

Re-annotation, improved large-scale assembly and establishment of a catalogue of noncoding loci for the genome of the model brown alga *Ectocarpus*

Alexandre Cormier, Komlan Avia, Lieven Sterck, Thomas Derrien, Valentin Wucher, Gwendoline Andres, Misharl Monsoor, Olivier Godfroy, Agnieszka Lipinska, Marie-Mathilde Perrineau, et al.

► **To cite this version:**

Alexandre Cormier, Komlan Avia, Lieven Sterck, Thomas Derrien, Valentin Wucher, et al.. Re-annotation, improved large-scale assembly and establishment of a catalogue of noncoding loci for the genome of the model brown alga *Ectocarpus*. *New Phytologist*, Wiley, 2016, 214 (1), pp.219-232. <10.1111/nph.14321>. <hal-01402123>

HAL Id: hal-01402123

<http://hal.upmc.fr/hal-01402123>

Submitted on 24 Nov 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Re-annotation, improved large-scale assembly and establishment**
2 **of a catalogue of non-coding loci for the genome of the model**
3 **brown alga *Ectocarpus***

4
5 **Alexandre Cormier¹, Komlan Avia¹, Lieven Sterck^{2,3,4}, Thomas Derrien⁵, Valentin**
6 **Wucher⁵, Gwendoline Andres⁶, Misharl Monsoor⁶, Olivier Godfroy¹, Agnieszka**
7 **Lipinska¹, Marie-Mathilde Perrineau¹, Yves Van De Peer^{2,3,4,7}, Christophe Hitte⁵,**
8 **Erwan Corre⁶, Susana M. Coelho¹, J. Mark Cock^{1*}**

9
10 ¹Sorbonne Université, UPMC Univ Paris 06, CNRS, Algal Genetics Group, UMR 8227,
11 Integrative Biology of Marine Models, Station Biologique de Roscoff, CS 90074, F-29688,
12 Roscoff, France, ²Department of Plant Systems Biology, VIB, Ghent, Belgium, ³Department of
13 Plant Biotechnology and Bioinformatics, Ghent University, Ghent, Belgium, ⁴Bioinformatics
14 Institute Ghent, Technologiepark 927, 9052 Ghent, Belgium, ⁵IGDR CNRS-UMR6290 –
15 Université Rennes 1, Rennes, France, ⁶Abims Platform, CNRS-UPMC, FR2424, Station
16 Biologique de Roscoff, CS 90074, 29688 Roscoff, France, ⁷Department of Genetics, Genomics
17 Research Institute, University of Pretoria, Pretoria, South Africa.

18
19 *Author for correspondence: Tel: 33 (0)2 98 29 23 60; Email: cock@sb-roscoff.fr

20
21 **Key words:** Alternative splicing, Brown algae, *Ectocarpus*, Genetic markers, Genome
22 reannotation, Long non-coding RNAs, *Saccharina japonica*, Stramenopile

23
24 **Summary**

25 • The genome of the filamentous brown alga *Ectocarpus* was the first to be completely
26 sequenced from within the brown algal group and has served as a key reference genome both
27 for this lineage and for the stramenopiles.

28 • We present a complete structural and functional reannotation of the *Ectocarpus* genome.

29 • The large-scale assembly of the *Ectocarpus* genome was significantly improved and genome-
30 wide gene re-annotation using extensive RNA-seq data improved the structure of 11,108
31 existing protein-coding genes and added 2,030 new loci. A genome-wide analysis of splicing
32 isoforms identified an average of 1.6 transcripts per locus. A large number of previously

33 undescribed non-coding genes were identified and annotated, including 717 loci that produce
34 long non-coding RNAs. Conservation of lncRNAs between *Ectocarpus* and another brown
35 alga, the kelp *Saccharina japonica*, suggests that at least a proportion of these loci serve a
36 function. Finally, a large collection of SNP-based markers was developed for genetic analyses.
37 These resources are available through an updated and improved genome database.

- 38 • This study significantly improves the utility of the *Ectocarpus* genome as a high-quality
39 reference for the study of many important aspects of brown algal biology and as a reference for
40 genomic analyses across the stramenopiles.

