



Royal Netherlands Institute for Sea Research

This is a postprint version of:

Villanueva, L., Rijpstra, W. I. C., Schouten, S., & Sinninghe Damsté, J. S. (2014). Genetic biomarkers of the sterol-biosynthetic pathway in microalgae. *Environmental Microbiology Reports*, 6(1), 35-44.

Published version: <http://dx.doi.org/10.1111/1758-2229.12106>

Link NIOZ Repository: [www.vliz.be/nl/imis?module=ref&refid=239837](http://www.vliz.be/nl/imis?module=ref&refid=239837)

[Article begins on next page]

The NIOZ Repository gives free access to the digital collection of the work of the Royal Netherlands Institute for Sea Research. This archive is managed according to the principles of the [Open Access Movement](#), and the [Open Archive Initiative](#). Each publication should be cited to its original source - please use the reference as presented. When using parts of, or whole publications in your own work, permission from the author(s) or copyright holder(s) is always needed.

# Genetic biomarkers of the sterol-biosynthetic pathway in microalgae

Laura Villanueva<sup>\*</sup>, W. Irene C. Rijpstra, Stefan Schouten, and  
Jaap S. Sinninghe Damsté

NIOZ Royal Netherlands Institute for Sea Research, Department of Marine Organic  
Biogeochemistry, PO Box 59, 179AB Den Burg, The Netherlands

To be submitted to *Environmental Microbiology*

<sup>\*</sup>To whom correspondence should be addressed. Royal Netherlands Institute for Sea  
Research P.O. Box 59, NL-1790 AB Den Burg, The Netherlands. E-mail:  
laura.villanueva@nioz.nl, phone number: +31 (0)222-369-428, fax number: +31 (0)222-  
319-674.

**Running title:** Biomarkers of sterol biosynthesis in microalgae

**Keywords:** Cycloartenol synthase, sterols, diatom, phytoplankton, phylogenomics, lipid  
biomarker, microalgae

## Summary

Sterols are isoprenoid lipids present in all eukaryotes. These compounds have been used to determine the composition of algal communities in marine and lake environments, and because of their preservation potential, have been used to reconstruct past eukaryotic presence and diversity in the geological record. In the last years there have been major advances in understanding the sterol biosynthetic pathways and the enzymes involved. Here, we have explored the diversity and phylogenetic distribution of the gene coding the cycloartenol synthase protein, a key enzyme of the phytosterol biosynthetic pathway. The cycloartenol synthase gene (CSG) was annotated in genomes of diatoms and other microalgae using protein homology with previously annotated CSG sequences. Based on this, primers for the detection of CSG sequences were designed and evaluated in cultures and environmental samples. A comparison of the phylogeny of the recovered CSG sequences in combination with sequence data of Rubisco gene sequences demonstrates the potential of CSG sequences as phylogenetic marker, as well as an indicator for the identity of sterol-producing organisms in the environment. The proposed gene-based approach can be used to assess the sterol-forming potential of algal groups independent of physiological conditions.

## Introduction

Biomarker lipids have been extensively used for determining the composition and function of microbial communities in past and modern environments (e.g. Hinrichs et al. 1999; Sinninghe Damsté et al., 2002; Kuypers et al., 2003; Talbot et al., 2003; Brocks et al., 2005, Volkman et al., 1998; Volkman 2003, etc). Lipids make excellent molecular fossils because of their relative resistance to degradation, and because some have structures unique to certain taxonomic groups. The combination of DNA-based diversity studies (mainly based on ribosomal rRNA gene taxonomy) and chemotaxonomic characterization of lipids has been shown to be a powerful approach to constrain the diversity of microbial communities (e.g. Stefen et al., 1999; Sinninghe Damsté et al., 2004; Villanueva et al., 2004; Rampen et al., 2010). Some studies have also compared biomarker lipids with functional/metabolic genes to assess both the diversity of certain microbial groups as well as their potential ability to perform an activity (e.g. Ertefai et al., 2008, Pitcher et al., 2011).

Sterols are important lipid biomarkers and are present in all eukaryotic organisms. These lipids have been considered as important tools for molecular paleontologists because sterols can be preserved as e.g. steranes in the fossil record for billions of years (Summons et al., 1999; Peters et al., 2005; Brocks and Pearson, 2005, among others). Steranes are thus molecular fossils for the presence sterol-producing organisms and their distribution can give taxonomic information (Moldowan et al., 1990; Brocks et al., 1999; Peters et al., 2005; Kodner et al., 2008). Furthermore, the cyclization of squalene to sterols and some of the following steps in the sterol biosynthetic pathway require

molecular oxygen, and the presence of steroids in the fossil record are thus indicators of oxygenation of the atmosphere and oceans (Summons et al., 1999; 2006).

The diversity of sterols and their synthetic pathways has been studied extensively and revealed a wide variety of structures (A-ring and side chain alkylation, cyclopropane rings, unsaturations, etc) some of which can be specific for certain eukaryotic groups (Volkman et al., 1998; Volkman, 2003; Rampen et al., 2010). However, most sterols are usually not exclusive of a specific group or genera and, in addition, a change in sterol distributions may also result from changes in environmental and growing conditions rather than community composition changes (Fabregas et al., 1997; Shifrin & Chisholm, 1980; Rampen et al., 2009b). Furthermore, the taxonomic distribution of sterols in microalgae is fully based on culture analysis and thus may not reflect environmental diversity since large numbers of environmental gene sequences are different from those of cultivated species. Finally, the algal taxonomy has been mostly based on morphology, or more recently on genetic data based on the 18S rRNA gene and plastid-encoded large subunit of Rubisco (ribulose 1,5-bisphosphate carboxylase, *rbcL*) enzyme, that do not necessarily reflect the structural diversity of sterols (Moniz & Kaczmarek 2009; Rampen et al., 2010).

One approach to solve the above issues is to examine the presence and diversity of genes involved in the biosynthesis of biomarker lipids as an evidence of the potential ability to biosynthesize the compound of interest, as well as a phylogenetic marker. For example, Pearson and collaborators (2007, 2009) investigated the phylogeny of the producers of hopanoids, isoprenoid bacterial lipids, by analyzing the sequence diversity and distribution of the squalene-hopane cyclase (*sqhC*) gene, concluding that the ability

of hopanoid production is not as widespread among bacteria as previously thought. Following the same approach, a recent study by Welander and collaborators (2010) investigated the genes involved in the synthesis of 2-methylhopanoid and showed that the gene required for the C-2 methylation in hopanoids was found in bacterial taxa other than cyanobacteria, invalidating the use of 2-methylhopanes as biomarkers of the appearance of oxygenic photosynthesis on Earth (Welander et al., 2010).

In this study, we have made use of recent advances on the phylogenomics of the sterol biosynthetic pathway (Desmond & Gribaldo, 2009) and the growing availability of complete or draft genomes of microalgae, allowing the identification of key genes of the phytosterol biosynthetic pathway. Among all the enzymes of the sterol pathway, oxidosqualene cyclases (OSCs) are one of the most conserved at the sequence level and homologues have been detected in all species capable of sterol synthesis (Desmond & Gribaldo, 2009). There are two main types of OSCs based on the end product of the cyclization: lanosterol synthases (found in animals, fungi, choanozoa, trypanosomatids and dinoflagellates), and cycloartenol synthases (found in higher plants, red and green algae, amoebzoa, diatoms, euglenids and heterolobosea). Previous studies have also identified conserved active sites and specific amino acid residues responsible for particular steps in the cyclization cascade (see Summons et al, 2006 for a review).

We targeted the gene encoding the cycloartenol synthase enzyme (CSG) because it is the first specific step in the phytosterol biosynthetic pathway (Fig. 1) and because it is possible to detect homologues of this gene in different organisms due to its conservation at the sequence level (Summons et al., 2006). We have focused on the characterization of the CSG of diatoms as these unicellular algae are thought to be the most common group

of eukaryotic phytoplankton in modern oceans and responsible for approximately 40% of marine primary productivity (Falkowski et al., 1998; Moniz&Kaczmarska 2009). Thus, they are likely one of the most important steroid- producing organisms in marine environments.

We searched for conserved areas of the cycloartenol synthase amino acid sequences in diatom representatives and designed specific primers to recover a fragment sequence with phylogenetic value (i.e. a variable sequence comprised between conserved motifs). We also explored the diversity and distribution of the key CSG of the phytosterol biosynthetic pathway, and evaluated the potential to link the analysis of lipid biosynthetic genes to patterns of distribution and abundance of their specific biomarker (i.e. sterol) in natural environments.

## **Results and Discussion**

### *Annotation and evolutionary analysis of CSG sequences in diatoms*

Only three genomes of the phylum Bacillariophyta (diatoms) are currently available in public databases, either annotated or in draft, i.e. *Thalassiosira pseudonana*, *Phaeodactylum tricornutum*, and *Fragilariopsis cylindricus*. *T. pseudonana* protein GI:223995517 has been previously annotated as cycloartenol synthase;-2,3-epoxysqualene mutase-like protein (Armbrust et al., 2004). A protein blast against non-redundant protein sequences revealed a 71% identity with the *P.tricornutum* protein GI:219120893 formerly annotated as acetyl-coenzyme A synthetase (Bowler et al., 2008). A protein alignment by ClustalW (multiple alignment Gap penalty 10, gap extension 0.2, and Gonnet protein weight matrix) of the two sequences was submitted to

the JGI Genome portal (<http://genome.jgi-psf.org>) protein blast tool using as a reference database the draft genome of the diatom *F. cylindricus*. A protein ortholog to the putative cycloartenol synthase protein (CSP) detected in *T. pseudonana* and *P. tricornutum* was assigned in the scaffold 9 of the draft genome of *F. cylindricus* based on a 63% and 61% protein identity, respectively. In this way, putative CSG sequences in all three diatom whole-genome sequences available in public databases were assigned.

The evolutionary divergence between the *T. pseudonana*, *F. cylindricus* and *P. tricornutum* open reading frames (ORFs) of the CSG was computed by using the Jukes-Cantor model (codon positions included 1st+2nd+3rd+Noncoding and all ambiguous positions were removed for each sequence pair) showing between 0.45 to 0.49 base substitutions per site and a sequence identity between 25–31%, which suggests important sequence divergence (diversity) even between CSG sequences of diatoms.

#### *CSG sequences in other diatoms*

We developed primers for detection of CSG sequences in other diatom genera. These primers were designed to match conserved amino acid sites of the CSP (see Table 1 and supplementary Fig. 1–2). The areas of conserved amino acid positions in the three ORF of the CSG were investigated in an alignment performed by ClustalW (Thompson et al., 1994). For the reverse primer Cycloart\_R and CycloR\_TPF (Table 1), an area comprising the amino acid motif GYNGSQC was chosen because it is conserved across different phyla following the cycloartenol branch of the sterol biosynthesis, and also because it includes the tyrosine (Y) amino acid position 381 (Y381 motif) that is one of the conserved amino acid residues responsible for particular steps in the cyclization cascade (Summons et al., 2006).



The designed primer pairs (Table 1) were tested on DNA extracted from 14 cultures of diatoms, including the three diatoms, *T. pseudonana*, *P. tricornutum*, and *F. cylindricus*, for which whole genome data were available and based on which the primers were designed (supplementary table 1). The primer pairs Cycloart\_F/R and CycloF\_TPF F&R gave the expected size PCR product. Half of the diatom cultures tested showed positive amplification, including *Skeletonema costatum* CCMP 1281, *Skeletonema subsalsum* CCAP 1077/8, *Pseudo-nitzschia seriata* CCMP1309, *Extubocellulus spinifer* CCMP 393, and, as expected, the three diatom strains that were used to design the primers, *T. pseudonana* CCMP 1335, *P. tricornutum*, and *F. cylindricus* CCMP 1102. Encouragingly, different genera showed positive results, suggesting that the primers may be generally applicable within the diatom group. However, in some other cases the designed primers failed to amplify members of the same genus (*Thalassiosira* and *Extubocellulus*). This suggests that the sequence of the CSG is more diverse than inferred based on the three diatom whole genome sequences available and our primers may not be generally applicable to all diatoms.

