

1 **Short title: Integrated molecular phylogeny of diatom FCPs**

2

3 **Molecular phylogeny of fucoxanthin-chlorophyll *a/c* proteins from *Chaetoceros gracilis* and**  
4 **Lhcq/Lhcf diversity**

5

6 Minoru Kumazawa<sup>a</sup>, Hiroyo Nishide<sup>b</sup>, Ryo Nagao<sup>c</sup>, Natsuko Inoue-Kashino<sup>d</sup>, Jian-Ren Shen<sup>c</sup>, Takeshi  
7 Nakano<sup>a</sup>, Ikuo Uchiyama<sup>b</sup>, Yasuhiro Kashino<sup>d</sup>, Kentaro Ifuku<sup>e</sup>

8

9 <sup>a</sup>Graduate School of Biostudies, Kyoto University, Kyoto 606-8502, Japan

10 <sup>b</sup>National Institute for Basic Biology, National Institutes of Natural Sciences, Aichi 444-8585, Japan

11 <sup>c</sup>Research Institute for Interdisciplinary Science and Graduate School of Natural Science and  
12 Technology, Okayama University, Okayama 700-8530, Japan

13 <sup>d</sup>Graduate School of Life Science, University of Hyogo, Hyogo 678-1297, Japan

14 <sup>e</sup>Graduate School of Agriculture, Kyoto University, Kyoto 606-8502, Japan

15

16 ORCID IDs: 0000-0003-1911-7308 (M.K.); 0000-0001-8212-3001 (R.N.); 0000-0001-8282-5249  
17 (N.I.); 0000-0002-9197-5000 (I.U.); 0000-0002-4278-9359 (Y.K.); 0000-0003-4471-8797 (J.-R.S.);  
18 0000-0002-3188-7770 (T.N.); 0000-0003-0241-8008 (K.I.)

19

20 **One sentence summary:** Phylogenetic analysis of fucoxanthin-chlorophyll *a/c* proteins in *C. gracilis*  
21 revealed five major subfamilies and one minor subfamily, providing insights into the diversification  
22 of light-harvesting systems in red algae.

23

24 **List of author contributions:**

25 M.K., T.N., and K.I. conceived the project; M.K. performed the phylogenetic analysis; N.I. and Y.K.  
26 obtained the NGS data; M.K., H.N., and I.U. refined the genome information of *C. gracilis* and  
27 provided the amino acid sequences of FCPs; R.N. and J.R.S. provided the structural models of diatom  
28 photosystems; M.K. and K.I. wrote the manuscript; and all authors contributed to the interpretation of  
29 the results and improvement of the manuscript.

30

31 The author responsible for distribution of materials integral to the findings presented in this article in  
32 accordance with the policy described in the Instructions for Authors is: Kentaro Ifuku  
33 (ifuku.kentaro.2m@kyoto-u.ac.jp)

34

35

36

1    **Funding information**

2    This work was supported in part by the JST ALCA (grant nos. JPMJAL1105, JPMJAL1608 [K. I. and  
3    Y. K.]), by the JSPS KAKENHI (grant nos. JP20H031160 [K.I.], JP20K06528, JP21K19085,  
4    JP20H02914 [R.N.], and JP17H06433 [J.R.S.]), and by a Collaborative Research Program from  
5    National Institute for Basic Biology (grant no. 21-306 [K.I., Y.K., and I.U.]).

6

7    **Corresponding author:** Kentaro Ifuku; E-mail, [ifuku.kentaro2m@kyoto-u.ac.jp](mailto:ifuku.kentaro2m@kyoto-u.ac.jp); Tel., +81-75-753-

8    6109

1 **Abstract**

2 Diatoms adapt to various aquatic light environments and play major roles in the global carbon cycle  
3 using their unique light-harvesting system, i.e., fucoxanthin chlorophyll *a/c* binding proteins (FCPs).  
4 Structural analyses of photosystem II (PSII)-FCPII and photosystem I (PSI)-FCPI complexes from the  
5 diatom *Chaetoceros gracilis* have revealed the localization and interactions of many FCPs; however,  
6 the entire set of FCPs has not been characterized. Here, we identified 46 FCPs in the newly assembled  
7 genome and transcriptome of *C. gracilis*. Phylogenetic analyses suggested that these FCPs could be  
8 classified into five subfamilies: Lhcr, Lhcf, Lhcx, Lhcz, and novel Lhcq, in addition to a distinct type  
9 of Lhcr, CgLhcr9. The FCPs in Lhcr, including CgLhcr9 and some Lhcqs, had orthologous proteins  
10 in other diatoms, particularly those found in the PSI-FCPI structure. By contrast, the Lhcf subfamily,  
11 some of which were found in the PSII-FCPII complex, seemed to be diversified in each diatom species,  
12 and the number of Lhcqs differed among species, indicating that their diversification may contribute  
13 to species-specific adaptations to light. Further phylogenetic analyses of FCPs/light-harvesting  
14 complex (LHC) proteins using genome data and assembled transcriptomes of other diatoms and  
15 microalgae in public databases suggest that our proposed classification of FCPs was common among  
16 various red-lineage algae derived from secondary endosymbiosis of red algae, including Haptophyta.  
17 These results provided insights into the loss and gain of FCP/LHC subfamilies during the evolutionary  
18 history of the red algal lineage.

19

## 1 Introduction

2 Diatoms are a group of photosynthetic Stramenopiles (or Heterokonts), which are red-lineage  
3 secondary symbiotic algae with plastids derived from Rhodophyta (red algae). Diatoms are major  
4 primary producers in modern oceans (José et al., 2019). Unlike the green lineage, diatoms have a  
5 brown color owing to the presence of photosynthetic pigments different from those of the green lineage  
6 (e.g., chlorophyll [Chl] *a*, Chl *c*, fucoxanthin, and diadinoxanthin) in the light-harvesting pigment  
7 protein complex (LHC) surrounding their photosystems. This LHC in diatoms is called fucoxanthin  
8 Chl *a/c* binding protein (FCP), which absorbs light with blue-green wavelengths and thus captures  
9 more light in aqueous environments. In addition, LHC/FCP can function in non-photochemical  
10 quenching (NPQ), which dissipates the excitation energy of excessively absorbed light as heat (Niyogi  
11 and Truong 2013; Ruban 2018; Goss and Lepetit 2015; Wobbe et al. 2016; Giovagnetti and Ruban  
12 2018). The core subunits of photosynthetic protein complexes are highly conserved among oxygenic  
13 photosynthetic organisms; however, in most eukaryotic photosynthetic organisms, the LHC shows  
14 diversified sequences and pigment compositions to adapt to the living environment (Büchel, 2015;  
15 Büchel, 2020).

16 The first X-ray crystal structure of photosystem II of cyanobacteria was reported at the atomic  
17 level (Umena et al., 2011), and the structures of photosystems in green-lineage plants have been  
18 resolved by both X-ray crystallography and cryo-electron microscopy (EM) (Mazor et al., 2015; Qin  
19 et al., 2015; Wei et al., 2016; Mazor et al., 2017; Su et al., 2017; Shen et al., 2019). In addition,  
20 structures of photosystems in red-lineage plants have also been reported; photosystem II (PSII) of  
21 Rhodophyta *Cyanidium caldarium* was resolved by X-ray crystal structural analysis (Ago et al., 2016),  
22 whereas photosystem I of Rhodophyta *Cyanidioschyzon merolae* was resolved by cryo-EM (Pi et al.,  
23 2018). The structures of the PSII-FCPII supercomplex of the centric diatom *Chaetoceros gracilis* was  
24 reported as the first photosystem structure in secondary symbiotic algae (Nagao et al., 2019; Pi et al.,  
25 2019), followed by that of PSI-FCPI of *Chaetoceros gracilis* (Nagao et al., 2020). The structure of  
26 *Chaetoceros gracilis* PSI-FCPI, which has an increased number of FCPs, has also been reported (Xu  
27 et al., 2020). Accordingly, the molecular phylogeny of diatom FCPs can be interpreted on the basis of  
28 structural information and mass spectrometric identification of FCPs in the complexes separated by  
29 sucrose density gradient or native polyacrylamide gel electrophoresis (PAGE).

30 In our report on diatom PSI-FCPI (Nagao et al. 2020), we argued that FCPs in the outer  
31 edge of PSI-FCPI should belong to a novel group, Lhcq, a phylogenetic group different from that of  
32 Lhcr, which is commonly found in red-algal PSI (Nagao et al., 2020). Hoffman et al., (2011) reported  
33 the LHC/FCP phylogeny using expressed sequence tags of red-lineage species and the genome  
34 sequences of the red alga *Cyanidioschyzon merolae*, pennate diatom *Phaeodactylum tricorutum*, and  
35 centric diatom *Thalassiosira pseudonana*.

36 Here, we performed a more comprehensive analysis of the draft genome of *Chaetoceros*



1 *gracilis*, its FCP sequences, and the LHC/FCP sequences obtained from the genomes of other diatoms  
2 and algae in the red lineage. Overall, our results suggest that the diversified subfamilies of LHC/FCP,  
3 particularly those of Lhcf and Lhcq, had occurred in the common ancestral origin of red lineage algae,  
4 contributing to their high adaptability and prosperity in the ocean.

## 5 6 **Results**

### 7 *Assembly, gene prediction, and genome completeness*

8 Next-generation sequencing (NGS) data suggested that the estimated size of the *Chaetoceros gracilis*  
9 genome was 35.4 Mbp. The assembled draft nuclear genome contained 791 scaffolds and 3,408  
10 contigs with an N50 of 180 kbp and GC content of 37.3% (**Fig. 1A, B**). In total, 15,484 genes were  
11 predicted as nuclear-coded genes by BRAKER2 (Hoff et al., 2016). The assembly included chloroplast  
12 scaffolds with gene prediction performed using DOGMA. Benchmarking Universal Single-Copy  
13 Orthologs (BUSCO) suggested that 96% of conserved single-copy genes of *Stramenopiles\_odb10*  
14 were included in the predicted genes in the *Chaetoceros gracilis* nuclear genome (**Fig. 1C**), similar to  
15 values for *Thalassiosira pseudonana* (97%) and *Phaeodactylum tricornutum* (97%; **Supplemental**  
16 **Table S1**), indicating that our *Chaetoceros gracilis* draft genome had adequate coverage of essential  
17 genes. Using OrthoFinder, the predicted nuclear-coded genes in the *Chaetoceros gracilis* genome were  
18 classified into 7,320 orthogroups, including 5,563 orthogroups (77%) common with *Thalassiosira*  
19 *pseudonana* and 5,451 orthogroups (74.5%) common with *Phaeodactylum tricornutum* (**Fig. 1B**).

### 21 *Molecular phylogeny of FCPs obtained from genomes and transcriptomes*

22 *Chaetoceros gracilis* FCPs (CgFCPs) were exhaustively searched using the draft genome and long-  
23 read transcriptome data. Forty-four CgFCPs were obtained from the *Chaetoceros gracilis* draft  
24 genome using the FCP genes of *Thalassiosira pseudonana* (TpFCPs, 30 genes) and *Phaeodactylum*  
25 *tricornutum* (PtFCPs, 39 genes) in the NCBI RefSeq database as queries. The CgFCPs were further  
26 complemented by a long-read transcriptome, IsoSeq, for *Chaetoceros gracilis* from two culture  
27 conditions. Transcriptomes were refined using IsoSeq3 and isONclust. Refined transcriptomes by  
28 IsoSeq3 had BUSCO scores of 73% and 78%, respectively, and those from isONclust had scores of  
29 78% and 77%, respectively. Two CgFCPs, CgLhcf13 and CgLhcf14, which were not found in the draft  
30 genome, were detected by BLASTP search of transcriptomes using the same query set. Additionally,  
31 44 TpFCPs and 42 PtFCPs were exhaustively extracted from RefSeq genomes (Armbrust et al., 2004;  
32 Bowler et al., 2008; Rastogi et al., 2018) using BLASTP similarity search with 30 TpFCPs, 39 PtFCPs,  
33 and 46 CgFCPs as a query set. This exhaustive FCP search revealed that *Chaetoceros gracilis*,  
34 *Thalassiosira pseudonana*, and *Phaeodactylum tricornutum* had 46, 44, and 42 FCPs, respectively  
35 (**Supplemental Tables S2–4**).

36 Phylogenetic analyses of CgFCPs with curated FCPs from *Thalassiosira pseudonana* (**Fig.**

1 **2A)** or *Phaeodactylum tricornutum* (**Fig. 2B**) suggested that CgFCPs could be systemically named  
2 using the four major types: Lhcr, Lhcf, and Lhcx annotated in previous studies (Koziol *et al.*, 2007;  
3 Dittami *et al.*, 2010; Hoffman *et al.*, 2011) and the new subfamily named Lhcq. The Lhcr type included  
4 CgLhcr4 and CgLhcr9, as well as the Lhcz subfamily. Although CgLhcr4 and CgLhcr9 were not  
5 branched into the Lhcr clade in our phylogenetic analysis, they were included in “Lhcr” because of  
6 their locations in the PSI-FCPI complex (Nagao *et al.*, 2020).

7 The Lhcr subfamily is a red-algal-type LHC shared among both red algae and red-lineage  
8 secondary symbiotic algae. The Lhcr subfamily consists of LHCI in red algae (Pi *et al.*, 2018). The  
9 Lhcf subfamily was named after fucoxanthin, while other FCP subfamily proteins also bind it. Lhcf  
10 was also named as Fcp; However, this nomenclature is confusing and should be avoided. The Lhcf  
11 clade contained a branch of CgLhcf9, in which PtLhcf15 was included as a red-light-induced FCP  
12 (**Fig. 2B**). The unique functions of PtLhcf15 under red light conditions have been suggested  
13 (Herbstová *et al.*, 2017). Lhcx subfamily proteins are involved in photoprotection through NPQ in  
14 diatoms (Buck *et al.*, 2019). This Lhcx subfamily is homologous to Lhcsr, which is also responsible  
15 for energy-dependent NPQ (qE) (Tokutsu and Minagawa, 2013; Giovagnetti and Ruban, 2018). The  
16 Lhcz subfamily was found in Cryptophyceae, Haptophyta, and Chlorarachniophyta, although its  
17 expression, function, and localization are unknown (Koziol *et al.*, 2007). The Lhcz subfamily in  
18 diatoms has also been reported by Dittami *et al.* (2010). This Lhcz subfamily was assigned to the Lhcr  
19 clade or as a sister clade of the Lhcr subfamily in our phylogenetic trees. Therefore, the systematic  
20 names of the Lhcz subfamily are described as Lhcr herein.

21 The fifth subfamily is Lhcq, a novel FCP subfamily proposed in our previous study (Nagao  
22 *et al.*, 2020). The functions of the Lhcq subfamily are unknown. Although Lhcq proteins were not  
23 annotated in the model diatoms, Lhcqs were partially annotated as Lhcy proteins by Nymark *et al.*  
24 (2013) and Clade V by Hoffman *et al.* (2011). The Lhcq clade was distinguished by high support  
25 values (e.g., 93.9/1/95 for SH-aLRT support [%]/aBayes support/ultrafast bootstrap support [%]; **Fig.**  
26 **2A**). The Lhcq subfamily was more similar to the Lhcf subfamily than to the Lhcr and Lhcx  
27 subfamilies based on likelihood mapping analysis (Strimmer and von Haeseler, 1997) (**Supplemental**  
28 **Fig. S1**).

29 In addition to the five subfamilies, a minor number of FCPs comprised a monophyletic clade  
30 containing CgLhcr9. CgLhcr9 homologs also included the protein Pt17531 (protein ID 17531 in the  
31 JGI database; *Phaeodactylum tricornutum* CCAP 1055/1 v2.0,  
32 <https://phycocosm.jgi.doe.gov/Phatr2/Phatr2.home.html>). The functions of CgLhcr9 homologs were  
33 unknown until their localization in the PSI-FCPI complex was reported (Nagao *et al.*, 2020). CgLhcr9  
34 homologs in *Thalassiosira pseudonana* and *Phaeodactylum tricornutum* were named Lhcq because  
35 they were phylogenetically independent of the typical Lhcr subfamily; however, CgLhcr9 itself was  
36 still considered a member of the Lhcr subfamily because of its structural composition in PSI-FCPI.

1           Based on the above phylogenetic analysis, we proposed that 44 TpFCPs and 42 PtFCPs  
2 could be renamed into the four subfamily names (Lhcr, Lhcf, Lhcx, and Lhcq; **Supplemental Tables**  
3 **S2, S3**). In particular, TpFCPs and PtFCPs belonging to the Lhcq clade were renamed as Lhcq using  
4 our new annotations. Some Lhcrs, previously considered Lhcas (e.g., RefSeq ID: XP\_002287377.1  
5 and XP\_002289005.1) were renamed as Lhcrs. Consequently, there were nine Lhcrs, 14 Lhcfs, three  
6 Lhcxs, six Lhczs, 13 Lhcqs including CgLhcr4, and CgLhcr9 in *Chaetoceros gracilis*; 11 Lhcrs, 12  
7 Lhcfs, six Lhcxs, five Lhczs, nine Lhcqs, and one CgLhcr9 homolog in *Thalassiosira pseudonana*;  
8 and nine Lhcrs, 17 Lhcfs, four Lhcxs, seven Lhczs, four Lhcqs, and one CgLhcr9 homolog in  
9 *Phaeodactylum tricornutum* (**Fig. 2A, B**).