41

42 Introduction

43 *Ectocarpus* has been studied since the nineteenth century and work on this organism has
44 provided many insights into novel aspects of brown algal biology (Müller, 1967; Charrier *et al.*,
45 2008). This long research history, together with several features of the organism that make
46 it well adapted for genetic and genomic approaches (Coelho *et al.*, 2012a), led to it being
47 proposed as a general model organism for the brown algae in 2004 (Peters *et al.*, 2004) and to
48 the initiation of a genome sequencing project that produced a first complete genome assembly
49 in 2010 (Cock *et al.*, 2010). The publication of the genomic sequence was followed up with
50 the development of many additional tools and resources including a genetic map (Heesch *et al.*,
51 2010), gene mapping techniques, microarrays (Dittami *et al.*, 2009; Coelho *et al.*, 2011),
52 transcriptomic data (Ahmed *et al.*, 2014; Lipinska *et al.*, 2015), proteomic techniques (Ritter
53 *et al.*, 2008) and bioinformatics tools (Gschloessl *et al.*, 2008; Prigent *et al.*, 2014). These
54 genomic resources are currently being exploited to further our understanding of a broad range
55 of processes, including life cycle regulation (Coelho *et al.*, 2011), sex determination (Lipinska
56 *et al.*, 2013, 2015; Ahmed *et al.*, 2014), development and morphology (Le Bail *et al.*, 2011),
57 interactions with pathogens (Zambounis *et al.*, 2012) and metabolism (Meslet-Cladière *et al.*,
58 2013; Prigent *et al.*, 2014).

59 The brown algae are an important taxonomic group for several reasons; they are key primary
60 producers in many coastal ecosystems and have a major influence on marine biodiversity and
61 ecology (Dayton, 1985; Steneck *et al.*, 2002; Bartsch *et al.*, 2008; Klinger, 2015; Wahl *et al.*,
62 2015). Brown algae also represent an important resource of considerable commercial value
63 (Kijjoo & Sawangwong, 2004; Smit, 2004; Hughes *et al.*, 2012) and industrial exploitation of
64 these organisms has increased markedly in recent years with the expansion of aquaculture
65 activities, particularly in Asia (Tseng, 2001). Finally, brown algae are also of phylogenetic

66 interest because they are very distantly related to well-studied groups such as the animals, fungi
67 and land plants and, moreover, have evolved complex multicellularity independently of these
68 other lineages (Cock *et al.*, 2010; Cock & Collén, 2015). Comparative analyses between brown
69 algae and members of these other eukaryotic supergroups therefore potentially provide a means
70 to explore deep evolutionary events of broad, general importance.

71 A high-quality genome resource is essential if these important features of the brown algae
72 are to be investigated effectively. The version of the *Ectocarpus* genome that was published in
73 2010 (Cock *et al.*, 2010) included detailed manual annotations of many of the genes but gene
74 structure predictions were based on a limited amount of transcriptomic data (Sanger expressed
75 sequence tags) and the large-scale assembly of the sequence contigs only associated about 70%
76 of the genome sequence with linkage groups. Moreover, annotation efforts had focused almost
77 exclusively on protein-coding genes, largely ignoring the non-coding component of the
78 genome. The study described here set out to address these shortfalls, exploiting the large
79 amount of transcriptomic data now available and using recently developed genetic and
80 bioinformatic approaches to improve both the assembly and annotation of the genome. A high-
81 density, RAD-seq-based genetic map was used to anchor sequence scaffolds onto the
82 chromosomes, considerably improving the large-scale assembly of the genome. In addition, a
83 complete reannotation of the genome was carried out based on extensive RNA-seq data. This
84 updated version of the genome annotation includes information about transcript isoforms and
85 integrates non-coding loci such as microRNAs (miRNAs) and long non-coding RNAs
86 (lncRNAs). Finally, we report additional resources including a genome-wide set of single
87 nucleotide polymorphisms for genetic mapping and improvements to the genome database
88 such as the addition of a JBrowse-based genome browser that allows multiple types of genome-
89 wide data to be visualised simultaneously.

90

91 **Materials and Methods**

92 **Biological material**

93 *Ectocarpus* strains were cultured as described previously (Coelho *et al.*, 2012b). The male
94 genome sequenced strain Ec32 (reference CCAP 1310/4 in the Culture Collection of Algae and
95 Protozoa, Oban, Scotland) is a meiotic offspring of a field sporophyte, Ec17, collected in 1988
96 in San Juan de Marcona, Peru (Peters *et al.*, 2008). Ec722 is a UV-mutagenised descendant of
97 Ec32. The female outcrossing line Ec568 is derived from a sporophyte collected in Arica in
98 northern Chile (Heesch *et al.*, 2010).