To check the identity of the partial CSG sequences obtained from the diatom cultures, other than those of *T. pseudonana*, *P. tricornutum*, and *F. cylindricus*, we performed a translated nucleotide query against protein database blast (xblast). The first subject homolog entry was *T. pseudonana* protein GI:223995517 for all the partial CSG sequences obtained from the diatom cultures, confirming a high homology with annotated CSG sequences and thus their identity as putative CSG sequences as well. The translated CSG sequences were aligned with *T. pseudonana*, *F. cylindricus* and *P. tricornutum* putative CSP sequences in order to investigate the phylogenetic diversity

between protein sequences (supplementary figure 2). The alignment revealed several amino acid (aa) changes between diatom genera, also between closely related species, such as in the case of *S. costatum* and *S. subsalsum* with a 96% identity and 3 aa changes in the 193 aa-fragment sequences analyzed (aa residues 91, 105 and 144 of the alignment). It is also important to highlight that the CSP sequences of *Pseudo-nitzschia seriata* and *F. cylindricus* were the sequences with a lower percentage of identity in comparison to the rest of the diatom sequences but relatively close compared to each other (82% identity; Suppl Fig 2). The same aa change was observed in these two sequences with respect to the others in several positions of the fragment, but only in two of the positions the amino acid was replaced by another of the same nature (see supplementary Figure 2): i.e. position aa 113 in the alignment, threonine/serine (T/S) in *P. seriata* and *F. cylindricus* respectively vs. cysteine (C) in the rest (all of them nucleophilic amino acids); position 118, lysine (K) in *P. seriata* and *F. cylindricus* respectively vs. arginine (R) in the rest (all of them basic amino acids). Generally, the conservation of specific motifs of the CSP has allowed its annotation in different phyla (see Summons et al., 2006 for a review). However, a close look at the amino acid sequence of the three putative CSP sequences annotated in *T. pseudonana*, *P. tricornutum*, and *F. cylindricus* demonstrates some flexibility in the protein sequence that would suggest that CSP has been influenced by evolutionary pressures even within the phylum Bacillariophyta.

#### *Phylogeny and evolution of CSG sequences in other microalgal groups*

In order to construct the phylogenetic diversity and evolution of CSG with respect to other phytoplanktonic groups, we searched for related microalgal sequences in genomic

databases. To this end, the putative CSP sequences of *T. pseudonana*, *F. cylindricus* and *P. tricornutum* were considered query sequences in a protein blast in NCBI (BlastP; Altschul et al., 1990) against non-redundant protein sequences and complete/draft genomes available of microalgae. Draft genomes of microalgae available in JGI DOE, such as *Emiliania huxleyi*, were also screened for protein orthologs by pBLAST.

The obtained annotated CSP sequences in microalgae (Table 2) were used to construct a maximum likelihood tree (Figure 2A) based on the entire nucleotide translated sequence. In addition, the partial CSG sequences comprised by the primers used in this study were translated and the obtained protein sequences were used to build a CSP tree (Figure 2B). This tree included the partial CSP sequences deduced by xblast obtained from the diatom cultures as well as the annotated sequences are listed in Table 2. The phylogeny based on CSP in microalgae was compared with a phylogeny by maximum likelihood method based on the *rbcL* protein sequence (Figure 2C), as well as the 18S rRNA gene sequence (supplementary Fig. 3) of the same algae, as previous studies have suggested that the higher variability of the *rbcL* of the chloroplast with respect to the 18S rRNA might make it more suited for phylogenetic studies (Evans et al., 2007).

The topology of the CSP, *rbcL* protein and the 18S rRNA gene sequence trees (Figure 2, Suppl Fig 3) was similar compared to each other, with distinctive clustering between Heterokontophyta (Bacillariophyceae, Pelagophyceae and Phaeophyceae) and Chlorophyta (Chlorophyceae, Trebouxiophyceae and Mamiellophyceae). However, both 18S rRNA gene and *rbcL* protein trees showed a more clear separation of the Bacillariophyceae group (diatoms) from the other microalgae, while the CSG tree did not indicate a clear

divergence of Bacillariophyceae and the other two sequences of Heterokontophyta (*Aureococcus anophagefferens* and *Ectocarpus siliculosus*). Other groups such as Prymnesiophyceae (*Emiliana huxley*) are clearly separated in the tree topology of the CSP and 18S rRNA gene trees from the other microalgae, while in the *rbcL* protein tree the *E.huxleyi* protein sequence is clustered with *A. anophagefferens* and *E. siliculosus* (Figure 2C). In both the CSP and the 18S rRNA gene trees the divergence of the Chlorophyceae and Trebouxiophyceae groups (both Chlorophyta) looks similar in comparison with the *rbcL* protein tree. In conclusion, the CSP-based phylogeny follows the same distribution of groups of the *rbcL* protein and 18S rRNA gene-based reconstruction, being able to cluster major groups (e.g. Haptophyta, Chlorophyta and Heterokontophyta), as well as families (e.g. Bacillariophyceae), and genera. This demonstrates the potential of the CSG sequences to be used as a phylogenetic marker. However, the different clustering generated by CSP, *rbcL* protein and 18S rRNA gene observed in some cases, such as in *E. huxleyi*, *E. siliculosus* and *A. anophagefferens* requires further attention as it might help to clarify the evolutionary history of these three groups.

When we focus on the diatom sequences clustering, a comparison between the protein tree based on the entire CSP and the studied fragment displays the same topology (Figure 2A and 2B), which supports the phylogenetic value of the CSG fragment amplified by the primers introduced in this study. On the other hand, the phylogenetic clustering of the diatom group by the 18S rRNA gene (Suppl. Fig. 3) and *rbcL* protein sequences (Figure 2C) separated the two main groups of centric and pennate diatoms: order Thalassiosirales (*Thalassiosira/Skeletonema*) plus order Cymatosyrales (*Extubocellulus*), and order

Bacillariales (*P. seriata*/*F. cylindricus*) plus order Naviculales (*Phaeodactylum*), respectively. These differences in the phylogenetic distribution of diatoms based on these three genetic markers (i.e. CSP, *rbcL* protein, and 18S rRNA gene) will be more clear once other CSG sequences become available.

The alignment of the putative CSP sequences of *T. pseudonana*, *P. tricornutum*, and *F. cylindricus* revealed amino acid changes, suggesting that the CSG is more diverse than previously thought even between members of the same phyla (see below). Thus, in order to test the occurrence of evolutionary events of positive and purifying selection in the codifying region of the CSP sequence we applied the Nei-Gojobori method (based on computing the numbers of synonymous and non-synonymous substitutions and the numbers of potentially synonymous and non-synonymous sites; Nei and Gojobori 1986) to the annotated sequences in Bacillariophyceae (diatoms; Table 3) and Chlorophyta (green algae; Table 4). The identity matrix for the Bacillariophyceae class (Table 3A) revealed the highest percentage of identity between CSP sequences of the same genus (*S. costatum* and *S. subsalsum*, 92%). CSP sequences of diatoms belonging to the same order such as Thalassiosirales (*Skeletonema* & *T. pseudonana*) and Bacillariales (*Pseudo-nitzschia seriata* & *F. cylindricus*) had a percentage of identity of 81% and 88%, respectively. The lowest percentage of identity was observed between *P. tricornutum* (order Naviculales) and *P. seriata*/*F. cylindricus* (Bacillariales; 60%). In the Z-test of selection for Bacillariophyceae (Table 3B), the test statistic ( $ds-dN$ ; number of synonymous substitutions per synonymous site minus number of non-synonymous substitutions per non-synonymous site) or probability of rejecting the null hypothesis of neutrality ( $ds=dN$ ) in favor of the alternative hypothesis ( $ds>dN$ , purifying selection) is

shown above the diagonal. In this case, the test of selection clearly stated the role of purifying selection with a significant  $p$ -value in all cases. As expected, the lower value of the test statistic is found between species of the same genus (*Skeletonema*) as the members of the same genus have very similar CSP sequences and thus, any amino acid change in the protein would result in purifying selection (favoring synonymous substitutions that code for the same amino acid).

For Chlorophyta (green algae), CSP sequences of members of the same class, e.g. *Chlamydomonas/Volvox* (class Chlorophyceae) and *Ostreococcus/Micromonas* (class Mamiellophyceae) had a protein percent identity 61–71% (Table 4A). For the class Bacillariophyceae the following CSP sequence identity percentages can be found in the different taxonomic divisions: Phylum (40–50%) < Class (60–70%) < Order (80–90%) < Genus (more than 95%). Thus, although there is a high degree of conservation of CSP sequences (see Summons et al, 2006 for a review), the CSG sequence is diverse enough to distinguish species of the same genus supporting its value as a phylogenetic marker.

#### *Detection of CSG sequences in environmental samples*

We tested the CSG primers (Table 1) on several environmental samples for the detection of diatom-related CSG sequences. These samples were a microbial mat from the island Schiermonnikoog that is characterized by the presence of pennate diatoms (e.g. *Navicula*, *Diploneis*, *Amphora* and *Cylindrotheca*) as shown using microscopic methods (Dijkman et al., 2010). Furthermore, we also analyzed suspended particulate matter (SPM) from North Sea surface water which has high contents of diatom pigments (diatoxanthin, diadinoxanthin, data not shown) and in which *Thalassiosira*, *Chaetoceros*

and *Skeletonema* species were previously observed (Cadee and Hegeman, 2002; Brandsma et al., 2012).

Primer pairs cycloart\_F/R and cycloF\_TPF F&R designed based on the diatom sequences available gave positive results on the environmental samples tested in this study. Thus, CSG sequences were amplified and sequenced from the microbial mat and North Sea SPM and analyzed phylogenetically by the maximum likelihood method. The sequences recovered from the North Sea water column mainly clustered with the order Thalassiosirales (*Skeletonema* and *Thalassiosira pseudonana*), while the majority of the diatom mat CSG sequences were closer to the representative sequence of *Extubocellulus spinifer* (order Cymatosyrales; centric diatom) with some also closer to the Thalassiosirales order cluster (Figure 3). As a comparison we also sequenced the *rbcL* gene: its distribution was more diverse than the CSG sequence-based tree (Figure 4). The majority of the *rbcL* gene sequences retrieved from the North Sea SPM clustered with the order Thalassiosirales but also to sequences such as *Phaeodactylum tricornutum* (order Naviculales). *rbcL* gene sequences retrieved from the microbial mat were also more diverse and clustered with *Extubocellulus spinifer* as well as with the orders Thalassiosirales, Bacillariales (*Pseudo-nitzschia* & *Fragilariopsis*) and Naviculales.

In addition, we also characterized the sterols in these samples to obtain an idea of the sterol diversity and potentially already identify specific diatom sterols. The main sterols in the Schiermonnikoog mat were cholest-5-en-3 $\beta$ -ol (cholesterol), cholesta-5,24-dien-3 $\beta$ -ol (desmosterol), 24-methylcholest-5,24(28)-dien-3 $\beta$ -ol, 24-methylcholest-5-en-3 $\beta$ -ol (campesterol), and 24-ethylcholest-5-en-3 $\beta$ -ol (sitosterol) (see Suppl Table 2 for details). These sterols have been reported for diatom cultures (Rampen et al. 2010), specifically,

24-methylcholesta-5,24(28)-dien-3 $\beta$ -ol is found in high abundances in some centric diatoms such as *Thalassiosira* and *Skeletonema* genera (Rampen et al., 2010). In the North Sea water SPM the main sterols were cholesterol, desmosterol and 24-methylcholest-5,24(28)-dien-3 $\beta$ -ol. The diversity of phytosterols detected in the North Sea SPM was thus lower than in the diatom mat and only one sterol, 5 $\alpha$ (H)-24-methylcholest-22en-3 $\beta$ -ol, was uniquely found in the North Sea water but not in the diatom mat. The abundant presence of 24-methylcholesta-5,24(28)-dien-3 $\beta$ -ol in both environments may indicate the importance of diatoms in these two systems. In addition, other sterols such as brassicasterol, also known as diatomsterol, 24-methylcholest-5,22E-dien-3 $\beta$ -ol, and 24-ethylcholesta-5,22E-dien-3 $\beta$ -ol (stigmasterol) have also been associated with diatoms (Rontani & Volkman, 2005). However, the presence of these sterols does not provide conclusive for their origin of diatoms as all these sterols are also found in other algae while the sterol composition also does not indicate which diatom genera are present (Volkman et al., 1998; Volkman, 2003; Rampen et al., 2010 and references cited therein).

A comparison of the CSG tree with those of the sterol distributions support the idea that diatoms from the order Thalassiosirales in both samples may be the source of the 24-methylcholest-5,24(28)-dien-3 $\beta$ -ol. In addition, the higher diversity of diatom CSG sequences detected in the microbial mat compared to the North Sea SPM also corresponds to a higher diversity of detected sterols. The detection of sequences homologous to the Bacillariales and Naviculales orders in the microbial mat suggests that cholesta-5,22-dien-3 $\beta$ -ol and 24-ethylcholesta-5,22-dien-3 $\beta$ -ol may be sourced by these diatoms as they are also dominant sterols in cultivated relatives (Rampen et al., 2010).



