1 Supplemental Methods for Traller et al "Genome and methylome of the oleaginous
2 diatom *Cyclotella cryptica* reveal genetic flexibility toward a high lipid phenotype"

3

4 *DNA extraction and purification*

5 　　　Liquid cell cultures concentrated either by filtration using a 3.0μm polycarbonate

6 filter (DNA for genome sequencing) or by centrifugation in 50mL conical tubes at 4000

7 xg for 8 minutes (methylome DNA) using in an Eppendorf 5810R centrifuge. DNA for

8 genome and bisulfite sequencing was purified using CsCl as described in [1]. To remove

9 RNA contamination, DNA for bisulfite sequencing was treated with 10mg/mL stock of

10 RNase A for 15 minutes at 37ºC.

11

12 *Genome Library Construction and Sequencing*

13 　　Three libraries were prepared following Illumina's standard genomic DNA paired end

14 construction: "PE-short" (post-assembly empirical insert lengths ~90-235 nucleotides (nt)

15 exclusive of adapters [preparation target 180-230 nt], with mode 123 nt), "PE-medium"

16 (~155-330 nt [preparation target 230- 330 nt], mode 221 nt), and "PE-long" (~225-460 nt

17 [preparation target 330-430 nt] with mode 305 nt) for ~58% of inserts, with the balance

18 in a second mode ~60-225 nt peaking at ~105 nt). These were sequenced as paired end

19 76-mer + 76-mer reads on an Illumina GA-IIx 120-tiles/lane run ("TP003") at the UCLA

20 BSCRC Core Sequencing facility, using two dedicated lanes for each of PE-short and PE-

21 medium, three dedicated lanes for PE-long, a dedicated PhiX control lane for RTA image

22 analysis autocalibration, and spiking in ≈1% Illumina PhiX into each non-control lane.

23 Two genomic DNA mate pair libraries (aiming for 10K nt effective inserts) were

24 prepared and run by Illumina service on a 48-tile/lane v3 HiSeq flow cell, each library on

25  a single dedicated lane: "MP-short" (effective ~2,100-3,320 nt with mode ~2,625 nt for

26  ~58% of inserts, with the balance in a PE-orientation [rather than MP-orientation] second

27  mode ~170-480 nt peaking at ~205 nt), and "MP-long" (~1,740-2,730 nt with mode

28  ~2,260 nt for ~63% of inserts, with the balance in a PE-oriented second mode ~176-535

29  nt peaking at ~223 nt). These were physically sequenced as paired end 101-mer + 10-mer

30  index + 101-mer, with the index reads and the last base of each main end discarded (in

31  the usual way so that the last retained base has bidirectional RTA phasing corrections).

32  The number of raw read pairs for PE-short/medium/long/control is found in Additional

33  File 1, Figure S1b. Only read pairs with no 'N' basecalls were retained; due to the pattern

34  of 'N' basecalls in PE-long, the first four bases of each end of its lanes were discarded

35  before this filter (and subsequent uses of this library).

36      The number of raw read pairs for MP-short and long libraries was 142,455,072 and

37  154,107,079, respectively, and only RTA PF=1-passing pairs were retained (Additional

38  File 1, Figure S1b). Although not used as a filter, relative to PF=1 read pairs, the 'N'-free

39  read pairs for MP-short/long were ~98.9% / 98.6%.

40      The PE and MP libraries contributed ~23.4G nt and ~57.4G nt, respectively, for a

41  total of ~80.8G nt (≈461x coverage of a 175Mbp genome; for 65-mers: ~3.4G and

42  ~20.7G, total ~24.0G and ≈137x).

43      Using many iterations of a variety of standard and *ad hoc* assemblers and alignment

44  tools, with extensive inspection of intermediate stages and judgment calls made by hand,

45  read pairs from PE-short/medium/long and MP-short/long were used to construct high-

46  quality best assemblies from the available data for the chloroplast ("chrC") and

47  mitochondrial ("chrM") genomes in *C. cryptica*. In each case, a single complete circular

48    sequence of pure A/C/G/T's without gaps was formed (chrC 129,320 bp, chrM 58,021

49    bp). This was greatly assisted by the presence in NCBI of related genomes: KJ958480.1

50    for *Cyclotella* strain L04_2 chrC (~96% identity; also useful: KJ958481.1 for *Cyclotella*

51    strain WC03_2 chrC), and NC_007405.1 for *Thalassiosira pseudonana* chrM (more

52    distant; even on homologous stretches, overall percent identity ≈80%). The *C. cryptica*

53    chrM is estimated with 17,880 bp (~31%) being 120 exact copies of the 149 bp sequence

54    TTATCGGCCTCAAATCAAGCAGTGTTTAAGCTGGAAT

55    CTATCGGCCTCAAATCGAAACAGTGTTTTAGCCTGAAT

56    TTATCGGCCTCAAATCAAGCAGTGTTTAAGCTGGAAT

57    CTATCGGCCTCAAATCGAAACAGTGTTTTGCCTGAAT  (which is itself four

58    approximate tandem copies of a smaller unit). Given current data, this region cannot be

59    completely resolved, and there is likely additional variation here, and the included

60    number of copies is an estimate informed in part by depth of coverage relative to other

61    chrM sequence.

62        The main genome assembly was performed with an ABySS 1.3.1 SE+PE+MP

63    pipeline using $k$=65-mers, $t$=65, $q$=3, $e$=3, $E$=0, $c$=3, $m$=30, $p$=0.9, no scaffolding at

64    PopBubbles, $s$=200, $n$=10, overlap min=5 with scaffolding and join masking at simple

65    repeats, SimpleGraph $d$=6 and scaffolding, greedy MergePaths, $a$=4, and abyss-scaffold

66    min-gap=100. The SE stage used PE-short+medium+long, MP-short+long, the PE stage

67    was applied to PE-short+medium+long, and the MP stage was applied to MP-short+long.

68    The assembly was taken as the final scaffolds of nt length $\geq 130 = 2k$. Based on

69    alignments, scaffolds apparently consisting of PhiX, chrC, or chrM were removed.

70    Genome annotation revealed several nearly identical genomic contigs with nearly

71    identical fold coverage, which appear to be due to a failure of the assembler to collapse

72    identical contigs into one contig. In an effort to resolve this artifact from the assembler,

73    we performed all versus all BLASTn for each genomic contig against the whole genome

74    and removed contigs which contained a >95% threshold for query coverage and

75    nucleotide identity. This reduced the total number of contigs from 199,501 to 116,817

76    (41.3% reduction) and total genome size from 182,854,974bp to 161,759,242bp

77    (excluding the mitochondrial and chloroplast genomes).

78        Notation for contigs is in the form of gXXXXXX_YYYYY, where XXXXXX

79    gives contig length and YYYYY gives the approximate average coverage for that contig.

80    Contigs beginning with a 'g' are genomic sequences and contigs beginning with an 'r' are

81    mRNA sequences.

