

# 1 **Telomere-to-telomere genome assembly of** 2 ***Phaeodactylum tricornutum***

3 **Daniel J. Giguere<sup>1</sup>, Alexander T. Bahcheli<sup>1</sup>, Samuel S. Slattery<sup>1</sup>, Rushali**  
4 **R. Patel<sup>1</sup>, Martin Flatley<sup>2</sup>, Bogumil J. Karas<sup>1</sup>, David R. Edgell<sup>1</sup>, and**  
5 **Gregory B. Gloor<sup>1</sup>**

6 <sup>1</sup>**Department of Biochemistry, Schulich School of Medicine and Dentistry, Western**  
7 **University, London, ON N6A 5C1, Canada**

8 <sup>2</sup>**Suncor Energy, Sarnia, Ontario, Canada**

9 Corresponding author:

10 Gregory B. Gloor<sup>1</sup>

11 Email address: [ggloor@uwo.ca](mailto:ggloor@uwo.ca)

## 12 **ABSTRACT**

13 *Phaeodactylum tricornutum* is a marine diatom with a growing genetic toolbox available and is being  
14 used in many synthetic biology applications. While most of the genome has been assembled, the  
15 currently available genome assembly is not a completed telomere-to-telomere assembly. Here, we used  
16 Oxford Nanopore long reads to build a telomere-to-telomere genome for *Phaeodactylum tricornutum*.  
17 We developed a graph-based approach to extract all unique telomeres, and used this information  
18 to manually correct assembly errors. In total, we found 25 nuclear chromosomes that comprise all  
19 previously assembled fragments, in addition to the chloroplast and mitochondrial genomes. We found  
20 that chromosome 19 has filtered long-read coverage and a quality estimate that suggests significantly  
21 less haplotype sequence variation than the other chromosomes. This work improves upon the previous  
22 genome assembly and provides new opportunities for genetic engineering of this species, including  
23 creating designer synthetic chromosomes.

## 25 **INTRODUCTION**

26 *Phaeodactylum tricornutum* is a marine diatom that is described as a "diatom cell factory" (Butler et al.,  
27 2020) because it can be used to manufacture valuable commercial products. Recent genetic toolbox  
28 expansions, such as delivering episomes by bacterial conjugation (Karas et al., 2015) CRISPR-editing  
29 tools (Russo et al., 2018; Moosburner et al., 2020; Sharma et al., 2018; Stukenberg et al., 2018; Slattery  
30 et al., 2018; Serif et al., 2018), generation of auxotrophic strains (Zaslavskaja et al., 2000; Buck et al., 2018;  
31 Slattery et al., 2020), and the identification of highly active endogenous promoters (Erdene-Ochir et al.,  
32 2019) that are enabling rapid implementation of new product designs into commercial-scale production.

33 The genome of *P. tricornutum* CCAP 1055/1 was sequenced in 2008, and resulted in a scaffold-level  
34 assembly predicting 33 chromosomes (NCBI assembly ASM15095v2) (Bowler et al., 2008). Chloroplast  
35 and mitochondrial genomes have also been published (Oudot-Le Secq et al., 2007; Oudot-Le Secq and  
36 Green, 2011), and have previously identified as targets for genetic engineering (Cochrane et al., 2020), as  
37 well as other chromosomes (Karas et al., 2013). Although the Bowler *et al.*, assembly contains several  
38 telomere-to-telomere chromosomes, many scaffolds have only zero or one telomere, suggesting they  
39 are either incomplete or fragments of another chromosome. More recent work identifying centromeric  
40 sequences (Diner et al., 2017) in *P. tricornutum* has suggested that there may be less than 33 chromosomes.

41 To generate a telomere-to-telomere assembly of *P. tricornutum* CCAP 1055/1, we used a hybrid  
42 approach with ultra-long reads from the Oxford Nanopore minION platform and highly accurate short  
43 reads from the Illumina NextSeq platform. We also introduce a novel graph-based approach to manually  
44 resolve telomere-related assembly errors. This approach identifies all unique telomere sequences and we  
45 demonstrate how it can be applied to manually correct assembly errors adjacent to chromosome ends.  
46 The full structural context of the *P. tricornutum* genome provides additional information for potential  
47 synthetic biology applications to manipulate the genome of this diatom cell factory.

## 48 METHODS

### 49 Growth

50 *Phaeodactylum tricornutum* (Culture Collection of Algae and Protozoa CCAP 1055/1) was grown in L1  
51 medium without silica at 18° C under cool white fluorescent lights (75 mE m<sup>-2</sup> s<sup>-1</sup>) and a photoperiod of  
52 16 h light:8 h dark as described previously (Slattery et al., 2018).

### 53 DNA extraction

54 200 mL of culture (approximately 5 x 10<sup>8</sup> cells) was spun at 3000 X g for 10 minutes at 4° C. The  
55 pellet was resuspended in 1 mL TE (pH 8.0) and added dropwise to a mortar (pre-cooled at -80° C)  
56 pre-filled with liquid nitrogen. The frozen droplets were ground into a fine powder with a mortar and  
57 pestle, being careful to keep the cells from thawing by adding more liquid nitrogen as necessary. The  
58 frozen powder was transferred to a 15 mL Falcon tube where 2 mL of lysis buffer was added (1.4 M NaCl,  
59 200 mM Tris-HCl pH 8.0, 50 mM EDTA, 2% (w/v) CTAB, RNase A (250 µg/mL) and proteinase K (100  
60 µg/mL)). The solution was mixed very slowly by inversion, incubated for 30 minutes at 37° C (mixed very  
61 slowly halfway through incubation). Cellular debris was pelleted at 6000 X g for 5 minutes. Lysate was  
62 transferred to a new 15 mL Falcon tube. One volume of 25:24:1 phenol:chloroform:isoamyl alcohol was  
63 added, mixing slowly by inversion. The sample was centrifuged at 6000 X g for 5 minutes. The aqueous  
64 phase was transferred as slow as possible to a new Falcon tube. One volume of 24:1 chloroform:isoamyl  
65 alcohol was added, and mixed slowly with end-over-end inversion. The sample was centrifuged at 6000 X  
66 g for 5 minutes. Approximately 450 µL of the aqueous phase was transferred into new 1.5 mL Eppendorf  
67 tubes. To the Eppendorf tube, 1/10 volume of 3 M NaAc pH 5.2 and 2 volumes (final volume) of ice-cold  
68 100% ethanol were added, mixing slowly by end-over-end inversion. The sample was centrifuged at 16  
69 000 X g for 5 minutes, and washed twice with 500 µL 70% ethanol. Ethanol was decanted, and the pellet  
70 was dried for approximately 10 minutes by inverting on a paper towel. The pellet was resuspended in 100  
71 µL 10 mM Tris-HCl pH 8.0, 0.1 mM EDTA pH 8.0. After resuspending overnight at 4° C, short DNA  
72 fragments were then selectively removed using the Short Read Eliminator (SRE) kit from Circulomics  
73 (Baltimore). DNA from the same extraction was used for sequencing on both the Oxford Nanopore  
74 minION and Illumina NextSeq 550 platform.