The fact that CSP appear to be less conserved than other phylogenetic markers has the disadvantage that this protein might be more difficult to annotate based on protein homology. On the other hand, it might provide clues about the evolutionary placement of certain organisms better than more conserved proteins. In general, CSG and *rbcL* gene analysis in environmental samples has proven effective for surveying the sterol-forming diatom community. However, the fact that *rbcL* gene analysis has elucidated more diversity indicates that CSG primers are still limited in their diversity coverage.

### *Implications*

The approach presented here based on specific gene searching in whole-genome and metagenomic databases has allowed the design of effective primers for the detection of a key gene in the sterol biosynthetic pathway (i.e. CSG) in microalgae. The comparison between the phylogenetic reconstruction of microalgae 18S rRNA gene, *rbcL* protein and CSP sequences supports the value of CSG sequences as marker of the presence and phylogeny of sterol-producing microalgae. However, further studies are needed to improve the diversity coverage of the CSG marker by sequencing more whole genomes of diatoms and other microalgae. Through this it will be possible to redesign the developed primers and assign a more refined taxonomic identification, which will lead to a more accurate the association between CSG sequences, their producer, and the sterol composition in environmental samples .

The genomic characterization of enzymes involved in lipid biosynthetic pathways opens a new chapter in organic geochemistry studies. Genomic approaches provide an independent assessment of the organism ability to produce a molecule of interest without extensive screening of cultures. Our study has expanded the range of CSG sequences

available, introduced a quick screening of environmental samples for the diversity of sterol-forming microorganisms and may, once the sequence coverage has increased, provide a link between sterols and their main sources. This is also the starting point of other studies involving determination of the abundance and expression of this key gene of the phytosterol biosynthetic pathway in environmental samples. Ultimately, the CSG has potential to elucidate the origin of certain microalgae groups and as a molecular clock to track the appearance of the sterol biosynthetic pathway.

## Material and Methods

### *Sampling*

A diatom-dominated microbial mat was sampled in the sandy beach of the Dutch barrier island Schiermonnikoog (53°29'N and 6°08'E; for a more detailed description see Stal et al., 1985) in January 2010, transported to the lab at 4°C and then stored at –80°C until further analysis. North Sea surface water was sampled at a Jetty platform at the NIOZ at the western entrance of the North Sea into the Wadden Sea at the Island Texel (53°0'2''N, 4°7'2''E). With each incoming-tide, water from the coastal North Sea moves as far as 25 km into the Wadden Sea (Potsma, 1954). At high tide, water collected represents Dutch coastal North Sea waters since the estuarine influence is minimal. Strong tidal currents assure that the water is vertically mixed. Therefore, surface water samples taken during high tide are representative of the entire water column. Suspended particulate matter (SPM) sample was taken on 24<sup>th</sup> March 2010. For DNA analysis, measured volumes (ca. 1 L) of water were filtered through a 142 mm diameter, 0.2 µm pore size polycarbonate filter (Millipore, Billerica, MA) and stored at –80°C until extraction. For lipid analyses, a measured volume (ca. 20 L) of water was filtered sequentially through pre-ashed 3 µm, and 0.7 µm-pore-size, glass fiber filters (GF/F, Pall, 142 mm filter diameter). GF/F filters were stored at –40°C until extraction. Diatom cultures were obtained from culture collection, grown in batch cultures and harvested at the end of the logarithmic growth phase by filtration on pre-ashed 47 mm, 0.7 µm GF/F filters (see Rampen et al., 2010).

417 *DNA extraction*

418 Approximately 0.2 g of wet weight of diatomaceous-mat was homogenized by a DNA-  
419 clean spatula and extracted by the Power Biofilm DNA extraction kit from MoBio  
420 (Carlsbad, CA) according to the manufacturer's instructions. North Sea water samples on  
421 0.2 µm PCC filters were extracted by a bead-beating protocol followed by DNeasy  
422 extraction kit of Qiagen (Valencia, CA). Diatom cultures were extracted as described by  
423 Rampen et al. (2009a). DNA quality and concentration were evaluated by gel agarose  
424 electrophoresis and Nanodrop (Wilmington, DE).

425 *PCR amplification, cloning and sequencing*

426 Partial CSG sequences were amplified by using the primers listed in Table 1 on diatom  
427 pure cultures DNA extracts and environmental samples. PCR reaction mixture was the  
428 following (final concentration): Q-solution 1×; PCR buffer 1×; BSA (200 µg/ml); dNTPs  
429 (20 µM); primers (0.2 pmol/µl); MgCl<sub>2</sub> (1.5 mM); 1.25 U Taq polymerase (Qiagen,  
430 Valencia, CA, USA) or BioThermD Taq DNA polymerase (Semiramis Genetics Ltd.,  
431 Manchester, UK). PCR conditions for these amplifications were the following: 95°C, 5  
432 min; 40× [95°C, 1 min; T<sub>m</sub>, 1 min; 72°C, 1 min]; final extension 72°C, 5 min. A gradient  
433 PCR cycle was performed for each set of primers and samples from 48 to 55°C melting  
434 temperature. Amplification of Rubisco gene was performed as described by Rampen et  
435 al., 2009a. Positive amplification bands were excised from agarose gel and gel or PCR  
436 purified (QIAquick gel/PCR purification kit, Qiagen) and cloned in the TOPO-TA  
437 cloning® kit from Invitrogen (Carlsbad, CA, USA) and transformed in *E. coli* TOP10  
438 cells following the manufacturer's recommendations. Recombinant clones plasmid DNAs  
439 were purified by Qiagen Miniprep kit and screening by sequencing using M13F (-20) (5'-

GTA AAA CGA CGG CCA G-3') and M13R (5'-CAG GAA ACA GCT ATG AC-3') primers with BigDye® v1.1 sequencing kit in house on a ABI PRISM® 310 Genetic analyzer (Applied Biosystems, Foster city, CA, USA) or sequenced in Macrogen Europe Inc.

#### *Alignments, tree reconstruction and evolutionary analyses*

Putative CSG partial sequences obtained from diatom pure cultures and environmental samples were translated to protein by submitting them as query sequences in translated blast (xblast: Find similar proteins to translated query in a protein database) and reviewed by manual annotation. DNA/Protein alignments were performed by ClustalW (multiple alignment Gap penalty 15, gap extension 6.66, and IUB DNA weight matrix) (Thompson et al., 1994). Mega5 software (Tamura et al., 2011) was used to estimate the best DNA/protein models for maximum likelihood analysis (automatic neighbor joining tree; statistical method, maximum likelihood; substitution type, nucleotide/amino acid; use all sites). In case of protein alignments, the choice of protein model by Mega5 (with models ranked by Bayesian information criterion, BIC) was contrasted with the choice of model of evolution for protein phylogeny given by ProtTest 2.4 (Abascal et al., 2005) with model selection criterion of Akaike Information Criterion (AIC) and (Bayesian Information Criterion) BIC. Maximum likelihood phylogenetic reconstruction of the 18S rRNA gene, partial CSG sequences and partial *rbcL* gene sequences was performed by Mega5 using the model with higher ranking (see figure legends for details). For the 18S rDNA, partial CSG sequences and partial *rbcL* gene sequences, the genera; time reversible model plus gamma distribution (GTR + G) was used (plus invariant sites, + I for the case of the *rbcL* gene sequences tree). *rbcL* protein and CSP sequence trees were

generated with the WAG+G+F model. Bootstrap analysis was performed in all cases with 1000 replicates. All sites were considered for the calculations and the maximum likelihood heuristic method chosen was nearest neighbor interchange (NNI). Evolutionary analyses for annotated partial CSP sequences were performed with MEGA5. Codon-based Z-test Test of Purifying Selection was also performed. Analyses were conducted using the Nei-Gojobori method and bootstrap methods with 1000 replicates.

#### *Data submission*

Partial CSG sequences were deposited in GenBank under the accession numbers (accession number pending to be assigned). Partial *rbcL* gene sequences were deposited under accession numbers (accession number pending to be assigned).

#### *Lipid extraction, separation and GC/MS detection*

Freeze-dried microbial mats were homogenized by lipid-free mortar and pestle and the GF/F filters were freeze-dried and cut into small pieces with sterile scissors before being ultrasonically extracted four times using dichloromethane (DCM)/methanol (MeOH) (1:1, v/v) as described in Rampen et al 2009a. An aliquot of the total lipid extracts was separated over a pipette-column filled with Al<sub>2</sub>O<sub>3</sub>, using hexane:dichloromethane (DCM; 9:1, v:v) and DCM: methanol (MeOH; 1:1, v:v) to elute the apolar and sterol fractions, respectively. Prior to analysis by gas chromatography (GC) and gas chromatography–mass spectrometry (GC–MS), the sterol fractions were silylated by adding 25 µl BSTFA [N,O-bis(trimethylsilyl)trifluoroacetamine] and 25 µl pyridine and heating the mixtures at 60°C for 20 min. GC and GC–MS analyses were performed as described by Rampen et

al. (2009a). Sterols (as their TMS derivatives) were identified based on their mass spectra and relative retention times in comparison with literature data.

## References

- Abascal, F., Zardoya R., and Posada D. (2005) ProtTest: Selection of best-fit models of protein evolution. *Bioinformatics* 21: 2104-2105.
- Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990) Basic local alignment search tool. *J Mol Biol* 215: 403-410.
- Armbrust, E.V., Berges, J.A., Bowler, C., Green, B.R., Martinez, D., Putnam, N.H., Zhou, S., Allen, A.E., Apt, K.E., Bechner, M., and Brzezinski, M.A. et al. (2004) The genome of the diatom *Thalassiosira pseudonana*: ecology, evolution, and metabolism. *Science* 306: 79-86.
- Bowler, C. et al. (2008) The Phaeodactylum genome reveals the evolutionary history of diatom genomes. *Nature* 456:239-244.
- Brandsma, J., Hopmans, E.C., Phillipart, K.J.M., Vedhuis, M.J.W., Schouten, S., and Sinninghe Damste, J.S. (2012) Low temporal variations in the intact polar lipid composition of North Sea coastal marine water reveals limited chemotaxonomic value. *Biogeosciences* 9: 1073-1084.
- Brocks, J.J., Logan, G.A., Buick, R., and Summons, R.E. (1999) Archean molecular fossils and early rise of eukaryotes. *Science* 285: 1033-1036.
- Brocks, J.J., and Pearson, A.P. (2005) Building the biomarker tree of life. *Rev Mineral Geochem* 59: 233-258.
- Brocks, J.J., Love, G.D., Summons, R.E., Knoll, A.H., Logan, G.A., and Bowden, S.A. (2005) Biomarker evidence for green and purple sulphur bacteria in a stratified Palaeoproterozoic sea. *Nature* 437: 866-870.



528 Cadée, G.C., and Hegeman, J. (2002) Phytoplankton in the Marsdiep at the end of the  
 529 20th century; 30 years monitoring biomass, primary production, and Phaeocystis  
 530 blooms. *J Sea Res* 48: 97-110.

531 Desmond, E., and Gribaldo, S. (2009) Phylogenomics of sterol biosynthesis: Insights into  
 532 the origin, evolution and diversity of a key eukaryotic feature. *Genome Biol Evol* 1:  
 533 364-381.

534 Dijkman, N.A., Boschker, H.T.S., Stal, L.J., and Kromkamp, J.C. (2010) Composition  
 535 and heterogeneity of the microbial community in a coastal microbial mat as revealed  
 536 by the analysis of pigments and phospholipid-derived fatty acids. *J Sea Res* 63: 62-70.

537 Ertefai, T.F., Fisher, M.C., Fredricks, H.F., Lipp, J.S., Pearson, A., Birgel, D., Udert,  
 538 K.M., Cavanaugh, C.M., Gschwend, P.M., Hinrichs, K.U. (2008) Vertical distribution  
 539 of microbial lipids and functional genes in chemically distinct layers of a highly  
 540 polluted meromictic lake. *Org Geochem* 39: 1572-1588.

541 Evans, K.M., Wortley, A.H., and Mann, D.G. (2007) An assessment of potential diatom  
 542 “barcode” genes (cox1, rbcL, 18S, and ITS rDNA) and their reffectiveness in  
 543 determining relationships in Sellaphora (Bacillariophyta). *Protist* 158: 349-364.