82

83    *DNA Bisulfite sequencing and analysis*

84        1 μg of *C. cryptica* nuclear DNA for bisulfite treatment was resuspended in 50 μl of

85    EB buffer (QIAGEN) and sonicated in AFA-fiber microTubes using a Covaris S2

86    machine (Duty Cycle = 10%; Intensity = 5; Cycles/Burst = 200; for 6 minutes) to obtain

87    100-300 bp fragments. The DNA was then subjected to End-Repair, A-tailing and

88    Adapter Ligation using Illumina TruSeq DNA Sample Prep kit v2 according to

89    manufacturer's instructions. The Adapter-ligated DNA was purified and size-selected

90    using AMPure XP beads. DNA was then bisulfite treated using EpiTect kit (QIAGEN)

91    using the following conversion protocol: 95°C 5min, 60°C 25min, 95°C 5min, 60°C

92    85min, 95°C 5min, 60°C 175min, 95°C 5min, 60°C 25min, 95°C 5min, 60°C 85min,

93   95°C 5min, 60°C 175min. Bisulfite-treated DNA was then desulphonated according to

94   manufacturer's protocol ("Small Amount of Fragmented DNA" variant) and DNA eluted

95   twice with EB. Converted DNA was amplified with MyTaq 2x mix (Bioline): 98°C 2

96   min; 12 cycles of 98°C 15 sec, 60°C 30 sec, 72°C 30 sec; 72°C 5 min. Amplified DNA

97   was diluted to 10 nM and sequenced using Illumina HiSeq2000 (100 single end reads).

98      Bisulfite converted reads were inspected for sequencing quality using FastQC 0.10.1.

99   Reads passing the Illumina quality filter ('PF' value equal to 1) were retained and aligned

100  to the genome assembly using BS-Seeker2 in local alignment mode with Bowtie2 as the

101  aligner [2]. For the purpose of this study, a 'methylated' base pair is defined as a cytosine

102  which has >= 4 methylated reads, similar to as described in [3]. Therefore an

103  'unmethylated' cytosine is that which has >= 4 unmethylated reads. Those cytosines

104  which do not have at least 4 aligned reads from the bisulfite sequencing are considered to

105  not have enough data sufficiently conclude whether that site is methylated or not. Genes

106  that are 'methylated' are those which are defined as having >=50% of callable CG sites

107  (the most common motif for methylation) containing a 'methylated' cytosine.

108  Unmethylated genes are <50%.

109  *RNA sequencing*

110     Total RNA was purified from cultures of *C. cryptica* in log-phase growth under

111  conditions of either silicon starvation or nitrogen starvation. For RNA isolation, 750mL

112  of liquid cell culture for each time point was harvested and treated with 20mg/mL

113  cycloheximide and concentrated by filtration. Cells were stored in -80ºC prior to

114  extraction. Total RNA was extracted according to [4,5].

115     Five µg of total RNA per sample was treated with Turbo DNase (Ambion) for 30 min

116 at 37°C according to the manufacturer's instructions in order to remove any

117 contaminating DNA. The resulting RNA was purified by ethanol precipitation, and RNA

118 quality was evaluated on a BioAnalyzer RNA Nano kit (Agilent). RNAseq libraries were

119 prepared using the Illumina TruSeq mRNA Sample Prep kit (Illumina) according to

120 manufacturer's protocols (Rev. A). Sequencing was performed at the UCLA Broad Stem

121 Cell Research Center sequencing core on a HiSeq 2000 sequencer (Illumina) using a

122 mixture of 50+50 nt paired end reads and 100 nt single end reads. 11 libraries were

123 sequenced on a single lane each, 2 libraries were multiplexed onto a 1 lane. Pooled raw

124 reads were demultiplexed and all reads mapped to transcriptome data, with remaining

125 unmapped reads mapped to reference *Cyclotella* genome version cycCry0dot2 using

126 TopHat 2.0.4 allowing for two mismatches, reporting only unique mappings [6]. Bam

127 files were processed through HTSeq 0.5.3 using the "intersection-nonempty" method

128 which assigns a read to a gene only if the read overlaps with only one gene in its entirety.

129 Single-end and paired-end data were combined and all data was run through DESeq 1.8.3

130 to obtain FPKM counts [7]. A combination of 50+50 nt paired end reads and 100 nt

131 single end reads were pooled to facilitate accurate mapping of the genes.

132

133 *Genome Annotation*

134     Gene model predictions were generated from several pipelines as follows: (1)

135 FGENESH using the built-in diatom training set; (2) standalone AUGUSTUS using the

136 built-in '*Chlamydomonas reinhardtii'* training parameters; (3) standalone AUGUSTUS

137 using the 100 longest FGENESH predictions as a training set; (4) web-based

138  AUGUSTUS trained on the de-novo RNA assembly; and (5) MAKER with FGENESH,

139  AUGUSTUS, GeneMarkES analyses enabled [8-11].  All prediction software was run

140  using default settings except where noted.  Several genes were selected where

141  intron/exon boundaries were well characterized in *T. pseudonana* and used to test the

142  accuracy of the gene model predictions. AUGUSTUS predictions from set (4) that

143  overlap MAKER predictions from set (5) constitute the 'high-confidence' collection of

144  gene predictions.  For gene structure, it was determined that (4) more accurately

145  predicted start/stop sites and intron/exon boundaries. Gene set (4 and 5) were functionally

146  annotated [12].

147      Repetitive elements were identified using RepeatMasker 4.0.1 with RMBLASTN

148  2.2.27+, using the Bacillariophyta repeat library from RepBase and default settings.

149  Additional repeat masking by RepeatMasker was performed using a custom de-novo

150  library generated using RepeatModeler 1.0.7 with the NCBI BLAST engine.

151      The diatom genomes used in Figure 9 and S2, OrthoMCL, and RBH analyses are

152  *T. pseudonana* v3.0 [13], *P. tricornutum* v2.0 [14], *F. cylindrus* v1.0 [15], *P. multiseries*

153  v1.0 [16], and *T. oceanica* v.3.0 (NCBI accession numbers JP288099-JP2977110,

154  http://www.ncbi.nlm.nih.gov)

155  For reciprocal best BLAST hit (RBH) analysis, gene models from diatom genomes

156  (listed above) were aligned using BLASTp 2.2.28+ [17], with an e-value cutoff score of

157  1e-5 and a query coverage of at least 70%. RBH pairs were detected using python script

158  [18].

159      Predicted proteins from *T. pseudonana* and *C. cryptica* genomes were evaluated for

160  phylogenetic relatedness to sequences in NCBI GenBank nr (accessed October 16, 2015)

161    using the DarkHorse program version 1.5 with a threshold filter setting of 0.1 [24].

162    BLASTP alignments to GenBank nr sequences were required to cover at least 70% of

163    total query length and have e-value scores of $1e^{-5}$ or better for inclusion in this analysis.

164        For phylogenetic analysis in Figure 9 and S2, and to further investigate proteins of

165    interest, amino acid sequences were aligned using MUSCLE with 10 maximum number

166    of iterations and default parameters [19]. Trees were generated using default parameters

167    in RAxML_GUI v1.3 for 100 iterations and visualized using FigTree v1.4.0 [20]. For

168    each tree shown, bootstrap values are listed and have been midpoint rooted.