### 75 Sequencing

76 All sequencing reads are publically available on the European Nucleotide Archive under project PR-  
77 JEB42700. All raw .fast5 files are available on under accession number ERR5858460.

78 An Oxford Nanopore minION flow cell R9.4.1 was used with the SQK-LSK109 Kit according to  
79 the manufacturer's protocol version GDE\_9063\_v109\_revK\_14Aug2019, with one alteration: for DNA  
80 repair and end-prep, the reaction mixture was incubated for 15 minutes at 20° C and 15 minutes at 65° C.  
81 Basecalling was performed after the run with Guppy (Version 3.6). NanoPlot (De Coster et al., 2018) was  
82 used to generate Q-score versus length plots and summary statistics. The read N50 of the unfiltered reads  
83 was approximately 35 kb (Fig. S1). Nanopore reads are available under accession ERR5207170.

84 For Illumina sequencing, the Nextera XT kit was used to prepare 2X75 paired-end mid-output NextSeq  
85 550 run at the London Regional Genomics Center (lrgc.ca). Reads were trimmed using Trimmomatic  
86 v0.36 (Bolger et al., 2014) in paired end mode with the following settings: AVGQUAL:30 CROP:75  
87 SLIDINGWINDOW:4:25 MINLEN:50 TRAILING:15. SLIDINGWINDOW AND TRAILING were  
88 added to remove poor quality base calls. Only paired end reads were retained. Illumina reads are available  
89 under accession ERR5198869.

### 90 Assembly

#### 91 *Telomere identification*

92 We first obtained sequences for the end of every linear chromosome. The sequence of the telomere repeats  
93 for *P. tricornutum* are known from the previous assembly (Bowler et al., 2008) to be repeats of AACCT.  
94 All long reads larger than 50 kilobases with 3 or more consecutive telomeric repeats (or the reverse  
95 complement) were extracted by filtering using NanoFilt (De Coster et al., 2018) and by string matching  
96 using grep. All-versus-all mapping of the telomeric reads was performed using minimap2 (Li, 2016).  
97 Only overlapping reads with a minimum query coverage of 95 % were retained.

98 To determine the sequence of unique telomeres for each chromosome, a network graph was generated  
99 with iGraph (Csardi and InterJournal, 2006). Each read name was used as a vertex, and edges were

100 generated between each overlapping read with more than 95% query coverage. Noise was filtered by  
101 removing any group of overlaps with less than 5X coverage. There were 95 vertices that had greater than  
102 20X coverage; that is, there are 95 unique telomere sequence groups. Most groups had approximately 40X  
103 coverage (number of long reads per group), however, several outliers had with more than 60X coverage.  
104 These represent duplicated regions in the telomeres that are not unique (i.e., more than one haplotype or  
105 chromosome contains this sequence). We estimate that ninety-five unique telomere groups indicates that  
106 there are at least 24 or more chromosomes (95 chromosomes divided by 2 ends per chromosome divided  
107 by 2 haplotypes). There are at least 88 groups of telomere sequences that are unique and can be used to  
108 improve the assembly, with the remaining 7 possibly duplicated. The longest read of each telomere group,  
109 typically greater than 100 kb in length, was retained as a representative telomere sequence for correction.  
110 Example code for this is available in Code S7.

### 111 **Assembly**

112 Several recent assembly algorithms and with multiple parameters were attempted, but we found overlap-  
113 layout consensus to provide the most contiguous assembly as a starting point. Sequencing reads longer  
114 than 75 kilobases were used for initial assembly with miniasm, (Li, 2016) using the parameters -s 30000  
115 -m 10000 -c 5 -d 100000. From this initial assembly, unitigs were manually completed with the following  
116 approach:

117 1) Mapping of telomeric reads against the unitig. If no telomere was present on the unitig and a high  
118 query coverage alignment was found, the unitig was extended to the telomere sequence of the mapped  
119 telomere. 2) After telomere extension (or confirmation), reads longer than 50 kb were mapped to the  
120 unitig to confirm overlapping coverage over the entire chromosome. Coverage was evaluated using  
121 only reads larger than 50 kb and with higher than 60% query coverage, with an alignment score:length  
122 ratio less than 2 (similar to previous validation methods)(Giguere et al., 2020). A query coverage of  
123 only 60% was chosen to allow for potential haplotype divergence. 3) Telomere-to-telomere unitigs with  
124 overlapping ultra-long read coverage and no gaps were deemed validated and brought forward to improve  
125 base accuracy by read polishing.

126 The chloroplast and mitochondrial genomes were assembled using a reference based approach by first  
127 extracting all reads that aligned to the reference chloroplast and mitochondria with high query coverage.  
128 Reads were then *de-novo* assembled using miniasm.

