

 Open access • Journal Article • DOI:10.1111/NPH.14321

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

— [Source link](#) 

[Alexandre Cormier](#), [Komlan Avia](#), [Lieven Sterck](#), [Thomas Derrien](#) ...+11 more authors

Institutions: [University of Paris](#), [Ghent University](#), [University of Rennes](#), [Centre national de la recherche scientifique](#)

Published on: 01 Apr 2017 - [New Phytologist](#) (Wiley)

Related papers:

- [The *Ectocarpus* genome and the independent evolution of multicellularity in brown algae](#)
- [Saccharina genomes provide novel insight into kelp biology](#)
- [A sequence-tagged genetic map for the brown alga *Ectocarpus siliculosus* provides large-scale assembly of the genome sequence](#)
- [A Haploid System of Sex Determination in the Brown Alga *Ectocarpus* sp.](#)
- [OUROBOROS is a master regulator of the gametophyte to sporophyte life cycle transition in the brown alga *Ectocarpus*](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/re-annotation-improved-large-scale-assembly-and-bwivvugxyw>



HAL
open science

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, 2017, 214 (1), pp.219-232. 10.1111/nph.14321 . hal-01402123

HAL Id: hal-01402123

<https://hal.sorbonne-universite.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*
45 *al.*, 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*
51 *al.*, 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

263 cytogenetic analysis of another *Ectocarpus* species, *E. siliculosus* indicated the presence of
264 approximately 25 chromosomes (Müller, 1966, 1967).

265

266 **Reannotation of gene structure based on RNA-seq data**

267 The initial set of *Ectocarpus* gene models (referred to hereafter as the v1 annotation) was
268 generated using EuGene (Foissac *et al.*, 2008) based on a limited amount of transcriptomic
269 information (91,041 Sanger expressed sequence tags, ESTs; Cock *et al.*, 2010) and therefore
270 involved a significant amount of *de novo* prediction. The v1 annotation has been gradually
271 improved since 2010 by the addition of 325 functional and 410 structural annotations for
272 individual genes through the Orcae database (Sterck *et al.*, 2012). This gene-by-gene approach
273 improved the quality of the annotation of a number of selected genes but it was necessary to
274 extend the approach to improve annotation quality across the whole genome.

275 A genome-wide reannotation, hereafter referred to as the v2 annotation, was therefore
276 carried out based on the analysis of 642 million reads of RNA-seq data from ten different
277 libraries (Ahmed *et al.*, 2014; Lipinska *et al.*, 2015 and this study; Table S1). This data was
278 assembled into 34,551 *de novo* transcripts using Trinity (Grabherr *et al.*, 2011).
279 GenomeThreader (Gremme *et al.*, 2005) was able to align 91% of these transcripts to the
280 genome. Gene prediction for the v2 annotation was then carried out using EuGene and the
281 34,551 *de novo* transcripts, along with 83,502 Sanger ESTs and SpliceMachine (Degroeve *et*
282 *al.*, 2005) splice site predictions. The 21,958 preliminary gene models generated by this
283 prediction were then compared with the 16,256 genes of the v1 annotation (Cock *et al.*, 2010)
284 using AEGeAn (Standage & Brendel, 2012) and a combination of automatic and manual
285 criteria were used to evaluate the predictions and select the optimal gene model for each locus.
286 This genome-wide reannotation integrated the results of the manual gene-by-gene annotation
287 carried out since publication of the v1 annotation by preferentially retaining high quality, expert
288 functional and structural annotations.

289 The 21,958 preliminary gene predictions included 1) genes that were identical to the v1
290 prediction (10,426 genes), 2) genes that were structurally different to their v1 counterpart
291 (6,295 genes) and 3) novel loci that were not predicted by the v1 annotation (5,237 genes). For
292 the first set of genes, the v1 gene models were replaced with the RNA-seq-based models,
293 providing considerable additional information about the UTR structure of the genes (e.g. Fig.
294 2A). When the RNA-seq-based prediction differed from the v1 model, manual inspection was
295 used to select the optimal model for each locus (e.g. Fig. 2B; see Methods and Materials for

296 details). This second set of genes also included predictions which indicated that v1 annotation
 297 genes needed to be fused (e.g. Fig. 2C) or split (e.g. Fig. 2D). Novel RNA-seq-based
 298 predictions, not present in the v1 annotation, were filtered to remove probable false positives.
 299 Predictions were retained only if 1) their transcripts had an abundance of >1 RPKM across the
 300 entire (merged) set of RNA-seq data, 2) the start codon of the gene was not located in a repeated
 301 region (to exclude transposon-derived ORFs; Yandell & Ence, 2012) and 3) their coding region
 302 was >100 bp. After applying these filters, 2,030 of the new predictions were retained and
 303 integrated into the v2 annotation.

304 Overall, the addition of these new genes and updates to the existing genes (fusing or splitting
 305 existing gene models) brought the total number of genes in the v2 annotated genome to 17,418
 306 (Table 1). The transition from the v1 to the v2 version of the genome annotation involved the
 307 modification of 11,108 of the v1 gene models, of which 5,336 were altered within their coding
 308 regions (Table 2). Of the former, 784 involved gene fusions (to produce 404 genes in the v2
 309 annotation), 19 involved splitting v1 annotation gene predictions (to create 38 genes in the v2
 310 annotation) and 123 genes were removed. The v2 annotation now includes coordinates for at
 311 least one of the UTR regions for 78.7% of the 17,418 genes (compared to 52.6% for the v1
 312 annotation; Fig. 3, Table 1). The v2 annotation is publically available through the ORCAE
 313 database (<http://bioinformatics.psb.ugent.be/orcae/overview/Ectsi>; Sterck *et al.*, 2012).