169

170    *Subcellular Localization Prediction*

171        All open reading frames were analyzed using several computational tools for

172    predicting the likely location of the protein product within the cell.  Nucpred [21] and

173    NetNES [22] provided estimations of nuclear localization for each putative protein.

174    PREDOTAR [23], PSORT [24] , HECTAR [25], CELLO [26] , PredAlgo [27], SignalP

175    3.0 [28] TargetP, ChloroP  [29] and ASAFind [30] provided possible localizations to

176    mitochondria, chloroplast, endoplasmic reticulum, or plastid. The presence of N-terminal

177    signal peptides and transmembrane regions were assessed using PHOBIUS [31,32], the

178    SignalP 3.0 portion of TargetP and PredAlgo. Specific sequences and predicted cleavage

179    sites for signal peptides were taken from SignalP 3.0.  All open reading frames from *C.*

180    *cryptica* were submitted to web-servers in an automated fashion using scripts for html or

181    webmail submission.

182        Methods for targeting predictions shown in Figure 4 are as follows: for plastid

183    targeting, all proteins with predicted positive ASAfind plastid targeting were surveyed for

184     proper cleavage site using SignalP 3.0 [28, 30]. Predicted SignalP 3.0 cleavage sites were

185     then defined by the guidelines addressed in Figure 5 of [33] and split into canonical

186     plastid (AF, GF, and SF cleavage sites), noncanonical plastid (AW, AY, AL, GW, GY,

187     GL, SW, SY, SL), periplastid and other (AH, AI, AM, AR, AE, AG, SH, SI, SM, SR, SE,

188     SG, GH, GI, GM, GR, GE, GG). Periplastid predicted cleavage sites which did not

189     contain a positive ChloroP prediction were then categorized as ER/secreted. Any negative

190     plastid ASAfind result but SignalP 3.0 positive with a predicted periplastid cleavage site,

191     and positive ChloroP were also classified as periplastid. ER/secreted proteins were then

192     defined as proteins with negative plastid ASAfind prediction, positive SignalP 3.0, with a

193     cleavage site not one that is specified in [33].

194          Mitochondrially-targeted proteins were classified according to any of the

195     following prediction combinations for HECTAR, Predotar, and TargetP (listed in

196     respective order): Type II, Mito, Mito; SigP, Mito, Mito; Mito, Plastid, Mito; Mito, Mito,

197     SigP; Mito, Mito, Mito; Mito, Mito, Plastid; Mito, Mito, no prediction; Mito, no

198     prediction, Mito; no prediction, Mito, Mito. MitoProt was also used for all predicted

199     proteins shown in Figure 5-8 [34]. All other proteins that did not fall into the above

200     criterion were classified as 'cytosol and other.'

201

202     *Vector Construction and Diatom Transformation*

203          The vector utilizing the *C. cryptica* rpL41 promoter and terminator sequences was

204     assembled using the GeneArt® Seamless Cloning and Assembly Kit (Invitrogen), with a

205     Gateway™ vector (Invitrogen) as the backbone. The construction of the *T. pseudonana*

206     *fcp* [35] and nitrate reductase (NR) [36,37]  vectors is described elsewhere.

207        Each vector was co-transformed with another vector expressing the *nat1* gene

208        (received from N. Kroger, Germany), which confers resistance to the antibiotic

209        nourseothricin, under control of the acetyl coenzyme A carboxylase promoter. *C. cryptica*

210        was transformed using tungsten microparticle bombardment as described elsewhere

211        (Shrestha and Hildebrand, 2014), with the following changes: $1 \times 10^7$ cells in exponential

212        phase were spread onto artificial sea water (ASW) containing 1.5% Bacto Agar (Becton

213        Dickinson, USA) and no antibiotics [38]. Immediately after bombardment, 10ml ASW

214        was added to the plates, which were incubated in low light for 18 hours.  Following the

215        incubation, cells were washed off the plates using the ASW they were incubated with,

216        and the entire volume was transferred to 125ml Erlenmyer flasks with 50ml ASW with

217        100µg/ml nourseothricin (clonNAT; Werner BioAgents, Germany). After 7 days of

218        growth with shaking, 10ml of the cultures were transferred to 125ml Erlenmeyer flasks

219        with fresh 50ml ASW plus 100µg/ml nourseothricin and grown to exponential phase. The

220        highest GFP expressors were selected using a sorting flow cytometer (Influx, Becton

221        Dickinson, USA), recovered in liquid ASW for 2 days, then plated on ASW agar plates

222        with 100µg/ml nourseothricin. Colonies were picked and transferred to 24 well plates,

223        each well containing 2 ml ASW with 100µg/ml nourseothricin.

224

225        *Assessment of Conditional Expression Using the Nitrate Reductase Promoter*

226        Three independent clones were grown to exponential phase in 24 well plates, as

227        described above. 2µl of culture was transferred to two new wells per clone. One of the

228        wells contained 2ml ASW, which would allow for expression of GFP under the NR

229        promoter. The other well contained 2ml modified ASW, in which ammonia replaced

230 nitrate as the nitrogen source, which would repress the expression of GFP under the NR

231 promoter [36]. The cultures were allowed to reach exponential phase and imaged using

232 Zeiss Axio Observer Z1 Inverted Microscope (Zeiss Axioscope, Carl Zeiss Microimaging

233 Inc., USA). The filter set for fluorescent imaging was Emission LP 515 nm for

234 chlorophyll autofluorescence, and Zeiss #38HE Excitation BP 470/40 nm, Dichromatic

235 mirror FT 495 nm, Emission BP 525/50 nm for GFP.

236

237

238

239

240
241
242 1. Jacobs JD, Ludwig JR, Hildebrand M, Kukel A, Feng T-Y, Ord RW, et al.
243 Characterization of two circular plasmids from the marine diatom Cylindrotheca
244 fusiformis: plasmids hybridize to chloroplast and nuclear DNA. Molec. Gen. Genet.
245 1992;233:302–10.

246 2. Guo W, Fiziev P, Yan W, Cokus S, Sun X, Zhang MQ, et al. BS-Seeker2: a versatile
247 aligning pipeline for bisulfite sequencing data. BMC Genomics. 2013;14:774.

248 3. Lopez DA, Hamaji T, Kropat J, De Hoff P, Morselli M, Rubbi L, et al. Dynamic
249 changes in the transcriptome and methylome of Chlamydomonas reinhardtii throughout
250 its life cycle. plantphysiol.org. American Society of Plant Biologists;
251 2015;169:00861.2015–743.

252 4. Smith SR, Glé C, Abbriano RM, Traller JC, Davis A, Trentacoste E, et al. Transcript
253 level coordination of carbon pathways during silicon starvation-induced lipid
254 accumulation in the diatom Thalassiosira pseudonana. New Phytol. 2016;210:810-904.

255 5. Hildebrand M, Dahlin K. Nitrate transporter genes from the diatom Cylindrotheca
256 fusiformis (Bacillariophyceae): mRNA levels controlled by nitrogen source and by the
257 cell cycle. J. Phycol. 2000;36:702-13.

258 6. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R. TopHat2: accurate alignment of
259 transcriptomes in the presence of insertions, deletions and gene fusions. Genome Biol.
260 2013;14:R36.

261 7. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for
262 RNA-seq data with DESeq2. Genome Biol. 2014;15:550.

263 8. Stanke M, Morgenstern B. AUGUSTUS: a web server for gene prediction in
264 eukaryotes that allows user-defined constraints. Nucleic Acids Res. 2005;33:W465–7.