### 129 **Polishing**

130 Due the repetitive nature of the genome and the diploid nature of *P. tricornutum*, raw assemblies were  
131 polished using an iterative approach with racon (Vaser et al., 2017), medaka (Oxford Nanopore)  
132 (Walker et al., 2014) as described in the Methods S6. Briefly, after each polishing iteration, we corrected  
133 errors that were introduced by the polishing algorithms as described in Methods S6, and modified the  
134 medaka polishing by filtering reads using a minimum of 50% query coverage. The assembly was first  
135 polished by nanopore reads only, followed by Illumina read polishing using Pilon. For the chloroplast and  
136 mitochondria, the subset of reads identified as either chloroplast or mitochondria were used for polishing.

### 137 **Methylation**

138 5mC methylation sites were predicted using Megalodon v2.2.1 (Oxford Nanopore Technologies) using  
139 the model `res_dna_r941_min_modbases_5mC_CpG_v001.cfg` from the Rerio repository (Oxford Nanopore  
140 Technologies) with Guppy 4.5.2. A default threshold of 0.75 was used as a minimum score for modified  
141 base aggregation (probability of modified/canonical base) to produce the final aggregated output. The  
142 percentage of reads methylated at the predicted are plotted in Fig. S2.

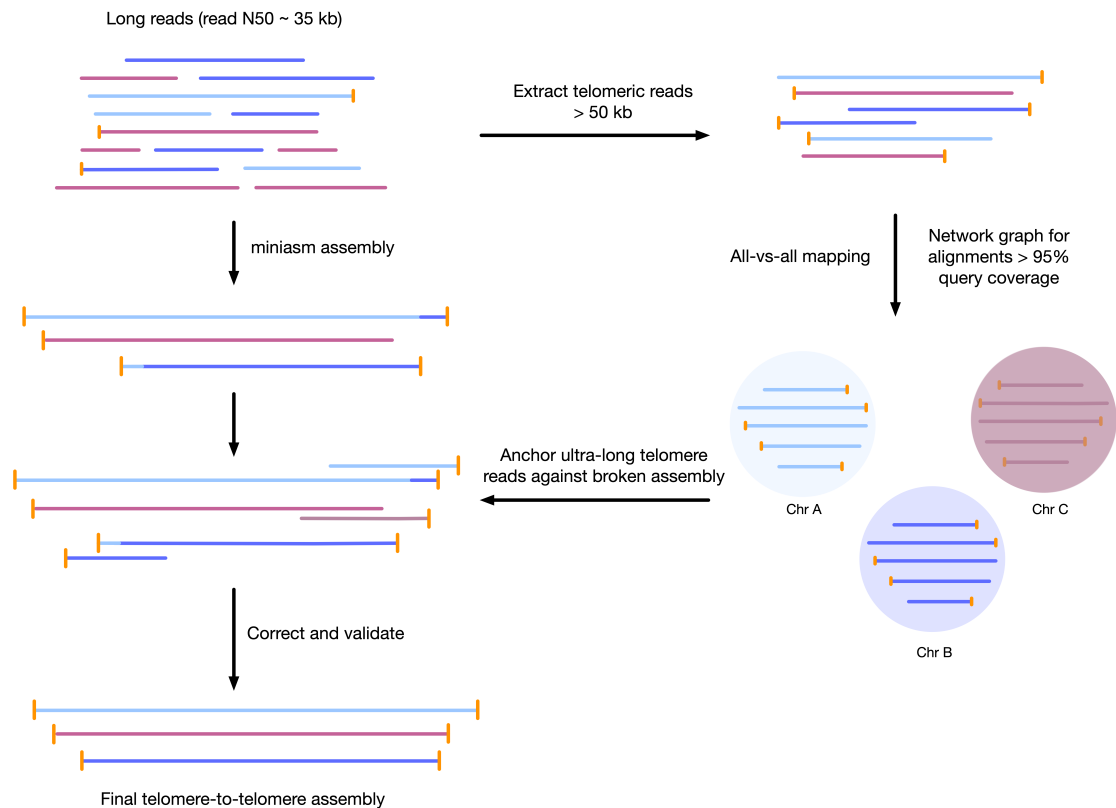
143 All data and code used to generate figures are available from DOI: 10.5281/zenodo.4731048.

## 144 **RESULTS**

### 145 **Workflow**

146 We developed a sample and library preparation protocol that provided exceptionally long reads. We  
147 observed a read N50 of 35 kilobases, with the longest reads just over 300 kb, following sequencing with  
148 the Oxford Nanopore minION platform. Of the 7.8 gigabases of raw sequence data, approximately 2.5  
149 gigabases were from reads longer than 50 kilobases (Figure S1). We found that chromosomes assembled  
150 using standard approaches were often mis-assembled around telomeres, or were fragmented and only

151 contained 1 telomere. To correct each unitig, we leveraged the unique ultra-long telomere reads as  
152 as described in the Methods S6 and in Figure 1. This approach was used to manually identify a tiling path  
153 for each chromosome until each chromosome was contiguous from telomere to telomere and validated by  
154 a tiling overlapping read path.



**Figure 1.** Workflow for telomere-to-telomere genome assembly. Telomere-containing nanopore reads larger than 50 kb are extracted, and are mapped in all-vs-all mode using minimap2. The resulting alignments are filtered by 95 % query coverage, and a network graph is created using iGraph using read names as vertices, and alignments between reads as edges. Each resulting cluster represents one end of a chromosome. The initial reads are used for overlap-consensus assembly using miniasm. On a chromosome-by-chromosome basis, ultra-long read coverage is plotted. If an assembled chromosome is missing a telomere or has an assembly error revealed by a lack of overlapping read coverage, the longest read from each telomere cluster is mapped against the chromosome, and the resulting telomere is used to manually correct the assembly and extend to the telomere using an overlap-layout consensus approach.