10           The centric diatoms *Chaetoceros gracilis* and *Thalassiosira pseudonana* had orthologous  
11 gene sets of Lhcr, Lhcz, Lhcq, and CgLhcr9 homologs, whereas some gene duplications and a minor  
12 exception, i.e., CgLhcr13 (Lhcz), were absent in *Thalassiosira pseudonana* (**Fig. 2A**). Notably, Lhcf-  
13 and Lhcx-type FCPs formed branches within each species, suggesting that Lhcr, Lhcz, Lhcq, and  
14 CgLhcr9 homologs may have conserved functions in both species, whereas Lhcf and Lhcx may have  
15 been differentiated within each species. A similar tendency was observed between *Chaetoceros*  
16 *gracilis* and *Phaeodactylum tricornutum* (**Fig. 2B**); however, *Phaeodactylum tricornutum* had a  
17 smaller number of Lhcq genes compared with *Chaetoceros gracilis* and *Thalassiosira pseudonana*.  
18 All PtLhcqs had putative orthologous FCPs in *Chaetoceros gracilis*, although several CgLhcq  
19 homologs was missing in *Phaeodactylum tricornutum*. We further extended our phylogenetic analysis  
20 to FCPs from other diatoms, such as *Thalassiosira oceanica* (Lommer et al., 2012), *Fistulifera solaris*  
21 (Tanaka et al., 2015), *Fragilariopsis cylindrus* CCMP1102 (Mock et al., 2017), and *Pseudo-nitzschia*  
22 *multistriata* (**Supplemental Table S5**). The relatively conserved set of Lhcr, Lhcz, and CgLhcr9  
23 homologs was found among all species, except *Pseudo-nitzschia multistriata*, which has only three  
24 Lhcr-type FCPs (**Supplemental Fig. S2A–C**), and the completely conserved set of Lhcrs among six  
25 species corresponding to those of *Chaetoceros gracilis* is listed in **Supplemental Table S6**. Diverged  
26 sets of Lhcf, Lhcq, and Lhcx were observed among other diatoms.

### 27 28 *Localization of five major FCP subfamilies in PSI-FCPI and PSII-FCPII structures of Chaetoceros* 29 *gracilis*

30 Nagao et al. (2020) named FCPI proteins of *Chaetoceros gracilis* based on their localization in the  
31 PSI-FCPI structure. Internal FCPs that formed a ring-like structure around the PSI core were named  
32 CgLhcr1–10, and peripheral FCPs, which bound the above internal FCPs, were named CgLhcq1–6  
33 (**Fig. 3A**). Among these FCPs, CgLhcr1–3, CgLhcr5–8, and CgLhcr10 belonged to the Lhcr  
34 subfamily; CgLhcr4 and CgLhcq1–6 branched into the Lhcq clade; and CgLhcr9 branched into an  
35 independent clade (**Fig. 2A, B**). The larger *Chaetoceros gracilis* PSI-FCPI supercomplex reported by  
36 Xu et al. (2020) also contained CgLhcq9 (FCPI-19), CgLhcq12 (FCPI-2), CgLhcf3 (FCPI-12), another

1 CgLhcq6 (FCPI-20), and CgLhcq5 (FCPI-21; **Fig. 3B**). Moreover, this structure contained three  
2 additional FCPs (FCPI-1, -17, and -18); however, the amino acid sequences used for the structural  
3 modeling were those from *Phaeodactylum tricornutum* (PtLhcf3/4:  
4 XP\_002177868.1/XP\_002177869.1) and *Fragilariopsis cylindrus* (OEU13194.1, Fracy1:210193 in  
5 JGI, and A0A1E7F4Y9 in UniProtKB; OEU18584.1 had a partial sequence of OEU13194.1).  
6 OEU13194.1 is a red algal lineage chlorophyll *a/b*-binding-like protein (redCAP) and is different from  
7 typical LHC proteins, such as FCP (Sturm et al., 2013). Because of the low sequence similarity,  
8 redCAP proteins were not obtained in our BLASTP search. If the sequences used for structural  
9 modeling were relevant to the actual CgFCP sequences, two FCPs (FCPI-17 and FCPI-18) modeled  
10 by PtLhcf3/4 sequences may belong to the Lhcf subfamily. Thus, eight Lhcr-, 11 Lhcq-, one CgLhcr9-,  
11 and three Lhcf-type FCPs as well as one unknown FCP may function as light-harvesting antennae for  
12 PSI.

13 The positions of the five Lhcrs in red algal PSI-LHCI were similar to those of CgLhcr1,  
14 CgLhcr5, CgLhcr6, CgLhcr7, and the unknown FCPI-1 (**Fig. 3B**); however, their orthologous  
15 relationships were not supported by our phylogenetic analysis (**Supplementary Fig. S3A, B**).  
16 Interestingly, CgLhcr9 bound to the PSI core in an orientation opposite that of other endogenous FCPI  
17 proteins. The unique binding mode of CgLhcr9 in PSI-FCPI is related to its separation from other  
18 medial FCPI proteins in the phylogenetic tree (**Fig. 2A, B**). The position occupied by the Lhcq protein  
19 in the PSI-FCPI of *Chaetoceros gracilis* is completely absent in the PSI-LHCI of the red alga  
20 *Cyanidioschyzon merolae* (Nagao et al., 2020; Pi et al., 2018). Therefore, the Lhcq subfamily,  
21 including Lhcr4, is likely a new addition from secondary endosymbiosis.

22 *Chaetoceros gracilis* PSII-FCPII formed dimers and had two tetramers and three monomers  
23 of FCPs per PSII core (Nagao et al., 2019; Pi et al., 2019) (**Fig. 3C**). Nagao et al. (2019) revealed that  
24 two tetramers consisted of CgLhcf1, and Pi et al. (2019) reported that the center monomer of three  
25 monomers was Lhca2, which was renamed CgLhcr17 based on the systematic nomenclature in this  
26 study. The presence of the Lhcr-type FCP in the PSII-FCPII complex suggests its special function in  
27 light energy transfer. Because of the resolution limit, the molecular identities of the other two FCP  
28 monomers in the PSII-FCPII complex are still unknown.

### 29 30 *Putative organization of FCPs surrounding photosystems in other diatoms*

31 The detailed structures of photosystems from other diatoms have not been elucidated. However, an  
32 orthologous set of Lhcr-type FCPs, including CgLhcr4 from the Lhcq subfamily, except for CgLhcr17  
33 homologs, was detected in purified PSI complexes from both centric *Thalassiosira pseudonana* and  
34 pennate *Phaeodactylum tricornutum* using mass spectrometry (Lepetit et al., 2010; Grouneva et al.,  
35 2011; Ikeda et al., 2013; Calvaruso et al., 2020) (**Figs. 4, 5**). Notably, the PSI-FCPI complex of  
36 *Thalassiosira pseudonana* reported by Calvaruso et al. (2020) lacks TpLhcr18 and TpLhcr20,

1 corresponding to CgLhcr3 and CgLhcr10, respectively (**Fig. 4**). These FCPs would be detached during  
2 the isolation process. TpLhcr17, an ortholog of CgLhcr17, was detected in a PSII-FCPII fraction  
3 (Calvaruso *et al.*, 2020). Thus, most Lhcrs have specific and conserved functions as antennae for PSI,  
4 with the exception of CgLhcr17 homologs for PSII.

5 In the centric diatom *Thalassiosira pseudonana*, almost full sets of Lhcq subfamily proteins,  
6 except for TpLhcq9 and TpLhcq6, were detected in the PSI-FCPI band separated by native PAGE  
7 (Ikeda *et al.*, 2013), whereas only TpLhcq7 and TpLhcq8, corresponding to CgLhcr4 and CgLhcq12  
8 located in the inner-ring PSI-FCPI, were detected in two other studies (Grouneva *et al.*, 2011;  
9 Calvaruso *et al.*, 2020) (**Fig. 4**). In the pennate diatom *Phaeodactylum tricornutum* (**Fig. 5**), only  
10 PtLhcq2, orthologous to CgLhcr4, was detected (Lepetit *et al.*, 2010; Grouneva *et al.*, 2011),  
11 suggesting a conserved role of CgLhcr4 homologs in PSI-FCPI. By contrast, PtLhcq1 and PtLhcq4,  
12 corresponding to CgLhcq12 and CgLhcq10, respectively, were detected in PSI-FCPI in one study  
13 (Grouneva *et al.*, 2011). The discrepancy between the two studies may be related to differences in  
14 cultivation conditions or purification processes. TpLhcq10 and PtLhcq5—CgLhcr9 homologs—were  
15 detected in both *Thalassiosira pseudonana* and *Phaeodactylum tricornutum* PSI-FCPI (Lepetit *et al.*,  
16 2010; Grouneva *et al.*, 2011; Ikeda *et al.*, 2013). Taken together, the organization of FCPs in the inner-  
17 ring FCPI surrounding PSI was similar among *Chaetoceros gracilis*, *Thalassiosira pseudonana*, and  
18 *Phaeodactylum tricornutum*, whereas the composition of outwardly bound FCPs was diverse among  
19 diatom species.

20 CgLhcf1 forms homotetramers and serves as the main antennae in the PSII-FCPII of  
21 *Chaetoceros gracilis* (Nagao *et al.*, 2020). In *Thalassiosira pseudonana* (**Fig. 4**), TpLhcf1-7 and  
22 TpLhcf11 were detected in the PSII-FCPII fraction (Calvaruso *et al.*, 2020), consistent with the  
23 function of Lhcf-type FCPs as the main antennae for PSII. Additionally, TpLhcf1-7 and TpLhcf11  
24 were also detected as “FCP trimers” in other studies (Grouneva *et al.*, 2011; Nagao *et al.*, 2013),  
25 indicating that FCPII consisting of Lhcf s was loosely attached to PSII and easily detached during  
26 isolation. In *Phaeodactylum tricornutum* (**Fig. 5**), isolation of the PSII-FCPII complex has not been  
27 reported, although free “FCP trimers” have been reported in several studies (Lepetit *et al.*, 2010;  
28 Grouneva *et al.*, 2011; Gundermann *et al.*, 2013; Nagao *et al.*, 2013), potentially representing  
29 detachment of FCPII from PSII. Notably, the exact oligomeric state of the freely isolated “FCP trimer”  
30 is unknown, although the structures of FCP tetramers (Nagao *et al.*, 2019; Pi *et al.*, 2019) and dimers  
31 (Wang *et al.*, 2019) have been elucidated, and the trimeric form has been observed in cryo-EM single  
32 particle analysis (Arshad *et al.*, 2021).

33 As described previously, CgLhcf3 was putatively assigned to *Chaetoceros gracilis* PSI-  
34 FCPI (Xu *et al.*, 2020). TpLhcf10 was detected in the *Thalassiosira pseudonana* PSI-FCPI fraction  
35 (Calvaruso *et al.*, 2020), but not in PSII-FCPII nor the “trimer” (Grouneva *et al.*, 2011; Nagao *et al.*,  
36 2013; Calvaruso *et al.*, 2020). Therefore, some Lhcf-type proteins may serve as FCPI in both species;



1 nevertheless, the diversification of Lhcf-type seems to have occurred independently (**Figs. 2A, 4**). In  
2 *Phaeodactylum tricornerutum*, PtLhcf2, PtLhcf3/4, PtLhcf9, PtLhcf14, and PtLhcf17 were detected in  
3 the PSI-FCPI fraction (Lepetit *et al.*, 2010), whereas PtLhcf8 and PtLhcf17 were detected elsewhere  
4 (Grouneva *et al.*, 2011) (**Fig. 5**). These Lhcf-type FCPs may compensate for the smaller number of  
5 Lhcqs in PSI-FCPI of *Phaeodactylum tricornerutum*. Further structural analysis of PSI-FCPI of  
6 *Phaeodactylum tricornerutum* is required.

7

#### 8 *Motif analysis of FCPs*

9 LHC/FCP is a pigment-protein complex with three transmembrane alpha helices ( $\alpha 1$ ,  $\alpha 2$ , and  $\alpha 3$ )  
10 (Engelken *et al.*, 2010) or helices B, C, and A (Kühlbrandt *et al.*, 1994; Bassi *et al.*, 1999) from the N-  
11 terminus. In all FCPs, the  $\alpha 1$  and  $\alpha 3$  helices have sequence similarity and highly conserved glutamate  
12 (E64 and E163 in CgLhcf1) and arginine (R69 and R168 in CgLhcf1) residues that interact in an  
13 interhelix manner (E64-R168 and E163-R69); thus, the interaction between the  $\alpha 1$  and  $\alpha 3$  helices  
14 seems to be stabilized, as indicated in green plant LHCII (Engelken *et al.*, 2010) (**Fig. 6**). The highly  
15 conserved glutamates of the  $\alpha 1$  and  $\alpha 3$  helices also coordinate Chls in the conserved composition.

16 For further motif analysis, multiple expectation maximization for motif elicitation (MEME,  
17 version 5.3.0) (Bailey *et al.*, 2009) was performed using translated sequences of FCP genes from  
18 *Chaetoceros gracilis* and *Thalassiosira pseudonana* (**Supplemental Fig. S4**). LHC/FCP had two  
19 common carotenoid (Car)-binding motifs in the extended sequence of the N-terminal side of the  $\alpha 1$   
20 helix ( $\alpha 1$  extension) and in the loop region between helices  $\alpha 2$  and  $\alpha 3$  ( $\alpha 2$ - $\alpha 3$  loop), with GFDPLG  
21 or similar sequences in some varieties (**Fig. 6**) (Bassi *et al.*, 1999; Engelken, Brinkmann, and Adamska  
22 2010). These conserved Car-binding motifs were detected by MEME as motif-2, -5, or -9 in the  $\alpha 1$   
23 extension and in the  $\alpha 2$ - $\alpha 3$  loop of Lhcr, Lhcx, Lhcz, and Lhcq subfamilies (**Supplemental Figs. S4,**  
24 **S5**). The typical GFDPLG sequence was conserved in the  $\alpha 1$  extension of Lhcr and in the  $\alpha 2$ - $\alpha 3$  loop  
25 of Lhcr and Lhcx (**Fig. 6**). Although the motif in the  $\alpha 2$ - $\alpha 3$  loop of Lhcr was less conserved, only  
26 Lhcr, the most ancestral FCP subfamily in diatoms, had a typical Car-binding sequence in both regions.

27 The Car-binding motif in the  $\alpha 1$  extension of Lhcf and Lhex contained GFFDPLG, with  
28 addition of phenylalanine to the typical sequence. The corresponding region of Lhcq (14/22 Lhcq  
29 sequences from *Chaetoceros gracilis* and *Thalassiosira pseudonana*) was [K/G]X[F/Y]DPLN, which  
30 had the same number of amino acids as Lhcf and Lhex. By contrast, CgLhcq9, CgLhcq12, CgLhcr4,  
31 and its homolog TpLhcq7 had various deletions or insertions before conserved aspartic acid residues  
32 in the Car-binding motif, and CgLhcq1, CgLhcq3, and their homologs had different residues between  
33 [K/G] and X (**Supplemental Fig. S6A**). The motif in the  $\alpha 2$ - $\alpha 3$  loop of Lhcq showed minor variations  
34 with X[F/Y]DP[F/L]G, whereas that of Lhcf was largely different from the typical Car-binding motif  
35 represented by GXXDFG, lacking proline and having different locations of aspartic acid residues.  
36 Thus, Lhcf was identified as the newest subfamily among the five major FCP subfamilies.

1           The conserved Chl-coordinating motif of SX[S/A][L/M]P, a part of MEME motif-6  
2 (**Supplemental Fig. S4**) was located on the stromal side (**Fig. 6**) and contained only the N-terminal  
3 sequence of Lhcr proteins, except for those of CgLhcr3 and CgLhcr10. Interestingly, FCPs of Lhcf,  
4 Lhcx, Lhcz, and Lhcq subfamilies, except for CgLhcx3 and TpLhcx4–6, completely lacked this motif  
5 in their N-terminal region, indicating that this motif was lost during the functional differentiation of  
6 FCPs in diatoms. In the Lhcr subfamily, this motif coordinates Mg<sup>2+</sup> of Chl *a* via the carbonyl oxygen  
7 in the peptide bond of S/A at 2.1–2.6 Å distance (Nagao *et al.*, 2020; Xu *et al.*, 2020). These Chls are  
8 located between CgLhcr1/2, -2/3, -5/6, -6/7, -7/8, and -8/10 and may contribute to the assembly and  
9 stabilization of FCPI and to energy transfer to adjacent Chls. CgLhcr10 had a PEPPI sequence instead  
10 of SX[S/A][L/M]P, and its glutamate residue coordinated Chl *a* at 2.4 Å via its side chain. The N-  
11 terminal region coordinating Mg<sup>2+</sup> of Chl *a* may be an ancestral property of red algal LHC because  
12 this motif was also observed in Lhcr1 of *Cyanidioschyzon merolae* (Pi *et al.*, 2018). Similarly, another  
13 sequence motif coordinated Mg<sup>2+</sup> of Chl *a* in the LHC proteins of green-lineage plants: Lhcas (PDB  
14 ID: 6JO5), CP26, and CP29 (PDB ID: 6KAD) had PX[W/F]LP in *Chlamydomonas reinhardtii* (Sheng  
15 *et al.*, 2019; Su *et al.*, 2019; Suga *et al.*, 2019), and LhcbMs (PDB ID: 6KAD) had  
16 [P/A][K/L][F/W]LGP (Sheng *et al.*, 2019). Each motif in the LHC proteins of *Chlamydomonas*  
17 *reinhardtii* coordinated Mg<sup>2+</sup> of Chl *a* via its carbonyl oxygen in the main chain of W/F at a 2.1–2.8  
18 Å distance.