265 9. Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, et al. MAKER: An easy-
266 to-use annotation pipeline designed for emerging model organism genomes. Genome
267 Res. 2007;18:188–96.

268 10. Salamov AA. Ab initio Gene Finding in Drosophila Genomic DNA. Genome Res.
269 2000;10:516–22.

270 11. Borodovsky M, Lomsadze A. Eukaryotic gene prediction using GeneMark.hmm-E
271 and GeneMark-ES. Curr Protoc Bioinformatics. 2011;Chapter 4:Unit4.6.1–10.

272 12. Lopez D, Casero D, Cokus SJ, Merchant SS, Pellegrini M. Algal Functional
273 Annotation Tool: a web-based analysis suite to functionally interpret large gene lists
274 using integrated annotation and expression data. BMC Bioinformatics. 2011;12:282.

275 13. Joint Genome Institute. http://genome.jgi.doe.gov/Thaps3/Thaps3.home.html

276 14. Joint Genome Institute. http://genome.jgi.doe.gov/Phatr2/Phatr2.home.html

277 15. Joint Genome Institute. http://genome.jgi.doe.gov/Fracy1/Fracy1.home.html

278 16. Joint Genome Institute. http://genome.jgi.doe.gov/Psemu1/Psemu1.home.html

279 17. Altschul SF, Gish W, Miller W, Myers EW. Basic local alignment search tool. J. Mol.
280 Biol. 1990;215:403-10.

281 18. Qin H. hongqin/Simple-reciprocal-best-blast-hit-pairs.

282 19. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high
283 throughput. Nucleic Acids Research. 2004;32:1792–7.

284 20. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses
285 with thousands of taxa and mixed models. Bioinformatics. 2006;22:2688–90.

286 21. Brameier M, Krings A, MacCallum RM. NucPred Predicting nuclear localization of
287 proteins. Bioinformatics. 2007;23:1159–60.

288 22. La Cour T, Kiemer L, Mølgaard A, Gupta R, Skriver K, Brunak S. Analysis and
289 prediction of leucine-rich nuclear export signals. Protein Engineering Design and
290 Selection. 2004;17:527–36.

291 23. Small I, Peeters N, Legeai F, Lurin C. Predotar: A tool for rapidly screening
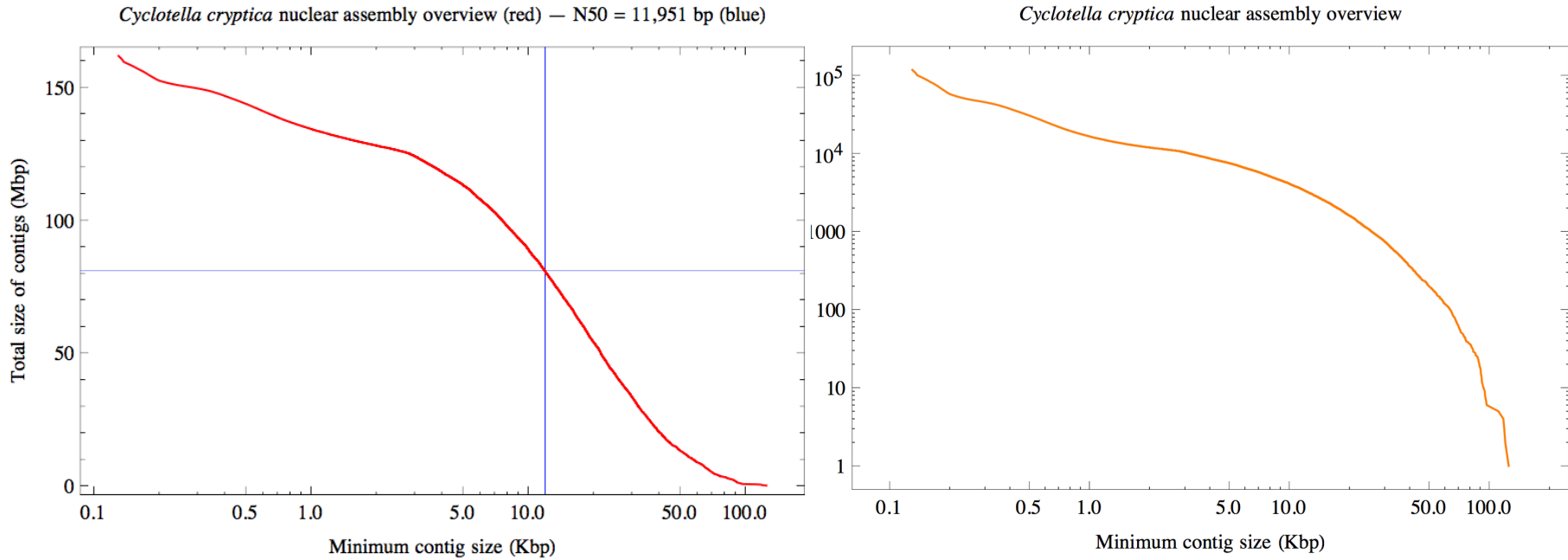292 proteomes forN-terminal targeting sequences. Proteomics. 2004;4:1581–90.

293    24. Horton P, Park KJ, Obayashi T, Fujita N, Harada H, Adams-Collier CJ, et al. WoLF
294    PSORT: protein localization predictor. Nucleic Acids Res. 2007;35:W585–7.

295    25. Gschloessl B, Guermeur Y, Cock JM. HECTAR: A method to predict subcellular
296    targeting in heterokonts. BMC Bioinformatics. 2008;9:393.

297    26. Yu C-S, Lin CJ, Hwang J-K. Predicting subcellular localization of proteins for Gram-
298    negative bacteria by support vector machines based on n-peptide compositions. Protein
299    Science. Cold Spring Harbor Laboratory Press; 2004;13:1402–6.

300    27. Tardif M, Atteia A, Specht M, Cogne G, Rolland N, Brugiere S, et al. PredAlgo: A
301    New Subcellular Localization Prediction Tool Dedicated to Green Algae. Mol Biol Evol.
302    2012;29:3625–39.

303    28. Dyrløv Bendtsen J, Nielsen H, Heijne von G, Brunak S. Improved Prediction of
304    Signal Peptides: SignalP 3.0. J Mol Biol. 2004;340:783–95.

305    29. Emanuelsson O, Brunak S, Heijne von G, Nielsen H. Locating proteins in the cell
306    using TargetP, SignalP and related tools. Nat Protoc. 2007;2:953–71.

307    30. Gruber A, Rocap G, Kroth PG, Armbrust EV, Mock T. Plastid proteome prediction
308    for diatoms and other algae with secondary plastids of the red lineage. Plant J.
309    2015;81:519–28.

310    31. Laws EA, Taguchi S, Hirata J, Pang L. Optimization of microalgal production in a
311    shallow outdoor flume. Biotechnol. Bioeng. 1988;32:140–7.

312    32. Käll L, Krogh A, Sonnhammer E. A Combined Transmembrane Topology and Signal
313    Peptide Prediction Method. J Mol Biol. 2004;338:1027–36.

314    33. Gruber A, Vugrinec S, Hempel F, Gould SB, Maier U-G, Kroth PG. Protein targeting
315    into complex diatom plastids: functional characterisation of a specific targeting motif.
316    Plant Mol Biol. 2007;64:519–30.