155 ***Tiling path of overlapping reads verify contiguity***

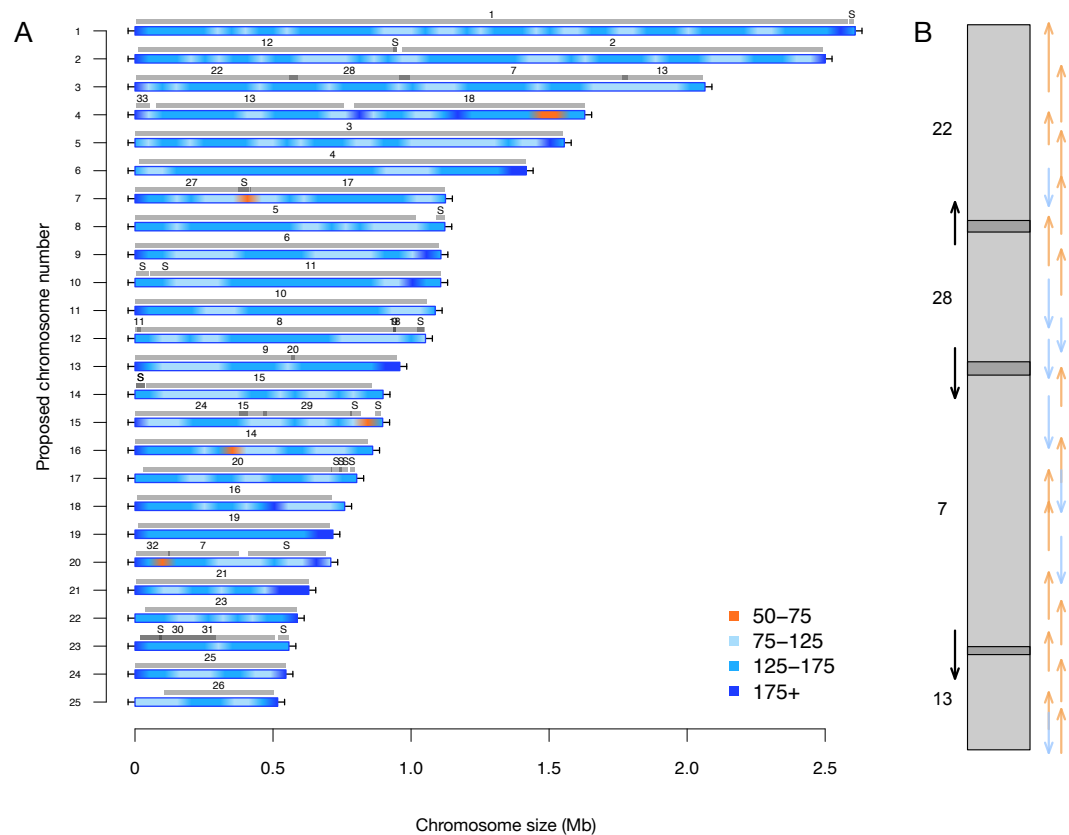
156 To ensure our genome assembly is contiguous, we generated multiple independent minimum tiling paths  
157 of overlapping long reads (Data S5, Fig. S2). Reads larger than 50 kb were mapped against the assembly  
158 using minimap2. To ensure no incorrect alignments were retained, any reads with less than 90% of the  
159 read aligned to the assembly were removed. From this subset, 5 independent minimum tiling paths that  
160 required at least 10 kb of overlap between each read was generated. All chromosomes have multiple  
161 independent (i.e., no common reads) tiling path of reads with a minimum overlap of 10 kb in the final  
162 validated assembly (5 independent paths shown in .paf format in Data S5), indicating that all chromosomes  
163 are contiguous. Chromosomes were manually corrected to meet this standard if necessary.

164 In addition to overlapping reads, Fig. S2 also shows the GC content for each chromosome. A previous  
165 study has proposed that centromeres could be identified by low GC content (Diner et al., 2017). The  
166 100 base window with the minimum GC content are shown in Fig. S2 and are highlighted between red  
167 dotted lines. Putative centromere sequences characterized by the lowest GC content window as previously  
168 described Diner et al. (2017) are highlighted between in red.

169 **Telomere-to-telomere assembly comprises previous scaffolds**

170 We ultimately obtained 25 telomere-to-telomere chromosome assemblies that recruit 98% of long reads,  
 171 and these chromosomes comprised all previously proposed chromosomes from Bowler *et al.*, (2008), as  
 172 well as circularized chloroplast and mitochondrial genomes. The median coverage for unfiltered long  
 173 reads across the nuclear genome was 202X, while median coverage for the chloroplast and mitochondrion  
 174 were approximately 6201X and 528X, respectively. This was calculated in 1000 base windows using  
 175 mosdepth (Pedersen and Quinlan, 2018).

176 A key feature of this updated assembly is the consistency with previous sequencing efforts (Bowler  
 177 *et al.*, 2008). Previously, 25 centromere sequences were identified (Diner *et al.*, 2017), suggesting that  
 178 there were fewer than the proposed 33 chromosomes. This agrees with our conclusion of 25 nuclear  
 179 chromosomes. In addition, we independently resolved the relative location for all of the previously  
 180 proposed partial chromosomes without internal inconsistencies in Figure 2 (i.e., scaffolds with only 1  
 181 telomere were resolved into a telomere-to-telomere chromosome).



**Figure 2.** A) Filtered long-read coverage and comparison to previous assembly. Reads longer than 20 kb were mapped against the assembly, filtered (minimum 20000 base alignment and 50 % query coverage), and genome coverage was calculated in 50 kb windows using mosdepth. Newly proposed chromosome names are left (by length). Scaffolds from the previous genome assembly (ASM15095v2) are overlaid as grey bars, aligned using minimap2 in asm5 mode and filtered to retain minimum 10 kb alignments. Numbers on top of grey bars indicate which previous scaffold number, with S representing small unassembled scaffolds. Horizontal "T" bars on each end indicate telomere-repeat presence. B) Visualization of proposed chromosome 3 with alignments to previous chromosomes. Dark gray regions indicate overlap. Coloured arrows on the right indicate minimum overlapping read path (orange = negative strand, blue = positive strand, black arrows on left show ultra-long reads that completely span regions where previous assembly could not assemble through).