19           CgLhcr3 and CgLhcr10 are unique among the diatom Lhcr subfamily because of differences  
20 in the N-terminal Chl-coordinating motif and their positions in the genome. Both CgLhcr3 and  
21 CgLhcr10 were encoded in scaffold00008 in the *Chaetoceros gracilis* draft genome, with their 5' ends  
22 placed head-to-head. Homologous genes for CgLhcr3 and CgLhcr10 (*TpLhcr18* and *TpLhcr20*,  
23 EJK71515.1 and RJK71517.1, and *PtLhcr14* and *PtLhcr13* in *Thalassiosira pseudonana*,  
24 *Thalassiosira oceanica*, and *Phaeodactylum tricornutum*, respectively) were also arranged in a head-  
25 to-head manner, suggesting differentiation from other Lhcr subfamily genes at an early stage of diatom  
26 diversification.

27           The C-terminal conserved motif PGSVP, a part of MEME motif-11 (**Supplemental Fig.**  
28 **S6B, C**), on the luminal side of the Lhcq subfamily coordinates Mg<sup>2+</sup> of Chl *a* via the carbonyl oxygen  
29 of the peptide bond from the serine residue (Nagao *et al.*, 2020) (**Fig. 6**). This PGSVP motif is  
30 conserved in the Lhcq subfamilies of *Chaetoceros gracilis*, *Thalassiosira pseudonana*, and  
31 *Phaeodactylum tricornutum*. The MEME motif-11, including this PGSVP motif, was also assigned to  
32 some Lhcf subfamily proteins (**Supplemental Fig. S6B**). However, the region of Lhcf proteins did  
33 not contain the first proline. Chl 316 with adjacent Chls and carotenoids in Lhcq proteins is involved  
34 in the excitation energy transfer between FCPs toward the circumferential direction in PSI-FCPI  
35 (Nagao *et al.*, 2020). The PGSVP motif corresponds to the TGKGP motif in LHCs of *Chlamydomonas*  
36 *reinhardtii* in multiple sequence alignment; however, the latter motif does not coordinate Mg<sup>2+</sup> of Chl

1 a. Therefore, the PGSVP motif specific to Lhcq was also obtained during FCP differentiation to retain  
2 additional Chl.

## 4 Discussion

### 5 *Distributions and functions of Lhcrs*

6 The Lhcr subfamily, which includes LHCI in red algae, is present in a wide variety of red algal lineages  
7 and contains an independent clade of the Lhcz subfamily (**Fig. 7A, Supplemental Figs. S3A, B, S7A–**  
8 **I**). However, LHCI of the two red algae analyzed herein did not show orthologous relationships with  
9 the Lhcrs of secondary endosymbiotic algae in the red algal lineage. By contrast, the secondary  
10 symbiotic algae in the red algal lineage had gene sets similar to those of Lhcrs with phylogenetic  
11 relevance to those in diatoms, suggesting that functional differentiation of the Lhcr subfamily occurred  
12 during secondary symbiotic events in the red algal lineage.

13 This Lhcr differentiation may have occurred gradually during red algae evolution; most diatoms  
14 shared an orthologous gene set of Lhcrs (**Supplemental Table S6**), although homologs of CgLhcr2  
15 and CgLhcr6 were not clearly distinguished in the phylogenetic trees. Other Stramenopiles have  
16 homologs of CgLhcr10 in Phaeophyceae and Raphidophyceae; however, they do not have homologs  
17 of CgLhcr3, and those of CgLhcr2 and CgLhcr6 were not clearly separated in their trees. In  
18 Haptophyta, the presence or number of homologs for CgLhcr2/CgLhcr6, CgLhcr3/CgLhcr10, and  
19 CgLhcr7/CgLhcr8 differed between species; *Phaeocystis antarctica* lacked CgLhcr2/CgLhcr6,  
20 CgLhcr3/CgLhcr10, and CgLhcr7/CgLhcr8, whereas *Emiliana huxleyi* lacked CgLhcr3/CgLhcr10  
21 and *Chrysochromulina tobinii* had homologs for CgLhcr2/CgLhcr6, but only one gene for  
22 CgLhcr7/CgLhcr8 and CgLhcr3/CgLhcr10. These differences could be related to incompleteness of  
23 genome or transcriptome analyses and may cause structural differences in the inner ring of FCPI  
24 surrounding PSI among the secondary symbiotic algae in the red algal lineage.

25 Interestingly, one central monomer (Nagao *et al.*, 2019; Pi *et al.*, 2019) of the three monomers  
26 in diatom PSII-FCPII was identified as CgLhcr17, the only FCPII belonging to the Lhcr subfamily.  
27 CgLhcr17 homologs were conserved in Stramenopiles, including *Pseudo-nitzschia multistriata*,  
28 Phaeophyceae, and Raphidophyceae, and only *Phaeocystis antarctica* from Haptophyta had a  
29 CgLhcr17 homolog (**Fig. 7A, Supplemental Fig. S7**). In the Haptophyta analyzed in this study,  
30 *Chrysochromulina tobinii* and CgLhcr17 homologs were missing, and in the other Haptophyta  
31 *Emiliana huxleyi*, the phylogeny among CgLhcr1, CgLhcr17, and XP\_005769864.1 was unclear.  
32 CgLhcr17 homologs belong to clade II described by Hoffman *et al.* (2011), which includes several  
33 haptophyte Lhcrs. Therefore, photosynthetic Stramenopiles and some Haptophytes may have an Lhcr-  
34 type FCP for PSII-FCPII. The specific function of CgLhcr17 in PSII-FCPII is unknown.

35 CgLhcr9 was designated Lhcr based on its binding location in *Chaetoceros gracilis* PSI-FCPI  
36 (Nagao *et al.*, 2020) and belongs to the independent clade from Lhcr and other subfamilies in the



1 phylogenetic tree. Clade VI described by Hoffman *et al.* (2011) corresponded to CgLhcr9 homologs,  
2 which were conserved in diatoms (**Fig. 7A, Supplemental Table S6**), other Stramenopiles, and two  
3 Haptophytes; however, the other Haptophyta *Emiliania huxleyi* showed unclear phylogeny around  
4 CgLhcr9.

#### 5 6 *Diversification of the Lhcf subfamily within each species*

7 The Lhcf subfamily is a group of FCPs that accumulates abundantly and includes FCPs assigned to  
8 the diatom PSII-FCPII. Free or “trimeric” fractions of FCPs are mainly composed of Lhcf-type FCPs,  
9 and few Lhcf s may also attach to the larger diatom PSI-FCPI (Xu *et al.*, 2020). Phylogenetic analyses  
10 suggested that Lhcf subfamily proteins were diversified within each species, indicating that changes  
11 in the Lhcf subfamily may be essential for adapting to light environments. One of the subordinate  
12 groups of Lhcf subfamilies is a group that includes CgLhcf9 and PtLhcf15, which is independent from  
13 other subclades of Lhcf subfamilies. PtLhcf15 constitutes red-shifted FCPs induced by red light  
14 exposure (Herbstová *et al.*, 2017). In this clade, *Chaetoceros gracilis* had only CgLhcf9, which may  
15 be induced by red light. CgLhcf9 homologs were identified in both diatoms and Haptophytes (**Figs.**  
16 **2A, 2B, 7A, Supplemental Fig. S2A–D**), indicating that the Lhcf subfamily is also involved in  
17 chronic adaptation. However, CgLhcf9 homologs were not detected in the Dinophyceae Peridinales  
18 transcriptome.

#### 19 20 *Lhcx subfamily for energy-dependent NPQ*

21 The Lhcx subfamily is responsible for energy-dependent NPQ (qE) components in diatoms and other  
22 red-lineage secondary symbiotic algae (Giovagnetti and Ruban, 2018) and is widely conserved among  
23 secondary symbiotic red-lineage algae. In green algae and moss, the homologous subfamily is called  
24 Lhcsr or LI818 protein (Zhu and Green, 2008), which is also responsible for qE (Bailleul *et al.*, 2010).  
25 Lhcx and Lhcsr subfamilies are classified into Clade V by Hoffman *et al.* (2011). Vascular plants lack  
26 Lhcsr subfamily proteins but have PsbS, which belongs to the LHC superfamily as an NPQ entity (Li  
27 *et al.*, 2000). Generally, Lhcx subfamily proteins are upregulated with increasing light intensity at the  
28 mRNA level, contributing to light intensity-dependent induction of NPQ, whereas only *PtLhcx1* is  
29 expressed constitutively under dark conditions (Bailleul *et al.*, 2010). Calvaruso *et al.* (2020) identified  
30 TpLhcx6\_1 in both PSI and PSII fractions of thylakoid membranes separated in the centric diatom  
31 *Thalassiosira pseudonana* under both low and high light conditions, TpLhcx4 in PSI from samples  
32 treated with high light, and TpLhcx1/2 in free fractions from samples treated with high light. Grouneva  
33 *et al.* (2011) identified TpLhcx4 in PSI from samples treated with high light, TpLhcx1/2 in free  
34 fractions from samples treated with high light, and PtLhcx1 in PSI and PtLhcx2 in free fractions of  
35 thylakoid membranes from *Phaeodactylum tricornutum*.

36 Orange carotenoid-binding protein (OCP) is responsible for NPQ in cyanobacteria and is

1 activated under high light conditions and connects to phycobilisome, membrane anchored light-  
2 harvesting pigment-protein complexes (Joshua *et al.*, 2005; Thurotte *et al.*, 2015; Kirilovsky, 2007;  
3 Kirilovsky and Kerfeld, 2016). Primary symbiotic red-lineage algae have NPQ capacity (Schubert *et*  
4 *al.*, 2011; Wu, 2016; Álvarez-Gómez *et al.*, 2019), whereas both OCP homologs and LhcX/Lhcr  
5 subfamilies are absent in red algal genomes (Tanaka *et al.*, 2004; Bhattacharya *et al.*, 2013). Thus, the  
6 molecular entity of NPQ in red algae remains unknown.

#### 7 8 *The Lhcq subfamily*

9 The Lhcq subfamily is a new subfamily of FCPs comprising the outer belt of FCPs in the PSI-FCPI  
10 complex of *Chaetoceros gracilis* (Nagao *et al.*, 2020; Xu *et al.*, 2020). *Chaetoceros gracilis* and  
11 *Phaeodactylum tricornutum* have different features of excitation energy transfer from FCP to PSI with  
12 different amounts of low-energy Chls in PSI and/or FCPI (Nagao *et al.*, 2018; Nagao *et al.*, 2019b;  
13 Tanabe *et al.*, 2020); this may be related to the reduced number of Lhcq proteins in *Phaeodactylum*  
14 *tricornutum* compared with *Chaetoceros gracilis*.

15 Among Lhcqs in *Chaetoceros gracilis*, CgLhcr4 was considered Lhcr because of its location  
16 in the inner ring of FCPs in PSI-FCPI, interacting with PsaB, PsaF, and Psa28. Unlike other Lhcq  
17 subfamily proteins, homologs of CgLhcr4 were widely conserved in the secondary symbiotic algae of  
18 the red lineage, i.e., Stramenopiles and Haptophytes (**Fig. 7A**). Therefore, the FCP compositions in  
19 the inner-ring PSI-FCPI should be conserved in Stramenopiles and Haptophytes. Photosynthetic  
20 Stramenopiles other than diatoms also had Lhcq subfamily proteins in addition to homologs of  
21 CgLhcr4; however, these are not orthologous to diatom Lhcqs. Haptophyta had Lhcq homologs  
22 belonging to a large sister clade of CgLhcq4, CgLhcq5, CgLhcq7, CgLhcq8, and CgLhcq10, in  
23 addition to CgLhcr4 homologs (**Supplemental Fig. S7C, E, H**). This suggested that peripheral region  
24 PSI-FCPI supercomplexes of other Stramenopiles and Haptophyta may be different from those of  
25 diatoms, indicating that CgLhcr4 homologs may be the oldest members of the Lhcq subfamily.

#### 26 27 *Hypothesis of diversification of FCP/LHC subfamilies in the red lineage*

28 Secondary symbiotic algae in the red algal lineage have obtained genes targeted for or encoded in the  
29 plastids from the symbiont of an ancient red alga. Red algae harvest light mainly via Lhcrs as antennae  
30 for PSI and use phycobilisome as antennae for PSII. However, phycobilisome genes are absent in the  
31 genomes of all secondary symbiotic algae in the red algal lineage. In diatoms, PSI uses Lhcrs, Lhcqs,  
32 the CgLhcr9 homolog, and several Lhcfs, whereas PSII uses the CgLhcr17 homolog and Lhcfs.

33 The diversification of LHC/FCP subfamilies was coupled with the phylogenetic  
34 diversification of red algal lineages. However, their phylogeny is complicated by symbiotic gene  
35 transfer (SGT) via primary, secondary, and tertiary endosymbiosis and horizontal gene transfer (HGT)  
36 (Keeling, 2013). Indeed, Dorrell *et al.* (2017) reported that 25% of plastid-targeted genes of red-

1 lineage secondary symbiotic algae were derived from the green lineage. Thus, we constructed a  
2 phylogenetic tree of 65 single-copy orthologous genes encoded in plastid genomes (**Supplemental**  
3 **Table S7**) detected by OrthoFinder (**Fig. 7**). In this tree, secondary symbiotic algae in the red algal  
4 lineage were suggested as monophyletic groups, consistent with a report by Kim *et al.* (2017).  
5 Haptophyta was inferred to be a sister group of Stramenopiles, consistent with the phylogenetic tree  
6 of the nuclear genome (Burki *et al.*, 2016). The Dinophyceae Peridinales was located in a clade of  
7 diatoms, as suggested by Horiguchi and Takano (2006). Overall, our phylogenetic tree using plastid-  
8 encoded genes did not contradict previously presented trees.

9 Confusion regarding the phylogenetic relationships of FCP/LHC subfamilies has hindered our  
10 understanding of the diversification process of red-lineage FCP/LHC subfamilies. Our FCP/LHC  
11 phylogenetic analysis including *Chlamydomonas reinhardtii* and *Physcomitrella patens* from green-  
12 lineage plants revealed that all subfamilies from the red lineage, except for Lhcx (Lhcsr) subfamilies,  
13 were independent from green-lineage LHCs, such as Lhca and Lhcb subfamilies, with high support  
14 values. Our analysis also suggested that the Lhcr subfamily was the most relevant subfamily and that  
15 Lhcq and Lhcf subfamilies were not derived from green-lineage LHC genes through HGT or SGT. By  
16 contrast, the process of acquiring Lhcxs was not clarified.

17 Similarities between Lhcq and Lhcf subfamilies were supported by likelihood mapping  
18 using 46 CgFCPs and 44 TpFCPs, suggesting that Lhcr, Lhcf, Lhcx, and Lhcq subfamilies could be  
19 grouped as (Lhcr, Lhcx)-(Lhcq, Lhcf). These findings were also supported by similarities in the  
20 pigment-binding motifs of Lhcqs and Lhcfs; the N-terminal motif in the stromal side coordinating Chl  
21 *a* was conserved in Lhcr, Lhca, and Lhcb subfamilies but absent in both Lhcq and Lhcf subfamilies,  
22 whereas the C-terminal motif in the luminal side coordinating Chl *a* was conserved in Lhcq and some  
23 Lhcf proteins. There were also differences between Lhcq and Lhcf subfamilies; for example, the Car-  
24 binding motif in the  $\alpha 2$ - $\alpha 3$  loop of Lhcq had proline, similar to other subfamilies, whereas that of Lhcf  
25 subfamily did not have proline. Thus, Lhcq and Lhcf may have a common ancestor derived from an  
26 ancestral Lhcr subfamily protein, and the Lhcf subfamily may have been derived from the Lhcq  
27 subfamily.

28 Based on these findings, we propose the following process for LHC/FCP diversification.  
29 First, the common ancestor of secondary symbiotic algae (excluding Cryptophyceae) acquired Lhcr  
30 subfamily genes from the red algal symbiont and diversified not only Lhcr genes but also CgLhcr9  
31 homologs and Lhcq genes, including CgLhcr4 homologs. This diversification enlarged the antenna  
32 size of the PSII. During this process, CgLhcr17 was derived from one of the Lhcr genes to fit into the  
33 PSII core instead of phycobilisomes, and Lhcf subfamily proteins diverged from the Lhcq subfamily,  
34 generating several monomers and tetramers and attaching to PSII as a major light-harvesting antenna.  
35 This hypothesis was supported by analyses of LHC/FCP distributions and their functions in various  
36 secondary symbiotic algae in the red lineage, particularly in those other than diatoms, using many

1 genome and transcriptome sequencing results from various species.