314 The *Ectocarpus* genome database was modified to take into account the large-scale
 315 assembly of the sequence scaffolds. In particular, the sequentially numbered locusIDs were
 316 modified to indicate sequential position on the pseudochromosome. The correspondence
 317 between the LocusIDs of the v1 and v2 annotations is given in Table S2 and is also available
 318 as a download from the genome database
 319 (<https://bioinformatics.psb.ugent.be/gdb/ectocarpusV2/>).

320

321 **Prediction of gene function**

322 The final 17,418 genes of the v2 annotation were further analysed to improve the prediction of
 323 gene function by comparing protein sequences with the InterPro database using InterProScan
 324 (Jones *et al.*, 2014) and by using Blast2GO (Conesa *et al.*, 2005) to assign gene ontology (GO)
 325 categories. This process allowed functional annotations and GO categories to be assigned to
 326 10,688 and 7,383 of the 17,418 v2 annotation genes, respectively (compared with 5,583 and
 327 5,989, respectively, for the v1 annotation; Table 1). Of the 2,030 genes that were present in the

328 v2 annotation but not the v1 annotation, 212 had matches in the public databases and 135 and
329 79 were assigned functional annotations and GO categories, respectively.

330

331 **Alternative splicing**

332 A previous search for alternative gene transcripts based on the 91,041 Sanger ESTs detected
333 isoforms for only a small percentage (2.9%) of the *Ectocarpus* genes (Cock *et al.*, 2010). Here
334 we carried out an updated search for alternative transcripts using the available RNA-seq data
335 (Table S1). The analysis focused on transcript isoforms with alternatively spliced coding
336 regions because variants of this type are more likely to have biological roles through the
337 production of variant protein products. A total of 10,723 alternative transcripts of this type
338 were detected genome-wide, associated with 7,362 (42.3%) of the 17,418 protein-coding
339 genes. This corresponded to an average of 1.62 transcripts per locus.

340 Whilst alternative splicing of gene transcripts can potentially lead to the production of two
341 or more protein products with different biological activities from a single genetic locus, this is
342 not necessarily the case and alternative transcripts can also represent spliceosomal errors or
343 correspond to variants that do not differ significantly from the principal transcript in terms of
344 transcript functionality. To assess the extent to which alternative splicing has the potential to
345 impact gene function in *Ectocarpus*, we used Interproscan (Jones *et al.*, 2014) to compare the
346 domain structures of the predicted protein products of the principal and alternative gene
347 transcripts of the 7,362 genes that exhibited alternative splicing of their coding region. This
348 analysis indicated that, on average, each isoform lacked about 21% of the domains that were
349 detected in the principal transcript. These marked differences between the domain structures
350 of the protein products of principal and alternative transcripts are likely to significantly modify
351 the activities of the alternative protein products.

352 In addition to this genome-wide approach, a more detailed analysis was carried out for four
353 genes that encoded proteins with multiple, repeated copies of small protein domains. Fig. 4
354 shows that the alternative transcripts of these genes encode multiple protein variants in which
355 repeated domains are included or excluded from the protein product in different combinations,
356 producing proteins with markedly different domain structures. Together with the genome-wide
357 analysis described above, these analyses suggested that alternative splicing is used in
358 *Ectocarpus* to combine protein domain modules to generate multiple protein isoforms from
359 individual loci.

360 Analysis of the types of alternative splicing events that give rise to transcript isoforms in
 361 *Ectocarpus* using the program SplAdder (Kahles *et al.*, 2016) indicated that the most common
 362 event was the use of an alternative 3' acceptor site (Table 3). Intron retention events were
 363 relatively rare, representing less than 12% of the detected events.

364

365 **Identification and integration of non-protein-coding genes**

366 With the exception of tRNA loci (Cock *et al.*, 2010), the v1 annotation provided very little
 367 information about non-protein-coding genes. The v2 annotation includes considerably more
 368 information about this type of locus, in particular integrating 64 microRNA (miRNA) loci, nine
 369 ribosomal RNA loci (rRNA) and 610 of the small nucleolar RNA (snoRNA) loci recently
 370 predicted by Tarver *et al.* (2015). The rRNA and snoRNA loci are listed in Tables S3 and S4;
 371 information about the miRNA loci can be found in Tarver *et al.* (2015).

372 In vertebrates most snoRNAs are located in introns (Hoeppner & Poole, 2012) but this is
 373 not the case in all species and only about 30% of *Ectocarpus* snoRNAs are intronic. Work in
 374 other species has shown that the main function of snoRNAs is to direct chemical modification
 375 of other RNA molecules, particularly ribosomal RNAs (reviewed in Bratkovic & Rogelj,
 376 2014). The two major classes of snoRNA, C/D box and H/ACA box, are principally involved
 377 in methylation and pseudouridylation of RNA molecules, respectively, but several alternative
 378 functions have been reported (Kehr *et al.*, 2014). *Ectocarpus* is predicted to have 95 C/D box
 379 and 515 H/ACA box snoRNAs. Note that the *Ectocarpus* snoRNAs were detected using
 380 ACAseeker and CDseeker and should therefore be considered predictions until their functions
 381 have been investigated experimentally.

382 A search of the *Ectocarpus* genome indicated that the core protein components that associate
 383 with both C/D and H/ACA box snoRNAs to form of sno-ribonucleoproteins (snoRNPs) are
 384 highly conserved in *Ectocarpus* (Table S5).