544 Fabregas, J., Aran, J., Morales, E.D., Lamela, T., and Otero, A. (1997) Modification of  
 545 sterol concentration in marine microalgae. *Phytochemistry* 46: 1189-1191.

546 Falkowski, P.G., Barber, R.T., and Smetacek, V. (1998) Biogeochemical controls and  
 547 feedbacks on ocean primary production. *Science* 281: 200-206.

548 Fischer, W.W., and Pearson, A. (2007) Hypotheses for the origin and early evolution of  
 549 triterpenoid cyclases. *Geobiol* 5: 19-34.

550 Hinrichs, K.-U., Hayes, J.M., Sylva, S.P., Brewer, P.G., and DeLong, E.F. (1999)  
 551 Methane-consuming archaeobacteria in marine sediments. *Nature* 398: 802-805.  
 552 Kodner, R.B., Pearson, A., Summons, R.E., and Knoll, A.H. (2008) Sterols in red and  
 553 green algae: quantification, phylogeny, and relevance for the interpretation of geologic  
 554 steranes. *Geobiol* 6: 411-420.  
 555 Kuypers, M.M.M., Sliekers, A.E., Lavik, G., Schmid, M., Jorgensen, B.B., Kuenen, J.G.,  
 556 Sinninghe Damsté, J.S., Strous, M., and Jetten, M.S.M. (2003) Anaerobic ammonium  
 557 oxidation by anammox bacteria in the Black Sea. *Nature* 422: 608-611.  
 558 Kuypers, M.M.M., van Breugel, Y., Schoten, S., Erba, E., and Sinninghe Damsté, J.  
 559 (2004) N<sub>2</sub>-fixing cyanobacteria supplied nutrient N for Cretaceous oceanic anoxic  
 560 events. *Geology* 32: 853-856.  
 561 Moldowan, J.M., Fago, F.J., Lee, C.Y., Jacobson, S.R., Watt, D.A., Slougui,  
 562 N., Jeganathan, A., and Young, D.C. (1990) Sedimentary 24-n propylcholestanes,  
 563 molecular fossils diagnostic of marine algae. *Science* 247: 309-312.  
 564 Moniz, M.B.J., and Kaczmarek, I. (2009) Barcoding diatoms: Is there a good marker?  
 565 *Mol Ecol Resources* 9: 65-74.  
 566 Nei, M., and Gojobori, T. (1986) Simple methods for estimating the numbers of  
 567 synonymous and nonsynonymous nucleotide substitutions. *Mol Biol Evol* 3: 418-426.  
 568 Pearson, A., Flood, S.R., Jorgenson, T.L., Fischer, W.W., and Higgins, M.B. (2007)  
 569 Novel hopanoid cyclases from the environment. *Environ Microbiol* 9: 2175-2188.  
 570 Pearson, A., Leavitt, W.D., Saenz, J.P., Summons, R.E., Tam, M.C.-M., and Close, H.G.  
 571 (2009) Diversity of hopanoids and squalene-hopene cyclases across a tropical land-  
 572 sea gradient. *Environ Microbiol* 11: 1208-1223.

573 Peters, K.E., Walters, C.C., and Moldowan, J.M. (2005) The biomarker guide.  
 574 Cambridge University Press, Cambridge, UK.  
 575 Pitcher, A., Villanueva, L., Hopmans, E.C., Schouten, S., Reichart, G.J., and Sinninghe  
 576 Damsté, J.S. (2011) Niche segregation of ammonia-oxidizing archaea and anammox  
 577 bacteria in the Arabian Sea oxygen minimum zone. *ISME J* 5: 1896-1904.  
 578 Potsma, H. (1954) Hydrography of the Dutch Wadden Sea. *Archives Néerlandaises de*  
 579 *Zoologie* 10: 405-511.  
 580 Rampen, S.W., Abbas, B.A., Schouten, S., and Sinninghe Damsté, J.S. 2010. A  
 581 comprehensive study of sterols in marine diatoms (Bacillariophyta): Implications for  
 582 their use as tracers for diatom productivity. *Limnol Oceanogr* 55: 91-105.  
 583 Rampen, S.W., Schouten, S., Koning, E., Brummer, G-J.A., and Sinninghe Damsté, J.S.  
 584 (2008) A 90 kyr upwelling record from the northwestern Indian Ocean using a novel  
 585 long-chain diol index. *Earth Planet Sci Letters* 276: 207-213.  
 586 Rampen, S.W., Schouten, S., Panoto, E., Brink, M., Andersen, R.A., Muyzer, G., Abbas,  
 587 B., and Sinninghe Damsté, J.S. (2009a) Phylogenetic position of *Attheya longicornis*  
 588 and *Attheya septentrionalis* (Bacillariophyta). *J Phycol* 45: 444-453.  
 589 Rampen, S.W., Schouten, S., Schefuß, E., and Sinninghe Damsté, J.S. (2009b) Impact of  
 590 temperature on long chain diol and mid-chain hydroxyl methyl alkanoate composition  
 591 in Proboscidea diatoms: Results from culture and field studies. *Org Geochem* 40: 1124-  
 592 1131.  
 593 Rontani, J-F., and Volkman, J.K. (2005) Lipid characterization of coastal cyanobacterial  
 594 mats from the Camargue (France). *Org Geochem* 36: 251-272.

595 Shifrin, N.S., Chisholm, S.W. (1980) Phytoplankton lipids: Environmental influences on  
 596 production and possible commercial applications. In *Algae biomass*. Shelef, G., and  
 597 Soeder, C.J. (eds). Elsevier, Amsterdam, pp 623-645.  
 598 Sinninghe Damsté, J.S., Muyzer, G., Abbas, B., Rampen, S.W., Masse, G., Guyallard,  
 599 W., Belt, S.T., Robert, J-M., Rowland, S.J., Moldowan, J.M., Barbati, S.M., Fago,  
 600 F.J., Denisevich, P., Dahl, J., Trindade, L.A.F., and Schouten, S. (2004) The rise of  
 601 rhisolenid diatoms. *Nature* 304: 584-587.  
 602 Sinninghe Damsté, J.S., Schouten, S., Hopmans, E.C., van Duin, A.C.T., and  
 603 Geenevasen, J.A.J. (2002) Crenarchaeol: the characteristic core glycerol dibiphytanyl  
 604 glycerol tetraether membrane lipid of cosmopolitan pelagic crenarchaeota. *J Lipid Res*  
 605 43: 1641-1651.  
 606 Stal, L.J., van Gernerden, H., and Krumbein, W.E. (1985) Structure and development of  
 607 a benthic marine microbial mat. *FEMS Microbiol Ecol* 31: 111-125.  
 608 Stephen, J.R., Chang, Y-J., Gan, Y.D., Peacock, A., Pfiffner, S.M., Barcelona, M.J.,  
 609 White, D.C., and Macnaughton, S.J. (1999) Microbial characterization of a JP-4 fuel –  
 610 contaminated site using a combined lipid biomarker/polymerase chain reaction-  
 611 denaturing gradient gel electrophoresis (PCR-DGGE)-based approach. *Environ*  
 612 *Microbiol* 1: 231-241.  
 613 Summons, R.E., Jahnke, L.L., Hope, J.M., and Logan, G.A. (1999) 2-methylhopanoids as  
 614 biomarkers for cyanobacterial oxygenic photosynthesis. *Nature* 400: 554-557.  
 615 Sumons, R.E., Bradley, A.S., Jahnke, L.L., and Waldbauer, J.R. (2006) Steroids,  
 616 triterpenoids and molecular oxygen. *Phil Trans R Soc B* 361: 951-968.

617 Talbot, H.M., Watson, D.F., Pearson, E.J., and Farrimond, P.(2003) Diverse biohopanoid  
 618 compositions of non-marine sediments. *Org Geochem* 34: 1353-1371.

619 Tamura, K., Peterson, P., Peterson, N., Stecher, G., Nei, M., and Kumar, S . (2011)  
 620 MEGA5: Molecular Evolutionary Genetics Analysis Using Maximum Likelihood,  
 621 Evolutionary Distance, and Maximum Parsimony Methods. *Mol Biol Evol* 28: 2731-  
 622 2739.

623 Thompson, J.D., Higgins, D.G., and Gibson, T.J. (1994) CLUSTAL W: improving the  
 624 sensitivity of progressive multiple sequence alignment through sequence weighting,  
 625 position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22: 4673-  
 626 4680.

627 Villanueva, L., Navarrete, A., Urmeneta, J., White, D.C., and Guerrero, R. (2004)  
 628 Combined phospholipid biomarker-16S rRNA gene denaturing gradient gel  
 629 electrophoresis analysis of bacterial diversity and physiological status in an intertidal  
 630 microbial mat. *Appl Environ Microbiol* 70: 6920-6926.

631 Volkman, J.K., Barrett, S.M., Blackburn, S.I., Mansour, M.P., Sikes, E.L., Gelin, F.  
 632 (1998) Microalgal biomarkers: a review of recent research developments. *Org*  
 633 *Geochem* 29: 1163-1179.

634 Volkman, J.K. (2003) Sterols in microorganisms. *Appl Microbiol Biotechnol* 60: 495-  
 635 506.

636 Welander, P.V., Coleman, M.L., Sessions, A.L., Summons, R.E., and Newman, D.K.  
 637 (2010) Identification of a methylase required for 2-methylhopanoid production and  
 638 implications for the interpretation of sedimentary hopanes. *Proc Natl Acad Sci USA*  
 639 107: 8537-8542.

## Figure Legends

Figure 1. Sterol biosynthetic pathway. Isoprenoid precursor isopentenyl diphosphate (IPP) is the precursor of squalene. Hopanoids are synthesized from the cyclization of squalene by a squalene-hopene cyclase (SHC) in a process independent of oxygen. For sterol synthesis squalene is transformed by a squalene monooxygenase (SQMO) that requires O<sub>2</sub>. Squalene epoxide is then cyclized either to lanosterol or to cycloartenol by lanosterol synthase or cycloartenol synthase.

Figure 2. Phylogenetic tree of cycloartenol synthase protein (CSP) and *rbcL* (Rubisco) protein sequences.

(A) Sequences annotated in microalgal genomes inferred by using the Maximum Likelihood method based on the Whelan And Goldman + Freq. model. A discrete Gamma distribution was used to model evolutionary rate differences among sites (5 categories (+G, parameter = 0.9821). The scale indicates number of substitutions per site. The analysis involved 13 amino acid sequences. There were a total of 993 positions in the final dataset. Bootstrap values (1000 replicates) are indicated on the nodes.

(B) Phylogenetic tree of CSP fragment comprised by the primers applied in this study annotated in microalgal genomes inferred by using the Maximum Likelihood method based on the General Reverse Transcriptase + Freq. model. A discrete Gamma distribution was used to model evolutionary rate differences among sites (5 categories (+G, parameter = 0.9267). The rate variation model allowed for some sites to be evolutionarily invariable ([+I], 0.0000% sites). The analysis involved 17

amino acid sequences. All ambiguous positions were removed for each sequence pair. There were a total of 228 positions in the final dataset.

(C) Phylogenetic tree of *rbcL* protein sequences in the microalgae under study inferred by using the Maximum Likelihood method based on the General Reverse Transcriptase (GRT) + Freq. model. A discrete Gamma distribution was used to model evolutionary rate differences among sites (5 categories (+G, parameter = 0.8751). The analysis involved 15 amino acid sequences. There were a total of 491 positions in the final dataset.

Figure 3. Phylogenetic tree of CSG sequences obtained from the environmental samples under study and inferred by using the Maximum Likelihood method based on the GTR model. A discrete Gamma distribution was used to model evolutionary rate differences among sites (5 categories (+G, parameter = 1.3092). The analysis involved 45 nucleotide sequences. Codon positions included were 1st+2nd+3rd+Noncoding. There were a total of 721 positions in the final dataset.

Figure 4. Phylogenetic tree of *rbcL* gene sequences obtained from the environmental samples under study and inferred by using the Maximum Likelihood method based on the GTR model. Gamma distribution was used to model evolutionary rate differences among sites (5 categories (+G, parameter = 0.2904). The rate variation model allowed for some sites to be evolutionarily invariable ([+I], 59.5472% sites). The analysis involved 54 nucleotide sequences. Codon positions included were 1st+2nd+3rd+Noncoding. There were a total of 1427 positions in the final dataset.