2

### 3 **Conclusion**

4 Our draft genome and transcriptome analyses suggested that *Chaetoceros gracilis* had 46 FCP genes,  
5 classified into five major subfamilies, i.e., Lhcr, Lhcf, Lhcx, Lhcz, and Lhcq, and one minor subfamily,  
6 i.e., CgLhcr9. FCPs of the inner light-harvesting ring of the PSI-FCPI complex were composed of  
7 Lhcrs, including CgLhcr9 and several Lhcqs, and were highly conserved in other diatom species. By  
8 contrast, Lhcfs, some of which were found in the PSII-FCPII complex, seemed to be diversified in  
9 each diatom species, and the number of Lhcqs differed among species. This indicated that  
10 diversification of Lhcf and Lhcq contributed to species-specific adaptations to the light environment.  
11 Other algae in Stramenopiles and Haptophyta possess the five major subfamilies and CgLhcr9  
12 homologs. Therefore, FCP/LHC diversification would have occurred in the common ancestral origin  
13 of red lineage algae.

14

### 15 **Materials and Methods**

#### 16 *Diatom cultivation*

17 The marine centric diatom *Chaetoceros gracilis* (UTEX LB 2658) was used for all analyses. Cell  
18 cultures were prepared in f/2 artificial seawater (Guillard, 1975) under 30  $\mu\text{mol photons m}^{-2} \text{s}^{-1}$  at 20°C  
19 with continuous shaking at 100 rpm. Additional cell culture for IsoSeq analysis was performed in  
20 artificial seawater under 30  $\mu\text{mol photons m}^{-2} \text{s}^{-1}$  at 30°C with continuous bubbling of air containing  
21 3% (v/v) CO<sub>2</sub> (Nagao *et al.*, 2007).

22

#### 23 *Genome sequencing and draft genome assembly*

24 Genomic DNA was isolated as previously described (Fischer *et al.*, 1999) and analyzed using the  
25 Genome Sequencer FLX+ System (GS FLX+; Roche Diagnostics, Basel, Switzerland), Genome  
26 Analyzer GAIIX (Illumina, Inc., San Diego, CA, USA), and Hiseq (Illumina, Inc.).

27 The sequencing library for GS FLX+ was prepared using the Paired End Library Preparation  
28 Method Manual (20 kb and 8 kb Span). The obtained library was amplified by emulsion polymerase  
29 chain reaction (PCR) using a GS FLK Titanium SV/LV emPCR Kit (Lib-L; Roche Diagnostics), added  
30 to a GS FLK Titanium PicoTiterPlate (Roche Diagnostics), and sequenced using the GS FLX+ System  
31 with a GS FLK Titanium Sequencing Kit XL+. The sequences of 744,262 reads containing  
32 319,847,738 bases (454 BaseCaller 2.6 in GS FLX+ system software) were used for the assembly  
33 process.

34 GAIIX sequencing was based on TruSeq DNA Sample Preparation v2 Guide Rev. A using a  
35 TruSeq DNA Sample Preparation v2 Kit. Sequences were called based on Genome Analyzer IIX User  
36 Guide version A and TruSeq SBS Kit v5 Reagent Preparation Guide (for Genome Analyzer) version

1 C. Sequences were processed following the Consensus Assessment of Sequence and Variation  
2 (CASAVA) v1.8 User Guide version B; 22,251,716 reads were processed.

3 The DNA sequences used in Hiseq were prepared using TruSeq. In total, 26,854,816 reads  
4 of 101 paired bases were obtained. Adapter sequences were eliminated using Cutadapt v1.1 (Martin,  
5 2011). Low-quality bases were trimmed using Trimmomatic v0.32 (Bolger, Lohse, and Usadel, 2014);  
6 paired reads of every 5 bases with a higher average quality score of 30 and whose lengths were longer  
7 than 74 survived; 16,222,538 reads survived (average length: 100.1).

8 Genome assembly was performed with all sequence data obtained from GS FLX+, GAIIx,  
9 and Hiseq using GS *De Novo* Assembler version 2.8 (Roche Diagnostics) with the following assembly  
10 parameters: -nrm/-het/-a0/-ml80%/-mi90/-urt/-large.

11

### 12 *RNA extraction*

13 RNA extraction was performed using a RNeasy kit (Qiagen Inc., Valencia, CA, USA) with some  
14 modifications. Cells were centrifuged, and pellets were suspended in 600  $\mu$ L RLT buffer containing  
15 1% (v/v)  $\beta$ -mercaptoethanol. The suspension was sonicated 10 times for 0.2–0.3 s each time using  
16 Handy Sonic UP-21P (Tomy, Japan) and centrifuged for 3 min at 15,000 rpm at room temperature.  
17 The supernatant was transferred to a new 1.5-mL microtube and mixed with the same volume of 70%  
18 ethanol. The mixture was further processed using the standard RNeasy protocol.

19

### 20 *RNA sequencing*

21 The library for RNA sequences was prepared using a Directional mRNA-Seq Library Prep. Pre-  
22 Release Protocol Rev.A with TruSeq RNA Sample Prep Kit and TruSeq Small RNA Sample Prep Kit  
23 (Illumina Inc.). The library was reverse transcribed, amplified by PCR with primers containing indexes,  
24 and purified by 6% agarose gel electrophoresis. Clustering was performed using cBot User Guide  
25 version F and TruSeq SR Cluster Kit v2 Reagent Preparation Guide (for cBot) version C and analyzed  
26 with Genome Analyzer IIX User Guide version A and TruSeq SBS Kit v5 Reagent Preparation Guide  
27 (for Genome Analyzer) version C. Base calling and processing were performed based on CASAVA  
28 v1.8 version B. The sequences were single reads of 75 bases. In total, 110,067,972 reads were obtained.  
29 RNA sequences were mapped to the genome using Hisat2 (Kim *et al.*, 2019).

30

### 31 *Gene prediction*

32 Genes in the *Chaetoceros gracilis* draft genome were predicted using BRAKER2 (Hoff *et al.*, 2016)  
33 with AUGUSTUS trained with RNA-seq mapping data. FCP genes were manually curated using  
34 RNA-seq mapping. Gene prediction of the chloroplast genome was performed using DOGMA  
35 (<https://dogma.cccb.utexas.edu/>). The information for predicted genes is available at ChaetoBase  
36 (<https://chaetoceros.nibb.ac.jp/>).

1

## 2 *Iso-Seq*

3 The Iso-seq libraries from two samples prepared under different cultivation conditions were generated  
4 according to the protocol provided by Pacific Biosciences (PN 101-763-800 Version 1; CA, USA),  
5 using NEBNext Single Cell/Low Input cDNA Synthesis & Amplification Module (New England  
6 Biolabs, MA, USA), Iso-Seq Express Oligo Kit, SMRTbell Express Template Prep Kit 2.0, and  
7 Barcoded Overhang Adapter Kit. The 5' and 3' primers GCAATGAAGTCGCAGGGTTGGG and  
8 AAGCAGTGGTATCAACGCAGAGTAC were used. Each Iso-Seq library was sequenced using  
9 Pacific Bioscience Sequel II. From each library, 1,968,854 and 1,035,268 raw reads were obtained,  
10 with average lengths of 2,744 and 3,310 bases, respectively.

11 The raw sequences were processed to ccs reads using SMRT Link v8.0.0, according to the  
12 SMRT Link User Guide (v8.0) version 09. The ccs reads were refined, and FLNC reads were created  
13 using IsoSeq3 refine (Gordon et al., 2015). FLNC reads were then clustered using IsoSeq3 (Gordon  
14 et al., 2015) with the verbose option and “use qvs”. The refined FLNC reads were also clustered using  
15 isONclust v0.0.6 (Sahlin and Medvedev, 2019). Open reading frames were extracted using  
16 TransDecoder v5.5.0 (Haas et al., 2013) from respective clustered reads processed by IsoSeq3 or  
17 isONclust, and two sets of translated sequences were obtained from each of two sets of raw reads.

18

## 19 *Genome and transcriptome quality assessment*

20 Basic statistics were analyzed using Seqkit (Shen et al., 2016). To assess genome or transcriptome  
21 assembly and gene prediction completeness, BUSCO (v4.0.6) (Seppey *et al.*, 2019) was performed in  
22 protein mode. The BUSCO lineage datasets used in our analyses were selected based on their  
23 taxonomy. Stramenopile\_odb10 was selected for diatoms, Raphidophyceae, and Pheophyceae  
24 (brown alga).

25

## 26 *Acquisition of the FCP sequences of Chaetoceros gracilis and model diatoms*

27 Forty-four sequences of CgFCPs were collected from the draft genome data using BLASTP 2.10.0  
28 similarity search (Altschul et al., 1990). In each BLASTP search, 30 and 39 known TpFCPs and  
29 PtFCPs collected from RefSeq were used as queries, and the expectation value (E-value) threshold  
30 was set to 1e-05. BLASTP searches were also conducted for each set of IsoSeq translated sequences  
31 generated by IsoSeq3 or isONclust. *CgLhc13* (GenBank ID: LC647435) and *CgLhc14* (GenBank  
32 ID: LC647436) were identified from the Iso-Seq sequences. Using 46 CgFCPs, known TpFCPs, and  
33 PtFCPs as queries, BLASTP similarity searches were performed to obtain FCP sequences from  
34 *Thalassiosira pseudonana* and *Phaeodactylum tricoratum* genomes (Armbrust *et al.*, 2004; Bowler  
35 *et al.*, 2008) with an E-value threshold of 1e-5. The 44 TpFCPs and 42 PtFCPs were phylogenetically  
36 analyzed for classification of their genes into different subfamilies, revised based on their phylogeny.



1 The lists of gene names (**Supplemental Table S3, S4**) were created based on UniProtKB  
2 (<https://www.uniprot.org/help/uniprotkb>).

3

4 *Acquisition of FCP sequences from other diatoms and Haptophytes and LHC sequences from red algae*

5 The sets of translated sequences of other diatoms were obtained from the genome assemblies of  
6 *Thalassiosira oceanica* (Lommer et al., 2012), *Fistulifera solaris* (Tanaka et al. 2015), *Fragilariopsis*  
7 *cylindrus* CCMP1102 (Mock et al., 2017), and *Pseudo-nitzschia multistriata*. Sets of translated  
8 sequences of other red-lineage microalgae were obtained from the following genome assemblies:  
9 Phaeophyceae, *Ectocarpus siliculosus* (Cock et al., 2010); three Haptophyta; *Emiliana huxleyi* (Read  
10 et al., 2013), *Phaeocystis antarctica* CCMP1374, and *Chrysochromulina tobinii* (Hovde et al., 2015);  
11 Cryptophyceae, *Guillardia theta* CCMP2712 (Curtis et al., 2012); and two Rhodophyta, red alga  
12 *Porphyridium purpureum* (Bhattacharya et al., 2013) and *Cyanidioschyzon merolae* (Tanaka et al.,  
13 2004). Sets of translated sequences derived from the RNA-seq data of Raphidophyceae *Chattonella*  
14 *antiqua* and Dinophyceae Peridinales *Heterocapsa circularisquama* were obtained from the database  
15 for research in harmful algal blooms (Shikata et al., 2019). Genome assemblies of green-lineage  
16 organisms have also been used, including *Chlamydomonas reinhardtii* (Merchant et al., 2007) and  
17 *Physcomitrella patens* (Rensing et al., 2008). The accession numbers and URLs are summarized in  
18 **Supplemental Table S5**. BLASTP 2.10.0 similarity searches with 1e-5 as the E-value threshold, using  
19 46 CgFCPs, 44 TpFCPs, and 42 PtFCPs as a query set, were conducted for the translated sequences  
20 of each genome. The identical LHC/FCP sequences were removed using CD-HIT v4.8.1 (Fu et al.,  
21 2012) with an identity threshold of 1.0.

22

23 *Maximum likelihood phylogenetic analysis*

24 Multiple sequence alignments were constructed using MAFFT-LINSI v7.4 (Katoh and Standley, 2013).  
25 Alignments were trimmed using ClipKIT (Steenwyk et al., 2020) with the “kpic-gappy” method, and  
26 maximum likelihood phylogenetic trees were constructed using IQ-TREE2 v2.0.7 (Minh et al., 2020).  
27 Ultrafast bootstrap (UFBoot2) (Hoang et al., 2018) approximation based on the model selected by  
28 ModelFinder (Kalyaanamoorthy et al., 2017) and the SH-like approximate likelihood ratio test  
29 (Guindon et al., 2010) were performed with 1000 replications in IQ-TREE2. aBayes test (Anisimova  
30 et al., 2011) was also performed. All trees were rerooted with Lhcx clade and drawn using  
31 FigTree (v1.4.4, <http://tree.bio.ed.ac.uk/software/figtree/>) or iTOL (5.7, <https://itol.embl.de/>) (Letunic  
32 and Bork, 2019). Likelihood mapping of Lhcq, Lhcx, Lhcf, and Lhcr (excluding Lhcz) subfamilies  
33 was performed using 46 CgFCPs and 44 TpFCPs. CgLhcr9 homologs and Lhcz subfamily sequences  
34 were ignored in this analysis.

35

36 *Motif analysis of Chaetoceros gracilis FCPs*

1 MEME (v5.3.0) (Bailey *et al.*, 2009) was performed with translated sequences of 46 CgFCP and 44  
2 TpFCPs to search 20 motifs. In MEME, the distribution of motifs was not limited, motif lengths were  
3 limited from 6 to 50, and other parameters were set to default. Alignment of *Chaetoceros gracilis* and  
4 *Thalassiosira pseudonana* FCPs, also used in the phylogenetic analysis, was applied to generate the  
5 amino acid sequence logos of the Car-binding motifs and Chl-coordinating motifs in Lhcr, Lhcf, Lhcx,  
6 and Lhcq. The logos were visualized using WebLogo 3.7.4 (Crooks *et al.*, 2004). The logo of the Car-  
7 binding motifs in the N-terminal extension of the  $\alpha 1$  helices ( $\alpha 1$  extensions) of Lhcr, Lhcf, Lhcx, and  
8 Lhcq were generated without using CgLhcr3, CgLhcr4, CgLhcf11, TpLhcx5, CgLhcq1, CgLhcq3,  
9 CgLhcq9, CgLhcq12, TpLhcq1, TpLhcq3, TpLhcq7, TpLhcq9, and TpLhcr18. As a result, 14 out of  
10 22 Lhcqs were used to generate this logo. The logo of the Car-binding motif in the loop region between  
11 helices  $\alpha 2$  and  $\alpha 3$  ( $\alpha 2$ - $\alpha 3$  loop) was generated without using the CgLhcf9 homolog clade, CgLhcf3,  
12 CgLhcf4, TpLhcr4, TpLhcr7, TpLhcr14, TpLhcf6, TpLhcf10, TpLhcx6\_1, and TpLhcq10 because of  
13 sequence divergence. The logo of the Chl-coordinating motif in the N-terminal extension of Lhcr was  
14 generated using the alignment of the proximal Lhcr subfamily, excluding CgLhcr3 and its homolog  
15 TpLhcr18. The logo of the Chl-coordinating motif in the C-terminal sequence of Lhcq was generated  
16 with alignment of the Lhcq subfamily, excluding TpLhcq9.

17

#### 18 *Visualization of the PSI-FCPI and PSII-FCPII structures*

19 PSI-FCPI, PSII-FCPII supercomplexes and FCP structures were visualized using the PyMOL  
20 Molecular Graphics System (Version 2.3.0 or 2.4.0 Schrödinger, LLC, <https://pymol.org/2/>).

21

#### 22 *Phylogenetic analysis of chloroplast genomes*

23 The set of translated sequences from chloroplast genomes was used for species phylogenetic analysis.  
24 The genomes were selected from NCBI (<https://www.ncbi.nlm.nih.gov/>) RefSeq or GenBank. All  
25 chloroplast genomes used in each analysis are listed in **Supplemental Table S5**, with corresponding  
26 accession numbers. Single-copy orthologs among each chloroplast genome were extracted using  
27 OrthoFinder v2.5.1 (Emms and Kelly, 2019). Every single-copy ortholog set was aligned using  
28 MAFFT v7.4 with the auto option and then trimmed using TrimAl (Capella-Gutiérrez *et al.*, 2009)  
29 with the automated1 option. The chloroplast phylogenetic tree was inferred using IQ-TREE2 with  
30 models automatically selected for each partition of the trimmed alignments. Bootstrap resampling was  
31 performed internally using UFBoot with 1000 replicates. Each tree was drawn using FigTree software.  
32 Glaucocystophyceae *Cyanophora paradoxa* was used as an outgroup in the chloroplast tree.

33

#### 34 **Accession Numbers**

35 Sequence data from this article can be found in our in-house database  
36 (<https://chaetoceros.nibb.ac.jp/>) or in the DDBJ Sequence Read Archive (DRA) under accession



1 numbers DRA012660 (genome sequencing), DRA012661 (RNA-Seq), and DRA012662 (Iso-Seq).

2

### 3 **Supplemental Data**

4 **Supplemental Figure S1.** Likelihood mapping of Lhcr, Lhcq, Lhcf, and Lhcx subfamilies (Strimmer  
5 and von Haeseler, 1997).

6 **Supplemental Figure S2.** Maximum-likelihood trees of FCPs/LHCs from *Chaetoceros gracilis* and  
7 other diatoms.

8 **Supplemental Figure S3.** Maximum likelihood tree of FCPs/LHCs from *Chaetoceros gracilis* and  
9 red algae.

10 **Supplemental Figure S4.** Maximum-likelihood tree of FCPs from *Chaetoceros gracilis* and  
11 *Thalassiosira pseudonana* showing the localization of the motifs generated by MEME.