Unique reads recruited	Unique filtered reads recruited
98.12%	74.00%

**Table 1.** Long reads were filtered to remove all reads smaller than 1000 bases and below a Q-score of 8 using NanoFilt. Lambda spike-in reads were then removed using NanoLyse, and the total number of reads was calculated using NanoPlot. The number of reads recruited was calculated by aligning the reads against the assembly using minimap2, and unique read ideas were counted. The number of filtered reads was counted after removing reads with less than 90 percent query coverage.

## 182 **Assembly quality**

183 To assess the quality of the assembly, we used Merqury (Rhie et al., 2020) to estimate the base-level  
184 accuracy and completeness by k-mer frequency, shown in Data S3. We found that the estimated quality  
185 value (estimated log-scaled probability of error for the consensus base calls by Merqury) ranged from  
186 27 - 53, depending on the chromosome. The mean quality value for nuclear chromosomes was 28.86,  
187 with chromosome 19 as an outlier at 43. The QV for all nuclear genomes except for 19 are likely lower  
188 because the chromosomes were polished using reads that are heterozygous - this can likely be improved  
189 in the future by binning the reads into haplotypes before polishing. The chloroplast and mitochondrial  
190 genomes have a quality value of 53 and 42, respectively. Importantly, the k-mer completeness estimate  
191 of 80% suggests that many k-mers in the Illumina reads are not represented in this genome assembly,  
192 implying there is still additional genomic sequence missing from this assembly. This was also the case  
193 when using the Bowler assembly as input for Merqury. This is expected because this assembly is not yet  
194 haplotype-resolved.

195 We also show that nearly all the reads are recruited in Table 1. When reads are filtered by query  
196 coverage, the number of reads recruited drops to 74%, indicating there is still sequence information to be  
197 determined by resolving the haplotypes.

## 198 **Filtered long-read coverage for Chromosome 19 is inconsistent with diploid state**

199 We observed that chromosome 19 has remarkably consistent (i.e., no drops in coverage) filtered long-read  
200 coverage relative to the other linear chromosomes (Figure 2, Fig. S2). *Phaeodactylum tricoratum*  
201 is a diploid organism but has not been observed to undergo meiosis. Therefore, we expected to observe  
202 2 unique haplotypes per chromosome. These 2 haplotypes can be inferred by drops in filtered-long  
203 read coverage to half the total coverage, which indicates that the chromosome is a combination of both  
204 haplotypes (Fig. S2). This variation in coverage is observed for all nuclear chromosomes with the striking  
205 exception of chromosome 19. Chromosome 19 also has duplicated coverage near the telomeres. These  
206 observations suggest that there are not two haplotypes for chromosome 19, suggesting a different recent  
207 history for this chromosome in this strain.

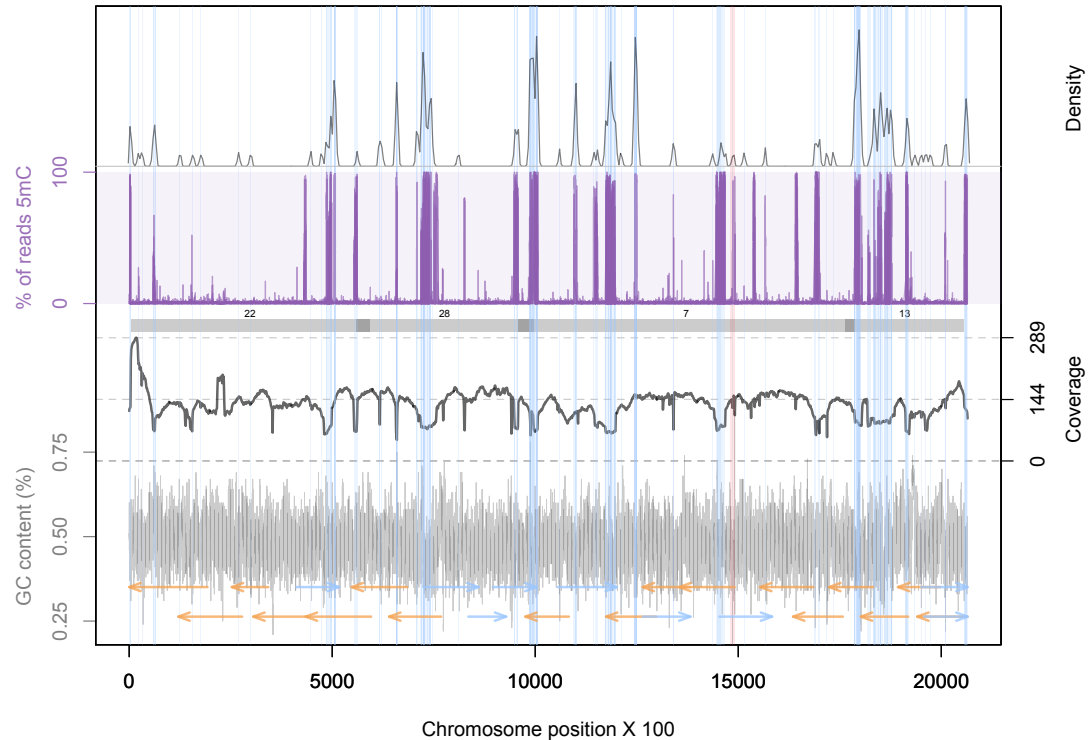
## 208 **5mC methylation and transposable elements**

209 The Extensive de-novo TE Annotator (EDTA) pipeline (Ou et al., 2019) was used to predict transposable  
210 elements in the genome. We found that the majority of transposable elements are long terminal repeat  
211 (LTR) retrotransposons (3.64% of the genome was found to be Copia-type, 5.65% were unknown - Data  
212 S4) while terminal inverted repeats were only 1.17% of the genome, and helitrons were 0.54% of the  
213 genome. Each LTR region is represented as shaded blue regions in Fig. S2 in blue, and density plots  
214 of the end locations are shown in the top quadrant. Chromosome 19 contained the fewest transposable  
215 elements at 50. The locations and density of LTR-retrotransposons are plotted in Figure 3 for proposed  
216 Chromosome 3 and Fig. S2 for all other chromosomes.

