2007). Aligned sequences containing

20    all species with available data (N = 1 to 4) were then analysed using phylogenetic

21    analyses. Divergence estimates were obtained using Bayesian Markov Chain Monte

22    Carlo (MCMC) analyse is implemented in Beast 2.5.0 (Bouckaert et al. 2019). The

23    analysis was carried out using a relaxed molecular clock approach with an uncorrected

24    log-normal distribution model of rate of variation, the HKY substitution model, four

1    gamma categories and a Yule model of speciation, for each gene family 5 runs were

2    carried out with 20 million MCMC generations sampled every 1,000[th] generation. The

3    results from all runs were combined and summarised using LOGcombiner v1.10.4,

4    these were checked for convergence using Tracer v1.7.1 (Rambaut, 2009). A maximum

5    clade credibility tree was generated using Tree Annotator v1.10.4 and was graphically

6    visualised using FigTreev 1.4.3 software (Rambaut, 2012). The Phylogenetic Maximum

7    Clade Credibility (MCC) trees were used to determine the rate of "speciation" (or more

8    accurately, the expansion) within a given gene family. The R package TESS (Höhna,

9    May and Moore, 2016) was used for this analysis with the mass extinction turned off,

10    measuring only for the speciation rate (expansion rate), using the phylogenetic trees

11    produced in the divergence time estimations.

12    Likelihood ratio tests ($2\Delta L$) were performed by PAML 4 (Yang 2007) between

13    3 Codeml site model pairs (M0/M3; M1a/M2a and M7/M8), with the first model

14    representing a site model without positive selection and the second a model that allows

15    a proportion of sites to be under positive selection. In all 3 cases the model with positive

16    selection was a significantly better fit to the data. For the M2a and M8 models, several

17    sites were predicted to be positively selected by Bayes Empirical Bayes analysis (Yang

18    et al. 2005). The sites 457 and 460 were identified by both models with high confidence

19    ($> 95\%$) and 457 by M8.

20

21    **Analysis of lipids, fatty acids, and quinones/quinols**

22    Crude lipids were extracted from the cells by the method of Bligh and Dyer (1959). The

23    lipid fraction was evaporated, and then the residue was heated at 90°C for 2 h with 2

24    mL of 5% (w/w) HCl-methanol to obtain fatty acid methyl esters. The methanol

1     solution was extracted with 2 mL of n-hexane twice. The layer of n-hexane was

2     concentrated to a minimum volume for use in gas-chromatographic analysis. Gas

3     chromatography was performed, according to Mitani et al. (2017), with a fused-silica

4     capillary column (0.25 mm internal diameter x 50 m; ChrompackCp-Sil 88, Agilent

5     Technologies Inc., USA) with the oven temperature was increased from 150 to 210˚C at

6     $5°C \ min^{-1}$. The fatty acid composition was calculated by a Chromatopac C-R8A data

7     processor (Shimadzu Corp., Kyoto, Japan). Each fatty acid was identified by comparing

8     the retention time with those of known standards. Thin layer chromatograph (silica gel

9     60) of the crude lipids developed by hexane/diethyl ether/acetate (80:20:1 v/v/v). Lipids

10     were visualized under UV light at 365 nm after spraying the plate with 0.01% (w/v)

11     primulin in 80% (v/v) acetone (Wright 1971). Two-dimensional thin layer

12     chromatograph of the polar lipids was performed according to Sato (1991) with the

13     following solvent systems: first dimension,acetone/benzene/methanol/water (8:3:2:1

14     v/v/v/v); second dimension, chloroform/methanol/28% ammonium (13:7:1 v/v/v). Each

15     spot of polar lipids was detected as described above. Each lipid was identified by

16     spraying the plate with a reagent specific to each lipid class (Allen and Good 1971).

17     Each spot on a plate was scraped off, and then subjected to gas-chromatographic

18     analysis and fatty acid identification as described above.

19        Quinone/quinol extraction and detection was performed as described in

20     Kayama et al. (2020). In this analyses, we treated the quinone/quinol extract with ferric

21     chloride (final concentration, 1.2 mM) before the analysis for oxidation of total

22     quinones and quinols, to convert them to quinone forms for higher detection sensitivity

23     according to Kayama et al. (2020). We could identified only ubiquinone but no

24     plastoquinone from the cells of *N. putrida* NIES-4239 cultivated as described above for

30

1     genome sequencing, supporting lack of the complete pathway of plastoquinol synthesis

2     (Fig. 1).

3

4

18

19    **Author Contributions**

20    RK, GT, TM, and YN conceived this research. RK, YT, TM, MS, GT, and YN

21    sequenced, assembled, and annotated the genome. RK identified plastid metabolisms.

22    RK and TN detected HGT candidates. RO, RK, and SM performed comparative

23    transcriptome analyses. UC, BH, and VL searched and annotated CAZymes. DM and

24    RK identified peroxisome metabolisms. SS, KS, AT, CO, and TM performed the

1    coalescence analyses. KO and MK detected and identified lipids and fatty acids. KK

2    and RK identified mitochondrial metabolisms. RK, KI, and SM isolated and identified

3    the diatom. MK, TA, KI, HM performed growth experiments. MK and YK performed

4    quinone detection. RK, JW, SK, TM, and MS identified secretomes. RK, CVO, and TM

5    wrote the manuscript, and all the authors commented and approved the final version.

6

7    **Competing Interests statement**

8    There is no competing interest.

9

10    **References**

11    Annunziata R., Ritter A., Emidio Fortunato A., Manzotti A., Cheminant-Navarro S.,

12        Agier N., Huysman M.J.J., Winge P., Bones A.M., Bouget F., Cosentino

13        Lagomarsino M., Bouly J., and Falciatore A. bHLH-PAS protein RITMO1

14        regulates diel biological rhythms in the marine diatom *Phaeodactylum*

15        *tricornutum*. Proc Natl Acad Sci USA 116, 13137-13142 (2019)

16    Allen A.E., Dupont C.L., Oborník M., Horák A., Nunes-Nesi A., McCrow J.P., Zheng

17        H., Johnson D.A., Hu H., Fernie A.R., and Bowler C. Evolution and metabolic

18        significance of the urea cycle in photosynthetic diatoms. Nature 473, 203–207

19        (2011)

20    Armbrust E.V., Berges J.A., Bowler C., Green B.R., Martinez D., Putnam N.H., Zhou

21        S., Allen A.E., Apt K.E., Bechner M., et al. The genome of the diatom

22        *Thalassiosira pseudonana*: ecology, evolution, and metabolism. Science 306,

23        5693, 79-86 (2004)

1    Ashworth J., Coesel S., Lee A., Armbrust E.V., Orellana M.V., and Baliga N.S.

2        Genome-wide diel growth state transitions in the diatom *Thalassiosira*

3        *pseudonana*. Proc Natl Acad Sci USA 110, 7518-7523 (2013)

4    Aylward F.O., Eppley J.M., Smith J.M., Chavez F.P., Scholin, C.A., and DeLong E.F.

5        Microbial community transcriptional networks are conserved in three domains at

6        ocean basin scales. Proc Natl Acad Sci USA 112, 5443–5448 (2015)

7    Bouckaert R., Vaughan T.G., Barido-Sottani J., Duchêne S., Fourment M.,

8        Gavryushkina A., Heled J., Jones G., Kühnert D., De Maio N., and Matschiner

9        M. BEAST 2.5: An advanced software platform for Bayesian evolutionary

10       analysis. PLoS Comput Biol 15, e1006650 (2019).