Table 1. List of primers for the detection of CSG sequences

Name	AA sequence	Primer (5'-3')	Name	AA sequence	Primer
Cycloart_F	WLLPNWF (146–152 aa)*	TGGCTKCTMCCMAACTGGTTT	Cycloart_R	GYNGSQC (310–316 aa)	G[CATTGGCTKCCGTTRTAGCC]
CycloF_TPF	PNW(F/I)PFHP (149–156 aa)	CCMAACTGGWTTTCCTTTYCATC	CycloR_TP F	GYNGSQC (310–316 aa)	[CAYTGGCTKCCGTTRTAKCC]†

\*aa (amino acid) positions of the cycloartenol synthase *T.pseudonana* (protein GI:223995517).

†Primer CycloR\_TPF has a nucleotide less than Cycloart\_R primer but we indicate the same aa motif as a reference.



Table 2. Cycloartenol synthase ortholog proteins deduced from database.

Organism	Ortholog	Phylogeny
<i>Ostreococcus tauri</i> OTH95 Green Algae	XM_003077949 XP_003077997.1*	Viridiplantae; Chlorophyta; Mamiellophyceae; Mamiellales
<i>Ostreococcus lucimarinus</i> CCE9901 Green Algae	XM_001416533 XP_001416570.1	Viridiplantae; Chlorophyta; Mamiellophyceae; Mamiellales
<i>Micromonas pusilla</i> CCMP1545 Green Algae	XM_003055917 XP_003055963.1	Viridiplantae; Chlorophyta; Mamiellophyceae; Mamiellales
<i>Micromonas sp.</i> RCC299 Green Algae	XM_002507604 XP_002507650.1	Viridiplantae; Chlorophyta; Mamiellophyceae; Mamiellales
<i>Chlamydomonas reinhardtii</i> Green Algae	XM_001689822 XP_001689874.1*	Viridiplantae; Chlorophyta; Chlorophyceae; Chlamydomonadales; Chlamydomonadaceae
<i>Chlorella variabilis</i> Green algae	GL433842 EFN56189.1	Viridiplantae; Chlorophyta; Trebouxiophyceae; Chlorellales; Chlorellaceae
<i>Volvox carteri</i> ft. <i>nagariensis</i> Green algae	XM_002955246 XP_002955292.1	Viridiplantae; Chlorophyta; Chlorophyceae; Chlamydomonadales; Volvocaceae; Volvox
<i>Thalassiosira pseudonana</i> Diatom	XM_002287396 XP_002287432.1*	Stramenopiles; Bacillariophyta; Coscinodiscophyceae; Thalassiosirophycidae; Thalassiosirales; Thalassiosiraceae
<i>Phaeodactylum tricornutum</i> Diatom	XM_002185642 XP_002185678.1†	Stramenopiles; Bacillariophyta; Bacillariophyceae; Bacillariophycidae; Naviculales; Phaeodactylaceae
<i>Fragilariopsis cylindricus</i> Diatom	Scaffold 9	Stramenopiles; Bacillariophyta; Bacillariophyceae; Bacillariophycidae; Bacillariales; Bacillariaceae
<i>Aureococcus anophagefferens</i>	GL833121.1 EGB12462.1	Stramenopiles; Pelagophyceae; Aureococcus
<i>Ectocarpus siliculosus</i>	Scaffold setg_148 CBN75619.1	Stramenopiles; PX clade; Phaeophyceae; Ectocarpales; Ectocarpaceae; Ectocarpus.
<i>Emiliana huxleyi</i> 1516	Scaffold 360	Haptophyceae; Isochrysidales; Noelaerhabdaceae

\*annotated as putative or cycloartenol synthase protein. †Annotated as acetyl-coenzyme A synthetase. The rest of the proteins were annotated as hypothetical/predicted protein and assigned as putative cycloartenol synthases in this study based on percentage of identity with annotated orthologs in other taxa. Peptide sequences matching query sequences with E-value  $<1e^{-150}$  and aligning over 55% of the query protein length were considered as significant.

Table 3. Evolutionary analysis of annotated partial cycloartenol synthase proteins in class Bacillariophyceae (diatoms)

(A)

Identity matrix						
<i>S. costatum</i>						
<i>S. subsalsum</i>	0.96					
<i>Extubocellulus</i>	0.76	0.75				
<i>Pseudo-nitzchia</i>	0.67	0.66	0.62			
<i>Fragilariopsis</i>	0.64	0.64	0.64	0.82		
<i>Thalassiosira</i>	0.81	0.81	0.73	0.64	0.63	
<i>Phaeodactylum</i>	0.70	0.69	0.71	0.6	0.61	0.67

(B)

Z test of selection	<i>S. costatum</i>	<i>S. subsalsum</i>	<i>Extubocellulus</i>	<i>Pseudo-nitzchia</i>	<i>Fragilariopsis</i>	<i>Thalassiosira</i>	<i>Phaeodactylum</i>
<i>S. costatum</i>		6.98	10.82	8.93	9.52	11.09	9.97
<i>S. subsalsum</i>	0.00		12.26	9.87	4.66	11.53	10.81
<i>Extubocellulus</i>	0.00	0.00		7.53	9.21	9.75	10.75
<i>Pseudo-nitzchia</i>	0.00	0.00	0.00		10.34	9.87	9.23
<i>Fragilariopsis</i>	0.00	0.00	0.00	0.00		9.34	9.16
<i>Thalassiosira</i>	0.00	0.00	0.00	0.00	0.00		11.12
<i>Phaeodactylum</i>	0.00	0.00	0.00	0.00	0.00	0.00	

(A) Identity matrix: Identity values between sequences being 1, 100% identical. (B) Codon-based Test of Purifying Selection for analysis between sequences: The probability of rejecting the null hypothesis of strict-neutrality ( $d_N = d_S$ ) in favor of the alternative hypothesis ( $d_N < d_S$ ) (above diagonal) is shown. Below the diagonal the p-values are shown (less than 0.05 are considered significant at the 5% level). The test statistic ( $d_S - d_N$ ) is shown below the diagonal.  $d_S$  and  $d_N$  are the numbers of synonymous and nonsynonymous substitutions per site, respectively. The variance of the difference was computed using the bootstrap method (1000 replicates). Analyses were conducted using the Nei-Gojobori method (Nei & Gojori, 1986). The analysis involved 7 nucleotide sequences. All ambiguous positions were removed for each sequence pair. There were a total of 193 positions in the final dataset.

Table 4. Evolutionary analysis of annotated cycloartenol synthase proteins in phylum Chlorophyta (green algae)

(A)

Identity matrix					
<i>Volvox</i>					
<i>Ostreococcus</i>	0.467				
<i>Chlamydomonas</i>	0.715	0.496			
<i>Micromonas</i>	0.510	0.607	0.527		
<i>Chlorella</i>	0.427	0.395	0.464	0.415	

(B)

Z test of selection	<i>Volvox</i>	<i>Ostreococcus</i>	<i>Chlamydomonas</i>	<i>Micromonas</i>	<i>Chlorella</i>
<i>Volvox</i>		12.04	13.11	17.15	10.47
<i>Ostreococcus</i>	0.00		11.33	19.91	11.38
<i>Chlamydomonas</i>	0.00	0.00		17.02	6.88
<i>Micromonas</i>	0.00	0.00	0.00		15.96
<i>Chlorella</i>	0.00	0.00	0.00	0.00	

(A) Identity matrix: Identity values between sequences being 1, 100% identical. (B) Codon-based Test of Purifying Selection for analysis between sequences: The probability of rejecting the null hypothesis of strict-neutrality ( $d_N = d_S$ ) in favor of the alternative hypothesis ( $d_N < d_S$ ) (above diagonal) is shown. Below the diagonal the p-values are shown (less than 0.05 are considered significant at the 5% level). The test statistic ( $d_S - d_N$ ) is shown below the diagonal.  $d_S$  and  $d_N$  are the numbers of synonymous and nonsynonymous substitutions per site, respectively. The variance of the difference was computed using the bootstrap method (1000 replicates). Analyses were conducted using the Nei-Gojobori method (Nei & Gojori, 1986). The analysis involved 7 nucleotide sequences. All ambiguous positions were removed for each sequence pair. There were a total of 826 positions in the final dataset.

Suppl Table 1. Diatom cultures investigated in this study.

Species	Sterol composition*	PCR result
<i>Thalassiosira pseudonana</i> CCMP 1335	C <sub>28</sub> 11(85), 15(6); C <sub>29</sub> 25(5), 26(4)	+
<i>Thalassiosira gravida</i> CCMP 986	C <sub>28</sub> 11(95); C <sub>29</sub> 25(5)	-
<i>Thalassiosira aff antarctica</i> CCAP 1085/9	C <sub>27</sub> 2(4); C <sub>28</sub> 11(83); C <sub>29</sub> 25(3)	-
<i>Fragilaropsis cylindricus</i> CCMP 1102	C <sub>27</sub> 2(68), 5(32)	+
<i>Pseudo-nitzschia seriata</i> CCMP 1309	C <sub>27</sub> 2(82), 5(18)	+
<i>Extubocellulus spinifer</i> CCMP 393	C <sub>27</sub> 2(79), 5(2); C <sub>28</sub> 10(13), 11(1), 13(1), 14(1); C <sub>29</sub> 29(3)	+
<i>Extubocellulus cribiger</i> CCAP 1026/1	C <sub>27</sub> 2(74), 5(4); C <sub>28</sub> 10(12), 11(4), 13(1), 15(2); C <sub>29</sub> 29(3)	-
<i>Phaeodactylum tricornutum</i>	C <sub>28</sub> 10(99), 15(1)	+
<i>Attheya septentrionalis</i> CS 425/03	C <sub>27</sub> 5(7), 5(2); C <sub>28</sub> 11(37), 15(19); C <sub>29</sub> 26(tr), 35(38)	-
<i>Skeletonema costatum</i> CCMP 1281	C <sub>27</sub> 3(5), 5(8), 8(1); C <sub>28</sub> 11(73), 18(11); C <sub>29</sub> 25(1), 26(1)	+
<i>Skeletonema subsalsum</i> CCAP 1077/8	C <sub>27</sub> 3(10), 5(44); C <sub>28</sub> 11(24), 15(5); C <sub>29</sub> 25(tr), 26(11), 35(6)	+
<i>Proboscia indica</i> CCMP 1896	C <sub>28</sub> 11(90), 18(10)	-
<i>Proboscia inermis</i> CCAP 1064/1	C <sub>27</sub> 5(44); C <sub>28</sub> 11(56)	-
<i>Proboscia alata</i>	C <sub>27</sub> 3(34), 5(29); C <sub>28</sub> 11(35); C <sub>29</sub> 24(2)	-

\*Data from Rampen et al., 2010.: Numbers in front of parentheses correspond to the sterol numbers as it follows: [2 (cholesta-5,22E-dien-3 $\beta$ -ol); 3 (cholesta-5,24-dien-3 $\beta$ -ol); 5 (cholest-5-en-3 $\beta$ -ol); 8 (5 $\alpha$ -cholestan-3 $\beta$ -ol); 10 (24-methylcholesta-5,22E-dien-3 $\beta$ -ol); 11 (24-methylcholesta-5,24(28)-dien-3 $\beta$ -ol); 13 (23-methylcholesta-5,22E-dien-3 $\beta$ -ol); 14 (23-methylcholesta-5,23(28)-dien-3 $\beta$ -ol); 15 (24-methylcholesta-5-en-3 $\beta$ -ol); 18 (24-methylcholesta-24(28)-en-3 $\beta$ -ol); 24 (24-ethylcholesta-5,22E-dien-3 $\beta$ -ol); 25 (24-ethylcholesta-5,24(28E)-dien-3 $\beta$ -ol); 26 (24-ethylcholesta-5,24(28Z)-dien-3 $\beta$ -ol); 29 (23,24-dimethylcholesta-5,22E-dien-3 $\beta$ -ol); 35 (24-ethylcholest-5-en-3 $\beta$ -ol)]. Values in parentheses represent the concentration of the individual sterol, as a percentage of the total sterols. tr indicates relative abundances of <0.5%.

Suppl Table 2. Sterol composition of North Sea SPM and microbial mat samples.