317    34. Claros MG, Vincens P. Computational method to predict mitochondrially imported
318    proteins and their targeting sequences. Eur J Biochem. 1996.

319    35. Shrestha RP, Hildebrand M. Evidence for a Regulatory Role of Diatom Silicon
320    Transporters in Cellular Silicon Responses. Eukaryot Cell. 2014;14:29–40.

321    36. Poulsen N, Chesley PM, Kröger N. Molecular genetic manipulation of the diatom
322    Thalassiosira pseudonana (Bacillariophyceae). J. Phycol. 2006;42:1059–65.

323    37. Trentacoste EM, Shrestha RP, Smith SR, Glé C, Hartmann AC, Hildebrand M, et al.
324    Metabolic engineering of lipid catabolism increases microalgal lipid accumulation
325    without compromising growth. Proc Nat Acad Sci. 2013;110:19748–53.

326    38. Darley WM, Volcani BE. Role of silicon in diatom metabolism. Exp Cell Res.
327    1969;58:334–42.

328

Figure S1: (a) *Cyclotella cryptica* nuclear assembly overview broken down by minimum contig size (Kbp) versus total Mbp of all contigs (left) and minimum contig size versus the total number of contigs (right) and (b) genomic sequencing data.

a.



b.

| Genomic Libraries | Number of Raw Reads for Paired End Libraries | | | | Percentage passing pairs | | | | Percentage passing pairs after removal of pairs with 'N' |
|---|---|---|---|---|---|---|---|---|---|
| | Lane 1 | Lane 2 | Lane 3 | Total Reads | Lane 1 | Lane 2 | Lane 3 | Average | Average |
| Paired End Short | 65270867 | 17592293 | | 82863160 | 4.8 | 57 | | 71.1 | 73.0 |
| Paired End Medium | 65333130 | 66235838 | | 131568968 | 68.5 | 69.5 | | 60.0 | 79.6 |
| Paired End Long | 69147985 | 69578501 | 65179466 | 203905952 | 27.8 | 34 | 65.6 | 42.5 | 57.9 |
| Paired End Control | 30301553 | | | | 92.8 | | | | 99.5 |
| Mate Pair Short | 142455072 | | | | 96.9 | | | | 98.9 |
| Mate Pair Long | 154107079 | | | | 96.7 | | | | 98.6 |

Supplementary table S1. Statistics from different gene model prediction pipelines. The pipelines are: (1) FGENESH Gene predictions, Diatom training set, (2) Augustus Gene predictions V1, *Chlamydomonas* training set, no RNAseq data, (3) Augustus Gene predictions V2, FGENESH100 training set, RNAseq data (4) Augustus Gene predictions V3, 'self' trained, RNAseq data (5) MAKER Gene predictions, (Augustus self trained + FGENESH + GeneMarkES). Details are presented in Supplementary Methods. Data presented includes all predicted models regardless of read counts from RNAseq data and prior to removing apparent duplicate contigs (Additional File 1: Supplementary Methods, Genome Library Construction and Sequencing).

| | Gene Model Prediction Pipeline | | | | |
| --- | --- | --- | --- | --- | --- |
| | CcFgenesh | CcAugustusV1 | CcAugustusV2 | CcAugustusV3 | CcMAKERV1 |
| Total Models | 50,288 | 6,295 | 24,819 | 33,682 | 9,049 |
| Average Model Length (bp) | 1,561.8 | 926.1 | 1,746.2 | 1,265 | 1,999.7 |
| Average Exons | 3.72 | 1.22 | 2.65 | 1.95 | 3.69 |
| Avg. Exon Length (bp) | 247.4 | 115.8 | 561.5 | 585.6 | 457.6 |
| Avg. # Introns | 2.72 | 0.22 | 1.65 | 0.95 | 2.69 |
| Avg. Intron Length (bp) | 129.8 | 542 | 153.9 | 128.1 | 128.1 |

# Table S2: Repeat sequences in *C. cryptica* identified using (a) RepBase data and (b) RepeatModeler Data.

## a.

Repeat sequences in *Cyclotella cryptica* using RepBase data
Sequences:        199501
total length:  182854974 bp  (174198679 bp excl N/X-runs)
GC level: 43.01%
Bases masked:    7037708 bp ( 3.85 %)

| | Number of elements* | Length occupied (bp) | Percentage of sequence |
|---|---|---|---|
| **Retroelements** | 12264 | 4694900 | 2.57 |
| **SINEs:** | 0 | 0 | 0 |
| **Penelope** | 13 | 1650 | 0 |
| **LINEs:** | 58 | 4854 | 0 |
| CRE/SLACS | 0 | 0 | 0 |
| L2/CR1/Rex | 0 | 0 | 0 |
| R1/LOA/Jockey | 0 | 0 | 0 |
| R2/R4/NeSL | 0 | 0 | 0 |
| RTE/Bov-B | 2 | 100 | 0 |
| L1/CIN4 | 0 | 0 | 0 |
| **LTR elements:** | 12206 | 4690046 | 2.56 |
| BEL/Pao | 0 | 0 | 0 |
| Ty1/Copia | 2901 | 1432661 | 0.78 |
| Gypsy/DIRS1 | 9305 | 3257385 | 1.78 |
| Retroviral | 0 | 0 | 0 |
| | | | |
| **DNA transposons** | 346 | 113716 | 0.06 |
| hobo-Activator | 0 | 0 | 0 |
| Tc1-IS630-Pogo | 0 | 0 | 0 |
| En-Spm | 0 | 0 | 0 |
| MuDR-IS905 | 0 | 0 | 0 |
| PiggyBac | 2 | 103 | 0 |
| Tourist/Harbinger | 305 | 105959 | 0.06 |
| Other (Mirage, P-element, Transib) | 0 | 0 | 0 |
| | | | |
| **Rolling-circles** | 0 | 0 | 0 |
| **Unclassified:** | 28 | 6197 | 0 |
| **Total interspersed repeats:** | | 4814813 | 2.63 |
| **Small RNA:** | 142 | 17166 | 0.01 |
| **Satellites:** | 0 | 0 | 0 |
| **Simple repeats:** | 32133 | 1928189 | 1.05 |
| **Low complexity:** | 3624 | 294614 | 0.16 |