11   Bowler C., Allen A.E., Badger J.H., Grimwood J., Jabbari K., Kuo A., Maheswari U.,

12       Martens C., Maumus F., Otillar R.P. et al. The *Phaeodactylum* genome reveals

13       the evolutionary history of diatom genomes. Nature 456, 239–244 (2008)

14   Bruckner C.G., Rehm C., Grossart H.P., and Kroth P.G. Growth and release of

15       extracellular organic compounds by benthic diatoms depend on interactions with

16       bacteria. Environ Microbiol 13, 1052–1063 (2011)

17   Buhmann M.T., Schulze B., F€orderer A., Schleheck D., and Kroth P.G. Bacteria may

18       induce the scretion of mucin-like proteins by the diatom *Phaeodactylum*

19       *tricornutum.* J Phycol 52, 463–474 (2016)

20   Chauton M.S., Winge P., Brembu T., Vadstein O., and Bones A.M. Gene regulation of

21       carbon fixation, storage, and utilization in the diatom *Phaeodactylum*

22       *tricornutum* acclimated to light/dark cycles. Plant Physiol 161, 1034–1048

23       (2013)

1  Chin C.S., Peluso P., Sedlazeck F.J., Nattestad M., Concepcion G.T, Clum A., Dunn C.,

2      O'Malley R., Figueroa-Balderas R., Morales-Cruz A. et al. Phased diploid

3      genome assembly with single-molecule real-time sequencing. Nature Methods

4      13, 1050–1054 (2016)

5  Coesel S., Mangogna M., Ishikawa T., Heijde M., Rogato A., Finazzi G., Todo T.,

6      Bowler C., and Falciatore A. Diatom PtCPF1 is a new cryptochrome/photolyase

7      family member with DNA repair and transcription regulation activity. EMBO

8      Rep 10, 655-661 (2009)

9  Dell'Aquila G., Zauner S, Heimerl T., Kahnt J., Samel-Gondesen V., Runge S., Hempel

10     F., and Maier U.G. Mobilization and cellulardistribution of phosphate in the

11     diatom *Phaeodactylum tricornutum.* Front Plant Sci 11, 579 (2020)

12 Dorrell R.G., Azuma T., Nomura M., Audren de Kerdrel G., Paoli L., Yang S., Bowler

13     C., Ishii K., Miyashita H., Gile G.H., and Kamikawa R. Principles of plastid

14     reductive evolution illuminated by nonphotosynthetic chrysophytes. Proc Natl

15     Acad Sci USA 116, 6914-6923 (2019)

16 Emms D.M., and Kelly S. OrthoFinder: solving fundamental biases in whole genome

17     comparisons dramatically improves orthogroup inference accuracy. Genome

18     Biol 16, 157 (2015)

19 Enright A.J., Van Dongen S., and Ouzounis C.A. An efficient algorithm for large-scale

20     detection of protein families. Nucleic Acids Res 30, 1575–1584 (2002)

21 Field C.B., Behrenfeld M.J., Randerson J.T., and Falkowski P. Primary production of

22     the biosphere: integrating terrestrial and oceanic components. Science 281, 237–

23     240 (1998)

1    Fortunato A.E., Jaubert M., Enomoto G., Bouly J., Raniello R., Thaler M., Malviya S.,

2        Bernardes J.S., Rappaport F., Gentili B., et al. Diatom phytochromes reveal the

3        existence of far-red-light-based sensing in the ocean. Plant Cell 28, 616–628

4        (2016)

5    Freese J.M., and Lane C.E. Parasitism finds many solutions to the same problems in red

6        algae (Florideophyceae, Rhodophyta). Mol Biochem Parasitol 214, 105-111

7        (2017)

8    Frischkorn K.R., Haley S.T., and Dyhrman S.T. Coordinated gene expression between

9        *Trichodesmium* and its microbiome over day–night cycles in the North Pacific

10       Subtropical Gyre. ISME J 12, 997–1007 (2018)

11   Fujiwara T., Hirooka S., Ohbayashi R., Onuma R., and Miyagishima S. Relationship

12       between cell cycle and diel transcriptomic changes in metabolism in a

13       unicellular red alga. Plant Physiol 183, 1484-1501 (2020)

14   Gruber A., Rocap G., Kroth P.G., Armbrust E.V., and Mock T. Plastid proteome

15       prediction for diatoms and other algae with secondary plastids of the red lineage.

16       Plant J 81, 519-528 (2015)

17   Hadariová L., Vesteg M., Hampl V., and Krajčovič J. Reductive evolution of

18       chloroplasts in non-photosynthetic plants, algae and protists. Curr Genet 64,

19       365–387 (2018)

20   Hernández Limón M.D., Hennon G.M.M., Harke M.J., Frischkorn K.R., Haley S.T.,

21       Dyhrman S.T. Transcriptional patterns of *Emiliania huxleyi* in the North Pacific

22       Subtropical Gyre reveal the daily rhythms of its metabolic potential. Environ

23       Microbiol 22, 381-396 (2020)

1    Hoff K.J., Lange S., Lomsadze A., Borodovsky M., Stanke M. BRAKER1:

2        unsupervised RNA-Seq-based genome annotation with GeneMark-ET and

3        AUGUSTUS. Bioinfomatics 32, 767-769 (2016)

4    Huysman M.J.J., Martens C., Vandepoele K., Gillard J., Rayko E., Heijde M., Bowler

5        C., Inzé D., Van de Peer Y., De Veylder L., and Vyverman W. Genome-wide

6        analysis of the diatom cell cycle unveils a novel type of cyclins involved in

7        environmental signalling. Genome Biol 11, R17 (2010)

8    Ishii K., and Kamikawa R. Growth characterization of non-photosynthetic diatoms,

9        *Nitzschia* spp., inhabiting estuarine mangrove forests of Ishigaki Island, Japan.

10       Plankton Benthos Res 12, 164-170 (2017)

11   Janouškovec J., Paskerova G.G., Miroliubova T.S., Mikhailov K.V., Birley T., Aleoshin

12       V.V., Simdyanov T.G. Apicomplexan-like parasites are polyphyletic and widely

13       but selectively dependent on cryptic plastid organelles. eLife 8:e49662 (2019)

14   Jones, P., Binns D., Chang H.Y., Fraser M., Li W., McAnulla C., McWilliam H.,

15       Maslen J., Mitchell A., Nuka G., Pesseat S., Quinn A.F., Sangrador-Vegas A.,

16       Scheremetjew M., Yong S.Y., Lopez R., and Hunter S. InterProScan 5: genome-

17       scale protein function classification. Bioinformatics 30, 1236-1240 (2014)

18   Kamikawa R., Yubuki N., Yoshida M., Taira M., Nakamura N., Ishida K., Leander

19       B.S., Miyashita H., Hashimoto T., Mayama S., et al. Multiple losses of

20       photosynthesis in *Nitzschia* (Bacillariophyceae). Phycol Res 63, 19–28 (2015a)

21   Kamikawa R., Tanifuji G., Ishikawa S.A., Ishii K., Matsuno Y., Onodera N.T., Ishida

22       K., Hashimoto T., Miyashita H., Mayama S., et al. Proposal of a twin arginine

23       translocator system-mediated constraint against loss of ATP synthase genes from

24       nonphotosynthetic plastid genomes. Mol Biol Evol 32, 2598–2604 (2015b)

1   Kamikawa R., Moog D., Zauner S., Tanifuji G., Ishida K., Miyashita H., Mayama H.,

2       Hashimoto T., Maier U.G., Archibald J.M. et al. A non-photosynthetic diatom

3       reveals early steps of reductive evolution in plastids. Mol Biol Evol 34, 2355–

4       2366 (2017)

5   Kayama M., Chen J.F., Nakada T., Nishimura Y., Shikanai T., Azuma T., Miyashita H.,

6       Takaichi S., Kashiyama Y., and Kamikawa R. A non-photosynthetic green alga

7       illuminates the reductive evolution of plastid electron transport systems. BMC

8       Biol 18, 126 (2020a)

9   Kayama M., Maciszewski K., Yabuki A., Miyashita H., Karnkowska A., and Kamikawa

10      R. Highly reduced plastid genomes of the non-photosynthetic

11      dictyochophyceans *Pteridomonas* spp. (Ochrophyta, SAR) are retained for

12      tRNA-Glu-based organellar heme biosynthesis. Front Plant Sci 11, 602455

13      (2020b)

14  Khatibi PA, Newmister SA, Rayment I, McCormick SP, Alexander NJ, Schmale DG

15      Bioprospecting for trichothecene 3-O-acetyltransferases in the fungal genus

16      Fusarium yields functional enzymes with different abilities to modify the

17      mycotoxin deoxynivalenol. Appl Environ Microbiol 77, 1162–1170 (2011)

18  Kim D., Langmead B., Salzberg S.L. HISAT: a fast spliced aligner with low memory

19      requirements. Nature Methods 12, 357–360 (2015).