99

100 RNA-seq

101 The analyses carried out in this study used RNA-seq data generated for biological replicate
102 (duplicate) samples of partheno-sporophytes and of both young and mature samples for both
103 male and female gametophytes (ten libraries in all). The production of the young (Lipinska *et*
104 *al.*, 2015) and mature (Ahmed *et al.*, 2014) gametophyte RNA-seq data has been described
105 previously. For each of the replicate partheno-sporophyte samples, total RNA was extracted
106 and used as a template by Fasteq (CH-1228 Plan-les-Ouates, Switzerland) to synthesise cDNA
107 using an oligo-dT primer. The cDNA libraries were sequenced with Illumina HiSeq 2000
108 technology to generate 100 bp single-end reads. Data quality was assessed using the FASTX
109 toolkit (http://hannonlab.cshl.edu/fastx_toolkit/index.html) and the reads were trimmed and
110 filtered using a quality threshold of 25 (base calling) and a minimal size of 60 bp. Only reads
111 in which more than 75% of nucleotides had a minimal quality threshold of 20 were retained.
112 Table S1 shows the number of raw reads generated per sample and the number of reads
113 remaining after trimming and filtering (cleaned reads). The cleaned reads were mapped to the
114 *Ectocarpus* sp. genome (Cock *et al.*, 2010) (available at Orcae; Sterck *et al.*, 2012) using
115 Tophat2 and the Bowtie2 aligner (Kim *et al.*, 2013). More than 90% of the sequencing reads
116 for each library mapped to the genome.

117 *De novo* assembly of the pooled RNA-seq data from the ten libraries was carried out using
118 Trinity (Grabherr *et al.*, 2011) in normalized mode with default parameters. Weakly expressed
119 transcripts (isoform percentage <1 and RPKM <1) were removed from the dataset. The
120 remaining transcripts were aligned against the *Ectocarpus* reference genome (Ec32) using
121 GenomeThreader (Gremme *et al.*, 2005) with a maximum intron length of 26,000 bp, a
122 minimum coverage of 75% and a minimum alignment score of 90%.

123

124 Gene prediction

125 Gene prediction was carried out using the EuGene program (Foissac *et al.*, 2008), as described
126 previously (Cock *et al.*, 2010). Alignments of the Trinity RNA-seq-derived transcripts against
127 the *Ectocarpus* sp. reference genome were added to the EuGene pipeline in addition to the data
128 used for the v1 annotation, which included splice site predictions generated by SpliceMachine
129 (Degroeve *et al.*, 2005) and *Ectocarpus* Sanger EST data. The new set of EuGene gene structure
130 predictions were compared with the gene structures of the v1 annotation using AEGeAn
131 (Standage & Brendel, 2012) and a combination of automated and manual approaches was used

132 to select the optimal gene structures. Briefly, automatic validation of new predictions was
133 applied for genes where there were modifications to the UTRs, where additional exons were
134 added or where there were modifications to the detailed structure of existing exons. In cases
135 where the new model predicted exon lost, the prediction was retained only if there was 65%
136 similarity between the reference and the new model. This threshold was reduced to 30%
137 similarity when the reference gene had only 4 exons or less. A subset of about one hundred
138 genes for each class was manually reviewed to validate the automatic selection of gene
139 structures. GenomeView (Abeel *et al.*, 2012) was used to visualise RNA-seq read mapping
140 information.

141

142 **Manual annotation**

143 The v2 annotation took into account the functional and structural annotation of 325 and 410
144 genes, respectively, carried out through the Orcae database (Sterck *et al.*, 2012) since the
145 publication of the v1 annotation. Many of the structural annotations were based on the same
146 set of RNA-seq data that was used for the genome-wide gene structure prediction but exploited
147 transcripts that had been generated using a reference-guided approach with Tophat2 and
148 Cufflinks2 (Trapnell *et al.*, 2010; Kim *et al.*, 2013). Tophat2 was able to map 92% of the
149 cleaned reads to the genome sequence and 36,565 transcripts were assembled by Cufflinks2
150 (including multiple transcripts for some loci) using the mapping information and the initial
151 gene models as guides.