Diatomaceous microbial mat sterol composition	North Sea water column sterol composition
<ul style="list-style-type: none"> <li>- Cholesta-5,22E-dien-3<math>\beta</math>-ol</li> <li>- 5<math>\alpha</math>-cholest-22-en-3<math>\beta</math>-ol</li> <li>- Cholest-5en-3<math>\beta</math>-ol (cholesterol)</li> <li>- 5<math>\alpha</math>-cholestan-3<math>\beta</math>-ol (cholestanol)</li> <li>- Cholesta-5,24-dien-3<math>\beta</math>-ol (desmosterol)</li> <li>- 24-methylcholest-5,22E-dien-3<math>\beta</math>-ol (brassicasterol/diatomsterol)</li> <li>- 5<math>\alpha</math>(H)-cholest-7en-3<math>\beta</math>-ol</li> <li>- 24-methylcholest-5,24(28)-dien-3<math>\beta</math>-ol</li> <li>- 24-methylcholest-5-en-3<math>\beta</math>-ol (campesterol)</li> <li>- 23,24-dimethylcholesta-5,22E-dien-3<math>\beta</math>-ol</li> <li>- 24-ethylcholesta-5,22E-dien-3<math>\beta</math>-ol (stigmasterol)</li> <li>- 24-ethylcholesta-5-en-3<math>\beta</math>-ol (<math>\beta</math>-sitosterol)</li> <li>- 5<math>\alpha</math>(H)-23,24-dimethylcholestanol (5<math>\alpha</math>H-C29:0)</li> <li>- 24-ethylcholesta-5,24Z-Zden-3<math>\beta</math>-ol (isofucosterol)</li> <li>- 4<math>\alpha</math>, 23, 24-trimethyl-5<math>\alpha</math>(H)cholest-22-en-3<math>\beta</math>-ol (dinosterol)</li> </ul>	<ul style="list-style-type: none"> <li>- Cholesta-5,22E-dien-3<math>\beta</math>-ol</li> <li>- Cholest-5en-3<math>\beta</math>-ol (cholesterol)</li> <li>- 5<math>\alpha</math>-cholestan-3<math>\beta</math>-ol (cholestanol)</li> <li>- Cholesta-5,24-dien-3<math>\beta</math>-ol (desmosterol)</li> <li>- 24-methylcholest-5,22E-dien-3<math>\beta</math>-ol (brassicasterol/diatomsterol)</li> <li>- 5<math>\alpha</math>(H)-cholest-7en-3<math>\beta</math>-ol</li> <li>- 5<math>\alpha</math>(H)-24-methyl-cholest-22en-3<math>\beta</math>-ol</li> <li>- 24-methylcholest-5,24(28)-dien-3<math>\beta</math>-ol</li> <li>- 24-methylcholest-5-en-3<math>\beta</math>-ol (campesterol)</li> <li>- 24-ethylcholesta-5-en-3<math>\beta</math>-ol (<math>\beta</math>-sitosterol)</li> <li>- 24-ethylcholesta-5,24Z-Zden-3<math>\beta</math>-ol (isofucosterol)</li> </ul>