385 A screen was also carried out for potential long non-coding RNAs (lncRNAs) using the
 386 FEELnc lncRNA prediction pipeline (<https://github.com/tderrien/FEELnc>) and the RNA-seq
 387 data listed in Table S1. This analysis predicted the presence of 717 lncRNA loci in the
 388 *Ectocarpus* genome (Table S6), corresponding to a total of 952 different transcripts (1.3
 389 isoforms per locus on average). The mean size of the lncRNA transcripts was 1,708 nucleotides
 390 and varied between 200 (the defined minimal size) and 7,988 nucleotides. The lncRNA loci
 391 were classified based on their configuration relative to the nearest protein-coding gene in the
 392 genome (referred to in the following text as the adjacent gene) and included both loci that were

393 located entirely in an intergenic region (i.e. long intergenic non-coding RNAs or lincRNAs)
394 and loci that overlapped with their adjacent gene (Fig. S1). About 45% of the lincRNAs were
395 classed as lincRNAs. Expression analysis indicated that lincRNA transcripts were about eight-
396 fold less abundant on average than those of protein-coding genes (Fig. 5). A similar difference
397 in mean expression level has been observed in animal and land plant systems (Ulitsky & Bartel,
398 2013; Chekanova, 2015 and references therein). The *Ectocarpus* lincRNA loci tend to occur in
399 regions of the genome of low gene density. The mean distance of lincRNA loci from flanking
400 protein-coding genes is 8,654 bp, which is significantly longer (Wilcoxon test $P < 2.2e-6$) than
401 the mean distance between protein-coding loci (4,154 bp).

402 To determine whether lincRNAs exhibited differential expression patterns in different
403 tissues, we compared abundances of lincRNA transcripts in replicate samples of two different
404 tissues of the sporophyte stage, the strongly adhering, prostrate filaments of the basal system
405 and the upright filaments of the apical system (Peters *et al.*, 2008). DESeq2 identified 219
406 lincRNA loci that were differentially expressed between these two tissues, and 4,019
407 differentially expressed protein-coding genes ($\text{padj} < 0.1$ and $|\log_2\text{fold-change}| \geq 1$ in both
408 cases).

409 To determine the extent to which the sequences of the *Ectocarpus* lincRNAs have been
410 conserved over evolutionary time, we carried out a search for lincRNA loci in a second brown
411 algal genome, that of the kelp *Saccharina japonica* (Ye *et al.*, 2015). The *Ectocarpus* sp. and
412 *S. japonica* lineages are thought to have diverged between 80 and 110 mya (Silberfeld *et al.*,
413 2010; Kawai *et al.*, 2015). Predicted lincRNA loci were compared between the two species
414 rather than simply searching for sequences related to *Ectocarpus* lincRNAs in the *S. japonica*
415 genome as the former approach is more likely to detect *bona fide* orthologues (Ulitsky & Bartel,
416 2013). *S. japonica* transcripts were predicted using Stringtie (Pertea *et al.*, 2015) based on the
417 mapping of 220,551,196 million reads of RNA-seq data (Ye *et al.*, 2015), corresponding to
418 several different tissues, to the assembled genome sequence. Based on these data, FEELnc
419 predicted the presence of 2,840 lincRNA loci in the *S. japonica* genome (Table S7),
420 corresponding to a total of 3,568 different transcripts (1.3 isoforms per locus on average). The
421 mean size of the *S. japonica* lincRNA transcripts was 2,036 nucleotides and varied between 200
422 (the defined minimal size) and 26,887 nucleotides. As with the *Ectocarpus* lincRNAs, the *S.*
423 *japonica* lincRNAs were found to be organised in a range of configurations relative to the
424 adjacent gene on the genome (Fig. S2). Comparison of the sets of predicted lincRNAs from
425 *Ectocarpus* and *S. japonica* using Blastn identified 64 pairs of loci that exhibited reciprocal
426 best Blast matches with E-values lower than 10^{-4} (Table S8). These loci are highly likely to be

427 orthologous. Note that Blast comparisons may underestimate the extent of similarity between
428 *Ectocarpus* and *S. japonica* lncRNAs because the program relies on the presence of short
429 regions of high sequence conservation to seed alignments.

430 Comparison of pairs of orthologous lncRNAs from *Ectocarpus* and *S. japonica* (e.g. Fig. 6)
431 indicated that they tended to contain both conserved and species-specific domains, with the
432 latter usually being located at the ends of the RNA molecules. This suggests that there may not
433 be strong selection pressure on the length of the lncRNA molecules nor on the precise sites of
434 initiation and termination of the mature transcripts.

435

436 **Impact of the updated large-scale assembly and gene annotation on large-scale genome** 437 **features including the sex chromosome and an integrated viral genome**

438 Linkage group 30 of the v1 assembly was recently shown to correspond to the sex chromosome
439 in *Ectocarpus* (Ahmed *et al.*, 2014). This linkage group consisted of 20 scaffolds in the v1
440 assembly but has been considerably extended in the v2 assembly (chromosome 13 in Fig. 1)
441 with the addition of a further 16 scaffolds, increasing the estimated physical length of the
442 chromosome (cumulative scaffold length) from 4,994 to 6,933 kbp. The non-recombining sex-
443 determining region was not affected by this update, as all the additional scaffolds are located
444 in the pseudoautosomal regions of the chromosome. However, as we have recently described
445 a number of unusual features of the pseudoautosomal regions (Luthringer *et al.*, 2015), we
446 verified that these observations were still valid for the updated version of the chromosome.
447 This analysis confirmed that the updated pseudoautosomal regions continue to exhibit a
448 number of structural features that are intermediate between those of the autosomes and the sex-
449 determining region. In particular, compared with the autosomes, the updated pseudoautosomal
450 regions still exhibit significantly reduced gene density, increased content of transposable
451 element sequences, lower %GC content and the genes had significantly smaller and fewer
452 exons (supplementary Fig. S3). The conclusions of the Luthringer *et al.* (2015) study therefore
453 remain valid for the updated version of the sex chromosome.