## b.

Repeat sequences in *Cyclotella cryptica* using RepeatModeler data
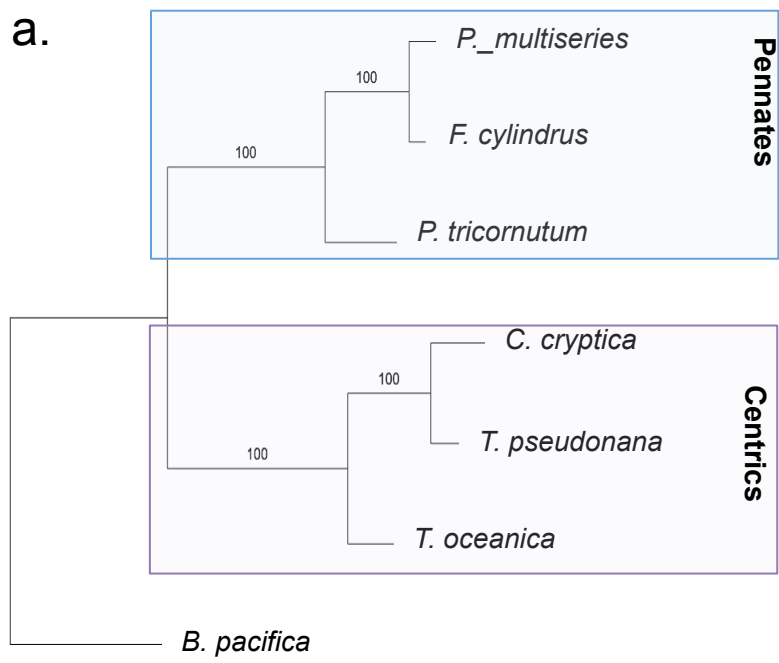Sequences:        199501
Total length:  182854974 bp  (174198679 bp excl N/X-runs)
GC level:        43.01 %
Bases masked:  98288109 bp ( 53.75 %)

| | Number of elements* | Length occupied (bp) | Percentage of sequence |
|---|---|---|---|
| **SINEs:** | 0 | 0 | 0 |
| ALUs | 0 | 0 | 0 |
| MIRs | 0 | 0 | 0 |
| **LINEs:** | 2626 | 1877047 | 1.03 |
| LINE1 | 0 | 0 | 0 |
| LINE2 | 0 | 0 | 0 |
| L3/CR1 | 0 | 0 | 0 |
| **LTR elements:** | 43176 | 15802661 | 8.64 |
| ERVL | 0 | 0 | 0 |
| ERVL-MaLRs | 0 | 0 | 0 |
| ERV_classI | 0 | 0 | 0 |
| ERV_classII | 0 | 0 | 0 |
| **DNA elements:** | 15402 | 5899536 | 3.23 |
| hAT-Charlie | 0 | 0 | 0 |
| TcMar-Tigger | 155 | 44023 | 0.02 |
| **Unclassified:** | 314059 | 73067346 | 39.96 |
| | | | |
| **Total interspersed repeats:** | | 96646590 | 52.85 |
| **Small RNA:** | 0 | 0 | 0 |
| **Satellites:** | 0 | 0 | 0 |
| **Simple repeats:** | 23628 | 1861243 | 1.02 |
| **Low complexity:** | 1672 | 129011 | 0.07 |

Figure S2: Phylogenetic comparison of diatom species with sequenced genomes. (a) 18S sequence comparison from [89], accession numbers are listed in Additional File 4 (b) Reciprocal best blast hits comparison. *C. cryptica* and *T. pseudonana* are the most similarly related centric diatoms with available genomic data.

a.



b.

| | Total Reciprocal BLAST Hits | Average % Identity |
|---|---|---|
| C. cryptica/T. pseudonana | 5498 | 63.6 |
| C. cryptica/T. oceanica | 4443 | 56.4 |
| T. pseudonana/T. oceanica | 3592 | 59.2 |
| P. tricornutum/F. cylindrus | 7060 | 55.5 |
| P. tricornutum/P. multiseries | 2878 | 54.2 |
| F. cylindrus/P. multiseries | 7773 | 66.0 |

Figure S3: Per-cytosine fraction methylation. (a) 0 hour, silicon-replete sample (b) 48 hour, silicon-deplete sample. All sites shown have a read coverage greater than or equal to 4.



a.

**Per-Cytosine Fraction Methylation
(Coverage >= 4) -- 0h**
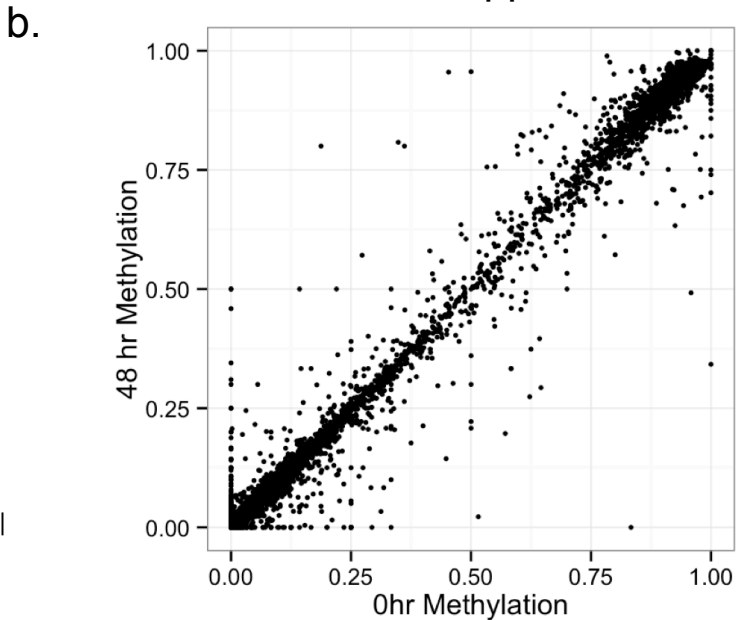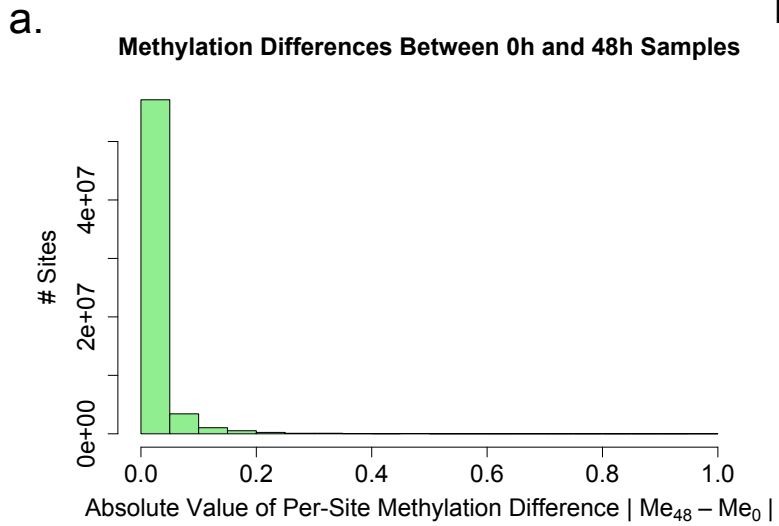
b.

**Per–Cytosine Fraction Methylation
(Coverage >= 4) –– 48h**

Figure S4: Per-site methylation differences between Silicon replete and deplete conditions. (a) Absolute value of the difference between 48h and 0h. (b) Slope between fraction methylation of genes in when comparing the two conditions is linear, outliers are apparent.

a.

**Methylation Differences Between 0h and 48h Samples**

b.