152

153 **Annotation of gene functions**

154 Putative functions were assigned to the v2 genes based on the identification of protein domains
155 using InterProScan, which carried out searches against all its component databases (Jones *et al.*
156 *et al.*, 2014). Gene ontology categories were assigned using Blast2GO (Conesa *et al.*, 2005). For
157 genes where a manually assigned function was already available (3,442 genes), the
158 InterProScan-based prediction was compared manually with the existing annotation and the
159 most relevant annotation retained.

160

161 **Detection of alternative transcripts**

162 To detect alternative transcripts of the set of 17,418 protein-coding loci, 507,634,855 million
163 reads of RNA-seq data corresponding to diverse tissues and life cycle stages (Table S1) were
164 mapped to the *Ectocarpus* genome using Bowtie2 (Langmead *et al.*, 2009) and transcripts were

165 predicted genome-wide using Stringtie (Pertea *et al.*, 2015) with default parameters, guided by
166 the annotation file from the v2 annotation. A Stringtie prediction was made for each library
167 based on TopHat2 mapping files. The results were merged using Cuffmerge (Trapnell *et al.*,
168 2010). Cuffcompare was used to assign the predicted transcripts to the reference genes.
169 Transcripts with 3' UTRs > 9300 bp and/or 5' UTRs > 2500 bp were discarded. Only potential
170 isoforms (class code = J, O and C) were retained. Prediction of the coding regions of the
171 alternative transcripts was carried out using Transdecoder (Haas *et al.*, 2013). ORF predictions
172 were filtered to retain complete coding sequences with both initiation and stop codons. The
173 longest ORF was retained for each transcript.

174 A global classification and quantification of the different types of alternative splicing that
175 generated the transcript isoforms was obtained using SplAdder (Kahles *et al.*, 2016) based on
176 the mapping of the pooled RNA-seq data.

177

178 **Detection of non-protein-coding genes**

179 The detection of microRNA, ribosomal RNA and snoRNA loci has been described previously
180 (Tarver *et al.*, 2015).

181 *Ectocarpus* lncRNA loci were detected using FEELnc
182 (<https://github.com/tderrien/FEELnc>) with default parameters and the output transcripts of the
183 Stringtie analysis described in the previous section. The same specificity threshold (0.97) was
184 used for both protein-coding and non-coding transcripts to predict lncRNA loci. Transcripts
185 overlapping annotated protein-coding genes (v2 annotation) were eliminated and a random
186 forest approach based on ORF coverage (i.e. length of the longest ORF / length of the lncRNA
187 transcript), transcript size and k-mer frequency was implemented to classify the remaining
188 transcripts as mRNAs or lncRNAs. Loci with mono-exonic transcripts were eliminated to limit
189 the inclusion of false positive loci due to read mapping ambiguity. An arbitrary minimum size
190 of 200 nt was applied to eliminate loci encoding small RNA transcripts. FEELnc also classifies
191 the predicted lncRNA loci by determining 1) if they overlap (genic) or not (intergenic) with
192 the nearest gene on the genome, designated the adjacent gene (and which can be a protein-
193 coding gene or small-RNA-encoding locus), 2) if genic lncRNAs overlap with intron or exon
194 regions of the adjacent gene and in which orientation, sense or antisense, and 3) how intergenic
195 lncRNAs are orientated with respect to the adjacent gene (within 10 kbp) on the chromosome
196 (same strand, convergent or divergent).

197 A similar approach was used to detect *S. japonica* lncRNA loci. For this genome, the
198 Stringtie transcript prediction used as input for FEELnc was based on mapping of 220,551,196
199 million RNA-seq reads to the *S. japonica* genome (Ye *et al.*, 2015). The RNA-seq data
200 corresponded to female gametes (127,607,414 reads, accession number SRR2064656), spores
201 (30,552,978 reads, accession number SRR2064654), thalli grown under blue light (11,981,830
202 reads, accession number SRR371552) or in the dark (12,657,652 reads, accession number
203 SRR371551), young sporophytes grown under blue (13,333,334 reads, accession number
204 SRR496757) or white (17,181,148 reads, accession number SRR496799) light and thalli
205 subjected to heat stress (7,236,840 reads, accession number SRR947066). Orthologous
206 *Ectocarpus* and *S. japonica* lncRNA loci were detected by carrying out reciprocal Blastn
207 searches (E-value < 10⁻⁴). Alignments of lncRNA sequences were carried out with SIM
208 (<http://web.expasy.org/sim/>) and visualised with Lalnview (Duret *et al.*, 1996).