454 The genome of *Ectocarpus* strain Ec32 contains an integrated copy of a large DNA virus,
455 closely related to the *Ectocarpus* phaeovirus EsV-1 (Cock *et al.*, 2010). Microarray analysis
456 had shown that all the viral genes were silent (Cock *et al.*, 2010) and the RNA-seq data analysed
457 here confirmed this observation, indicating complete silencing of this region of the
458 chromosome under all the conditions analysed (Fig. S4).

459

460 **A genome-wide variant resource for genetic analysis of brown algal gene function**

461 To create an additional genetic resource for gene mapping in *Ectocarpus*, a genome re-
462 sequencing approach was used to identify sequence variants (single nucleotide polymorphisms,
463 SNPs, and indels) across the entire genome. Hi-seq2500 Illumina technology was used to
464 generate 25,976,388,600 bp of paired-end, sequence reads (121x genome coverage) for the
465 female outcrossing line Ec568 (Heesch *et al.*, 2010). A total of 340,665 high quality sequence
466 variants (Table S9) were identified by comparing this data with the reference genome of the
467 male strain Ec32 (Cock *et al.*, 2010) plus the sex-determining region from the Ec32-related
468 female strain Ec597 (Ahmed *et al.*, 2014).

469 To further validate the sequence variants as potential genetic markers, we used a bulked
470 segregant approach to determine whether they behaved as Mendelian loci. Genomic DNA
471 extracts from a population of 180 segregating progeny derived from a cross between a UV-
472 mutagenised derivative of the reference genome strain Ec32 (strain Ec722) and the female
473 outcrossing line Ec568 were grouped into two bulked segregant pools (84 and 96 individuals)
474 and sequenced on an Illumina platform. Lists of SNP variants were then generated for the two
475 bulked segregant pools and the two lists compared to identify 390,804 shared SNPs that
476 exhibited a 1:1 segregation pattern in the progeny population and were therefore behaving as
477 Mendelian loci. Using this data, 237,839 of the 340,665 sequence variants obtained by mapping
478 the Ec568 DNA-seq data against the reference scaffolds (see above) were validated as
479 Mendelian genetic markers (Table S9). The average distance between adjacent pairs of the
480 genetic markers identified is 823 bp, providing a high-density resource for genetic analysis in
481 this species.

482

483 **Extension and improvement of the *Ectocarpus* genome database**

484 The v1 annotation of the *Ectocarpus* genome has been publically available on the Orcae
485 database (Sterck *et al.*, 2012) since its publication in 2010. We have updated the database by
486 adding the v2 annotation described in this study. In addition, a v2 annotation-based Jbrowse
487 genome browser has been created (<http://mmodev.sb-roscoff.fr/jbrowse/>) to allow
488 simultaneous visualisation of multiple types of data in a genome context. The Jbrowse genome
489 browser allows parallel visualisation of gene models for both coding and non-coding loci,
490 transcript predictions based on RNA-seq data, genetic markers including microsatellites and
491 SNP markers, raw RNA-seq data for both messenger RNAs and small RNAs, Sanger EST data,
492 micro-array data and tiling array data. The Jbrowse genome browser is complementary to the

493 Orcae database, providing an environment for the compilation and analysis of newly generated
494 data before information is definitively incorporated into Orcae, which is the reference database.
495 It is possible for registered users of the Jbrowse genome browser to create private versions in
496 order to upload unpublished and working datasets.

497

498 **Discussion**

499 The objective of the work reported here was to improve the utility of the *Ectocarpus* genome
500 sequence as a genomic resource.

501 A high-density, RAD-seq-based genetic map was exploited to significantly improve the
502 large-scale assembly of the genome. This approach allowed 90.5% of the genome sequence to
503 be assembled into 28 pseudo-chromosomes, providing a high quality reference genome for
504 future comparisons with other brown algal genomes focused on synteny and large-scale
505 organisation of chromosomal regions.

506 In addition, extensive RNA-seq data was used to improve 11,108 existing gene models and
507 to identify 2,030 new protein-coding genes. New data available in the public databases has
508 allowed the functional annotation associated with the protein-coding genes to be considerably
509 improved. Sixty-one percent of genes have now been assigned functional information,
510 compared with 34% in the v1 annotation.

511 The RNA-seq data was also exploited to evaluate the extent to which protein-coding genes
512 generate alternative transcripts. Wu et al. (2013) reported strong skews in codon usage at both
513 the 5' and 3' ends of *Ectocarpus* exons. Based on a preliminary analysis that indicated a low
514 level of alternative splicing compared with humans, these authors suggested that the skews
515 might reflect strong selection to preserve exon splicing enhancers to avoid mis-splicing of gene
516 transcripts. Our analysis, which was based on a significantly larger transcriptomic dataset,
517 detected a frequency of alternative splicing of about 1.62 transcripts per gene on average. It is
518 difficult to precisely evaluate whether *Ectocarpus* exhibits a particularly low level of
519 alternative splicing compared to other model organisms based on this value because estimates
520 for these other organisms are constantly being revised as more extensive transcriptomic
521 datasets become available. Based on current estimates, however, the frequency of alternative
522 splicing in *Ectocarpus* falls within the range of 1.2 to 3.4 transcripts per intron-containing gene
523 proposed for diverse model organisms with intron-rich genomes including humans, mouse,
524 *Drosophila melanogaster*, *Caenorhabditis elegans* and *Arabidopsis thaliana* (Kianianmomeni

525 *et al.*, 2014; Chen *et al.*, 2014; Lee & Rio, 2015; Zhang *et al.*, 2015), and therefore does not
526 appear to be exceptionally low.