Figure S5: Per-site methylation differences between silicon replete and deplete conditions. (a) Histogram showing absolute value of the difference between 48h and 0h. (b) Scatter plot showing the slope between fraction methylation of genes in when comparing the two conditions is linear, outliers are apparent.
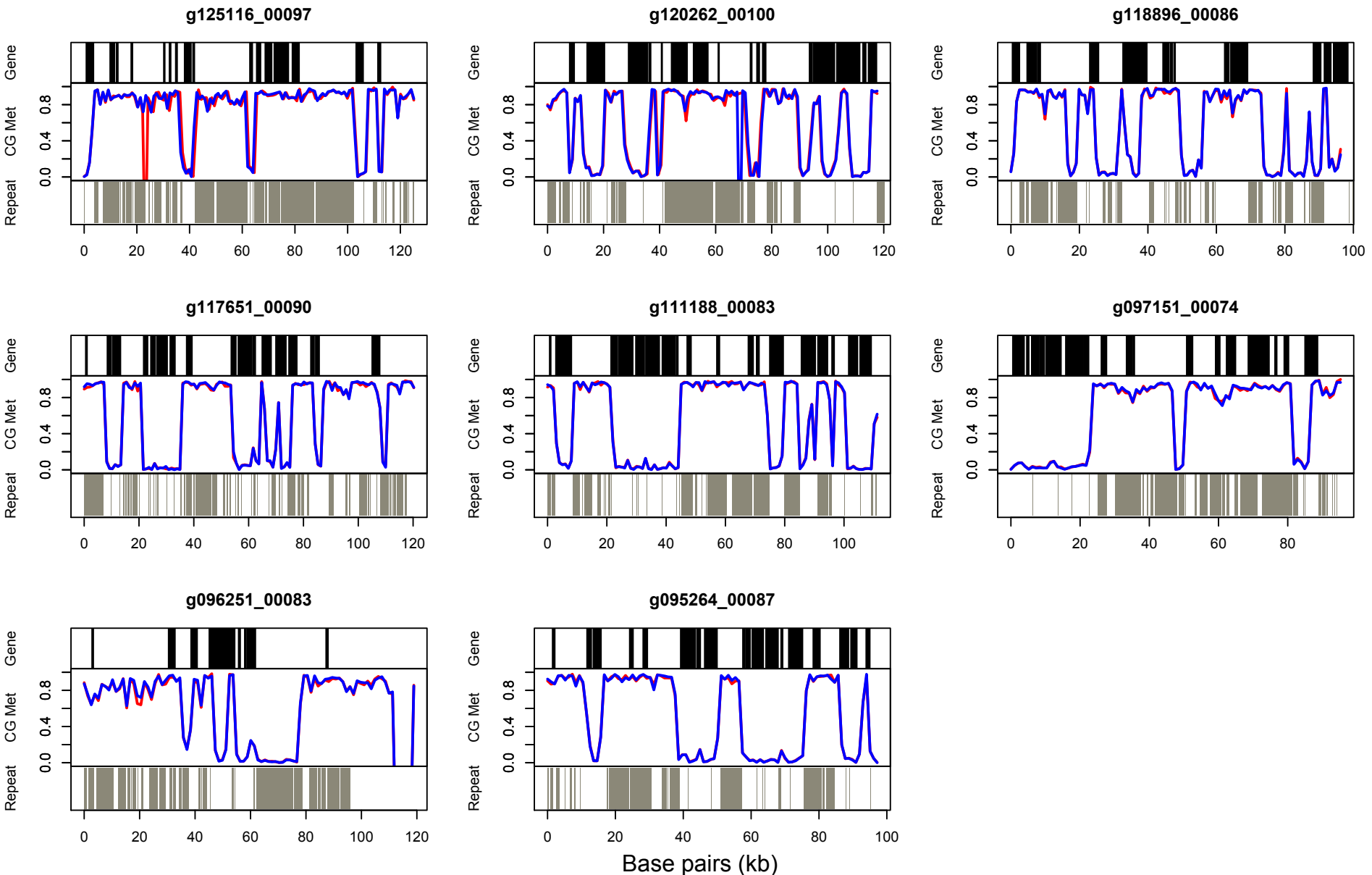
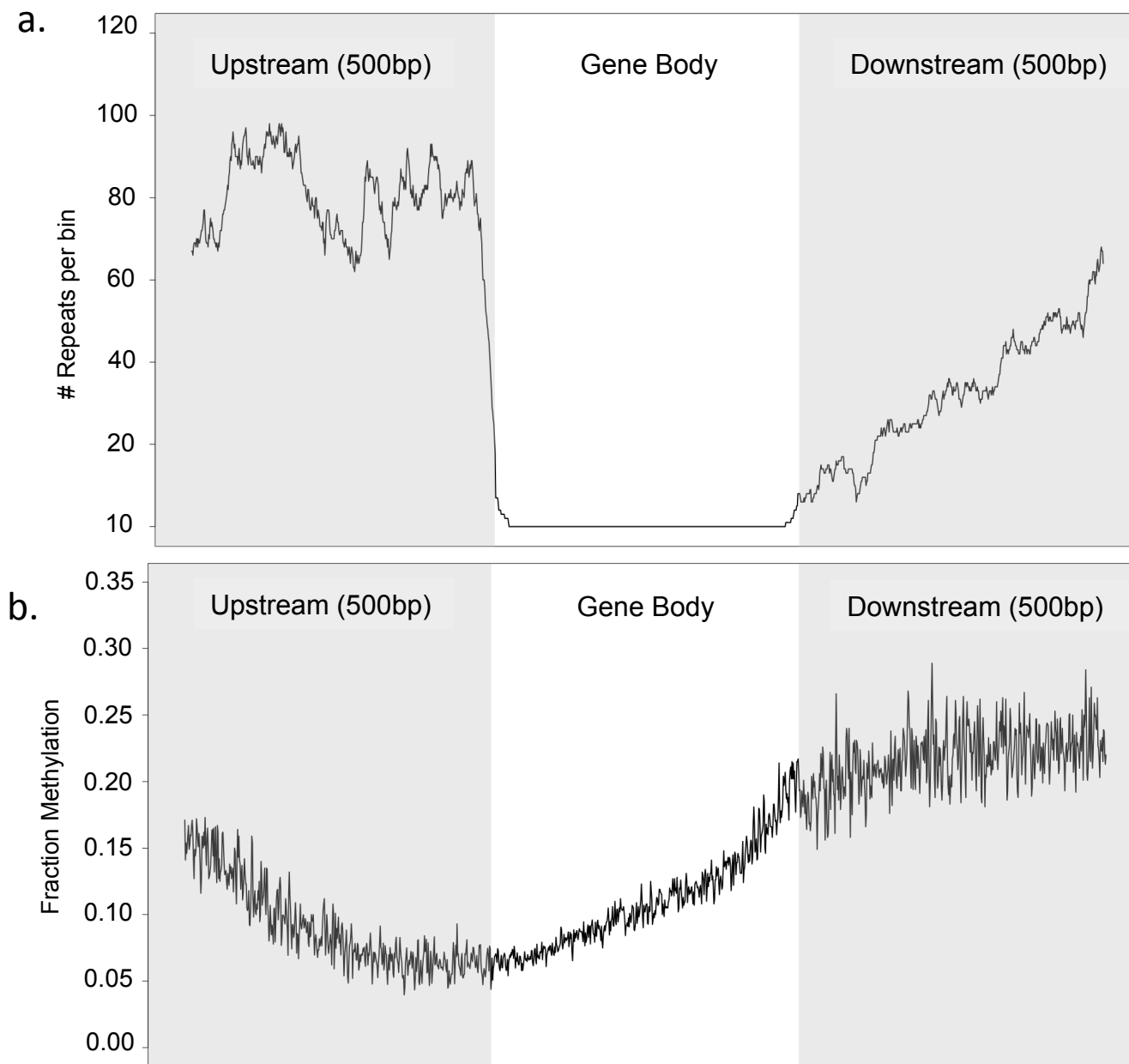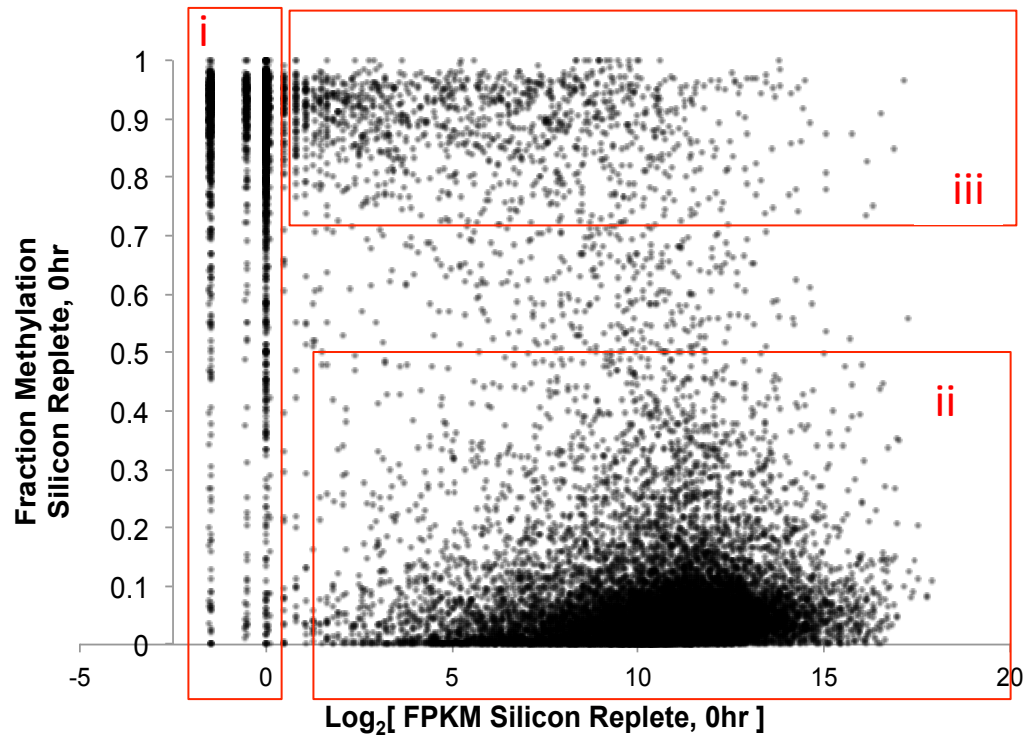Figure S6: Binned repeats (a) and fraction methylation across gene body (b) in *C. cryptica.*