20  Kissinger J.C., Brunk B.P., Crabtree J., Fraunholz M.J., Gajria B., Milgram A.J.,

21      Pearson D.S., Schug J., Bahl A., Diskin S.J., et al.. The *Plasmodium* genome

22      database. Nature 419, 490–492 (2002)

23  Kong G., Chen Y., Deng Y., Feng D., Jiang L., Wan L., Li M., Jiang Z., Xi P. The basic

24      leucine zipper transcription factor PlBZP32 associated with the oxidative stress

1    response is critical for pathogenicity of the lychee downy blight oomycete

2    *Peronophythora litchi*. mSphere 5, e00261-20 (2020).

3  Lachnit M, Buhmann MT, Klemm J, Kröger N, Poulsen N. Identification of proteins in

4    the adhesive trails of the diatom *Amphora coffeaeformis*. Phil Trans R Soc B

5    374: 20190196 (2019)

6  Li C.W., and Volcani B.E. Four new apochlorotic diatoms. Br Phycol J 22, 375-382

7    (1987)

8  Li H., Benedito V.A., Udvardi M.K., and Zhao P.X. TransportTP: A two-phase

9    classification approach for membrane transporter prediction and characterization

10    BMC Bioinformatics 10, 418 (2009)

11  Mann D.G. The species concept in diatoms. Phycologia 38, 437–495 (1999)

12  Marchler-Bauer A., and Bryant S.H. CD-Search: protein domain annotations on the fly.

13    Nucleic Acids Res 32, W327-W331 (2004)

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

15    Sanges R., Toseland A., Ward B.J. Evolutionary genomics of the cold-adapted

16    diatom *Fragilariopsis cylindrus*. Nature 541, 536–540 (2017)

17  Moog D., Nozawa A., Tozawa Y., and Kamikawa R. Substrate specificity of plastid

18    phosphate transporters in a non-photosynthetic diatom and its implication in

19    evolution of red alga-derived complex plastids. Sci Rep 10, 1167 (2020)

20  Moriya Y., Itoh M., Okuda S., Yoshizawa A.C., and Kanehisa M. KAAS: an automatic

21    genome annotation and pathway reconstruction server. Nucleic Acids Res 35,

22    W182-W185 (2007)

23  Morris J.J., Lenski R.E., Zinser E.R. The black queen hypothesis: evolution of

24    dependencies through adaptive gene loss. mBio 3, e00036-12 (2012)

Nei M, Kumar S: Molecular Evolution and Phylogenetics. New York: Oxford University Press, 17-203 (2000)

Ottesen E.A., Young C.R., Eppley J.M., Ryan J.P., Chavez F.P., Scholin C.A., DeLong E.F. Pattern and synchrony of gene expression among sympatric marine microbial populations. Proc Natl Acad Sci USA 110, E488-E497 (2013)

Ottesen E.A., Young C.R., Gifford S.M., Eppley J.M., Marin III R., Schuster S.C., Scholin C.A., DeLong E.F. Multispecies diel transcriptional oscillations in open ocean heterotrophic bacterial assemblages. Science 345, 207-212 (2014)

Rayko E., Maumus F., Maheswari U., Jabbari K., and Bowler C. Transcription factor families inferred from genome sequences of photosynthetic stramenopiles. New Phytol 188, 52–66 (2010)

Richards T.A., and Talbot N.J. Horizontal gene transfer in osmotrophs: playing with public goods. Nature Rev Microbiol 11, 720-727 (2013)

Richards T.A., and Talbot N.J. Osmotrophy. Curr Biol 28, R1179–R1180 (2018)

Smith S.R., Gillard J.T.F., Kustka A.B., McCrow J.P., Badger J.H., et al. Transcriptional orchestration of the global cellular response of a model pennate diatom to diel light cycling under iron limitation. PLOS Genetics 13, e1006688 (2017)

Smith S.R., Dupont C.L., McCarthy J.K., Broddrick J.T., Oborník M., Horák A., Füssy Z., Cihlář J., Kleessen S., Zheng H., et al. Evolution and regulation of nitrogen flux through compartmentalized metabolic networks in a marine diatom. Nature Commun 10, 4552 (2019)

39

1 Sun G., Xu Y., Liu H., Sun T., Zhang J., Hettenhausen C., Shen G., Qi J., Qin Y., Li J.,

2 et al.. Large-scale gene losses underlie the genome evolution of parasitic plant

3 *Cuscuta australis.* Nature Commun 9, 2683 (2018)

4 Vancaester E., Depuydt T., Osuna-Cruz C.M., and Vandepoele K. Comprehensive and

5 functional analysis of horizontal gene transfer events in diatoms. Mol Biol Evol.

6 37, 3243–3257 (2020)

7 Van Valen L. Molecular evolution as predicted by natural selection. J Mol Evol 3, 89–

8 101 (1974)

9 Vurture G.W., Sedlazeck F.J., Nattestad M., Underwood C.J., Fang H., Gurtowski J.,

10 and Schatz M.C. GenomeScope: fast reference-free genome profiling from short

11 reads. Bioinformatics 33, 2202–2204 (2017)

12 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An

13 integrated tool for comprehensive microbial variant detection and genome

14 assembly improvement. PLoS ONE 9, e112963 (2014)

15 Waterhouse R.M., Seppey M., Simão F.A, Manni M., Ioannidis P., Klioutchnikov G.,

16 Kriventseva E.V., Zdobnov E.M. BUSCO applications from quality assessments

17 to gene prediction and phylogenomics. Mol Biol Evol 35, 543–548 (2018)

18 Wolf Y.I., and Koonin E.V. Genome reduction as the dominant mode of evolution.

19 Bioessays 35, 829–837 (2013)

20

21 **References for Methods**

22 Allen C.F., and Good P. Acyl lipids in photosynthetic systems. Methods Enzymol 23,

23 523-547 (1971)

1   Altschul S.F., Gish W., Miller W., Myers E.W., and Lipman D.J. Basic local alignment

2       search tool. J Mol Biol. 215, 403-410 (1990)

3   Altschul S.F., Madden T.L., Schäffer A.A., Zhang J., Zhang Z., Miller W., Lipman D.J.

4       Gapped BLAST and PSI-BLAST: a new generation of protein database search

5       programs. Nucleic Acids Res 25, 3389-3402 (1997)

6   Aramaki T., Blanc-Mathieu R., Endo H., Ohkubo K., Kanehisa M., Goto S., Ogata H.

7       KofamKOALA: KEGG Ortholog assignment based on profile HMM and

8       adaptive score threshold. Bioinformatics 36, 2251-2252 (2020)

9   Bligh E.G., and Dyer W.J. A rapid method of total lipid extraction and purification. Can

10      J Biochem Physiol 37, 911–917 (1959)

11  Bolger A.M., Lohse M.A., and Usadel B. Trimmomatic: a flexible trimmer for Illumina

12      sequence data. Bioinformatics 30, 2114–2120 (2014)

13  Bray J.R., and Curtis J.T. An ordination of the upland forest communities of Southern

14      Wisconsin. Ecol Monogr 27, 325–349 (1957)

15  Camacho C., Coulouris G., Avagyan V., Ma N., Papadopoulos J., Bealer K., Madden

16      T.L. BLAST+: architecture and applications. BMC Bioinformatics 10, 421

17      (2009)

18  Capella-Gutiérrez S., Silla-Martínez J.M., and Gabaldón T. trimAl: a tool for automated

19      alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25,

20      1972-1973 (2009)

21  Castresana J. Selection of conserved blocks from multiple alignments for their use in

22      phylogenetic analysis. Mol Biol Evol 17, 540-552 (2000)

23  Cenci U., Sibbald S.J., Curtis B.A., Kamikawa R., Eme L., Moog D., Henrissat B.,

24      Maréchal E., Chabi M., Djemiel C., et al. Nuclear genome sequence of the

41

1      plastid-lacking cryptomonad *Goniomonas avonlea* provides insights into the

2      evolution of secondary plastids. BMC Biol 16, 137 (2018).