527 As far as the types of alternative splicing events are concerned, the *Ectocarpus* genome does
528 not show the same bias towards intron retention events that has been observed with members
529 of the green lineage such as *Arabidopsis* or *Volvox* (Reddy *et al.*, 2013; Kianianmomeni *et al.*,
530 2014). Instead, use of alternative 3' acceptor sites is very common (41% of events), a bias that
531 has not been observed in other genomes as far as we are aware. Analysis of the domain
532 composition of predicted protein products of alternative transcripts indicated that alternative
533 splicing is likely to contribute significantly to the complexity of the *Ectocarpus* proteome.

534 The initial v1 annotation focused on protein-coding genes. In this study a genome-wide
535 search was also carried out for non-coding genes, particularly lncRNA loci. Comparison of the
536 *Ectocarpus* lncRNAs with the lncRNA complement of the kelp *S. japonica* indicated that some
537 of the lncRNA loci were already present in the last common ancestor of these two species and
538 have been at least partially conserved, at the sequence level, over the period of about 80 and
539 110 mya (Silberfeld *et al.*, 2010; Kawai *et al.*, 2015) since the divergence of the two species.
540 Conserved regions were often associated within the same lncRNA with regions that had no
541 equivalent in the opposite species suggesting that brown algal lncRNAs may be organised in a
542 modular fashion and be relatively insensitive to the presence or absence of additional lengths
543 of sequence associated with functional modules. The catalogues of *Ectocarpus* and *S. japonica*
544 lncRNA loci are expected to serve as important reference sets for future analyses of lncRNA
545 function in the brown algae.

546 A genome-wide SNP resource was also developed as part of this study. This collection of
547 SNPs will be a valuable tool for future genetic analyses using *Ectocarpus* as a model system
548 (Cock *et al.*, 2011; Coelho *et al.*, 2012a). All of these new and updated resources have been
549 integrated into the *Ectocarpus* genome database, which has also been improved and extended
550 to facilitate exploitation of the genome data and associated information.

551 With the integration of the new information and resources described here, the *Ectocarpus*
552 genome represents one of the most extensively annotated genomes within the stramenopile
553 group and, as such, will serve as an important reference genome for future genome analysis
554 projects. Recently, the *Ectocarpus* genome provided a reference for the analysis of the larger
555 and more complex genome of the kelp *Saccharina japonica* (Ye *et al.*, 2015) and similar
556 comparisons are expected in the future as part of the many ongoing brown algal and
557 stramenopile genome projects.

558

559 **Acknowledgements**

560 We thank Toshiaki Uji for providing RNA-seq data, diverse members of the *Ectocarpus*
 561 Genome Consortium for manual annotation of genes through the Orcae database and an
 562 anonymous reviewer for comments that led to significant improvement of the manuscript. This
 563 work was supported by the Centre National de la Recherche Scientifique, the Agence Nationale
 564 de la Recherche (project Bi-cycle ANR-10-BLAN-1727, project Idealg ANR-10-BTBR-04-01
 565 and project Sexseaweed ANR-12-JSV7-0008), the University Pierre et Marie Curie and the
 566 European Research Council (grant agreement 638240). A.C. was supported by a grant from the
 567 Brittany Region.

568

569 **Author contributions**

570 AC reannotated the *Ectocarpus* genome, identified and characterised alternative transcripts and
 571 prepared data for database integration, LS and YVDP created the *Ectocarpus* v2 Orcae
 572 database, KA and AC constructed the pseudochromosomes using the genetic map, TD, VW,
 573 AC and CH identified the *Ectocarpus* and *S. japonica* lncRNAs, GA and MM created the
 574 JBrowse database, OG identified and catalogued the SNP markers, AL and MMP analysed
 575 data, JMC, EC and SMC designed and coordinated the research, JMC wrote the manuscript.
 576 All authors read and approved the final manuscript.

577

578 **References**

- 579 **Abeel T, Van Parys T, Saeys Y, Galagan J, Van de Peer Y. 2012.** GenomeView: a next-generation
 580 genome browser. *Nucleic Acids Res* **40**: e12.
- 581 **Ahmed S, Cock JM, Pessia E, Luthringer R, Cormier A, Robuchon M, Sterck L, Peters AF, Dittami SM,**
 582 **Corre E, et al. 2014.** A Haploid System of Sex Determination in the Brown Alga *Ectocarpus* sp. *Curr*
 583 *Biol* **24**: 1945–1957.
- 584 **Bartsch I, Wiencke C, Bischof K, Buchholz C, Buck B, Eggert A, Feuerpfeil P, Hanelt D, Jacobsen S,**
 585 **Karez R, et al. 2008.** The genus *Laminaria sensu lato*: recent insights and developments. *Eur J Phycol*
 586 **43**: 1–86.
- 587 **Bratkovic T, Rogelj B. 2014.** The many faces of small nucleolar RNAs. *Biochimica et biophysica acta*
 588 **1839**: 438–443.
- 589 **Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, Goodstein DM, Elsik CG, Lewis SE,**
 590 **Stein L, et al. 2016.** JBrowse: a dynamic web platform for genome visualization and analysis.
 591 *Genome Biology* **17**: 66.