Table S3: Summary of gene methylation (AUGUSTUS V3 models) in both experimental conditions.

| Silicon Replete 0hr | | | | | |
|---|---|---|---|---|---|
| | Count | Percent total of all genes | Percent total Excluding genes with insufficient coverage | Average Fraction Methylation | Average FPKM |
| Methylated genes | 4170 | 20 | 22 | 0.87 | 536 |
| Unmethylated genes | 14866 | 70 | 78 | 0.07 | 4593 |
| No Methylation Data | 2085 | 10 | - | ND | 1767 |
| Total genes | 21121 | | | 0.24 | 3520 |
| Silicon Deplete 48hr | | | | | |
| | Count | Percent total of all genes | Percent total Excluding genes with insufficient coverage | Average Fraction Methylation | Average FPKM |
| Methylated genes | 2627 | 12 | 24 | 0.88 | 523 |
| Unmethylated genes | 8453 | 40 | 76 | 0.07 | 4589 |
| No Methylation Data | 10041 | 48 | - | ND | 3641 |
| Total genes | 21121 | | | 0.26 | 3408 |

| | Methylated 48hr | Unmethylated 48hr |
|---|---|---|
| Methylated 0hr | 2584 | 31 |
| Unmethylated 0hr | 17 | 8165 |

Figure S7: Gene methylation relative to gene expression. (a) 3 populations emerged (Red boxes i-iii) when comparing average fraction methylation to Log2 FPKM. (b) Gene count binned according to FPKM and shaded according to methylation status. Line depicts the proportion of methylated genes per bin. Data shown is for silicon replete, 0 hour condition.
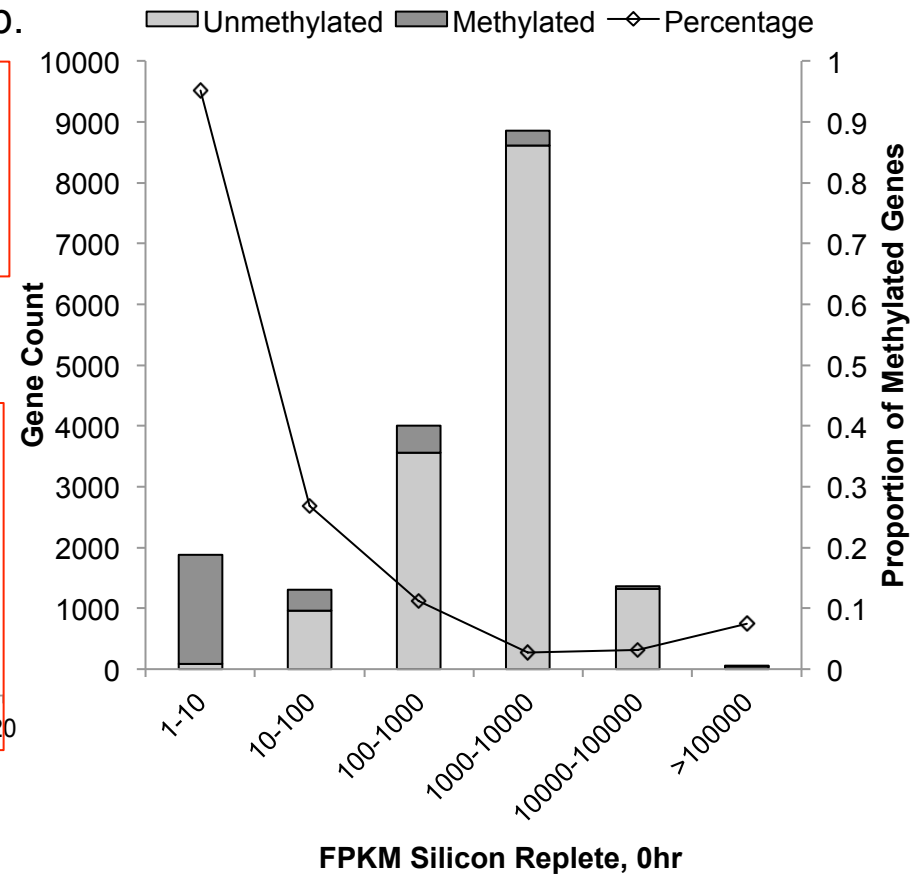
# Figure S8: Vector map for rpL41 construct for *C. cryptica.*