217 Previous studies have found that some transposable elements were hypermethylated (Veluchamy et al.,  
218 2013). Using chromosome scale nanopore methylation basecalling, we found a strong signal between  
219 many predicted LTR retrotransposons and methylation status (Figure 3, Fig. S2).

220 We examined the association between LTR transposon dense regions and the locations of low coverage  
221 regions and regions where the previous assembly failed to generate overlapping regions. We observed  
222 that scaffolds with overlapping regions (Fig. S2) generally were not assembled into full chromosomes  
223 because of ambiguity in the placement of the LTR-rich regions at the ends of the scaffolds. These are  
224 now resolved by the long-read assembly identified here. Additionally, many of the low-coverage regions

225 of our assembly overlap with the locations of the LTR-dense regions and this would be consistent with  
226 chromosomal rearrangements being more likely in these regions. Further investigation at these regions is  
227 required.



**Figure 3.** Genomic summary for proposed Chromosome 3. GC content was calculated in 100 base windows, and plotted in the bottom quadrant. An overlapping read tiling path, with a minimum overlap of 30 kb, is shown with orange indicating reads mapping to the negative strand and blue indicating reads mapping to the positive strand. Filtered long-read coverage is plotted in the second quadrant (minimum 20 kb length and 70% query coverage). Number scaffolds from the previous assembly are overlaid in gray bars, with dark grey representing overlapping regions. The third quadrant shows the proportion of reads that were called as methylated. The top quadrant shows the density of LTR-retrotransposon end points as predicted by the EDTA pipeline. Light blue box (full plot) shows the regions that are annotated at LTR-retrotransposons. Light red box indicates the 100 base window with the minimum GC content.

## 228 DISCUSSION

229 Here, we developed a graph-based approach to locate the unique telomere ends of all *P. tricornerutum*  
230 chromosomes, and applied this information to generate an telomere-to-telomere assembly. The new  
231 assembly incorporates all the previous chromosome fragments from Bowler (2008), and also includes  
232 circularized organelle genomes.

233 The chromosomes show marked variations in genome coverage in a two-fold range. This suggests that  
234 there are large regions of the chromosomes that have substantial haplotype differences. Strikingly, one  
235 chromosome has completely consistent coverage between the telomeres. While this needs to be further  
236 investigated, we speculate that this chromosome in this strain may have undergone a recent sequence  
237 homogenization event. Previous work has also found that chromosome 19 appears homozygous in the  
238 wild type strain (Russo et al., 2018). It has previously been speculated that *Phaeodactylum tricornerutum*  
239 may be capable of sexual reproduction (Mao et al., 2020; Patil et al., 2015), but there has yet to be

240 conclusive evidence.

241 We also demonstrate that nanopore sequencing can identify methylated regions associated with  
242 transposable elements (Fig. S2). While this strain of algae was not grown under stress conditions,  
243 we demonstrate that this new technology can be applied to to characterize the 5mC methylome of *P.*  
244 *tricornutum*, and can be applied to future differential methylation experiments.

245 Chromosome 19 has a high quality value of 43, while the other nuclear chromosomes have lower  
246 quality values around 28, representing an expected drop in per-base quality due to polishing the assembly  
247 with a heterozygous read set. Taken together with the consistent filtered-long read coverage, these data  
248 suggest that there are not highly divergent haplotypes of chromosome 19. In addition, the higher quality  
249 values for the organelle genomes and chromosome 19 indicate that this dataset is sufficient to generate  
250 highly accurate (quality value greater than 40) assemblies when haplotypes can be resolved. Furthermore,  
251 the k-mer completeness estimate suggests there is additional unknown genomic information (up to 20%  
252 of k-mers in Illumina reads), and we believe the remaining sequence can be identified by generating a  
253 haplotype-resolved genome assembly for *P. tricornutum*. We have deposited all raw short and long-read  
254 data publicly for use by the community as Project PRJEB42700 on the European Nucleotide Archive.

255 While this work improves on the previous genome assembly, we believe there is yet further room  
256 to improve our understanding of the *P. tricornutum* genome. First, we believe that haplotype-phasing  
257 can be completed with ultra-long read data. Furthermore, there modified bases can be further explored  
258 in relation to gene expression. We have therefore posted the raw files and encourage the community to  
259 investigate other genomic features with this data. We also recognize that the per-base accuracy of this  
260 assembly may be improved by binning the reads used for polishing by haplotype - as evidence by the  
261 high quality value estimate for chromosome 19 and the organelle genomes. The per-base quality of this  
262 genome assembly will likely benefit from resolving the haplotypes and re-analysis as improvements to  
263 basecalling and polishing algorithms become available in the future. This telomere-to-telomere genome  
264 assembly will make it possible to start designing and creating synthetic chromosomes in *Phaeodactylum*  
265 *tricornutum*.