12 **Supplemental Figure S5.** MEME motif logos contained the conserved carotenoid-binding motif  
13 “GFDPLG” with adjacent MEME motif logos.

14 **Supplemental Figure S6.** Specific motifs of Lhcq subfamily proteins: varieties of carotenoid-binding  
15 motifs and the novel C-terminal chlorophyll binding motif “PGSVP”.

16 **Supplemental Figure S7.** Maximum-likelihood phylogenetic tree of FCPs/LHCs from *Chaetoceros*  
17 *gracilis* and red- and green-lineage species.

18 **Supplemental Table S1.** List of assemblies used in FCP/LHC detection with BUSCO scores and  
19 lineages.

20 **Supplemental Table S2.** List of *Chaetoceros gracilis* FCPs with gene IDs or accession IDs.

21 **Supplemental Table S3.** List of all FCPs from *Thalassiosira pseudonana* with revised gene names.

22 **Supplemental Table S4.** List of all FCPs from *Phaeodactylum tricornutum* with revised gene names.

23 **Supplemental Table S5.** List of RefSeq or GenBank accession IDs or other references used to obtain  
24 the FCP/LHC sequences.

25 **Supplemental Table S6.** The conserved FCP set of diatoms, including the FCPs assigned to  
26 *Chaetoceros gracilis* photosystems.

27 **Supplemental Table S7.** List of RefSeq or GenBank accession IDs used to infer phylogenetic tree of  
28 chloroplast genes.

29

### 30 **Acknowledgements**

31 Computational resources were provided by the Data Integration and Analysis Facility, National  
32 Institute for Basic Biology. We would like to thank Editage ([www.editage.com](http://www.editage.com)) for English language  
33 editing.

34

### 35 **Competing interests**

36 The authors declare no competing interests.

1 **Figure Legends**

2 **Figure 1. Assessments of the *Chaetoceros gracilis* draft genome assembly. A,** General statistics of  
3 the *Chaetoceros gracilis* draft genome. **B,** Euler diagram of the orthogroups among *Chaetoceros*  
4 *gracilis* and two model diatom nuclear genomes, *Thalassiosira pseudonana* and *Phaeodactylum*  
5 *tricornutum*, with the draft genome size and the number of predicted genes. The diagram was  
6 generated using the Eulerr package (Wilkinson, 2012; Micallef and Rodgers, 2014) with R language.  
7 **C,** BUSCO scores for the predicted genes in the draft nuclear genome of *Chaetoceros gracilis* using  
8 the dataset Stramenopiles\_odb10.

9  
10 **Figure 2. Maximum-likelihood trees of FCPs from *Chaetoceros gracilis* (Cg) and *Thalassiosira***  
11 ***pseudonana* (Tp) and from *Chaetoceros gracilis* and *Phaeodactylum tricornutum* (Pt).** The trees  
12 were inferred using IQ-TREE 2 (Minh et al., 2020). The numbers of supporting values are SH-aLRT  
13 support (%)/aBayes support/ultrafast bootstrap support (%). Colors of clades are as follows: magenta,  
14 Lhcq subfamily; red, Lhcz subfamily; orange, Lhcr subfamily; brown, CgLhcr9 homologs; green, Lhcf  
15 subfamily (CgLhcf9 homolog clade is in gray); blue, Lhcx subfamily. Colors of gene names are as  
16 follows: red, *Chaetoceros gracilis* FCP; black, *Thalassiosira pseudonana* FCP. **A,** Maximum-  
17 likelihood tree of 46 CgFCPs and 44 TpFCPs. The tree was inferred using the LG+F+R4 model  
18 selected with ModelFinder (Kalyaanamoorthy et al., 2017). **B,** Maximum-likelihood tree of 46  
19 CgFCPs and 42 PtFCPs. The tree was inferred using the LG+F+R5 model selected with ModelFinder.

20  
21 **Figure 3. Structural arrangements of the photosystem I-FCPI supercomplex (A, PDB ID: 6L4U;**  
22 **B, PDB ID: 6LY5) and the photosystem II-FCPII supercomplex (C, PDB ID: 6J40) of**  
23 ***Chaetoceros gracilis*.** Top view of each supercomplex from the stromal side was depicted using  
24 PyMOL (Schrodinger LLC, 2015). The colors of FCPs are indicated as follows: magenta, Lhcq  
25 subfamily; red, Lhcz subfamily; orange, Lhcr subfamily; salmon pink, CgLhcr9 homologs; green,  
26 Lhcf subfamily. **A,** Sixteen FCPs were assigned in the PSI-FCPI supercomplex. **B,** Twenty FCPs were  
27 assigned, among which 24 FCPs were found in the larger PSI-FCPI supercomplex. Four unassigned  
28 FCPs are indicated as Xu *et al.* (2020). CgLhcq2 (q2\*) was assigned in **A:** 6L4U (Nagao et al., 2020),  
29 whereas CgLhcq6 was assigned in **B:** 6LY5 (Xu et al., 2020). **C,** CgLhcf1 tetramers were assigned in  
30 the dimeric PSII-FCPII supercomplex (Nagao *et al.*, 2019); CgLhcr17 (r17\*\*) was assigned in Pi *et al.*  
31 (2019). Two FCP monomers in each monomer of the PSII-FCPII were not assigned in both reports.  
32 The unassigned FCPs are shown in green.

33  
34 **Figure 4. Maximum-likelihood tree of FCPs from *Chaetoceros gracilis* (Cg) and *Thalassiosira***  
35 ***pseudonana* (Tp) combined with the table showing their previous detection in purified protein**  
36 **complexes.** The trees were inferred using IQ-TREE 2 (Minh et al., 2020) with the LG+F+R4 model

1 selected using ModelFinder (Kalyaanamoorthy et al., 2017). Numbers of supporting values are SH-  
2 aLRT support (%)/aBayes support/ultrafast bootstrap support (%). The tree was rerooted with the Lhcx  
3 subfamily. Detection of FCPs in each fraction or band is indicated by colored boxes as follows: red,  
4 PSI; blue, PSII; green, trimer; brown, free. Colors of clades are as follows: magenta, Lhcq subfamily;  
5 red, Lhcz subfamily; orange, Lhcr subfamily; brown, CgLhcr9 homologs; green, Lhcf subfamily  
6 (CgLhcf9 homolog clade is in gray); blue, Lhcx subfamily.

7

8 **Figure 5. Maximum-likelihood tree of FCPs from *Chaetoceros gracilis* (Cg) and *Phaeodactylum***  
9 ***tricornutum* (Pt) combined with the table showing their previous detection in purified protein**  
10 **complexes.** The trees were inferred using IQ-TREE 2 (Minh et al., 2020) with the LG+F+R5 model  
11 selected using ModelFinder (Kalyaanamoorthy et al., 2017). Numbers of supporting values are SH-  
12 aLRT support (%)/aBayes support/ultrafast bootstrap support (%). The tree was rerooted with the Lhcx  
13 subfamily. Detection of FCPs in each fraction or band is indicated by colored boxes as follows: red,  
14 PSI; blue, PSII; green, trimer; brown, free; purple, FCPs induced by red light. \*PtLhcr4, PtLhcr6,  
15 PtLhcr8, and PtLhcr10 proteins could be detected with a few peptides under HL, while they were  
16 completely missing under LL. Colors of clades are as follows: magenta, Lhcq subfamily; red, Lhcz  
17 subfamily; orange, Lhcr subfamily; brown, CgLhcr9 homologs; green, Lhcf subfamily (CgLhcf9  
18 homolog clade is in gray); blue, Lhcx subfamily.

19

20 **Figure 6. Structural localization and sequence logos of the pigment-binding motifs in FCPs.**  
21 CgLhcr5 and CgLhcq2 structures from *Chaetoceros gracilis* PSI-FCPI (PDB ID: 6L4U) were depicted  
22 using PyMOL (Schrodinger LLC, 2015). The cartoon model shows the side view of each FCP with  
23 the stromal side up. Not all Chls or carotenoids are shown. The amino acid residues and their  
24 coordinating or binding pigments are shown as stick models: carotenoid-binding motifs and  
25 carotenoids, purple; glutamate, orange; arginine, magenta; Lhcr N-terminal Chl-coordinating motif  
26 SX[S/A]X[L/M]P, yellow; Lhcq C-terminal Chl-coordinating motif PGSVP, cyan. Motif logos were  
27 created using WebLogo 3.7.4 (Crooks et al., 2004).

28

29 **Figure 7. Distribution of LHC/FCP subfamilies among red-lineage algae and hypothesis of these**  
30 **acquisitions based on the phylogenetic tree of chloroplast genes. A,** Numbers of FCP/LHC  
31 belonging to each subfamily detected from each species. **B,** Maximum-likelihood tree generated using  
32 chloroplast-encoded genes from various algal species indicating the estimated acquisition point of  
33 each LHC/FCP subfamily. The tree was constructed using models selected with ModelFinder  
34 (Kalyaanamoorthy et al., 2017) for each gene. The tree was rerooted with Graucocystophyceae.  
35 Numbers of supporting values are SH-aLRT support (%)/aBayes support/ultrafast bootstrap support  
36 (%).

1 **References**

- 2 **Ago H, Adachi H, Umena Y, Tashiro T, Kawakami K, Tian NKL, Han G, Kuang T, Liu Z, Wang F,**  
3 **et al** (2016) Novel features of eukaryotic photosystem II revealed by its crystal structure analysis  
4 from a red alga. *J Biol Chem* **291**: 5676–5687
- 5 **Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ** (1990) Basic local alignment search tool. *J Mol*  
6 *Biol* **215**: 403–410
- 7 **Álvarez-Gómez F, Korbee N, Figueroa FL** (2019) Effects of UV Radiation on Photosynthesis,  
8 Antioxidant Capacity and the Accumulation of Bioactive Compounds in *Gracilariopsis longissima*,  
9 *Hydropuntia cornea* and *Halopithys incurva* (Rhodophyta). *J Phycol* **55**: 1258–1273
- 10 **Anisimova M, Gil M, Dufayard JF, Dessimoz C, Gascuel O** (2011) Survey of branch support methods  
11 demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. *Syst*  
12 *Biol* **60**: 685–699
- 13 **Armbrust EV, Berges JA, Bowler C, Green BR, Martinez D, Putnam NH, Zhou S, Allen AE, Apt KE,**  
14 **Bechner M, et al** (2004) The genome of the diatom *Thalassiosira Pseudonana*: Ecology, evolution,  
15 and metabolism. *Science* (80- ) **306**: 79–86
- 16 **Arshad R, Calvaruso C, Boekema EJ, Büchel C, Kouřil R** (2021) Revealing the architecture of the  
17 photosynthetic apparatus in the diatom *Thalassiosira pseudonana*. *Plant Physiol* **186**: 2124–2136
- 18 **Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS** (2009)  
19 MEME Suite: Tools for motif discovery and searching. *Nucleic Acids Res* **37**: 202–208
- 20 **Bailleul B, Rogato A, De Martino A, Coesel S, Cardol P, Bowler C, Falciatore A, Finazzi G** (2010) An  
21 atypical member of the light-harvesting complex stress-related protein family modulates diatom  
22 responses to light. *Proc Natl Acad Sci U S A* **107**: 18214–18219
- 23 **Bassi R, Croce R, Cugini D, Sandonà D** (1999) Mutational analysis of a higher plant antenna protein  
24 provides identification of chromophores bound into multiple sites. *Proc Natl Acad Sci U S A* **96**:  
25 10056–10061
- 26 **Bhattacharya D, Price DC, Xin Chan C, Qiu H, Rose N, Ball S, Weber APM, Cecilia Arias M,**  
27 **Henrissat B, Coutinho PM, et al** (2013) Genome of the red alga *Porphyridium purpureum*. *Nat*  
28 *Commun.* doi: 10.1038/ncomms2931
- 29 **Bolger AM, Lohse M, Usadel B** (2014) Trimmomatic: A flexible trimmer for Illumina sequence data.  
30 *Bioinformatics* **30**: 2114–2120
- 31 **Bowler C, Allen AE, Badger JH, Grimwood J, Jabbari K, Kuo A, Maheswari U, Martens C, Maumus**  
32 **F, Otilar RP, et al** (2008) The *Phaeodactylum* genome reveals the evolutionary history of diatom  
33 genomes. *Nature* **456**: 239–244
- 34 **Büchel C** (2015) Evolution and function of light harvesting proteins. *J Plant Physiol* **172**: 62–75
- 35 **Büchel C** (2020) Light harvesting complexes in chlorophyll c-containing algae. *Biochim Biophys Acta -*  
36 *Bioenerg* **1861**: 148027

- 1 **Buck JM, Sherman J, Bártulos CR, Serif M, Halder M, Henkel J, Falciatore A, Lavaud J, Gorbunov**  
2 **MY, Kroth PG, et al** (2019) Lhc proteins provide photoprotection via thermal dissipation of  
3 absorbed light in the diatom *Phaeodactylum tricornutum*. *Nat Commun.* doi: 10.1038/s41467-019-  
4 12043-6
- 5 **Burki F, Kaplan M, Tikhonenkov D V., Zlatogursky V, Minh BQ, Radaykina L V., Smirnov A,**  
6 **Mylnikov AP, Keeling PJ** (2016) Untangling the early diversification of eukaryotes: A  
7 phylogenomic study of the evolutionary origins of centrohelida, haptophyta and cryptista. *Proc R Soc*  
8 *B Biol Sci.* doi: 10.1098/rspb.2015.2802
- 9 **Calvaruso C, Rokka A, Aro E-M, Büchel C** (2020) Specific Lhc proteins are bound to PSI or PSII  
10 supercomplexes in the diatom *Thalassiosira pseudonana*. *Plant Physiol* **183**: pp.00042.2020
- 11 **Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T** (2009) trimAl: A tool for automated alignment  
12 trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**: 1972–1973
- 13 **Cock JM, Sterck L, Rouzé P, Scornet D, Allen AE, Amoutzias G, Anthouard V, Artiguenave F, Aury**  
14 **JM, Badger JH, et al** (2010) The *Ectocarpus* genome and the independent evolution of  
15 multicellularity in brown algae. *Nature* **465**: 617–621
- 16 **Crooks GE, Hon G, Chandonia JM, Brenner SE** (2004) WebLogo: A sequence logo generator. *Genome*  
17 *Res* **14**: 1188–1190
- 18 **Curtis BA, Tanifuji G, Maruyama S, Gile GH, Hopkins JF, Eveleigh RJM, Nakayama T, Malik SB,**  
19 **Onodera NT, Slamovits CH, et al** (2012) Algal genomes reveal evolutionary mosaicism and the  
20 fate of nucleomorphs. *Nature* **492**: 59–65
- 21 **Dittami SM, Michel G, Collén J, Boyen C, Tonon T** (2010) Chlorophyll-binding proteins revisited - A  
22 multigenic family of light-harvesting and stress proteins from a brown algal perspective. *BMC Evol*  
23 *Biol* **10**: 365
- 24 **Dorrell RG, Gile G, McCallum G, Méheust R, Baptiste EP, Klinger CM, Brillet-Guéguen L,**  
25 **Freeman KD, Richter DJ, Bowler C** (2017) Chimeric origins of ochrophytes and haptophytes  
26 revealed through an ancient plastid proteome. *Elife* **6**: 1–45
- 27 **Emms DM, Kelly S** (2019) OrthoFinder: Phylogenetic orthology inference for comparative genomics.  
28 *Genome Biol* **20**: 1–14
- 29 **Engelken J, Brinkmann H, Adamska I** (2010) Taxonomic distribution and origins of the extended LHC  
30 (light-harvesting complex) antenna protein superfamily. *BMC Evol Biol.* doi: 10.1186/1471-2148-  
31 10-233
- 32 **Fischer H, Robl I, Sumper M, Kröger N** (1999) Targeting and covalent modification of cell wall and  
33 membrane proteins heterologously expressed in the diatom *Cylindrotheca fusiformis*  
34 (*Bacillariophyceae*). *J Phycol* **35**: 113–120
- 35 **Fu L, Niu B, Zhu Z, Wu S, Li W** (2012) CD-HIT: Accelerated for clustering the next-generation  
36 sequencing data. *Bioinformatics* **28**: 3150–3152