3  Criscuolo A., and Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a

4      new software for selection of phylogenetic informative regions from multiple

5      sequence alignments. BMC Evol Biol 10, 210 (2010)

6  Curtis B.A. et al. Algal genomes reveal evolutionary mosaicism and the fate of

7      nucleomorphs. Nature 492, 59–65 (2012)

8  Davis A., Abbriano R., Smith S.R., and Hildebrand M. Clarification of photorespiratory

9      processes and the role of malic enzyme in diatoms. Protist 168, 134-153 (2017)

10  de Hoon M.J.L., Imoto S., Nolan J., and Miyano S. Open source clustering software.

11      Bioinformatics 20, 1453–1454 (2004)

12  Emms D.M., and Kelly S. OrthoFinder: solving fundamental biases in whole genome

13      comparisons dramatically improves orthogroup inference accuracy. Genome

14      Biol 16, 157 (2015)

15  El-Gebali S., Mistry J., Bateman A., Eddy S.R., Luciani A., Potter S.C., Qureshi M.,

16      Richardson L.J., Salazar G.A., Smart A., et al. The Pfam protein families

17      database in 2019. Nucleic Acids Res 47, D427–D432 (2019)

18  Fabregat A., Korninger F., Viteri G., Sidiropoulos K., Marin-Garcia P., Ping P., Wu G.,

19      Stein L., D'Eustachio P., Hermjakob H. Reactome graph database: Efficient

20      access to complex pathway data. PLoS Comput Biol 14, e1005968 (2018a)

21  Fabregat A., Sidiropoulos K., Viteri G., Marin-Garcia P., Ping P., Stein L., D'Eustachio

22      P., Hermjakob H. Reactome diagram viewer: data structures and strategies to

23      boost performance. Bioinformatics 34, 1208-1214 (2018b)

1    Fukasawa Y., Tsuji J., Fu S.C., Tomii K., Horton P., and Imai K. MitoFates: improved
2        prediction of mitochondrial targeting sequences and their cleavage sites. Mol
3        Cell Proteom 14, 1113-1126 (2015)

4    Grabherr M.G., Haas B.J., Yassour M., Levin J.Z., Thompson D.A., Amit I., et al. Full-
5        length transcriptome assembly from RNA-Seq data without a reference genome.
6        Nature Biotechnol 29, 644–652 (2011)

7    Gonzalez N.H., Felsner G., Schramm F.D., Klingl A., Maier U.G., and Bolte K. A
8        single peroxisomal targeting signal mediates matrix protein import in diatoms.
9        PLoS ONE 6, e25316 (2011)

10    Haas B.J., et al. Improving the *Arabidopsis* genome annotation using maximal transcript
11        alignment assemblies. Nucleic Acids Res 31, 5654–5666 (2003)

12    Haas B.J., Papanicolaou A., Yassour M., Grabherr M., Blood P.D., Bowden J., et al. De
13        novo transcript sequence reconstruction from RNA-seq using the Trinity
14        platform for reference generation and analysis. Nature Protoc. 8, 1494–1512
15        (2013)

16    Hall TA. BioEdit: a user-friendly biological sequence alignment editor and analysis
17        program for windows 95/98/NT. Nucleic Acids Symp Ser 41, 95–98 (1999)

18    Höhna S., May M.R., and Moore B.R. TESS: an R package for efficiently simulating
19        phylogenetic trees and performing Bayesian inference of lineage diversification
20        rates. Bioinformatics 32, 789-791 (2016)

21    Ihaka R., and Gentleman R. R: A language for data analysis and graphics. J Comp
22        Graph Stat 5, 299-314 (1996)

1   Jassal B., Matthews L., Viteri G., Gong C., Lorente P., Fabregat A., Sidiropoulos K.,
2        Cook J., Gillespie M., Haw R., et al. The reactome pathway knowledgebase.
3        Nucleic Acids Res 48, D498-D503 (2020)

4   Johnson, L. K., Alexander, H., & Brown, C. T. (2019). Re-assembly, quality evaluation,
5        and annotation of 678 microbial eukaryotic reference transcriptomes.
6        GigaScience, 8(4), giy158.

7   Kamikawa R., Tanifuji G., Ishikawa S.A., Ishii K., Matsuno Y., Onodera N.T., Ishida
8        K., Hashimoto T., Miyashita H., Mayama S., et al. Proposal of a twin arginine
9        translocator system-mediated constraint against loss of ATP synthase genes from
10       nonphotosynthetic plastid genomes. Mol Biol Evol 32, 2598–2604 (2015b)

11  Kamikawa R., Azuma T., Ishii K., Matsuno Y., and Miyashita H. Diversity of
12       organellar genomes in non-photosynthetic diatoms. Protist 169, 351-361 (2018)

13  Kanehisa M., Goto S., Sato Y., Furumichi M., Tanabe M. KEGG for integration and
14       interpretation of large-scale molecular data sets. Nucleic Acids Res 40(Database
15       issue), D109-D114 (2012)

16  Katoh K., and Standley D.M. MAFFT multiple sequence alignment software version 7:
17       improvements in performance and usability. Mol Biol Evol 30, 772-780 (2013)

18  Kayama M., Chen J., Nakada T., Nishimura Y., Shikanai S., Azuma T., Miyashita H.,
19       Takaichi S., Kashiyama Y., and Kamikawa R. A non-photosynthetic green alga
20       illuminates the reductive evolution of plastid electron transport systems. BMC
21       Biol 18, 126 (2020)

22  Krogh A., et al. Prediction transmembrane protein topology with a hidden markov
23       model: application to complete genomes. J Mol Biol 305, 567–580 (2001)

1  Kume K., Amagasa T., Hashimoto T., and Kitagawa H. NommPred: Prediction of

2      mitochondrial and mitochondrion-related organelle proteins of nonmodel

3      organisms. Evol Bioinform Online. 14, 1176934318819835 (2018)

4  Langmead B., and Salzberg S.L. Fast gapped-read alignment with Bowtie 2. Nature

5      Methods 9, 357–359 (2012)

6  Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., et al. The Sequence

7      Alignment/Map format and SAMtools. Bioinformatics 25, 2078–2079 (2009)

8  Lombard V., Golaconda Ramulu H., Drula E., Coutinho P.M., and Henrissat B. The

9      carbohydrate-active enzymes database (CAZy) in 2013. Nucleic Acids Res 42,

10     D490–D495 (2014)

11 Löytynoja A. Phylogeny-aware alignment with PRANK. In Multiple sequence

12     alignment methods, pp. 155-170 (2014)

13 Mistry J., Finn R.D., Eddy S.R., Bateman A., and Punta M. Challenges in homology

14     search: HMMER3 and convergent evolution of coiled-coil regions. Nucleic

15     Acids Res 41, e121–e121 (2013)

16 Mitani E., Nakayama F., Matsuwaki I., Ichi I., Kawabata A., Kawachi M., and Kato M.

17     Fatty acid composition profiles of 235 strains of three microalgal divisions

18     within the NIES Microbial Culture Collection. Microb. Resour. Syst. 33, 19-29

19     (2017)

20 Mix A.K., Cenci U., Heimerl T., Marter P., Wirkner M.L., and Moog D. Identification

21     and localization of peroxisomal biogenesis proteins indicates the presence of

22     peroxisomes in the cryptophyte *Guillardia theta* and other "Chromalveolates".