## 266 FUNDING

267 DJG Mitacs number IT8360, Ontario Graduate Scholarship

268 SS: Mitacs number IT8360

269 GBG: Natural Sciences and Engineering Research Council of Canada (NSERC), grant number:  
270 RGPIN-03878-2015

271 BK: Natural Sciences and Engineering Research Council of Canada (NSERC), grant number: RGPIN-  
272 2018-06172

273 DE: Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant  
274 RGPIN-2015-04800

## 275 REFERENCES

- 276 Bolger, A. M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence  
277 data. *Bioinformatics*, 30(15):2114–2120.
- 278 Bowler, C., Allen, A. E., Badger, J. H., Grimwood, J., Jabbari, K., Kuo, A., Maheswari, U., Martens, C.,  
279 Maumus, F., Otiillar, R. P., Rayko, E., Salamov, A., Vandepoele, K., Beszteri, B., Gruber, A., Heijde,  
280 M., Katinka, M., Mock, T., Valentin, K., Verret, F., Berges, J. A., Brownlee, C., Cadoret, J.-P., Chiovitti,  
281 A., Choi, C. J., Coesel, S., De Martino, A., Detter, J. C., Durkin, C., Falciatore, A., Fournet, J., Haruta,  
282 M., Huysman, M. J. J., Jenkins, B. D., Jiroutova, K., Jorgensen, R. E., Joubert, Y., Kaplan, A., Kröger,  
283 N., Kroth, P. G., La Roche, J., Lindquist, E., Lommer, M., Martin-Jézéquel, V., Lopez, P. J., Lucas, S.,  
284 Mangogna, M., McGinnis, K., Medlin, L. K., Montsant, A., Oudot-Le Secq, M.-P., Napoli, C., Obornik,  
285 M., Parker, M. S., Petit, J.-L., Porcel, B. M., Poulsen, N., Robison, M., Rychlewski, L., Rynearson,  
286 T. A., Schmutz, J., Shapiro, H., Siaut, M., Stanley, M., Sussman, M. R., Taylor, A. R., Vardi, A., von  
287 Dassow, P., Vyverman, W., Willis, A., Wyrwicz, L. S., Rokhsar, D. S., Weissenbach, J., Armbrust,  
288 E. V., Green, B. R., Van de Peer, Y., and Grigoriev, I. V. (2008). The *Phaeodactylum* genome reveals  
289 the evolutionary history of diatom genomes. *Nature*, 456(7219):239–244.
- 290 Buck, J. M., Bártulos, C. R., Gruber, A., and Kroth, P. G. (2018). Blastocidin-S deaminase, a new selection  
291 marker for genetic transformation of the diatom *Phaeodactylum tricornutum*. *PeerJ*, 6:e5884.



- 292 Butler, T., Kapoore, R. V., and Vaidyanathan, S. (2020). Phaeodactylum tricornutum: A Diatom Cell  
293 Factory. *Trends in biotechnology*.
- 294 Cochrane, R. R., Brumwell, S. L., Soltysiak, M. P. M., Hamadache, S., Davis, J. G., Wang, J., Tholl, S. Q.,  
295 Janakirama, P., Edgell, D. R., and Karas, B. J. (2020). Rapid method for generating designer algal  
296 mitochondrial genomes. *Algal Research*, 50:102014.
- 297 Csardi, G. and InterJournal, T. N. (2006). The igraph software package for complex network research.  
298 *researchgate.net*.
- 299 De Coster, W., D'Hert, S., Schultz, D. T., Cruys, M., and Van Broeckhoven, C. (2018). NanoPack:  
300 visualizing and processing long-read sequencing data. *Bioinformatics*, 34(15):2666–2669.
- 301 Diner, R. E., Noddings, C. M., Lian, N. C., Kang, A. K., McQuaid, J. B., Jablanovic, J., Espinoza,  
302 J. L., Nguyen, N. A., Anzelmann, M. A., Jansson, J., Bielinski, V. A., Karas, B. J., Dupont, C. L.,  
303 Allen, A. E., and Weyman, P. D. (2017). Diatom centromeres suggest a mechanism for nuclear  
304 DNA acquisition. *Proceedings of the National Academy of Sciences of the United States of America*,  
305 114(29):E6015–E6024.
- 306 Erdene-Ochir, E., Shin, B.-K., Kwon, B., Jung, C., and Pan, C.-H. (2019). Identification and characterisa-  
307 tion of the novel endogenous promoter HASP1 and its signal peptide from Phaeodactylum tricornutum.  
308 *Scientific Reports*, 9(1):9941–10.
- 309 Giguere, D. J., Bahcheli, A. T., Joris, B. R., Paulssen, J. M., Gieg, L. M., Flatley, M. W., and Gloor, G. B.  
310 (2020). Complete and validated genomes from a metagenome. *bioRxiv*, 11(11):2020.04.08.032540.
- 311 Karas, B. J., Diner, R. E., Lefebvre, S. C., McQuaid, J., Phillips, A. P. R., Noddings, C. M., Brunson, J. K.,  
312 Valas, R. E., Deerinck, T. J., Jablanovic, J., Gillard, J. T. F., Beeri, K., Ellisman, M. H., Glass, J. I.,  
313 Hutchison, C. A., Smith, H. O., Venter, J. C., Allen, A. E., Dupont, C. L., and Weyman, P. D. (2015).  
314 Designer diatom episomes delivered by bacterial conjugation. *Nature Communications*, 6(1):6925–10.
- 315 Karas, B. J., Molparia, B., Jablanovic, J., Hermann, W. J., Lin, Y.-C., Dupont, C. L., Tagwerker, C.,  
316 Yonemoto, I. T., Noskov, V. N., Chuang, R.-Y., Allen, A. E., Glass, J. I., Hutchison, C. A., Smith, H. O.,  
317 Venter, J. C., and Weyman, P. D. (2013). Assembly of eukaryotic algal chromosomes in yeast. *Journal*  
318 *of Biological Engineering*, 7(1):1–12.
- 319 Li, H. (2016). Minimap and miniasm: fast mapping and de novo assembly for noisy long sequences.  
320 *Bioinformatics*, 32(14):2103–2110.
- 321 Mao, Y., Guo, L., Luo, Y., Tang, Z., Li, W., and Dong, W. (2020). Sexual reproduction potential implied  
322 by functional analysis of SPO11 in Phaeodactylum tricornutum. *Gene*, 757(7):144929.
- 323 Moosburner, M. A., Gholami, P., McCarthy, J. K., Tan, M., Bielinski, V. A., and Allen, A. E. (2020).  
324 Multiplexed Knockouts in the Model Diatom Phaeodactylum by Episomal Delivery of a Selectable  
325 Cas9. *Frontiers in microbiology*, 11.
- 326 Ou, S., Su, W., Liao, Y., Chougule, K., Agda, J. R. A., Hellinga, A. J., Lugo, C. S. B., Elliott, T. A., Ware,  
327 D., Peterson, T., Jiang, N., Hirsch, C. N., and Hufford, M. B. (2019). Benchmarking transposable  
328 element annotation methods for creation of a streamlined, comprehensive pipeline. *Genome biology*,  
329 20(1):1–18.
- 330 Oudot-Le Secq, M.-P. and Green, B. R. (2011). Complex repeat structures and novel features in the  
331 mitochondrial genomes of the diatoms Phaeodactylum tricornutum and Thalassiosira pseudonana.  
332 *Gene*, 476(1-2):20–26.
- 333 Oudot-Le Secq, M.-P., Grimwood, J., Shapiro, H., Armbrust, E. V., Bowler, C., and Green, B. R.  
334 (2007). Chloroplast genomes of the diatoms Phaeodactylum tricornutum and Thalassiosira pseudonana:  
335 comparison with other plastid genomes of the red lineage. *Molecular genetics and genomics : MGG*,  
336 277(4):427–439.
- 337 Patil, S., Moeys, S., von Dassow, P., Huysman, M. J. J., Mapleson, D., De Veylder, L., Sanges, R.,  
338 Vyverman, W., Montresor, M., and Ferrante, M. I. (2015). Identification of the meiotic toolkit in  
339 diatoms and exploration of meiosis-specific SPO11 and RAD51 homologs in the sexual species  
340 Pseudo-nitzschia multistriata and Seminavis robusta. *BMC Genomics*, 16(1):930–21.
- 341 Pedersen, B. S. and Quinlan, A. R. (2018). Mosdepth: quick coverage calculation for genomes and  
342 exomes. *Bioinformatics*, 34(5):867–868.
- 343 Rhie, A., Walenz, B. P., Koren, S., and Phillippy, A. M. (2020). Merqury: reference-free quality,  
344 completeness, and phasing assessment for genome assemblies. *Genome biology*, 21(1):245–27.
- 345 Russo, M. T., Cigliano, R. A., Sanseverino, W., and Ferrante, M. I. (2018). Assessment of genomic  
346 changes in a CRISPR/Cas9 Phaeodactylum tricornutum mutant through whole genome resequencing.