- 1 **Giovagnetti V, Ruban A V.** (2018) The evolution of the photoprotective antenna proteins in oxygenic  
2 photosynthetic eukaryotes. *Biochem Soc Trans* **46**: 1263–1277
- 3 **Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, Kang D, Underwood J, Grigoriev I V,**  
4 **Figueroa M, et al** (2015) Widespread Polycistronic Transcripts in Fungi Revealed by Single-  
5 Molecule mRNA Sequencing. *PLoS One* **10**: e0132628
- 6 **Goss R, Lepetit B** (2015) Biodiversity of NPQ. *J Plant Physiol* **172**: 13–32
- 7 **Grouneva I, Rokka A, Aro EM** (2011) The thylakoid membrane proteome of two marine diatoms outlines  
8 both diatom-specific and species-specific features of the photosynthetic machinery. *J Proteome Res*  
9 **10**: 5338–5353
- 10 **Guillard RRL** (1975) Culture of Phytoplankton for Feeding Marine Invertebrates BT - Culture of Marine  
11 Invertebrate Animals: Proceedings — 1st Conference on Culture of Marine Invertebrate Animals  
12 Greenport. In WL Smith, MH Chanley, eds, Springer US, Boston, MA, pp 29–60
- 13 **Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O** (2010) New algorithms and  
14 methods to estimate maximum-likelihood phylogenies: Assessing the performance of PhyML 3.0.  
15 *Syst Biol* **59**: 307–321
- 16 **Gundermann K, Schmidt M, Weisheit W, Mittag M, Büchel C** (2013) Identification of several sub-  
17 populations in the pool of light harvesting proteins in the pennate diatom *Phaeodactylum tricoratum*.  
18 *Biochim Biophys Acta - Bioenerg* **1827**: 303–310
- 19 **Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li**  
20 **B, Lieber M, et al** (2013) De novo transcript sequence reconstruction from RNA-seq using the  
21 Trinity platform for reference generation and analysis. *Nat Protoc* **8**: 1494–1512
- 22 **Herbstová M, Bína D, Kaňa R, Vácha F, Litvín R** (2017) Red-light phenotype in a marine diatom  
23 involves a specialized oligomeric red-shifted antenna and altered cell morphology. *Sci Rep* **7**: 1–10
- 24 **Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS** (2018) UFBoot2: Improving the  
25 Ultrafast Bootstrap Approximation. *Molecular biology and evolution*. *Mol Biol Evol* **35**: 518–522
- 26 **Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M** (2016) BRAKER1: Unsupervised RNA-Seq-  
27 based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics* **32**: 767–769
- 28 **Hoffman GE, Puerta MVS, Delwiche CF** (2011) Evolution of light-harvesting complex proteins from  
29 Chl c-containing algae. *BMC Evol Biol* **11**: 101
- 30 **Horiguchi T, Takano Y** (2006) Serial replacement of a diatom endosymbiont in the marine dinoflagellate  
31 *Peridinium quinquecorne* (Peridinales, Dinophyceae). *Phycol Res* **54**: 193–200
- 32 **Hovde BT, Deodato CR, Hunsperger HM, Ryken SA, Yost W, Jha RK, Patterson J, Monnat RJ,**  
33 **Barlow SB, Starkenburg SR, et al** (2015) Genome Sequence and Transcriptome Analyses of  
34 *Chrysochromulina tobin*: Metabolic Tools for Enhanced Algal Fitness in the Prominent Order  
35 Prymnesiales (Haptophyceae). *PLoS Genet* **11**: 1–31
- 36 **Ikeda Y, Yamagishi A, Komura M, Suzuki T, Dohmae N, Shibata Y, Itoh S, Koike H, Satoh K** (2013)



- 1 Two types of fucoxanthin-chlorophyll-binding proteins i tightly bound to the photosystem i core  
2 complex in marine centric diatoms. *Biochim Biophys Acta - Bioenerg* **1827**: 529–539
- 3 **José J, Karlusich P, Ibarbalz FM, Bowler C** (2019) Phytoplankton in the Tara Ocean. doi:  
4 10.1146/annurev-marine-010419
- 5 **Joshua S, Bailey S, Mann NH, Mullineaux CW** (2005) Involvement of phycobilisome diffusion in energy  
6 quenching in cyanobacteria. *Plant Physiol* **138**: 1577–1585
- 7 **Kalyaanamoorthy S, Minh BQ, Wong TKF, Von Haeseler A, Jermiin LS** (2017) ModelFinder: Fast  
8 model selection for accurate phylogenetic estimates. *Nat Methods* **14**: 587–589
- 9 **Katoh K, Standley DM** (2013) MAFFT multiple sequence alignment software version 7: Improvements  
10 in performance and usability. *Mol Biol Evol* **30**: 772–780
- 11 **Keeling PJ** (2013) The number, speed, and impact of plastid endosymbioses in eukaryotic evolution. *Annu*  
12 *Rev Plant Biol* **64**: 583–607
- 13 **Kim D, Paggi JM, Park C, Bennett C, Salzberg SL** (2019) Graph-based genome alignment and  
14 genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* **37**: 907–915
- 15 **Kim JI, Moore CE, Archibald JM, Bhattacharya D, Yi G, Yoon HS, Shin W** (2017) Evolutionary  
16 dynamics of cryptophyte plastid genomes. *Genome Biol Evol* **9**: 1859–1872
- 17 **Kirilovsky D** (2007) Photoprotection in cyanobacteria: The orange carotenoid protein (OCP)-related non-  
18 photochemical-quenching mechanism. *Photosynth Res* **93**: 7–16
- 19 **Kirilovsky D, Kerfeld CA** (2016) Cyanobacterial photoprotection by the orange carotenoid protein. *Nat*  
20 *Plants*. doi: 10.1038/nplants.2016.180
- 21 **Kozioł AG, Borza T, Ishida KI, Keeling P, Lee RW, Durnford DG** (2007) Tracing the evolution of the  
22 light-harvesting antennae in chlorophyll a/b-containing organisms. *Plant Physiol* **143**: 1802–1816
- 23 **Kühlbrandt W, Wang DN, Fujiyoshi Y** (1994) Atomic model of plant light-harvesting complex by  
24 electron crystallography. *Nature* **367**: 614–621
- 25 **Lepetit B, Volke D, Gilbert M, Wilhelm C, Goss R** (2010) Evidence for the existence of one antenna-  
26 associated, lipid-dissolved and two protein-bound pools of diadinoxanthin cycle pigments in diatoms.  
27 *Plant Physiol* **154**: 1905–1920
- 28 **Letunic I, Bork P** (2019) Interactive Tree of Life (iTOL) v4: Recent updates and new developments.  
29 *Nucleic Acids Res* **47**: 256–259
- 30 **Lommer M, Specht M, Roy AS, Kraemer L, Andreson R, Gutowska MA, Wolf J, Bergner S V.,**  
31 **Schilhabel MB, Klostermeier UC, et al** (2012) Genome and low-iron response of an oceanic diatom  
32 adapted to chronic iron limitation. *Genome Biol* **13**: R66
- 33 **Martin M** (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.  
34 *EMBnet.journal*; Vol 17, No 1 *Next Gener Seq Data Anal*. doi: 10.14806/ej.17.1.200
- 35 **Mazor Y, Borovikova A, Caspy I, Nelson N** (2017) Structure of the plant photosystem i supercomplex at  
36 2.6 Å resolution. *Nat Plants*. doi: 10.1038/nplants.2017.14

- 1 **Mazor Y, Borovikova A, Nelson N** (2015) The structure of plant photosystem I super-complex at 2.8 Å  
2 resolution. *Elife*. doi: 10.7554/eLife.07433
- 3 **Merchant SS, Prochnik SE, Vallon O, Harris EH, Karpowicz SJ, Witman GB, Terry A, Salamov A,**  
4 **Fritz-Laylin LK, Maréchal-Drouard L, et al** (2007) The *Chlamydomonas* Genome Reveals the  
5 Evolution of Key. *Natl institutes Heal* **318**: 245–250
- 6 **Micallef L, Rodgers P** (2014) eulerAPE: Drawing Area-Proportional 3-Venn Diagrams Using Ellipses.  
7 *PLoS One* **9**: e101717-
- 8 **Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R**  
9 (2020) IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic  
10 Era. *Mol Biol Evol* **37**: 1530–1534
- 11 **Mock T, Otilar RP, Strauss J, McMullan M, Pajananen P, Schmutz J, Salamov A, Sanges R, Toseland**  
12 **A, Ward BJ, et al** (2017) Evolutionary genomics of the cold-Adapted diatom *Fragilariopsis*  
13 *cylindrus*. *Nature* **541**: 536–540
- 14 **Nagao R, Kato K, Ifuku K, Suzuki T, Kumazawa M, Uchiyama I** (2020) Structural basis for assembly  
15 and function of a diatom photosystem I-light-harvesting supercomplex. *Nat Commun*. doi:  
16 10.1038/s41467-020-16324-3
- 17 **Nagao R, Kato K, Suzuki T, Ifuku K, Uchiyama I, Kashino Y, Dohmae N, Akimoto S, Shen JR,**  
18 **Miyazaki N, et al** (2019a) Structural basis for energy harvesting and dissipation in a diatom PSII–  
19 FCPII supercomplex. *Nat Plants* **5**: 890–901
- 20 **Nagao R, Takahashi S, Suzuki T, Dohmae N, Nakazato K, Tomo T** (2013) Comparison of oligomeric  
21 states and polypeptide compositions of fucoxanthin chlorophyll a/c-binding protein complexes  
22 among various diatom species. *Photosynth Res* **117**: 281–288
- 23 **Nagao R, Ueno Y, Yokono M, Shen JR, Akimoto S** (2019b) Effects of excess light energy on excitation-  
24 energy dynamics in a pennate diatom *Phaeodactylum tricornutum*. *Photosynth Res* **141**: 355–365
- 25 **Nagao R, Ueno Y, Yokono M, Shen JR, Akimoto S** (2018) Alterations of pigment composition and their  
26 interactions in response to different light conditions in the diatom *Chaetoceros gracilis* probed by  
27 time-resolved fluorescence spectroscopy. *Biochim Biophys Acta - Bioenerg* **1859**: 524–530
- 28 **Niyogi KK, Truong TB** (2013) Evolution of flexible non-photochemical quenching mechanisms that  
29 regulate light harvesting in oxygenic photosynthesis. *Curr Opin Plant Biol* **16**: 307–314
- 30 **Nymark M, Valle KC, Hancke K, Winge P, Andresen K, Johnsen G, Bones AM, Brembu T** (2013)  
31 Molecular and Photosynthetic Responses to Prolonged Darkness and Subsequent Acclimation to Re-  
32 Illumination in the Diatom *Phaeodactylum tricornutum*. *PLoS One*. doi:  
33 10.1371/journal.pone.0058722
- 34 **Pi X, Tian L, Dai HE, Qin X, Cheng L, Kuang T, Sui SF, Shen JR** (2018) Unique organization of  
35 photosystem I–light-harvesting supercomplex revealed by cryo-EM from a red alga. *Proc Natl Acad*  
36 *Sci U S A* **115**: 4423–4428



- 1 **Pi X, Zhao S, Wang W, Liu D, Xu C, Han G, Kuang T, Sui SF, Shen JR** (2019) The pigment-protein  
2 network of a diatom photosystem II–light-harvesting antenna supercomplex. *Science* (80- ). doi:  
3 10.1126/science.aax4406
- 4 **Qin X, Suga M, Kuang T, Shen JR** (2015) Structural basis for energy transfer pathways in the plant PSI-  
5 LHCI supercomplex. *Science* (80- ) **348**: 989–995
- 6 **Rastogi A, Maheswari U, Dorrell RG, Vieira FRJ, Maumus F, Kustka A, McCarthy J, Allen AE,**  
7 **Kersey P, Bowler C, et al** (2018) Integrative analysis of large scale transcriptome data draws a  
8 comprehensive landscape of *Phaeodactylum tricornutum* genome and evolutionary origin of diatoms.  
9 *Sci Rep*. doi: 10.1038/s41598-018-23106-x
- 10 **Read BA, Kegel J, Klute MJ, Kuo A, Lefebvre SC, Maumus F, Mayer C, Miller J, Monier A, Salamov**  
11 **A, et al** (2013) Pan genome of the phytoplankton *Emiliania* underpins its global distribution. *Nature*  
12 **499**: 209–213
- 13 **Rensing SA, Lang D, Zimmer AD, Terry A, Salamov A, Shapiro H, Nishiyama T, Perroud P-F,**  
14 **Lindquist EA, Kamisugi Y, et al** (2008) The *Physcomitrella* genome reveals evolutionary insights  
15 into the conquest of land by plants. *Science* **319**: 64–69
- 16 **Ruban A V.** (2018) Light harvesting control in plants. *FEBS Lett* **592**: 3030–3039
- 17 **Sahlin K, Medvedev P** (2019) De Novo Clustering of Long-Read Transcriptome Data Using a Greedy,  
18 Quality-Value Based Algorithm. *In* LJ Cowen, ed, *Res. Comput. Mol. Biol.* Springer International  
19 Publishing, Cham, pp 227–242
- 20 **Schrodinger LLC** (2015) The PyMOL Molecular Graphics System, Version 1.8.
- 21 **Schubert N, García-Mendoza E, Enríquez S** (2011) Is the photo-acclimatory response of Rhodophyta  
22 conditioned by the species carotenoid profile? *Limnol Oceanogr* **56**: 2347–2361
- 23 **Seppy M, Manni M, Zdobnov EM** (2019) BUSCO: Assessing Genome Assembly and Annotation  
24 Completeness BT - Gene Prediction: Methods and Protocols. *In* M Kollmar, ed, Springer New York,  
25 New York, NY, pp 227–245
- 26 **Shen L, Huang Z, Chang S, Wang W, Wang J, Kuang T, Han G, Shen JR, Zhang X** (2019) Structure  
27 of a C2S2M2N2-type PSII–LHCII supercomplex from the green alga *Chlamydomonas reinhardtii*.  
28 *Proc Natl Acad Sci U S A* **116**: 21246–21255
- 29 **Shen W, Le S, Li Y, Hu F** (2016) SeqKit: A cross-platform and ultrafast toolkit for FASTA/Q file  
30 manipulation. *PLoS One* **11**: 1–10
- 31 **Sheng X, Watanabe A, Li A, Kim E, Song C, Murata K, Song D, Minagawa J, Liu Z** (2019) Structural  
32 insight into light harvesting for photosystem II in green algae. *Nat Plants* **5**: 1320–1330
- 33 **Shikata T, Takahashi F, Nishide H, Shigenobu S, Kamei Y, Sakamoto S, Yuasa K, Nishiyama Y,**  
34 **Yamasaki Y, Uchiyama I** (2019) RNA-Seq Analysis Reveals Genes Related to Photoreception,  
35 Nutrient Uptake, and Toxicity in a Noxious Red-Tide Raphidophyte *Chattonella antiqua*. *Front*  
36 *Microbiol* **10**: 1–14

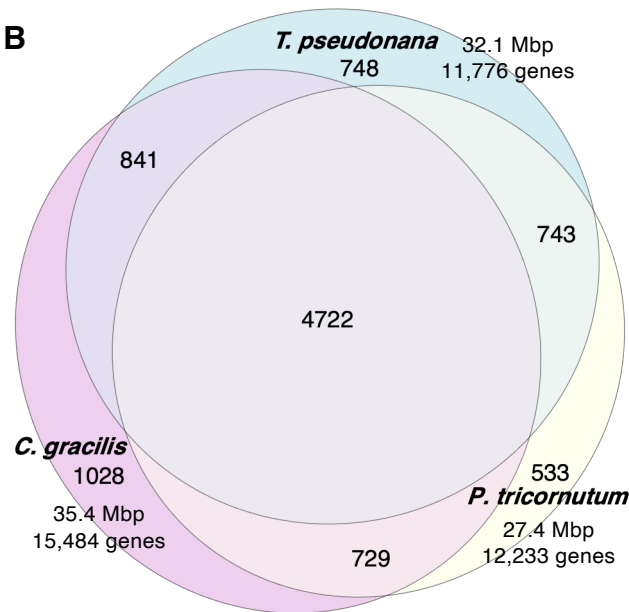
- 1 **Steenwyk J, Buida T, Li Y, Shen X-X, Rokas A** (2020) ClipKIT: a multiple sequence alignment-trimming  
2 algorithm for accurate phylogenomic inference. 1–17
- 3 **Strimmer K, von Haeseler A** (1997) Likelihood-mapping : A simple method to visualize phylogenetic.  
4 Proc Natl Acad Sci U S A **94**: 6815–6819
- 5 **Sturm S, Engelken J, Gruber A, Vugrinec S, G Kroth P, Adamska I, Lavaud J** (2013) A novel type of  
6 light-harvesting antenna protein of red algal origin in algae with secondary plastids. BMC Evol Biol.  
7 doi: 10.1186/1471-2148-13-159
- 8 **Su X, Ma J, Pan X, Zhao X, Chang W, Liu Z, Zhang X, Li M** (2019) Antenna arrangement and energy  
9 transfer pathways of a green algal photosystem-I-LHCI supercomplex. Nat Plants **5**: 273–281
- 10 **Su X, Ma J, Wei X, Cao P, Zhu D, Chang W, Liu Z, Zhang X, Li M** (2017) Structure and assembly  
11 mechanism of plant C2S2M2-type PSII-LHCII supercomplex. Science (80- ) **357**: 815–820
- 12 **Suga M, Ozawa SI, Yoshida-Motomura K, Akita F, Miyazaki N, Takahashi Y** (2019) Structure of the  
13 green algal photosystem I supercomplex with a decameric light-harvesting complex I. Nat Plants **5**:  
14 626–636
- 15 **Tanabe M, Ueno Y, Yokono M, Shen JR, Nagao R, Akimoto S** (2020) Changes in excitation relaxation  
16 of diatoms in response to fluctuating light, probed by fluorescence spectroscopies. Photosynth Res  
17 **146**: 143–150
- 18 **Tanaka K, Shimizu N, Sugano S, Sato N, Nozaki H, Ogasawara N, Kohara Y, Kuroiwa T** (2004)  
19 Genome sequence of the ultrasmall unicellular red alga Cyanidioschyzon merolae 10D. Nature **428**:  
20 653–7
- 21 **Tanaka T, Maeda Y, Veluchamy A, Tanaka M, Abida H, Maréchal E, Bowler C, Muto M, Sunaga  
22 Y, Tanaka M, et al** (2015) Oil accumulation by the oleaginous diatom *Fistulifera solaris* as revealed  
23 by the genome and transcriptome. Plant Cell **27**: 162–176
- 24 **Thurotte A, Igual RL, Wilson A, Comolet L, de Carbon CB, Xiao F, Kirilovsky D** (2015) Regulation  
25 of orange carotenoid protein activity in cyanobacterial photoprotection. Plant Physiol **169**: 737–747
- 26 **Tokutsu R, Minagawa J** (2013) Energy-dissipative supercomplex of photosystem II associated with  
27 LHCSR3 in *Chlamydomonas reinhardtii*. Proc Natl Acad Sci U S A **110**: 10016–10021
- 28 **Umena Y, Kawakami K, Shen JR, Kamiya N** (2011) Crystal structure of oxygen-evolving photosystem  
29 II at a resolution of 1.9Å. Nature **473**: 55–60
- 30 **Wang W, Yu LJ, Xu C, Tomizaki T, Zhao S, Umena Y, Chen X, Qin X, Xin Y, Suga M, et al** (2019)  
31 Structural basis for blue-green light harvesting and energy dissipation in diatoms. Science (80- ). doi:  
32 10.1126/science.aav0365
- 33 **Wei X, Su X, Cao P, Liu X, Chang W, Li M, Zhang X, Liu Z** (2016) Structure of spinach photosystem  
34 II-LHCII supercomplex at 3.2 Å resolution. Nature **534**: 69–74
- 35 **Wilkinson L** (2012) Exact and approximate area-proportional circular venn and euler diagrams. IEEE  
36 Trans Vis Comput Graph **18**: 321–331