23     Genome Biol Evol 10, 2834-2852 (2018)

1    Nguyen L.T., Schmidt H.A., Von Haeseler A., and Minh B.Q. IQ-TREE: a fast and

2        effective stochastic algorithm for estimating maximum-likelihood phylogenies.

3        Mol Biol Evol 32, 268-274 (2015)

4    Oksanen J. Multivariate Analysis of Ecological Communities in R: vegan tutorial. Cc

5        Oulu Fi Jarioksapopular Html (2014)

6    Pertea G., and Pertea M. GFF Utilities: GffRead and GffCompare. F1000Research 9,

7        304 (2020)

8    Petersen T.N., Brunak S., von Heijne G., and Nielsen H. SignalP 4.0: discriminating

9        signal peptides from transmembrane regions. Nature Methods 8, 785-786 (2011)

10   Quinlan A.R., and Hall I.M. BEDTools: A flexible suite of utilities for comparing

11       genomic features. Bioinformatics 26, 841–842 (2010)

12   Rambaut A. http://tree.bio.ed.ac.uk/software/tracer/ (2009).

13   Rambaut A. FigTree v1. 4. (2012).

14   Saldanha A.J. Java Treeview - Extensible visualization of microarray data.

15       Bioinformatics 20, 3246–3248 (2004)

16   Sato N. Lipids in *Cryptomonas* CR-1. I. Occurrence of betaine lipids. Plant Cell Physiol

17       32, 819-825 (1991)

18   Sidiropoulos K., Viteri G., Sevilla C., Jupe S., Webber M., Orlic-Milacic M., Jassal B.,

19       May B., Shamovsky V., Duenas C., et al. Reactome enhanced pathway

20       visualization. Bioinformatics 33, 3461-3467 (2017)

21   Simão F.A., Waterhouse R.M., Ioannidis P., Kriventseva E.V., Zdobnov E.M. BUSCO:

22       assessing genome assembly and annotation completeness with single-copy

23       orthologs. Bioinformatics 31, 3210–3212 (2015)

1  Sonnhammer E., et al. A hidden Markov model for predicting transmembrane helices in

2      protein sequences. Proc. ISMB 6, 175–182 (1998)

3  Talavera G., and Castresana J. Improvement of phylogenies after removing divergent

4      and ambiguously aligned blocks from protein sequence alignments. Syst Biol 56,

5      564-577 (2007)

6  UniProt Consortium. UniProt: a worldwide hub of protein knowledge. Nucleic Acids

7      Res. 47, D506-D515 (2019)

8  Wright R.S. A reagent for the non-destructive location of steroids and some other

9      lipophilic materials on a silica gel thinlayer chromatogram. J Chromat 59, 220-

10      221 (1971)

11  Yang Z, Wong W.S.W., and Nielsen R. Bayes empirical bayes inference of amino acid

12      sites under positive selection. Mol Biol Evol 22, 1107–1118 (2005)

13  Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. Mol Biol Evol 24,

14      1586-1591 (2007)

15

16  **Figure legends**

17  **Fig. 1 The heterotrophic diatom *Nitzschia putrida* and its plastid proteome. A.** The

18  frustule view of *N. putrida*. Bar = 20 µm. **B.** Estimated plastid proteome size in three

19  diatoms. Light and dark grey bars show low and high confident plastid-targeted proteins

20  identified by ASAFind, respectively. Data of two photosynthetic diatoms

21  *Phaeodactylum tricornutum* and *Thalassisira pseudonana* are derived from Gruber et al.

22  (2007). **C.** Unique and shared plastid-targeted orthogroups. High lighted in red is the

23  orthogroup exclusively shared by the two photosynthetic diatoms. **D.** Comparison of

24  KO ID numbers among plastids of the three diatoms. Each bar indicates numbers of

47

1    unique KO IDs in each functional category. **E.** Predicted metabolic map of the non-

2    photosynthetic plastid. Representative pathways found in photosynthetic diatom species

3    are shown. Green and light grey arrows show presence and absence of the responsible

4    protein sequences for the reactions in the genome, respectively. Amino acids are

5    highlighted in red. Glu: Glutamate, Gln: Glutamine, Orn: Ornithine, $NH^{4+}$: ammonium,

6    $NO_2$: nitrite, N-Acetyl-Glu: N-Acetylglutamate, N-Acetyl-Glu-5P: N-Acetylglutamyl-

7    5phosphate, N-Acetyl-Glu-semialdehyde: N-Acetylglutamyl-semialdehyde, a-

8    ketoglutarate: alpha-ketoglutarate, PPP-IX: protoporphyrin-IX, ALA: 5-aminolevulinic

9    acid, GSAH: Glutamate-1-semialdehyde, aa-tRNA-Glu: amino acyl-tRNA-glutamate,

10    tRNA-E: glutamyl tRNA, Chl *a*: chlorophyll *a*, α-TT: alpha-Tocotrierol, β-TT: beta-

11    tocotrierol, α-TP: alpha-tocopherol, β-TP: beta-tocopherol, PQ:

12    plastoquinone/plastoquinol, GA3P: glyceroaldehyde-3phosphate, Pyr: Pyruvate, IPP:

13    Isopentenylpyrophosphate, DPP: Dimethylallylpyrophosphate, Phy: Phytoene, β-Car:

14    beta-carotene, DHB4P: 3,4-Dihydroxy 2 butanone 4phosphate, DMR: 6,7-Dimethyl-8-

15    ribityl umazine, Asp: Aspartate, Lys: Lysine, Ile: Isoleucine, Leu: Leucine, Val: Valine,

16    MM: Methylmalate, ACL: Acetolactate, AcCoA: Acetyl-CoA, ACP: acyl-carrier

17    protein, 3KA-ACP: 3ketoacyl-ACP, A-ACP: acyl-ACP, E-ACP: enoyl-ACP, 3HA-

18    ACP: 3hydroxyacyl-ACP, PEP: phosphoenolpyruvate, G2P: 2 phosphoglycerate, G3P:

19    3 phosphoglycerate, DAH7P: 3-deoxy-7-phosphoheptulonate, DHQ: 3-dehydroquinic

20    acid, DHS: 3-dehydroshikimate, S3P: shikimate 3phosphate, EPS3P: 5-

21    enolpyruvylshikimate -3-phosphate, CM: Chorismate, Trp: Tryptophan, Tyr: Tyrosine,

22    Phe: Phenylalanine, G1,2P2: Glycerol 1,2 bisphosphate, R1,5P2: Ribulose 1,5

23    bisphosphate, Ru5P: Ribulose 5phosphate, Ri5P: Ribose 5phosphate, S7P:

24    sedoheptulose 7-phosphate, E4P: Erythrose 4phosphate, F6P: Fructose 6phosphate,

1  X5P: Xylulose 5phosphate, F1,6P2: Fructose 1,6 bisphosphate, Gly3P: Glycerate

2  3phosphate, LysoPA: Lysophosphatidic acid, PA: Phosphatidic acid, PG: Phosphatidyl

3  glycerol, G6P: Glucose 6phosphate, G1P: Glucose 1phosphate, UDPG: UDP-glucose,

4  UDPSQ: UDP-sulfoquinovose, SQDG:    Sulfoquinovosyl diacylglycerol, $SO_4^{2-}$:

5  sulfate, Cys: Cysteine, FeS: Iron-sulfur cluster, Ala: Alanine, Cyt: cytochrome *b6/f*

6  complex, PSI: photosystem I, PSII: Photosystem II.