- 592 **Charrier B, Coelho S, Le Bail A, Tonon T, Michel G, Potin P, Kloareg B, Boyen C, Peters A, Cock J.**
 593 **2008.** Development and physiology of the brown alga *Ectocarpus siliculosus*: two centuries of
 594 research. *New Phytol* **177**: 319–32.
- 595 **Chekanova JA. 2015.** Long non-coding RNAs and their functions in plants. *Curr Opin Plant Biol* **27**:
 596 207–16.
- 597 **Chen L, Bush SJ, Tovar-Corona JM, Castillo-Morales A, Urrutia AO. 2014.** Correcting for differential
 598 transcript coverage reveals a strong relationship between alternative splicing and organism
 599 complexity. *Molecular Biology and Evolution* **31**: 1402–1413.
- 600 **Cock JM, Collén J. 2015.** Independent emergence of complex multicellularity in the brown and red
 601 algae. In: Ruiz-Trillo I,, In: Nedelcu AM, eds. *Advances in Marine Genomics. Evolutionary transitions*
 602 *to multicellular life.* Springer Verlag, 335–361.
- 603 **Cock JM, Peters AF, Coelho SM. 2011.** Brown algae. *Curr Biol* **21**: R573–5.
- 604 **Cock JM, Sterck L, Rouzé P, Scornet D, Allen AE, Amoutzias G, Anthouard V, Artiguenave F, Aury J,**
 605 **Badger J, et al. 2010.** The *Ectocarpus* genome and the independent evolution of multicellularity in
 606 brown algae. *Nature* **465**: 617–21.
- 607 **Coelho SM, Godfroy O, Arun A, Le Corguillé G, Peters AF, Cock JM. 2011.** *OUROBOROS* is a master
 608 regulator of the gametophyte to sporophyte life cycle transition in the brown alga *Ectocarpus*. *Proc*
 609 *Natl Acad Sci U S A* **108**: 11518–11523.
- 610 **Coelho SM, Scornet D, Rousvoal S, Peters N, Dartevelle L, Peters AF, Cock JM. 2012a.** *Ectocarpus*: A
 611 model organism for the brown algae. *Cold Spring Harbor Protoc* **2012**: 193–198.
- 612 **Coelho SM, Scornet D, Rousvoal S, Peters NT, Dartevelle L, Peters AF, Cock JM. 2012b.** How to
 613 cultivate *Ectocarpus*. *Cold Spring Harb Protoc* **2012**: 258–261.
- 614 **Conesa A, Götz S, García-Gómez J, Terol J, Talón M, Robles M. 2005.** Blast2GO: a universal tool for
 615 annotation, visualization and analysis in functional genomics research. *Bioinformatics* **21**: 3674–6.
- 616 **Dayton P. 1985.** Ecology of Kelp Communities. *Annu Rev Ecol Syst* **16**: 215–245.
- 617 **Degroeve S, Saeys Y, De Baets B, Rouzé P, Van de Peer Y. 2005.** SpliceMachine: predicting splice
 618 sites from high-dimensional local context representations. *Bioinformatics* **21**: 1332–8.
- 619 **Dittami S, Scornet D, Petit J, Ségurens B, Da Silva C, Corre E, Dondrup M, Glatting K, König R, Sterck**
 620 **L, et al. 2009.** Global expression analysis of the brown alga *Ectocarpus siliculosus* (Phaeophyceae)
 621 reveals large-scale reprogramming of the transcriptome in response to abiotic stress. *Genome Biol*
 622 **10**: R66.
- 623 **Duret L, Gasteiger E, Perrière G. 1996.** LALNVIEW: a graphical viewer for pairwise sequence
 624 alignments. *Computer applications in the biosciences: CABIOS* **12**: 507–510.
- 625 **Foissac S, Gouzy JP, Rombauts S, Mathé C, Amselem J, Sterck L, Van de Peer Y, Rouzé P, Schiex T.**
 626 **2008.** Genome Annotation in Plants and Fungi: EuGene as a model platform. *Current Bioinformatics*
 627 **3**: 87–97.
- 628 **Godfroy O, Peters AF, Coelho SM, Cock JM. 2015.** Genome-wide comparison of ultraviolet and ethyl
 629 methanesulphonate mutagenesis methods for the brown alga *Ectocarpus*. *Mar Genomics*.