- 1    **Wobbe L, Bassi R, Kruse O** (2016) Multi-Level Light Capture Control in Plants and Green Algae. Trends  
2        Plant Sci **21**: 55–68
- 3    **Wu H** (2016) Effect of Different Light Qualities on Growth, Pigment Content, Chlorophyll Fluorescence,  
4        and Antioxidant Enzyme Activity in the Red Alga *Pyropia haitanensis* (Bangiales, Rhodophyta).  
5        Biomed Res Int. doi: 10.1155/2016/7383918
- 6    **Xiao-Ping Li C, Grossman AR, Rosenquist M, Jansson S, Niyogi KK, Li X, Bjo O** (2000) A pigment-  
7        binding protein essential for regulation of photosynthetic light harvesting. Nature **403**: 391–395
- 8    **Xu C, Pi X, Huang Y, Han G, Chen X, Qin X, Huang G, Zhao S, Yang Y, Kuang T, et al** (2020)  
9        Structural basis for energy transfer in a huge diatom PSI-FCPI supercomplex. Nat Commun **11**: 1–  
10       12
- 11   **Zhu S-H, Green BR** (2008) Light-Harvesting and Photoprotection in Diatoms: Identification and  
12        Expression of L818-Like Proteins. Photosynth Energy from Sun 261–264
- 13

**A**

Genome assembly	Statistics
Genome size	35.4 Mbp
Scaffolds	791
Contigs	3408
N50	180 kbp
GC content	37.3 %

**B**



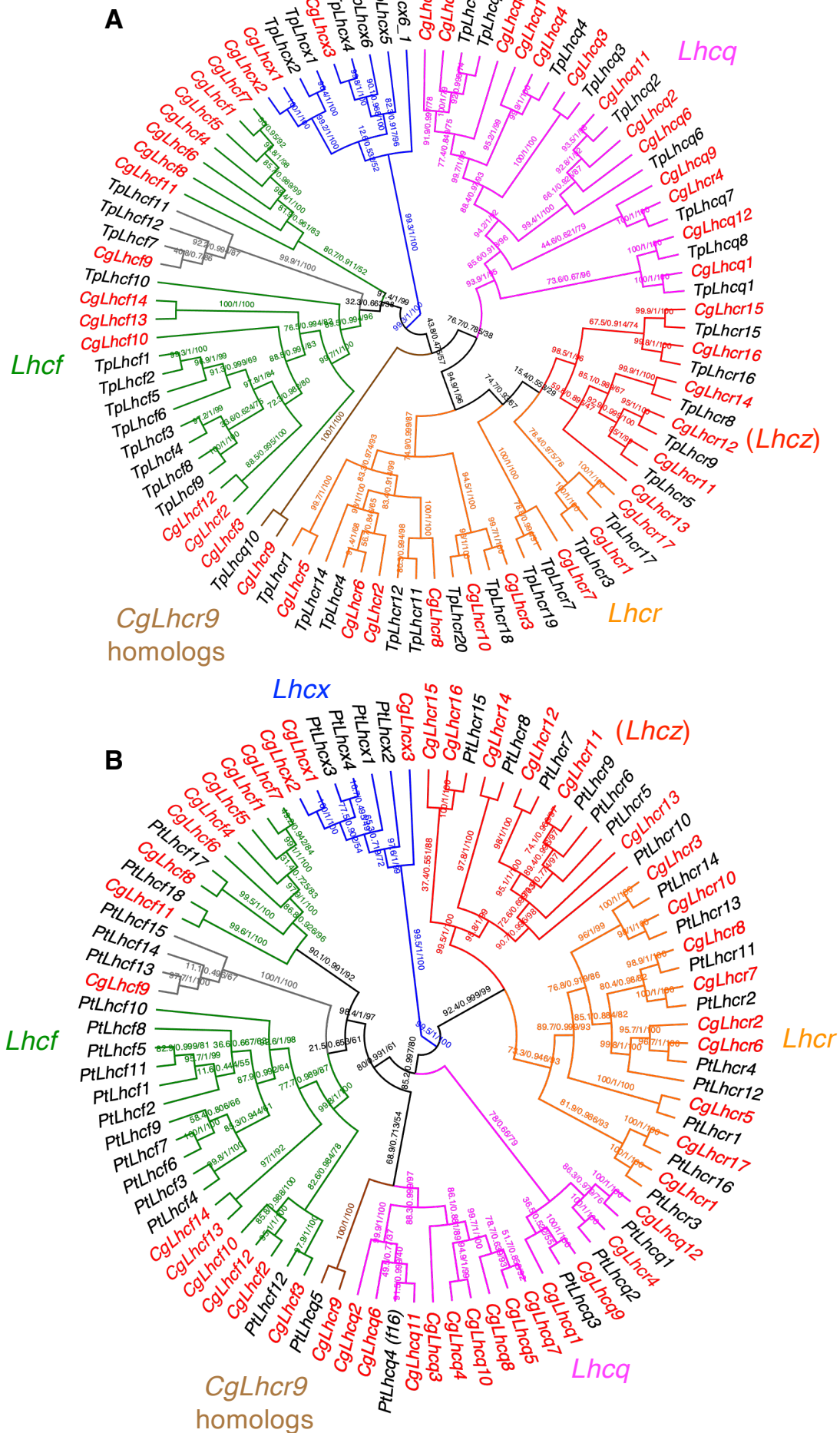
**C**

Results from dataset stramenopiles\_odb10

C:96.0%[S:90.0%,D:6.0%],F:3.0%,M:1.0%,n:100

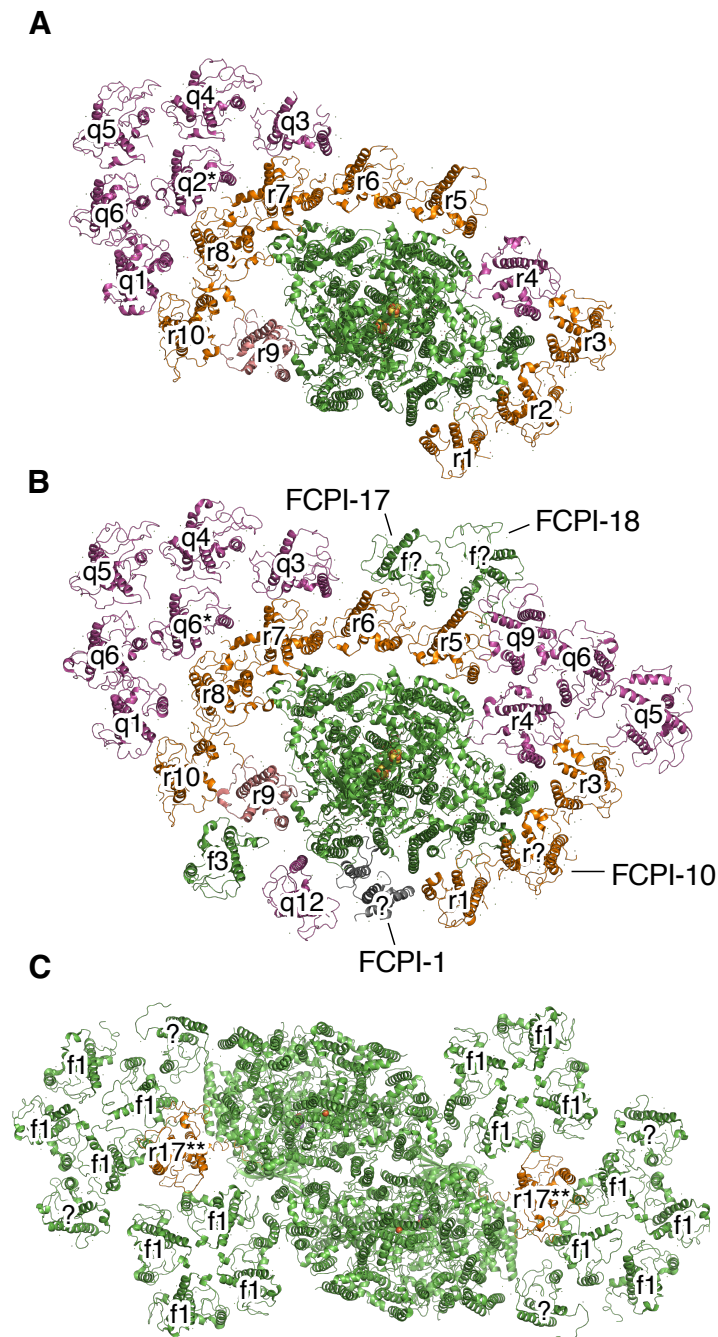
96	Complete BUSCOs (C)
90	Complete and single-copy BUSCOs (S)
6	Complete and duplicated BUSCOs (D)
3	Fragmented BUSCOs (F)
1	Missing BUSCOs (M)
100	Total BUSCO groups searched

**Figure 1. Assessments of the *Chaetoceros gracilis* draft genome assembly.** **A**, General statistics of the *Chaetoceros gracilis* draft genome. **B**, Euler diagram of the orthogroups among *Chaetoceros gracilis* and two model diatom nuclear genomes, *Thalassiosira pseudonana* and *Phaeodactylum tricornutum*, with the draft genome size and the number of predicted genes. The diagram was generated using the Eulerr package (Wilkinson, 2012; Micallef and Rodgers, 2014) with R language. **C**, BUSCO scores for the predicted genes in the draft nuclear genome of *Chaetoceros gracilis* using the dataset Stramenopiles\_odb10.

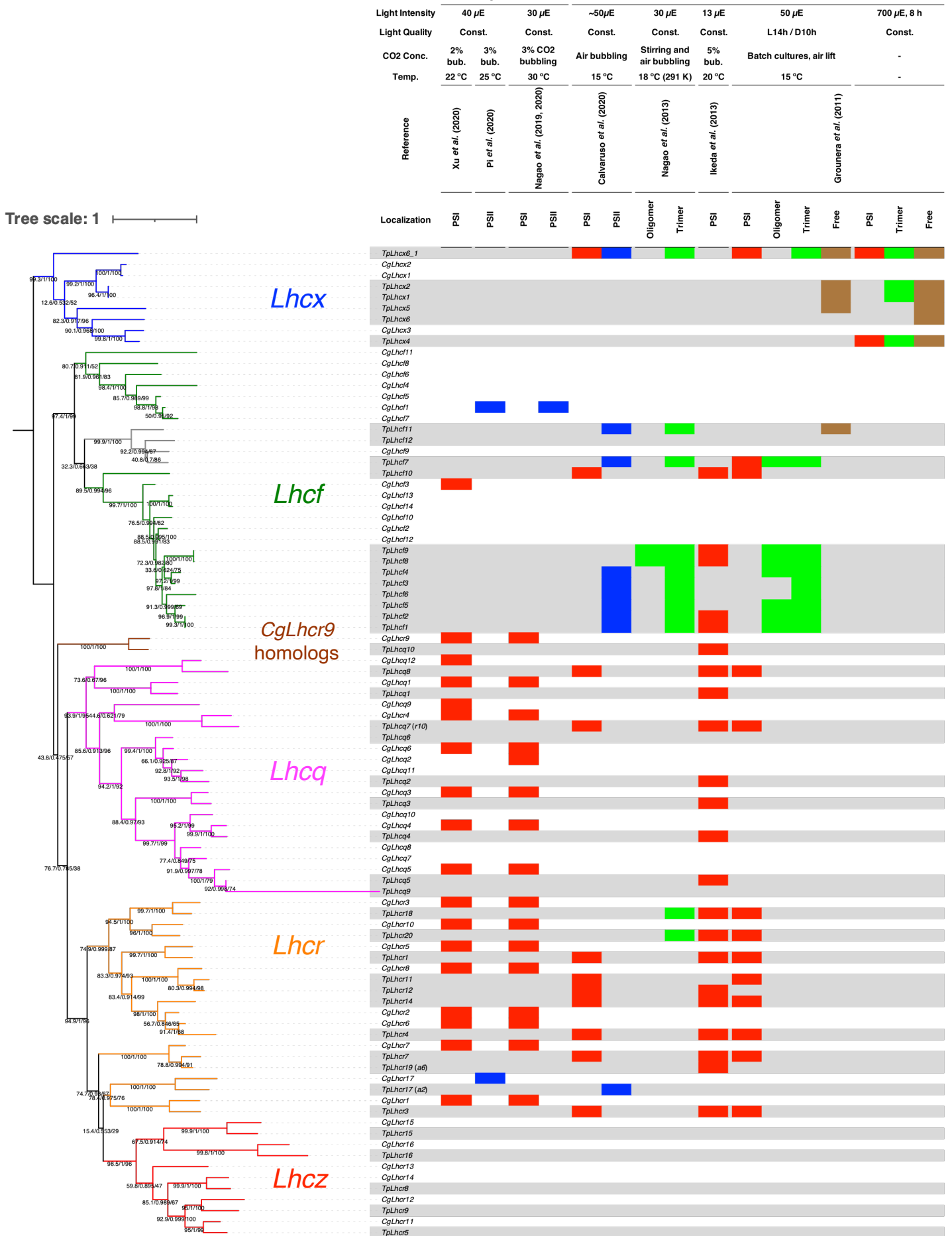


**Figure 2.** Maximum-likelihood trees of FCPs from *Chaetoceros gracilis* (Cg) and *Thalassiosira pseudonana* (Tp) and from *Chaetoceros gracilis* and *Phaeodactylum tricornutum* (Pt). The trees were inferred using IQ-TREE 2 (Minh *et al.*, 2020). The numbers of supporting values are SH-aLRT support (%) / aBayes support / ultrafast bootstrap support (%). Colors of clades are as follows: magenta, Lhcq subfamily; red, Lhcz subfamily; orange, Lhcr subfamily; brown, CgLhcr9 homologs; green, Lhcf subfamily; blue, Lhcx subfamily; grey, Lhcx subfamily. Colors of gene names are as follows: red, *Chaetoceros gracilis* FCP; black, *Thalassiosira pseudonana* FCP. **A**, Maximum-likelihood tree of 46 CgFCPs and 44 TpFCPs. The tree was inferred using the LG+F+R4 model selected with ModelFinder (Kalyaanamoorthy *et al.*, 2017). **B**, Maximum-likelihood tree of 46 CgFCPs and 42 PtFCPs. The tree was inferred using the LG+F+R5 model selected with ModelFinder.