209 DESeq2 with default parameters was used to detect *Ectocarpus* lncRNA and protein-coding
210 loci that were differently expressed in sporophyte basal versus upright filaments.

211

212 **Genome-wide identification of sequence variants**

213 Genome sequence data was generated for the female outcrossing line Ec568 using Illumina
214 HiSeq2500 technology (Fasteris, Switzerland), which produced 25,976,388,600 bp of 2x100
215 bp paired-end sequence. Sequence variants were detected as described previously (Godfroy *et al.*
216 *al.*, 2015).

217 To determine whether sequence variants behaved as Mendelian loci, a cross between a UV-
218 mutagenised derivative of the reference genome strain Ec32 (strain Ec722) and the female
219 outcrossing line Ec568 (Heesch *et al.*, 2010) was used to generate a population of 180 progeny
220 each corresponding to an independent meiotic event, segregating the two parental alleles of
221 each variant locus. Two libraries were constructed with pools of 84 and 96 haploid, partheno-
222 sporophyte individuals and sequenced using Illumina HiSeq2500 technology (Fasteris,
223 Switzerland) to generate 20,785,058,400 bp and 23,429,143,400 bp of 2x100 bp paired-end
224 sequence, respectively. Sequence variants were detected in each dataset as described previously
225 (Godfroy *et al.*, 2015) and VarScan was used to identify SNPs shared by the two pools of
226 haploid individuals. For each of these SNPs the sum of the variant frequencies observed in the
227 two pools was calculated, and only those for which this sum was between 0.8 and 1.2 were
228 retained. VarScan compare was then used to extract the Ec568 variants from the list of
229 Mendelian segregating SNPs.

230

231 **Database curation of the v2 annotation**

232 A Genome Browser was implemented based on Jbrowse (Buels *et al.*, 2016) using a Chado
233 database (Mungall & Emmert, 2007). The browser integrates both v1 and v2 reference gene
234 models, raw gene models predicted by EuGene, transcripts predicted by Cufflinks and EST and
235 RNA-seq read data.

236

237 **Accession numbers**

238 The accession numbers for the sequence data used in this article are given in supplementary
239 Table S1.

240

241 **Results**

242 **Improved chromosome-scale assembly of the *Ectocarpus* genome**

243 A microsatellite-based genetic map (Heesch *et al.*, 2010) was originally used to produce a
244 large-scale assembly of the *Ectocarpus* genome consisting of 34 pseudo-chromosomes (Cock
245 *et al.*, 2010) corresponding to the 34 linkage groups of the genetic map. The pseudo-
246 chromosomes were constructed by concatenating sequence scaffolds based on the genetic order
247 of sequence-anchored microsatellite markers on the genetic map (Cock *et al.*, 2010). However,
248 due to the low density of the markers, the large-scale assembly included only 325 of the 1,561
249 sequence scaffolds (70.1% of the total sequence length) and, moreover, only 40 (12%) of the
250 mapped scaffolds could be orientated relative to the chromosome (i.e. only 12% of the scaffolds
251 contained at least two microsatellite markers which recombined relative to each other).

252 To improve the large-scale assembly of the *Ectocarpus* genome, we took advantage of a
253 high-density, single nucleotide polymorphism (SNP)-based genetic map that has recently been
254 generated using a Restriction site associated DNA (RAD)-seq method (K. Avia, personal
255 communication). The 3,588 SNP markers used to construct the genetic map were mapped to
256 sequence scaffolds and the recombination information for these markers used to construct a
257 new set of pseudo-chromosomes (Fig. 1). The new large-scale assembly represents a significant
258 improvement because it integrates 531 of the 1,561 sequence scaffolds onto genetic linkage
259 groups (90.5% of the total sequence length) and 49% of these scaffolds have been orientated
260 with respect to their chromosome. Moreover, the high-density genetic map has allowed several
261 fragmented linkage groups / pseudo-chromosomes to be fused, reducing the total number from
262 34 to 28. The exact number of chromosomes in *Ectocarpus* sp. strain Ec32 is not known but