7

8  **Fig. 2 Loss of genes for the plastid-persoxisome metabolic flow and photoreceptors.**

9  **A.** Metabolic interactions between a mitochondrion and a non-photosynthetic plastid

10  and between a mitochondrion and a peroxisome. Black, orange, and blue arrows show

11  presence of responsible protein sequences for the reactions in a plastid, a

12  mitochondrion, and a peroxisome, respectively, while light grey arrows show absence

13  of responsible protein sequences. Dashed arrows show possible inter-organellar

14  metabolic flows. $NH4+$: ammonium, $CO2$: carbon dioxide, $HCO3-$: bicarbonate, Arg:

15  Arginine, Orn: Ornithine, Asp: Aspartate, βox: fatty acid β oxidation, Gln: Glutamine,

16  Glu: Glutamate, a-KG: alpha-ketoglutarate, Ser: Serine, Gly: Glycine, GCS: Glycine

17  cleavage system, Ala: Alanine, TCA: Tricarboxylic acid cycle, OAA: oxaloacetate,

18  BCAA: Branched chain amino acid synthesis, AAA: aromatic amino acid synthesis.

19  Other abbreviations are described in Fig. 1. B. Photoreceptor and cell cycle genes in the

20  *N. putrida* genome. The other genes are shown in Supplementary Fig. S5. Light green

21  and light grey boxes show presence and absence of corresponding genes, respectively.

22  **C.** Growth of the heterotrophic diatom under the different light conditions. Closed

23  boxes show growth in the continuous dark condition, while open boxes show growth in

24  the light-dark condition. Shaded in grey are the dark periods in the light-dark cultivation

1    conditions. **D.** (Left) Heatmap showing the reproducible expression patterns of genes

2    (Pearson's correlation coefficient < 0.9). *k*-means clustering was calculated for each

3    gene based on RPKM +1 values, which were transformed to log2 and centred by

4    median values. Yellow and blue indicate upregulation and downregulation of the gene,

5    respectively. (Right) The line graphs showing expression pattern of genes in each

6    cluster. The coloured line indicates the average value of the expression patterns.

7

8    **Fig. 3 Horizontal gene transfers and gene family diversification in the**

9    **heterotrophic diatom *Nitzschia putrida*. A.** Proportion of origins of laterally

10    transferred genes identified in the *N. putrida* genome. **B-E.** Phylogenetic trees of

11    glucose 1P-dehydrogenase (NAA00P12280; B), arginase family protein

12    (NAA12P00140;C), alpha-N-acetylgalactosidase (NAA03P01590; D), and

13    trichothecene 3-O-acetyltransferase, self-protection against trichothecenes

14    (NAA53P00570; E). The numbers on branches are maximum likelihood bootstrap

15    values higher than 80. All the trees showing possible *N. putrida*-specific LGTs are

16    provided in Supplementary Fig. 7. **F.** Venn diagram of KOG IDs shared by four

17    diatoms. *Nitzschia*: *Nitzschia putrida*, *Phaeodactylum*: *Phaeodactylum tricornutum*,

18    *Thalassiosira*: *Thalassiosira pseudonana*, *Fragilariopsis*: *Fragilariopsis cylindrus*. **G.**

19    Comparison of the number of KOG ID among diatoms. KOG categories are as follows:

20    A, RNA processing and modification; B, chromatin structure and dynamics; C, energy

21    production and conversion; D, cell cycle control, cell division and chromosome

22    partitioning; E, amino acid transport and metabolism; F, nucleotide transport and

23    metabolism; G, carbohydrate transport and metabolism; H, coenzyme transport and

24    metabolism; I, lipid transport and metabolism; J, translation, ribosomal structure and

1 biogenesis; K, transcription; L, replication, recombination and repair; M, cell wall,

2 membrane or envelope biogenesis; N, cell motility; O, post-translational modification,

3 protein turnover, chaperones; P, inorganic ion transport and metabolism; Q, secondary

4 metabolites biosynthesis, transport and catabolism; R, general function prediction only;

5 S, function unknown; T, signal transduction; U, intracellular trafficking, secretion and

6 vesicular transport; V, defence mechanisms; W, extracellular structures; Y, nuclear

7 structure; Z, cytoskeleton. **H.** Comparison of the number of genes assigned to each

8 KOG category. Other details are described above.

9

10 **Fig. 4 Diversity of transporters and carbohydrate active enzymes in *N. putrida*. A.**

11 Distribution of the number of transporters in each transporter family. There is no

12 significant difference detectable by Wilcoxon signed rank test. **B.** The gene number of

13 transporters in the 10 most abundant transporter families of *N. putrida*. 1. Resistance-

14 Nodulation-Cell Division (RND) Superfamily, 2. Solute:Sodium Symporter (SSS)

15 Family, 3. Voltage-gated Ion Channel (VIC) Superfamily, 4. Silicon Transporter (SIT)

16 Family, 5. Amino Acid/Auxin Permease (AAAP) Family, 6. P-type ATPase (P-ATPase)

17 Superfamily, 7. Drug/Metabolite Transporter (DMT) Superfamily, 8. Mitochondrial

18 Carrier (MC) Family, 9. Major Facilitator Superfamily (MFS), 10. ATP-binding

19 Cassette (ABC) Superfamily. **C.** Silicon transporter (SIT) genes tandemly located in the

20 contig 000000F. SIT genes are highlighted in light red with the gene IDs. **D.** Glycoside

21 Hydrolases (GH) families from the Carbohydrate Active enZyme (CAZy) database

22 focused on Bacillariophyta. The diagram shows a heatmap of CAZyme prevalence in

23 each taxon (number of a particular CAZyme family divided by the total number of

24 CAZyme genes in the organism); the white to blue color scheme indicates low to high

1    prevalence, respectively. Dendrograms (left and top) show respectively the relative taxa

2    proximity with respect co-occurrence of CAZyme families and the co-occurrence of

3    CAZyme families with one another within genomes (according to the method described

4    in Cenci et al. 2018). **E.** GH114 genes tandemly located in the contig 000022F. GH114

5    genes are highlighted in light green with the gene IDs.

6

7    **Fig. 5 Secretome of *N. putrida*. A.** The number of secretome tribes, including at least

8    four sequences, clustered by TribeMCL (Enright et al. 2002). Different colours

9    represent tribe categories as follows: 1. Species specific tribes, 2: tribes shared by two

10   species, 3: tribes shared by three species, and 4: tribes shared by all the four diatoms. **B.**

11   Proportion of each tribe category in diatoms. Details are described in **A. C.** Distribution

12   of the number of protein sequences in each secretome tribe. Each box represents the

13   interquartile range between the first and third quartiles (25th and 75th percentiles,

14   respectively), and the median is represented by the vertical line inside the box. The lines

15   protruding from either side of the box are the lowest and highest values within 1.5 times

16   the interquartile range from the first and third quartiles, respectively. Outliers were

17   omitted in the boxplot. The results of the Wilcoxon signed rank test was shown below

18   the box plot. *p*-values for three pairs of comparisons that include *N. putrida* as a

19   counterpart were adjusted by the Benjamini–Hochberg procedure. Asterisks show

20   significantly different pairs under $p < 0.01$. **D.** Sequence diversity of Leucin-rich repeat

21   protein sequences. Phylogenetic tree centred in the figure was reconstructed by IQ-

22   TREE. Domains contained in each sequence were predicted by NCBI Conserved

23   Domain Search and depicted next to the gene IDs; in only several sequences no domain

24   was predicted regardless of apparent homologies ($e^{-30}$ in TribeMCL) to LRR proteins.