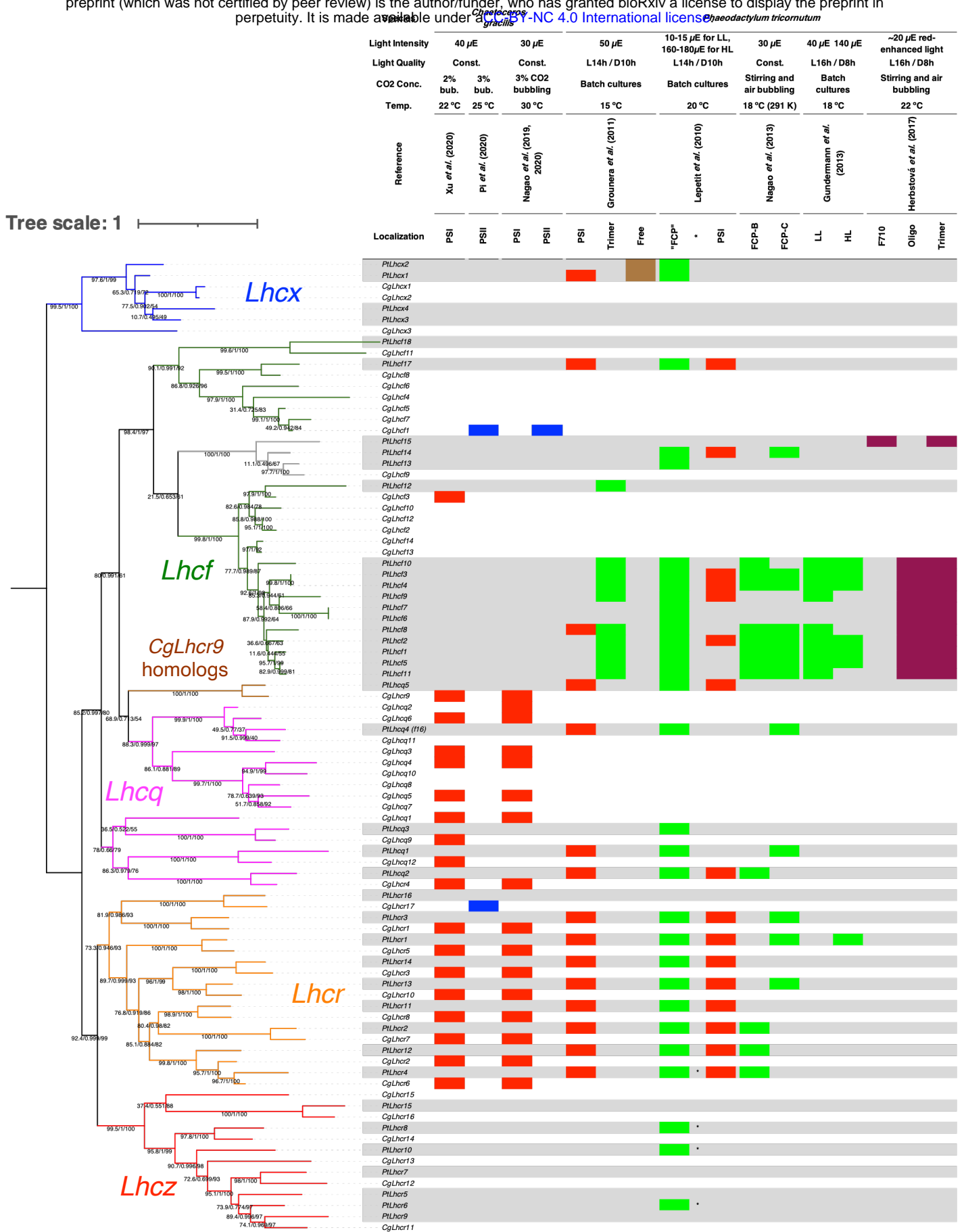




**Figure 3. Structural arrangements of the photosystem I-FCPI supercomplex (A, PDB ID: 6L4U; B, PDB ID: 6LY5) and the photosystem II-FCPII supercomplex (C, PDB ID: 6J40) of *Chaetoceros gracilis*.** Top view of each supercomplex from the stromal side was depicted using PyMOL (Schrodinger LLC, 2015). The colors of FCPs are indicated as follows: magenta, Lhcq subfamily; red, Lhcz subfamily; orange, Lhcr subfamily; salmon pink, CgLhcr9 homologs; green, Lhcf subfamily. **A**, Sixteen FCPs were assigned in the PSI-FCPI supercomplex. **B**, Twenty FCPs were assigned, among which 24 FCPs were found in the larger PSI-FCPI supercomplex. Four unassigned FCPIs are indicated as Xu *et al.* (2020). CgLhcr2 (q2\*) was assigned in **A**: 6L4U (Nagao *et al.*, 2020), whereas CgLhcr6 (q6) was assigned in **B**: 6LY5 (Xu *et al.*, 2020). **C**, CgLhcf1 tetramers were assigned in the dimeric PSII-FCPII supercomplex (Nagao *et al.*, 2019); CgLhcr17 (r17\*\*) was assigned in Pi *et al.* (2019). Two FCP monomers in each monomer of the PSII-FCPII were not assigned in both reports. The unassigned FCPs are shown in green.

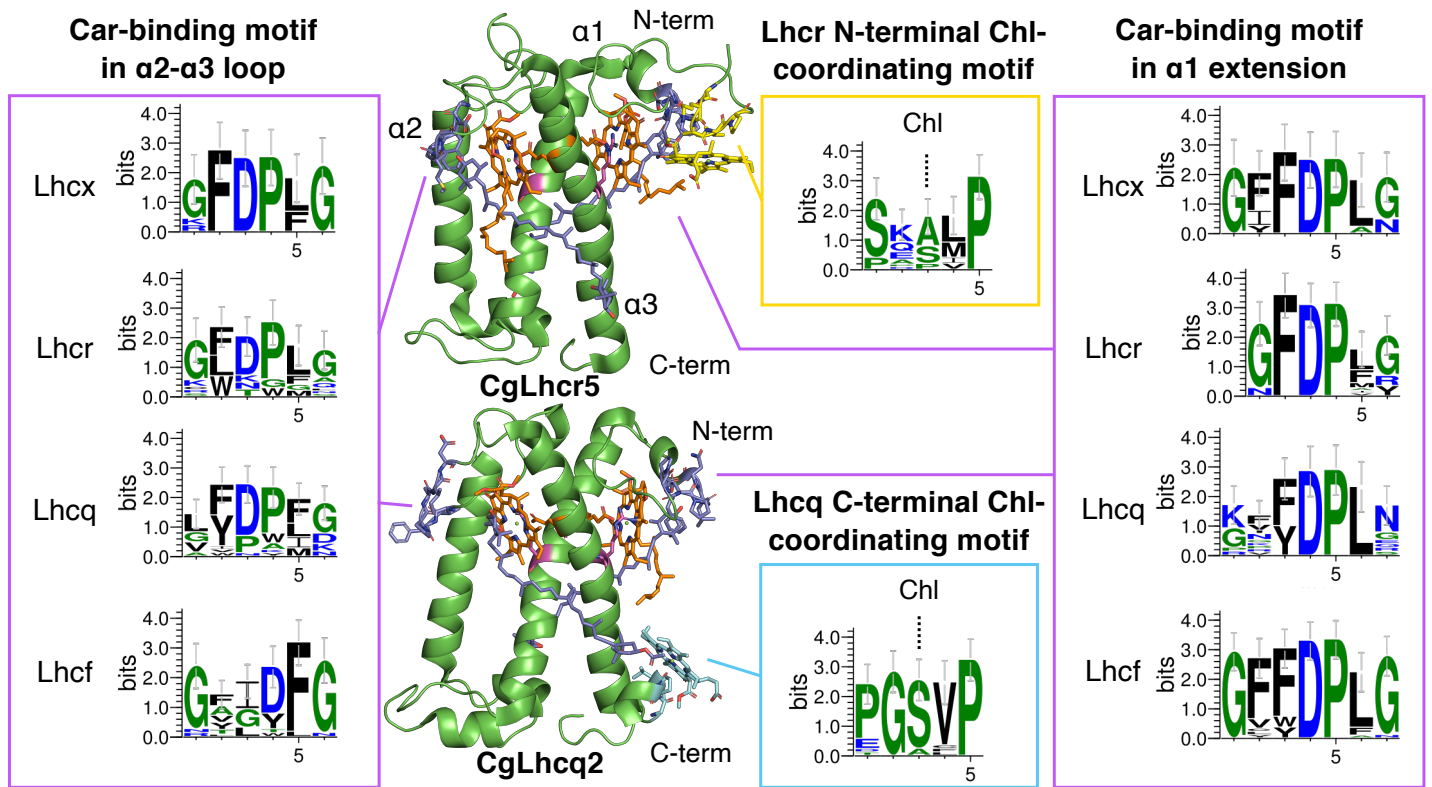


**Figure 4. Maximum-likelihood tree of FCPs from *Chaetoceros gracilis* (Cg) and *Thalassiosira pseudonana* (Tp) combined with the table showing their previous detection in purified protein complexes.** The trees were inferred using IQ-TREE 2 (Minh *et al.*, 2020) with the LG+F+R4 model selected using ModelFinder (Kalyaanamoorthy *et al.*, 2017). Numbers of supporting values are SH-aLRT support (%) / aBayes support / ultrafast bootstrap support (%). The tree was rooted with the Lhcx subfamily. Detection of FCPs in each fraction or band is indicated by colored boxes as follows: red, PSI; blue, PSII; green, trimer; brown, free. Colors of clades are as follows: magenta, Lhcq subfamily; red, Lhcz subfamily; orange, Lhcr subfamily; brown, CgLhcr9 homologs; green, Lhcf subfamily (CgLhcf9 homolog clade is in gray); blue, Lhcx subfamily.



**Figure 5. Maximum-likelihood tree of FCPs from *Chaetoceros gracilis* (Cg) and *Phaeodactylum tricornutum* (Pt) combined with the table showing their previous detection in purified protein complexes.** The trees were inferred using IQ-TREE 2 (Minh *et al.*, 2020) with the LG+F+R5 model selected using ModelFinder (Kalyaanamoorthy *et al.*, 2017). Numbers of supporting values are SH-aLRT support (%)/aBayes support/ultrafast bootstrap support (%). The tree was rerooted with the Lhcx subfamily. Detection of FCPs in each fraction or band is indicated by colored boxes as follows: red, PSI; blue, PSII; green, trimer; brown, free; purple, FCPs induced by red light. \*PtLhcr4, PtLhcr6, PtLhcr8, and PtLhcr10 proteins could be detected with a few peptides under HL, while they were completely missing under LL. Colors of clades are as follows: magenta, Lhcq subfamily; red, Lhcz subfamily; orange, Lhcr subfamily; brown, CgLhcr9 homologs; green, Lhcf subfamily (CgLhcf9 homolog clade is in gray); blue, Lhcx subfamily.

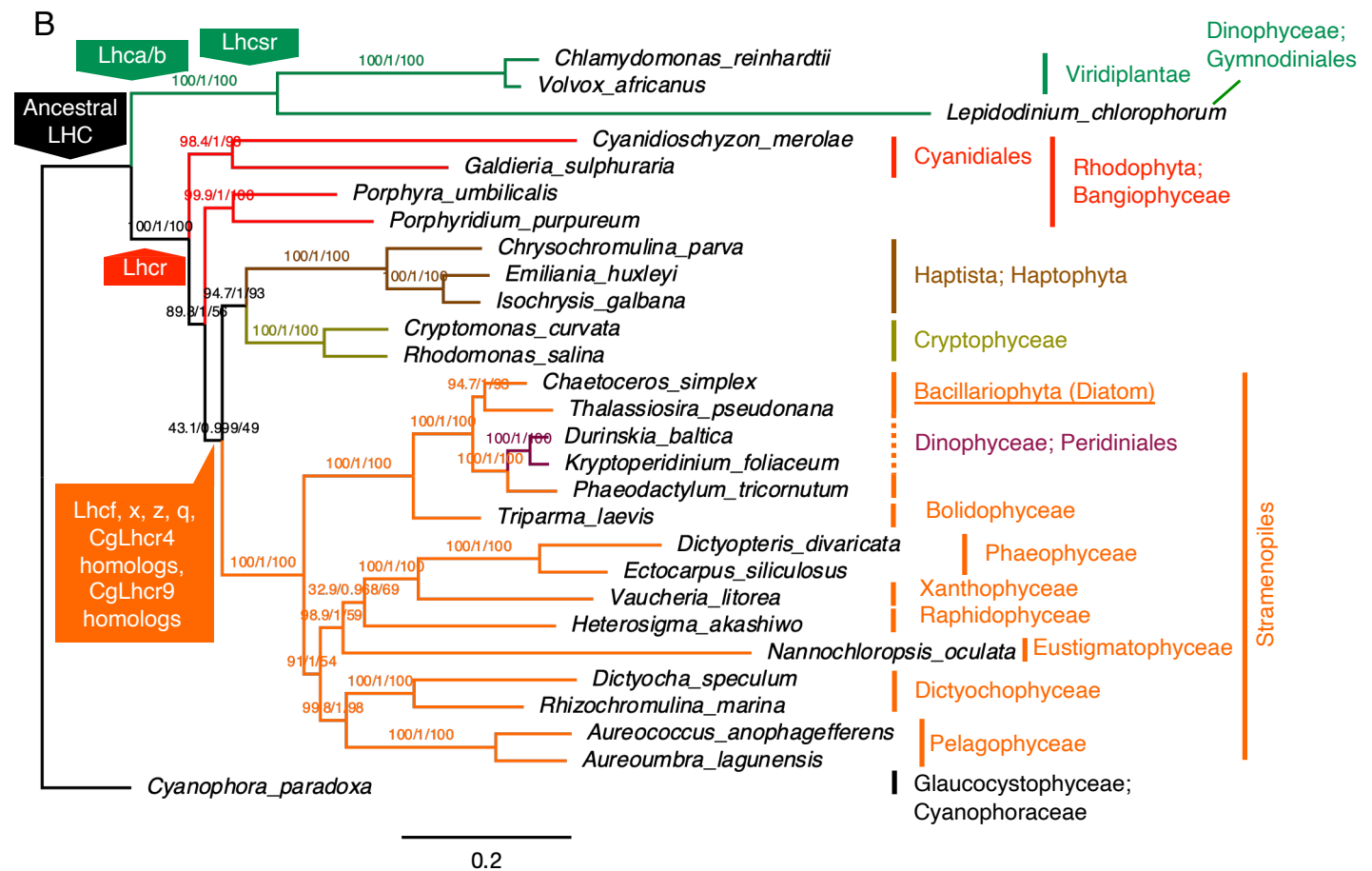




**Figure 6. Structural localization and sequence logos of the pigment-binding motifs in FCPs.** CgLhcr5 and CgLhcq2 structures from *Chaetoceros gracilis* PSI-FCPI (PDB ID: 6L4U) were depicted using PyMOL (Schrodinger LLC, 2015). The cartoon model shows the side view of each FCP with the stromal side up. Not all chlorophylls or carotenoids are shown. The amino acid residues and their coordinating or binding pigments are shown as stick models: carotenoid-binding motifs and carotenoids, purple; glutamate, orange; arginine, magenta; Lhcr N-terminal Chl-coordinating motif SX[S/A]X[L/M]P, yellow; Lhcq C-terminal Chl-coordinating motif PGSVP, cyan. Motif logos were created using WebLogo 3.7.4 (Crooks *et al.*, 2004).

**A**

Taxon	Species	Lhcr	CgLhcr17 homologs	Lhcz	Lhcx	Lhcf	CgLhcf9 homologs	Lhcq	CgLhcr4 homologs	CgLhcr9 homologs	Green-lineage	Other LHCs	Sum
Red alga (Rhodophyta)	<i>Cyanidioschyzon merolae</i>	2	0	0	0	0	0	0	0	0	0	0	2
Red alga (Rhodophyta)	<i>Porphyridium purpureum</i>	7	0	0	0	0	0	0	0	0	0	0	7
Brown alga (Phaeophyceae)	<i>Ectocarpus siliculosus</i>	9	1	1	14	20	0	2	1	1	0	4	53
Raphidophyceae	<i>Chattonella antiqua</i>	9	2	9	0	14	0	8	1	1	0	0	44
Dinophyceae Peridinales	<i>Heterocapsa circularisquama</i>	19	1	7	0	39	0	0	3	2	0	25	96
Diatom (Bacillariophyta)	<i>Pseudo-nitzschia multistriata</i>	3	1	8	4	12	1	9	0	0	0	0	38
Diatom (Bacillariophyta)	<i>Fragilariopsis cylindrus</i>	9	1	9	11	20	1	12	1	1	0	1	66
Diatom (Bacillariophyta)	<i>Fistulifera solaris</i>	12	2	10	6	22	5	8	2	2	0	0	69
Diatom (Bacillariophyta)	<i>Phaeodactylum tricornutum</i>	8	1	7	4	14	3	3	1	1	0	0	42
Diatom (Bacillariophyta)	<i>Chaetoceros gracilis</i>	8	1	6	3	13	1	12	1	1	0	0	46
Diatom (Bacillariophyta)	<i>Thalassiosira pseudonana</i>	10	1	5	6	9	3	8	1	1	0	0	44
Diatom (Bacillariophyta)	<i>Thalassiosira oceanica</i>	9	2	4	11	26	1	14	1	1	0	0	69
Haptophyta	<i>Emiliana huxleyi</i>	8	0	13	12	17	0	30	3	0	0	4	87
Haptophyta	<i>Chrysochromulina tobinii</i>	8	0	8	9	8	2	11	1	1	0	1	49
Haptophyta	<i>Phaeocystis antarctica</i>	14	1	17	28	13	0	35	3	1	0	2	114
Green alga (Chlorophyta)	<i>Chlamydomonas reinhardtii</i>	0	0	0	2	0	0	0	0	0	22	0	24
Land plant (Streptophyta)	<i>Physcomitrella patens</i>	0	0	0	2	0	0	0	0	0	45	0	47
Cryptophyceae	<i>Guillardia theta</i>	17	0	4	0	0	0	0	0	0	0	0	21



**Figure 7. Distribution of LHC/FCP subfamilies among red-lineage algae and hypothesis of these acquisitions based on the phylogenetic tree of chloroplast genes. A, Numbers of FCP/LHC belonging to each subfamily detected from each species. B, Maximum-likelihood tree generated using chloroplast-encoded genes from various algal species indicating the estimated acquisition point of each LHC/FCP subfamily. The tree was constructed using models selected with ModelFinder (Kalyaanamoorthy *et al.*, 2017) for each gene. The tree was rerooted with Glaucocystophyceae. Numbers of supporting values are SH-aLRT support (%) / aBayes support / ultrafast bootstrap support (%).**