Figure 1

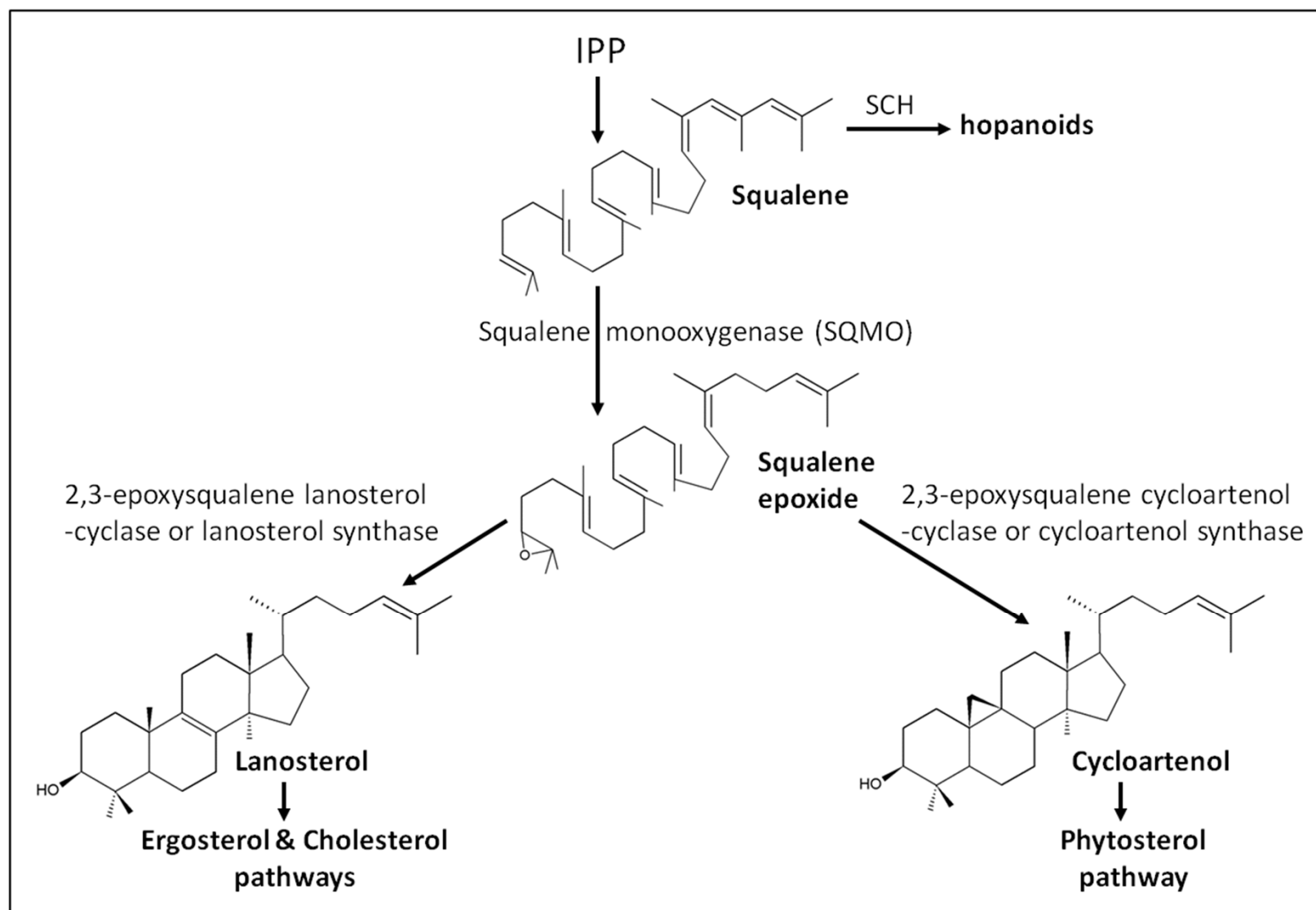


Figure 2A

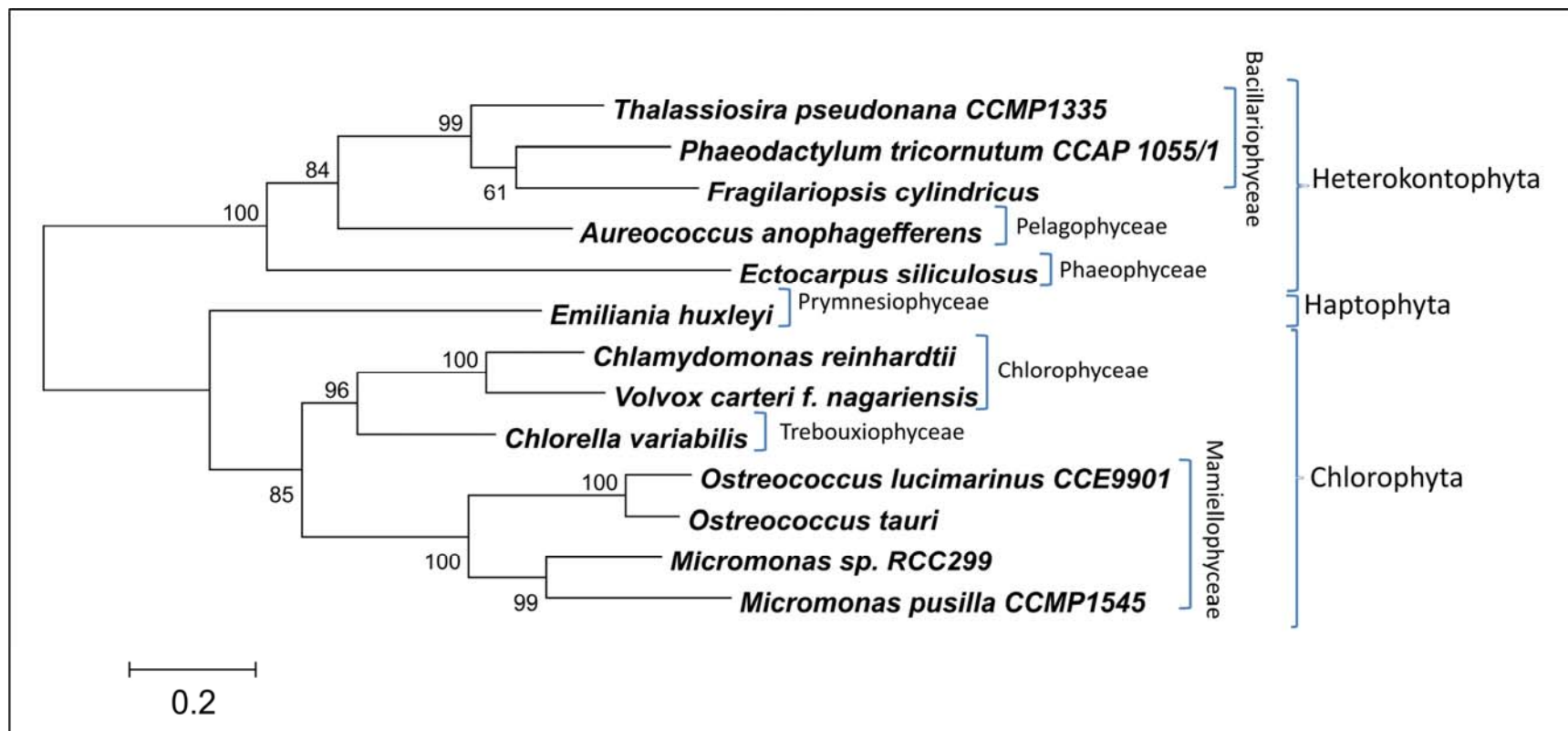


Figure 2B

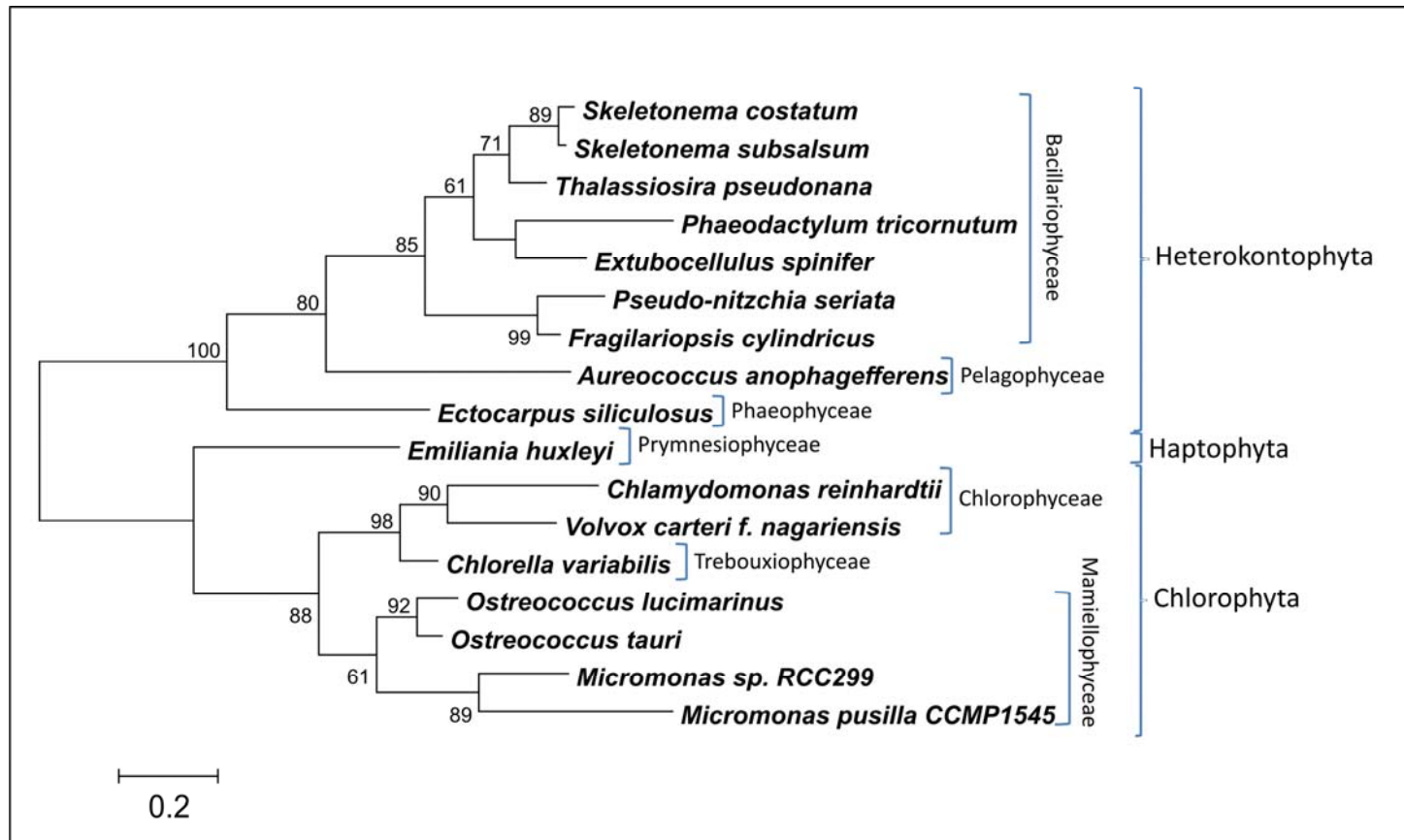




Figure 2C

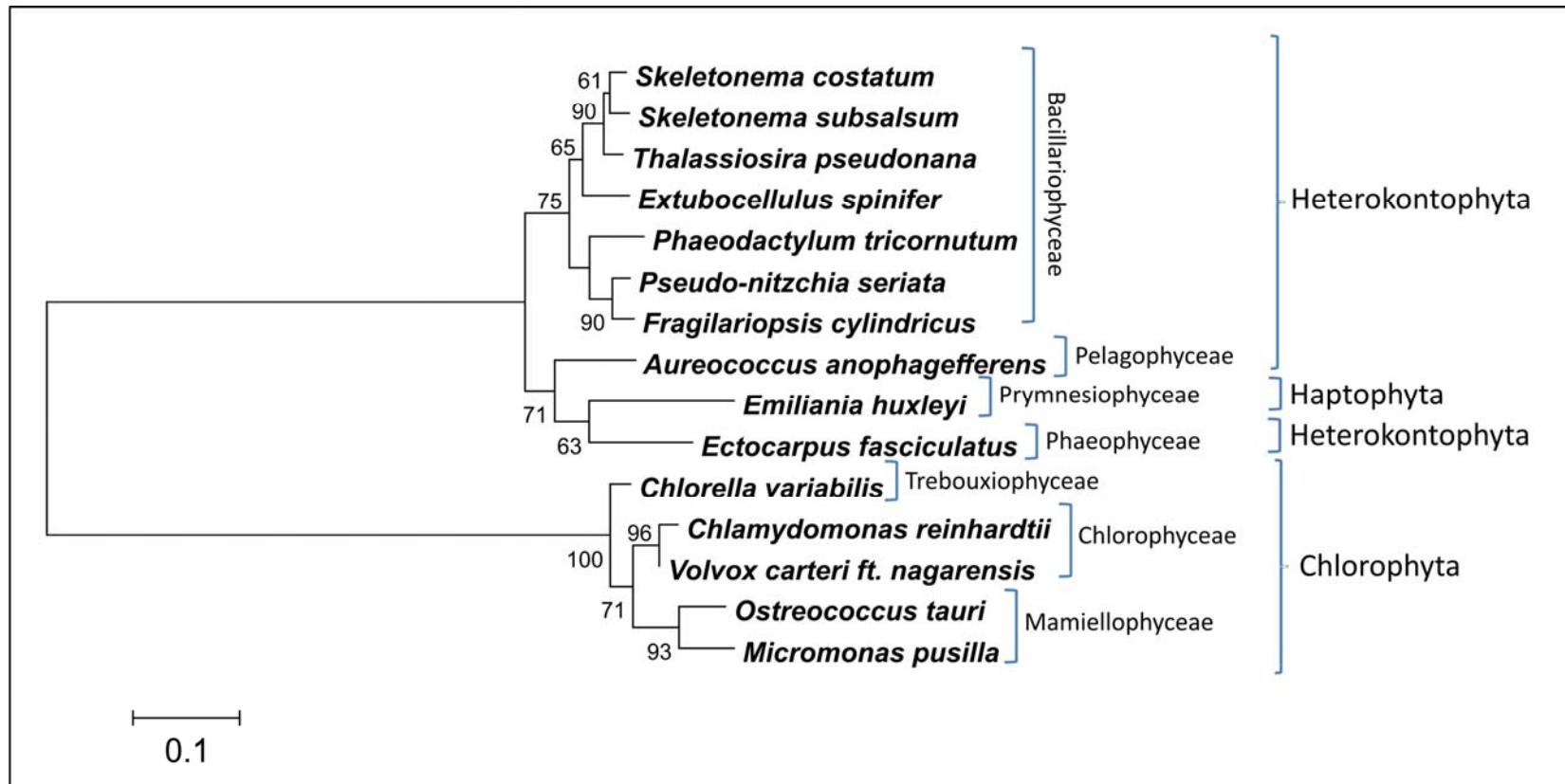


Figure 3

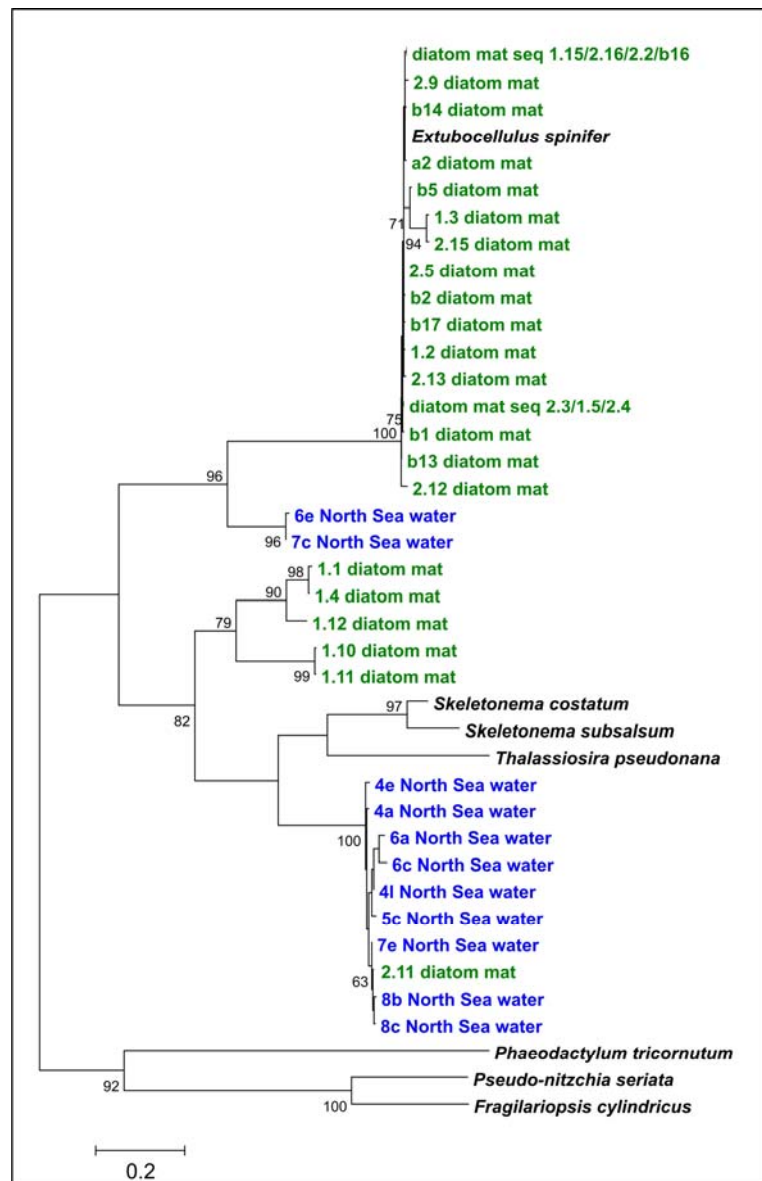
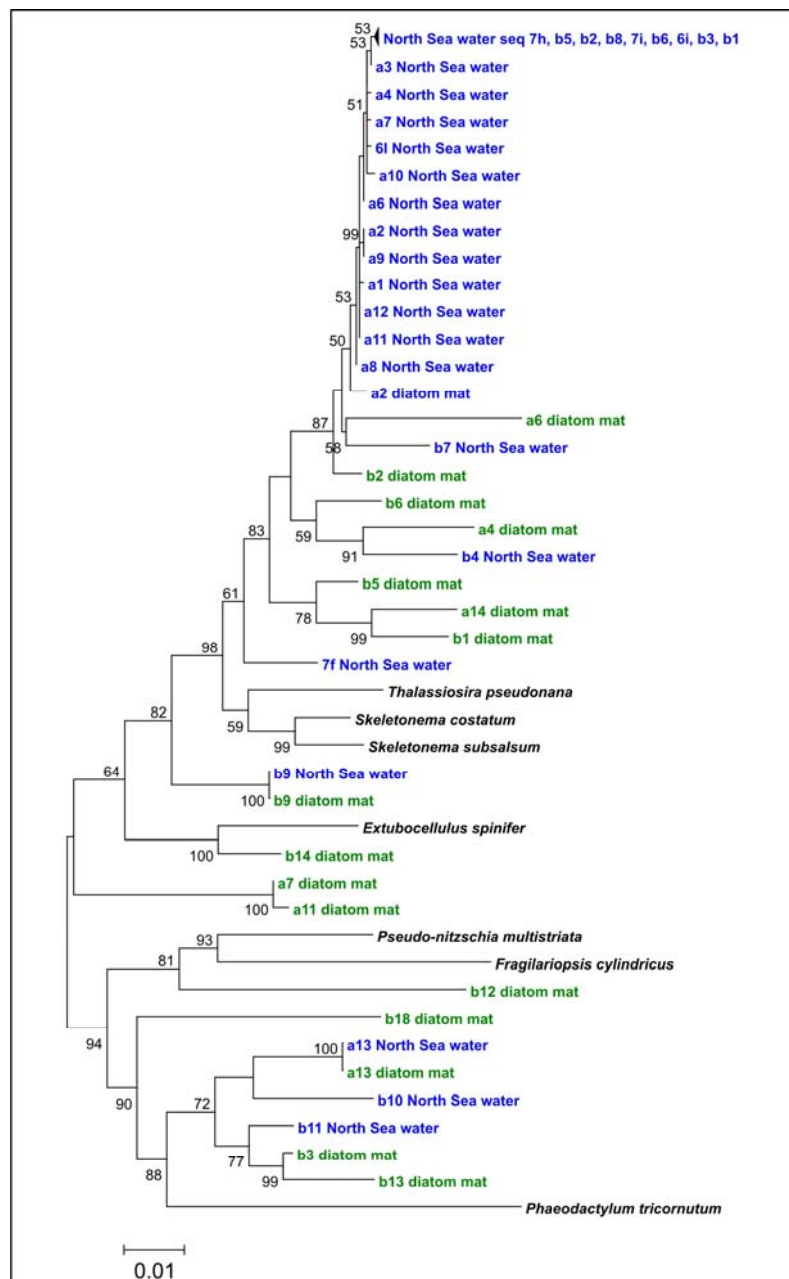
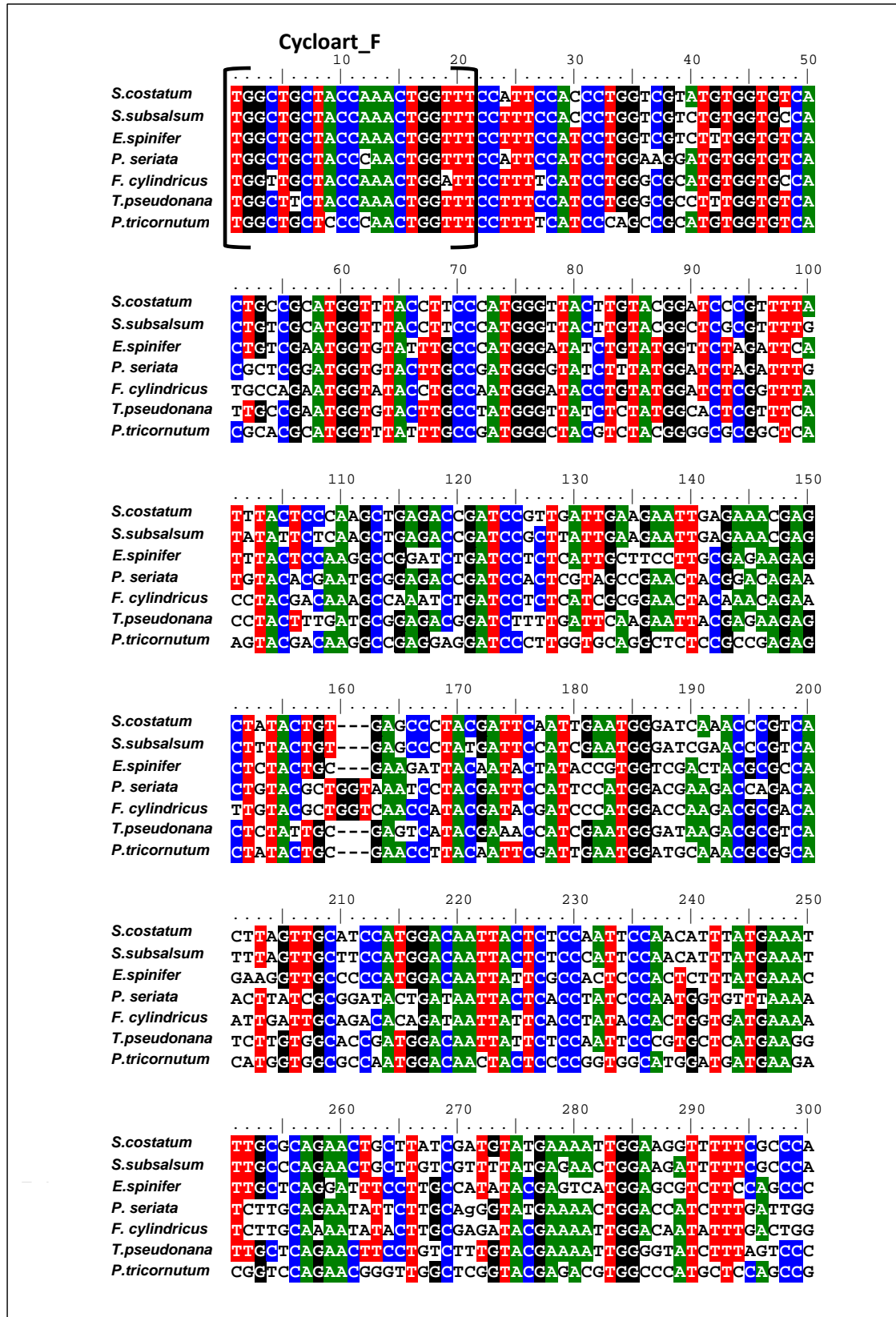