- 630 **Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,**
631 **Raychowdhury R, Zeng Q, et al. 2011.** Full-length transcriptome assembly from RNA-Seq data
632 without a reference genome. *Nat Biotechnol* **29**: 644–52.
- 633 **Gremme G, Brendel V, Sparks ME, Kurtz S. 2005.** Engineering a software tool for gene structure
634 prediction in higher organisms. *Information and Software Technology* **47**: 965–978.
- 635 **Gschloessl B, Guermeur Y, Cock J. 2008.** HECTAR: a method to predict subcellular targeting in
636 heterokonts. *BMC Bioinf* **9**: 393.
- 637 **Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B,**
638 **Lieber M, et al. 2013.** De novo transcript sequence reconstruction from RNA-seq using the Trinity
639 platform for reference generation and analysis. *Nature Protocols* **8**: 1494–1512.
- 640 **Heesch S, Cho GY, Peters AF, Le Corguillé G, Falentin C, Boutet G, Coëdel S, Jubin C, Samson G,**
641 **Corre E, et al. 2010.** A sequence-tagged genetic map for the brown alga *Ectocarpus siliculosus*
642 provides large-scale assembly of the genome sequence. *New Phytol* **188**: 42–51.
- 643 **Hoepfner MP, Poole AM. 2012.** Comparative genomics of eukaryotic small nucleolar RNAs reveals
644 deep evolutionary ancestry amidst ongoing intragenomic mobility. *BMC evolutionary biology* **12**:
645 183.
- 646 **Hughes AD, Kelly MS, Black KD, Stanley MS. 2012.** Biogas from Macroalgae: is it time to revisit the
647 idea? *Biotechnol Biofuels* **5**: 86.
- 648 **Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka**
649 **G, et al. 2014.** InterProScan 5: genome-scale protein function classification. *Bioinformatics (Oxford,*
650 *England)* **30**: 1236–1240.
- 651 **Kahles A, Ong CS, Zhong Y, Rättsch G. 2016.** SplAdder: Identification, quantification and testing of
652 alternative splicing events from RNA-Seq data. *Bioinformatics*.
- 653 **Kawai H, Hanyuda T, Draisma SGA, Wilce RT, Andersen RA. 2015.** Molecular phylogeny of two
654 unusual brown algae, *Phaeostrophion irregulare* and *Platysiphon glacialis*, proposal of the
655 Stschapoviales ord. nov. and Platysiphonaceae fam. nov., and a re-examination of divergence times
656 for brown algal orders. *Journal of Phycology* **51**: 918–928.
- 657 **Kehr S, Bartschat S, Tafer H, Stadler PF, Hertel J. 2014.** Matching of Soulmates: coevolution of
658 snoRNAs and their targets. *Molecular Biology and Evolution* **31**: 455–467.
- 659 **Kianianmomeni A, Ong CS, Rättsch G, Hallmann A. 2014.** Genome-wide analysis of alternative
660 splicing in *Volvox carteri*. *BMC genomics* **15**: 1117.
- 661 **Kijjoa A, Sawangwong P. 2004.** Drugs and Cosmetics from the Sea. *Mar Drugs* **2**: 73–82.
- 662 **Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. 2013.** TopHat2: accurate alignment of
663 transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol* **14**: R36.
- 664 **Klinger T. 2015.** The role of seaweeds in the modern ocean. *Perspect Phycol* **2**: 31–39.
- 665 **Langmead B, Trapnell C, Pop M, Salzberg SL. 2009.** Ultrafast and memory-efficient alignment of
666 short DNA sequences to the human genome. *Genome Biol* **10**: R25.

- 667 **Le Bail A, Billoud B, Le Panse S, Chenivesse S, Charrier B. 2011.** ETOILE Regulates Developmental
668 Patterning in the Filamentous Brown Alga *Ectocarpus siliculosus*. *Plant Cell* **23**: 1666–1678.
- 669 **Lee Y, Rio DC. 2015.** Mechanisms and Regulation of Alternative Pre-mRNA Splicing. *Annual Review of*
670 *Biochemistry* **84**: 291–323.
- 671 **Lipinska A, Cormier A, Luthringer R, Peters AF, Corre E, Gachon CMM, Cock JM, Coelho SM. 2015.**
672 Sexual dimorphism and the evolution of sex-biased gene expression in the brown alga *Ectocarpus*.
673 *Molecular Biology and Evolution* **32**: 1581–1597.
- 674 **Lipinska AP, D'hondt S, Van Damme EJM, De Clerck O. 2013.** Uncovering the genetic basis for early
675 isogamete differentiation: a case study of *Ectocarpus siliculosus*. *BMC genomics* **14**: 909.
- 676 **Luthringer R, Lipinska AP, Roze D, Cormier A, Macaisne N, Peters AF, Cock JM, Coelho SM. 2015.**
677 The Pseudoautosomal Regions of the U/V Sex Chromosomes of the Brown Alga *Ectocarpus* Exhibit
678 Unusual Features. *Molecular Biology and Evolution* **32**: 2973–2985.
- 679 **Meslet-Cladière L, Delage L, Leroux CJ, Goulitquer S, Leblanc C, Creis E, Gall EA, Stiger-Pouvreau V,**
680 **Czjzek M, Potin P. 2013.** Structure/Function Analysis of a Type III Polyketide Synthase in the Brown
681 Alga *Ectocarpus siliculosus* Reveals a Biochemical Pathway in Phlorotannin Monomer Biosynthesis.
682 *Plant Cell* **25**: 3089–103.
- 683 **Müller DG. 1966.** Untersuchungen zur Entwicklungsgeschichte der Braunalge *Ectocarpus siliculosus*
684 aus Neapel. *Planta* **68**: 57–68.
- 685 **Müller DG. 1967.** Generationswechsel, Kernphasenwechsel und Sexualität der Braunalge *Ectocarpus*
686 *siliculosus* im Kulturversuch. *Planta* **75**: 39–54.
- 687 **Mungall CJ, Emmert DB. 2007.** A Chado case study: an ontology-based modular schema for
688 representing genome-associated biological information. *Bioinformatics* **23**: i337–46.
- 689 **Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015.** StringTie enables
690 improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* **33**: 290–5.
- 691 **Peters AF, Marie D, Scornet D, Kloareg B, Cock JM. 2004.** Proposal of *Ectocarpus siliculosus*
692 (*Ectocarpales*, *Phaeophyceae*) as a model organism for brown algal genetics and genomics. *J Phycol*
693 **40**: 1079–1088.
- 694 **Peters AF, Scornet D, Ratin M, Charrier B, Monnier A, Merrien Y, Corre E, Coelho SM, Cock JM.**
695 **2008.** Life-cycle-generation-specific developmental processes are modified in the *immediate upright*
696 mutant of the brown alga *Ectocarpus siliculosus*. *Development* **135**: 1503–12.
- 697 **Prigent S, Collet G, Dittami SM, Delage L, Ethis de Corny F, Dameron O, Eveillard D, Thiele S,**
698 **Cambefort J, Boyen C, et al. 2014.** The genome-scale metabolic network of *Ectocarpus siliculosus*
699 (*EctoGEM*): a resource to study brown algal physiology and beyond. *Plant J* **80**: 367–81.
- 700 **Reddy ASN, Marquez Y, Kalyna M, Barta A. 2013.** Complexity of the alternative splicing landscape in
701 plants. *The Plant Cell* **25**: 3657–3683.
- 702 **Ritter A, Goulitquer S, Salaün J, Tonon T, Correa J, Potin P. 2008.** Copper stress induces biosynthesis
703 of octadecanoid and eicosanoid oxygenated derivatives in the brown algal kelp *Laminaria digitata*.
704 *New Phytol* **180**: 809–21.