- 347 *PeerJ*, 6:e5507.
- 348 Serif, M., Dubois, G., Finoux, A.-L., Teste, M.-A., Jallet, D., and Daboussi, F. (2018). One-step generation  
349 of multiple gene knock-outs in the diatom *Phaeodactylum tricornutum* by DNA-free genome editing.  
350 *Nature Communications*, 9(1):1–10.
- 351 Sharma, A. K., Nymark, M., Sparstad, T., Bones, A. M., and Winge, P. (2018). Transgene-free genome  
352 editing in marine algae by bacterial conjugation – comparison with biolistic CRISPR/Cas9 transforma-  
353 tion. *Scientific Reports*, 8(1):1–11.
- 354 Slattery, S. S., Diamond, A., Wang, H., Therrien, J. A., Lant, J. T., Jazey, T., Lee, K., Klassen, Z., Desgagné-  
355 Penix, I., Karas, B. J., and Edgell, D. R. (2018). An Expanded Plasmid-Based Genetic Toolbox Enables  
356 Cas9 Genome Editing and Stable Maintenance of Synthetic Pathways in *Phaeodactylum tricornutum*.  
357 *ACS synthetic biology*.
- 358 Slattery, S. S., Wang, H., Giguere, D. J., Kocsis, C., Urquhart, B. L., Karas, B. J., and Edgell, D. R. (2020).  
359 Plasmid-based complementation of large deletions in *Phaeodactylum tricornutum* biosynthetic genes  
360 generated by Cas9 editing. *Scientific Reports*, 10(1):1–12.
- 361 Stukenberg, D., Zauner, S., Dell’Aquila, G., and Maier, U. G. (2018). Optimizing CRISPR/Cas9 for the  
362 Diatom *Phaeodactylum tricornutum*. *Frontiers in Plant Science*, 9.
- 363 Vaser, R., Sović, I., Nagarajan, N., and Šikić, M. (2017). Fast and accurate de novo genome assembly  
364 from long uncorrected reads. *Genome Research*, 27(5):737–746.
- 365 Veluchamy, A., Lin, X., Maumus, F., Rivarola, M., Bhavsar, J., Creasy, T., O’Brien, K., Sengamalay,  
366 N. A., Tallon, L. J., Smith, A. D., Rayko, E., Ahmed, I., Le Crom, S., Farrant, G. K., Sgro, J.-Y.,  
367 Olson, S. A., Bondurant, S. S., Allen, A. E., Rabinowicz, P. D., Sussman, M. R., Bowler, C., and  
368 Tirichine, L. (2013). Insights into the role of DNA methylation in diatoms by genome-wide profiling in  
369 *Phaeodactylum tricornutum*. *Nature Communications*, 4(1):1–10.
- 370 Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C. A., Zeng,  
371 Q., Wortman, J., Young, S. K., and Earl, A. M. (2014). Pilon: an integrated tool for comprehensive  
372 microbial variant detection and genome assembly improvement. *PLoS ONE*, 9(11):e112963.
- 373 Zaslavskaja, L. A., Lippmeier, J. C., Kroth, P. G., Grossman, A. R., and Apt, K. E. (2000). Transformation  
374 of the diatom *Phaeodactylum tricornutum* (Bacillariophyceae) with a variety of selectable marker and  
375 reporter genes. *Journal of Phycology*, 36(2):379–386.