Figure 4



Suppl. Figure 1. Alignment cycloartenol synthase partial gene sequences involved in this analysis



310 320 330 340 350

*S.costatum* TTTTCGTGATGCATTCAGGAAGGCTGGACTTGATTTCTGTCTGGAATACAT  
*S.subsalsum* TTTTCGTGATGCAATTAGAAAGGCTGGACTTGACTTCTGTCTGGAGTATAT  
*E.spinifer* TTCAAGAATAGATTTAGGAAGATTGGTCTCAAGTTTTCGCTCGAGTACAT  
*P. seriata* TTCCGGAACTACTTTCCGCAACGAGGTTTGGATTTTACGATGGAATACAT  
*F. cylindricus* TTTTCGGAAATCACGTTTCGTCAACGAGGTTTAGAATTTTCAATGGAATACAT  
*T.pseudonana* TTCCGTAATGCTGTACGAAAGGCTGGGTGAAATACTGTCTCGAGTACAT  
*P.tricornutum* TTCAAAACGACGTCGCAAACTCGGCTTGGCCCTTCTGTGTCTGATACAT

360 370 380 390 400

*S.costatum* GCGAGCAGAAGATTTGCAAAACCAACTTTATTGATATTGGCCCTGTGAACA  
*S.subsalsum* GCGAGCGGAAGACTTTGCAAAACCAACTACATTGATATTGGCCCGGTGAACA  
*E.spinifer* GGCCGCTGAGGATTTGCAACGAAATTTATTGATATTGGCCAGTAAACA  
*P. seriata* GAAAGCAGAAGATTTGCAAAACAAATTTATTGATATTGGCCCGGTGAATA  
*F. cylindricus* GAAAGCTGAAGATTTTCAAAACAAATTTATTGATATTAGGTCCCGTTAACA  
*T.pseudonana* GCGAGCTGAAGACTTTGCAACCAATTTATTGATATTGGCCCGGTGAACA  
*P.tricornutum* GGCAGCGGAAGATCTGCAAAACCAACTTTATTGATATTAGGCCCGGTGAATA

410 420 430 440 450

*S.costatum* AAGCATTGAAATGGTCTCGGCGTTTCACTCTGCA-----  
*S.subsalsum* AAGCACTGAAATGGTCTCGGCTTTCCATCATGCA-----  
*E.spinifer* AGGCTCTCAATCTCGTTTCGGCTTATCAGCTGCG-----  
*P. seriata* AAATTCTGAATATGTTATCGATGATCACTATCAT-----AAAT  
*F. cylindricus* AAGTATTAAATATGTTGTCGATGATCATTATCATCATCATGGTGAACA  
*T.pseudonana* AAGCATTGAAATGGTCTCAGCGTTTCAT-----  
*P.tricornutum* AAGTGTCTCAATATGCTCTCGGCTTTTCATCAGCA-----

460 470 480 490 500

*S.costatum* -----AACAAATGATATCAACGACCCCGCAGTGCCTAG  
*S.subsalsum* -----AATAATGATATCAATGACCCCTGCAGTACGTAG  
*E.spinifer* -----GGTAACGATATCAACGCCACCCACGTACAGAA  
*P. seriata* AATAGTACAAAGGAAGATGAAAGCGATGACGGAGGGCTCTCATCGAAAA  
*F. cylindricus* AATGATATAAAGTTGGACAGCAACGACAATAATGAGCTTCTTATCGAGCG  
*T.pseudonana* -----GTCCGTAG  
*P.tricornutum* -----GGAAATGATTGTCATCATTCAACAGTGATGAA

510 520 530 540 550

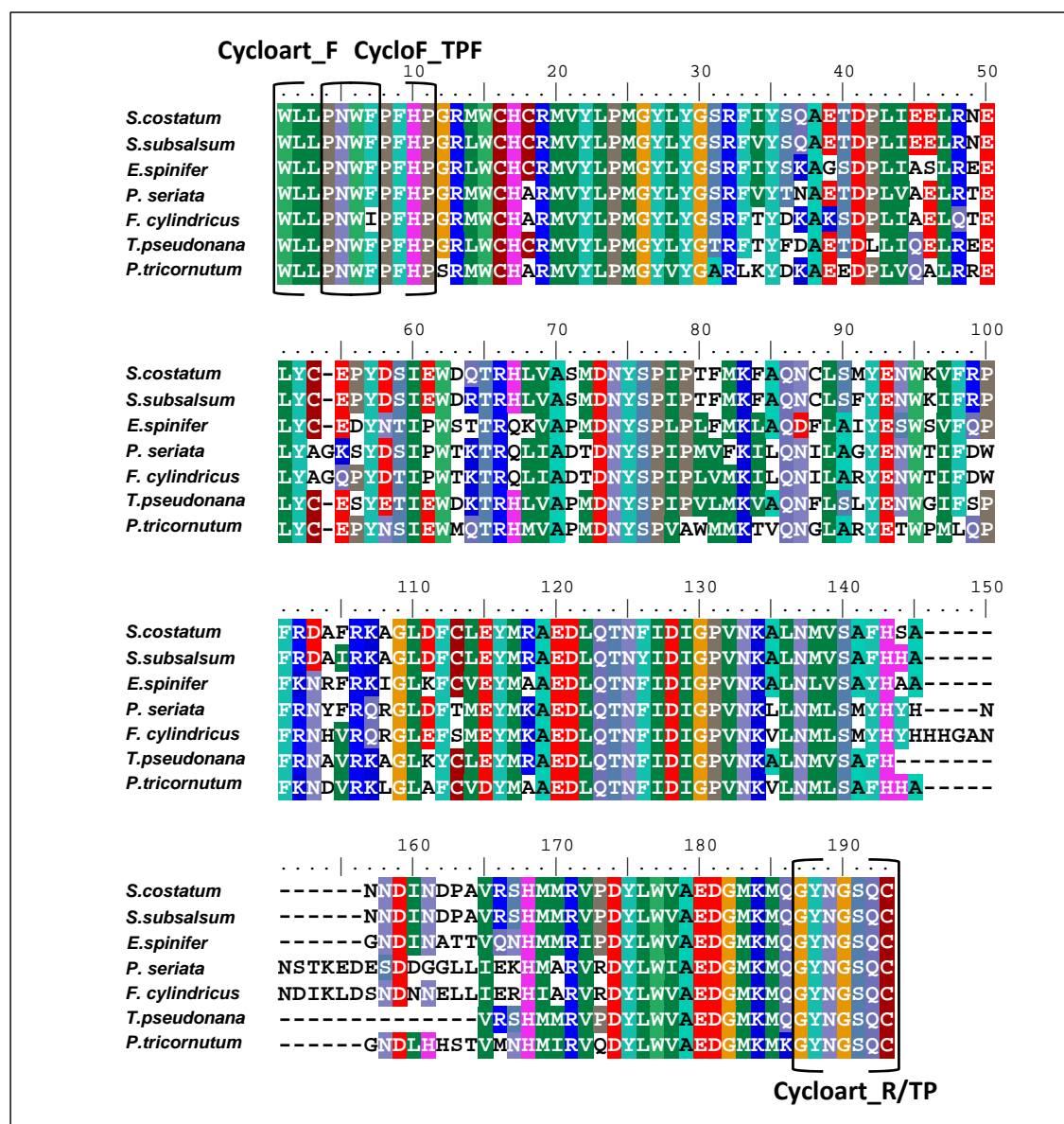
*S.costatum* TCACATGATGCGCTGCCGACTATCTTTGGGTGGCTGAGGATGGTATGA  
*S.subsalsum* TCACATGATGCGTGTGCCGATTATCTTTGGGTGGCTGAGGATGGTATGA  
*E.spinifer* TCACATGATGAGAATTCTTGATTATCTTTGGGTAGCGGAAGATGGAATGA  
*P. seriata* GCACATGGCACGGGTTCTGTGATTATTTGTGGATCGCCGAAGATGGGATGA  
*F. cylindricus* CCACATCGCCCGGTACGTGATTACCTCTGGGTAGCCGAAGATGGAATGA  
*T.pseudonana* TCACATGATGCGCTACCACTACCTCTGGGTAGCCGAAGATGGAATGA  
*P.tricornutum* CCACATGATTGAGTTGAGACTATTTATGGGTGCGGAAGATGGCATGA

#### Cycloart\_R

560 570

*S.costatum* AAATGCAAGGCTATAACGGCAGCCAATGC  
*S.subsalsum* AAATGCAAGGCTACAACGGAGCCAATGC  
*E.spinifer* AGATGCAAGGCTATAACGGCAGCCAATGC  
*P. seriata* AGATGCAAGGCTACAACGGCAGCCAATGC  
*F. cylindricus* AAATGCAAGGATATAACGGCAGCCAATGT  
*T.pseudonana* AAATGCAAGGCTACAACGGAGCCAATGC  
*P.tricornutum* AAATGCAAGGCTATAACGGCAGCCAATGC

Suppl. Figure 2. Protein alignment sequences of partial cycloartenol synthase (CSG) sequences from diatom pure cultures inferred by xblast



Suppl. Figure 3. Phylogenetic tree of 18S rRNA sequences (DNA) in the microalgae under study inferred by using the Maximum Likelihood method based on the GTR model. A discrete Gamma distribution was used to model evolutionary rate differences among sites (5 categories (+G, parameter = 0.2831). The analysis involved 15 nucleotide sequences. All ambiguous positions were removed for each sequence pair. There were a total of 3589 positions in the final dataset.