1    **E.** Expression of the 10 largest tribes in *N. putrida* during the 25 hours cultivation.

2    Heatmap showing the expression patterns of genes in two independent experiments. *k*-

3    means clustering was calculated for each gene based on RPKM + 1 values, which were

4    transformed to log2 and centred by median values. Yellow and blue indicate

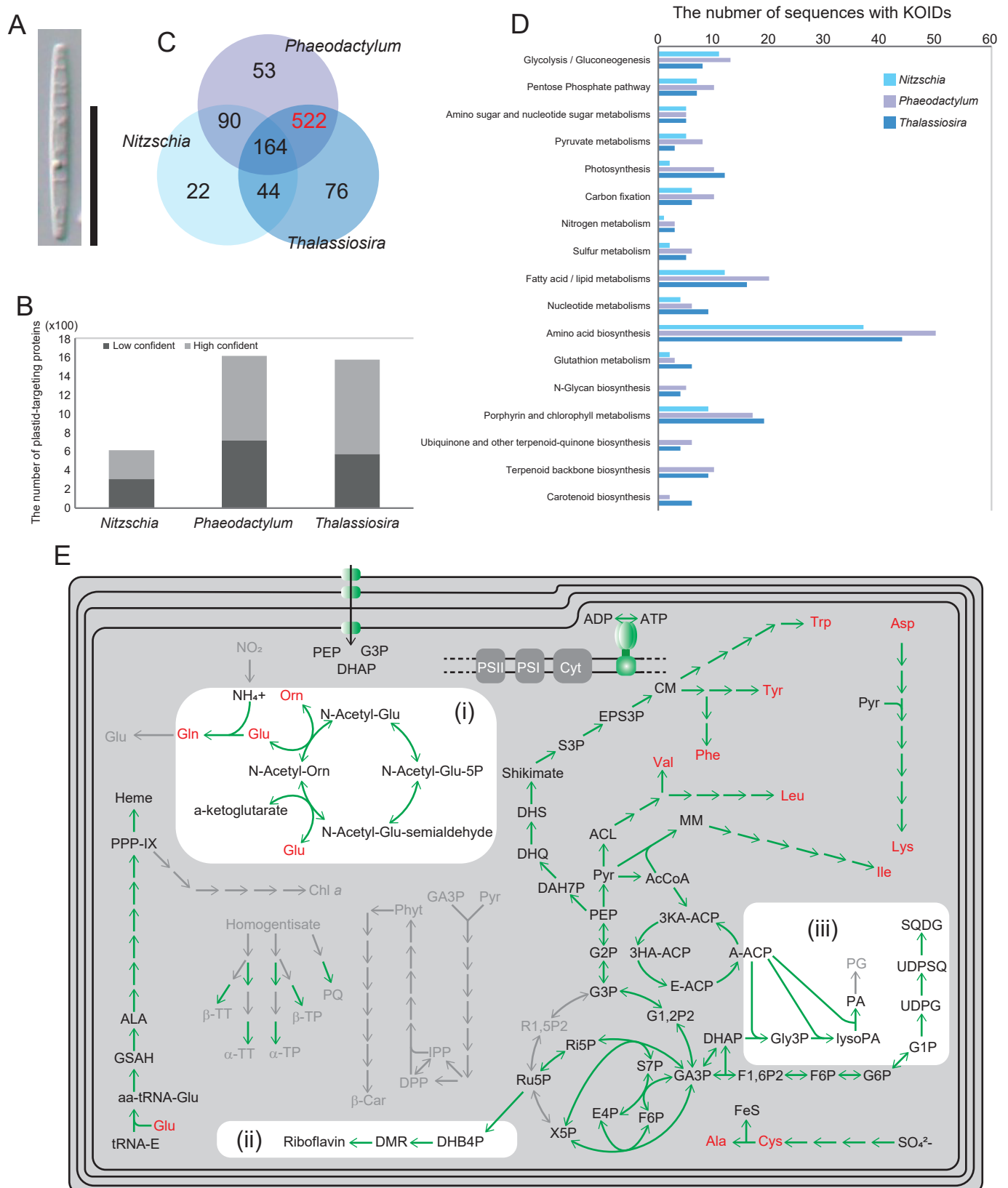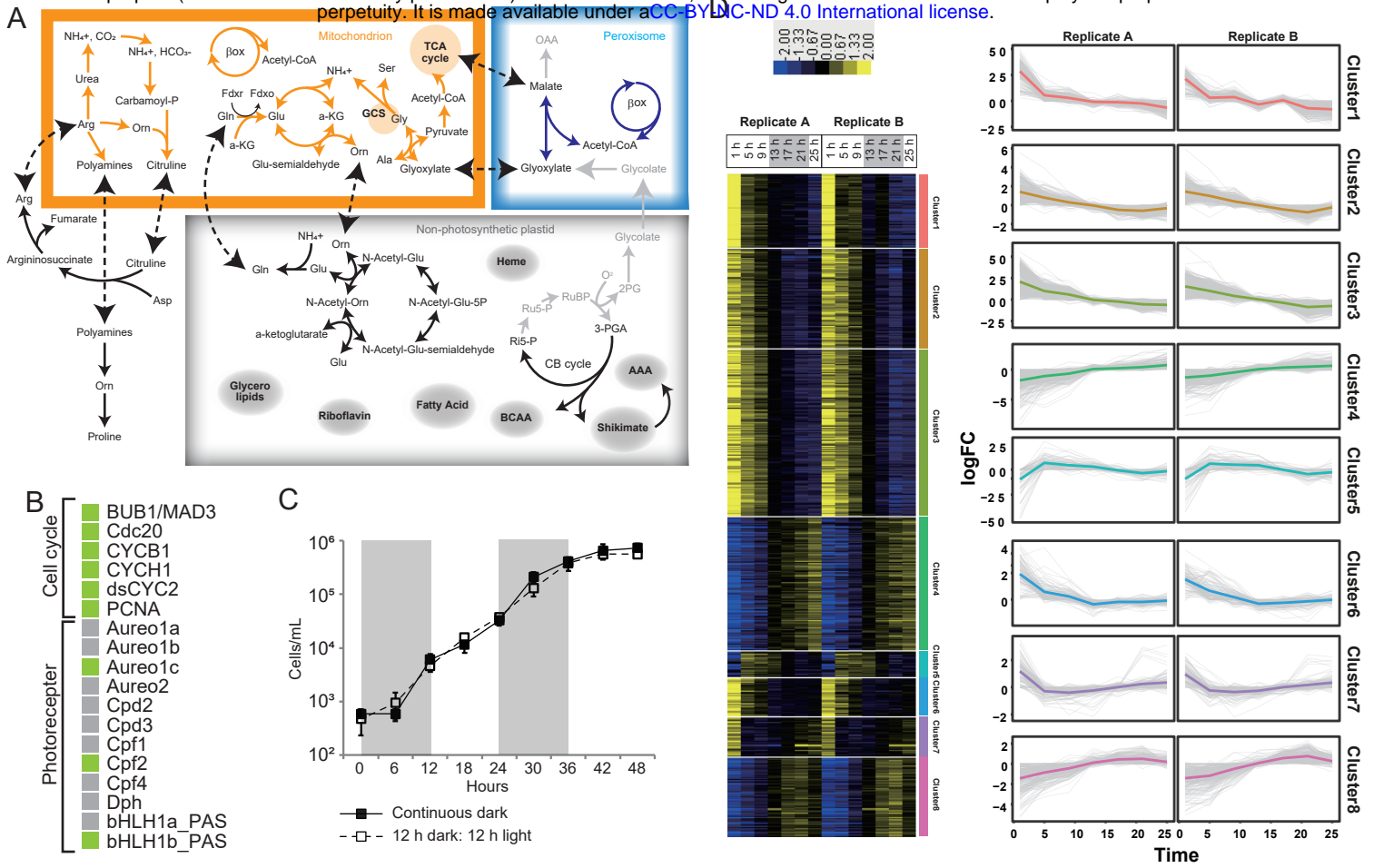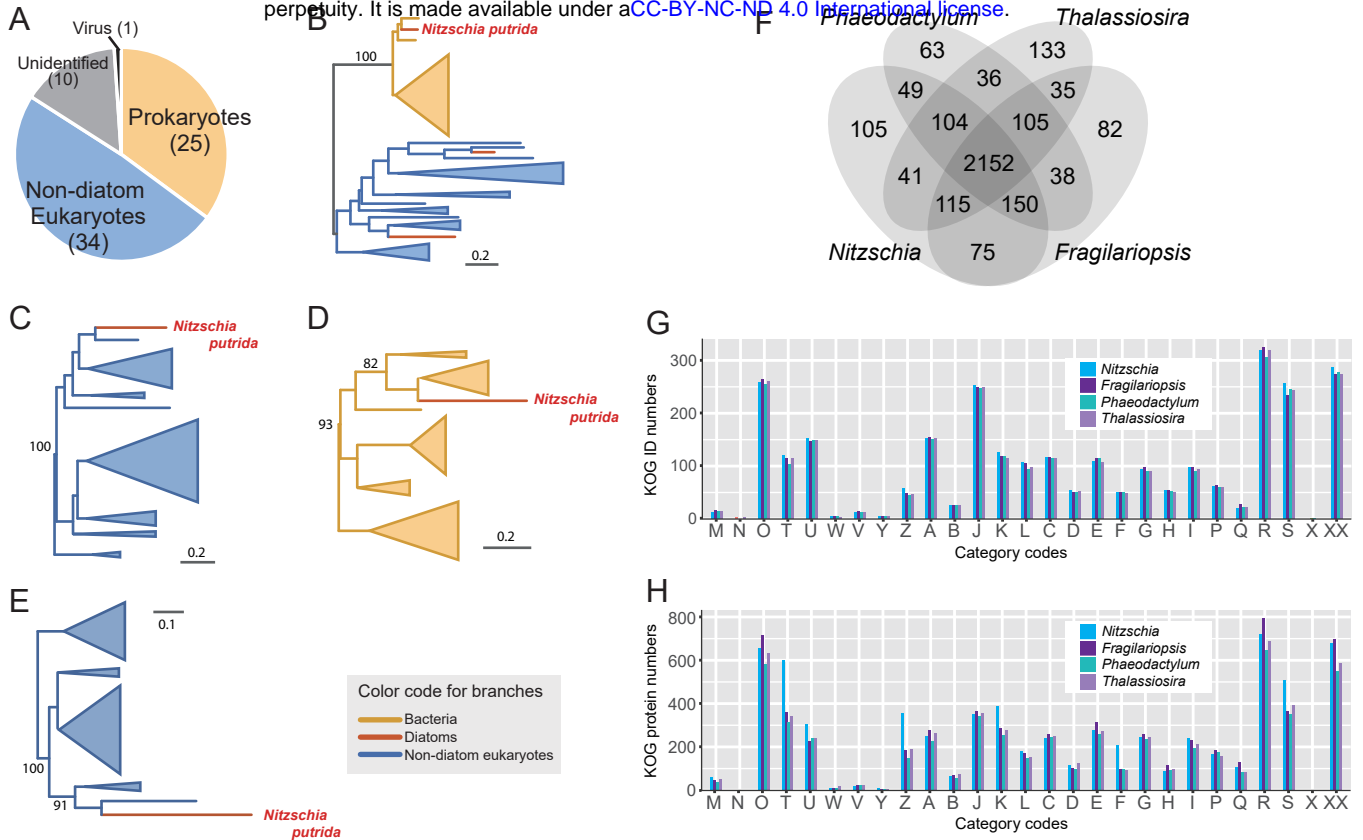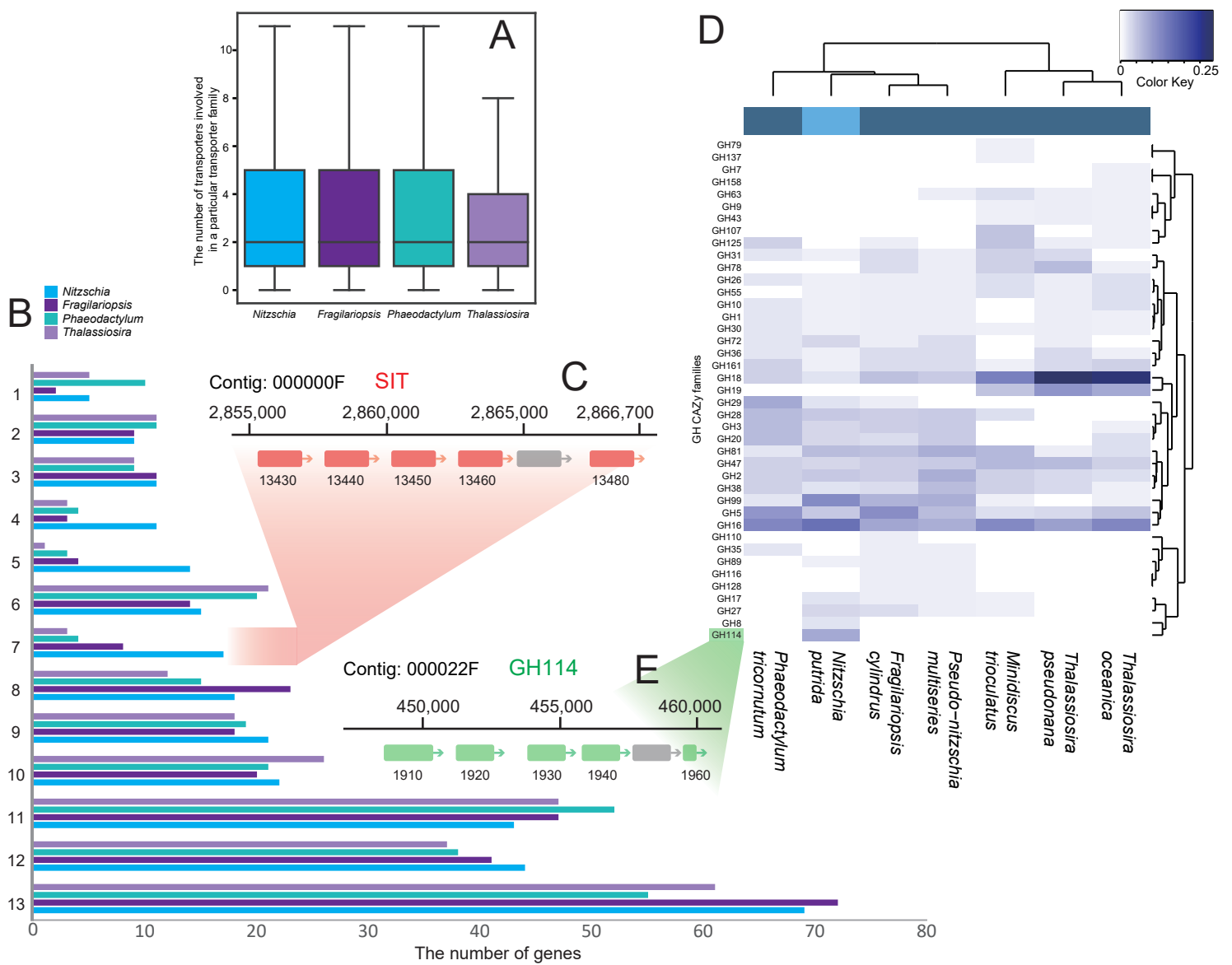5    upregulation and downregulation of the gene, respectively.

6

Fig. 1

Fig. 2

Fig. 3

Fig. 4

Fig. 5