- 705 **Silberfeld T, Leigh JW, Verbruggen H, Cruaud C, de Reviers B, Rousseau F. 2010.** A multi-locus time-
 706 calibrated phylogeny of the brown algae (Heterokonta, Ochrophyta, Phaeophyceae): Investigating
 707 the evolutionary nature of the 'brown algal crown radiation'. *Mol Phylogenet Evol* **56**: 659–74.
- 708 **Smit AJ. 2004.** Medicinal and pharmaceutical uses of seaweed natural products: A review. *J Appl*
 709 *Phycol* **16**: 245–262.
- 710 **Standage DS, Brendel VP. 2012.** ParsEval: parallel comparison and analysis of gene structure
 711 annotations. *BMC Bioinformatics* **13**: 187.
- 712 **Steneck RS, Graham MH, Bourque BJ, Corbett D, Erlandson JM, Estes JA, Tegner MJ. 2002.** Kelp
 713 forest ecosystems: biodiversity, stability, resilience and future. *Environ Conserv* **29**: 436–459.
- 714 **Sterck L, Billiau K, Abeel T, Rouzé P, Van de Peer Y. 2012.** ORCAE: online resource for community
 715 annotation of eukaryotes. *Nat Methods* **9**: 1041.
- 716 **Tarver JE, Cormier A, Pinzón N, Taylor RS, Carré W, Strittmatter M, Seitz H, Coelho SM, Cock JM.**
 717 **2015.** microRNAs and the evolution of complex multicellularity: identification of a large, diverse
 718 complement of microRNAs in the brown alga *Ectocarpus*. *Nucl Acids Res* **43**: 6384–6398.
- 719 **Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ,**
 720 **Pachter L. 2010.** Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts
 721 and isoform switching during cell differentiation. *Nat Biotechnol* **28**: 511–5.
- 722 **Tseng C. 2001.** Algal biotechnology industries and research activities in China. *J. Appl. Phycol.* **13**:
 723 375–380.
- 724 **Ulitsky I, Bartel DP. 2013.** lincRNAs: genomics, evolution, and mechanisms. *Cell* **154**: 26–46.
- 725 **Wahl M, Molis M, Hobday AJ, Dudgeon S, Neumann R, Steinberg P, Campbell AH, Marzinelli E,**
 726 **Connell S. 2015.** The responses of brown macroalgae to environmental change from local to global
 727 scales: direct versus ecologically mediated effects. *Perspect Phycol* **2**: 11 – 29.
- 728 **Wu X, Tronholm A, Caceres EF, Tovar-Corona JM, Chen L, Urrutia AO, Hurst LD. 2013.** Evidence for
 729 deep phylogenetic conservation of exonic splice-related constraints: splice-related skews at exonic
 730 ends in the brown alga *Ectocarpus* are common and resemble those seen in humans. *Genome Biol*
 731 *Evol* **5**: 1731–45.
- 732 **Yandell M, Ence D. 2012.** A beginner's guide to eukaryotic genome annotation. *Nat Rev Genet* **13**:
 733 329–42.
- 734 **Ye N, Zhang X, Miao M, Fan X, Zheng Y, Xu D, Wang J, Zhou L, Wang D, Gao Y, et al. 2015.**
 735 *Saccharina* genomes provide novel insight into kelp biology. *Nat Commun* **6**: 6986.
- 736 **Zambounis A, Elias M, Sterck L, Maumus F, Gachon CM. 2012.** Highly dynamic exon shuffling in
 737 candidate pathogen receptors... What if brown algae were capable of adaptive immunity? *Mol Biol*
 738 *Evol* **29**: 1263–1276.
- 739 **Zhang R, Calixto CPG, Tzioutziou NA, James AB, Simpson CG, Guo W, Marquez Y, Kalyna M, Patro**
 740 **R, Eyras E, et al. 2015.** AtRTD - a comprehensive reference transcript dataset resource for accurate
 741 quantification of transcript-specific expression in *Arabidopsis thaliana*. *The New phytologist* **208**: 96–
 742 101.

743

744 **Supporting information**

745 Additional supporting information may be found in the online version of this article.

746 **Fig. S1** Classification of *Ectocarpus* lncRNAs.747 **Fig. S2** Classification of *S. japonica* lncRNAs.748 **Fig. S3** Comparisons of structural characteristics of the sex-determining and pseudoautosomal
749 regions of the sex chromosome with both a representative autosome and with all autosomes for
750 both the v1 and v2 versions of the *Ectocarpus* genome annotation.751 **Fig. S4** Suppressed transcription from a viral genome inserted into chromosome 6.752 **Table S1** *Ectocarpus* RNA-seq data used in this study. Reads were cleaned using the Fastx
753 toolkit.754 **Table S2** Correspondences between v1 and v2 LocusIDs.755 **Table S3** List of the rRNA loci in the assembled *Ectocarpus* genome.756 **Table S4** List of predicted snoRNA loci in the *Ectocarpus* genome.757 **Table S5** *Ectocarpus* orthologues of core protein components of snoRNPs.758 **Table S6** List of predicted lncRNA loci in the *Ectocarpus* genome.759 **Table S7** List of predicted lncRNA loci in the *S. japonica* genome.760 **Table S8** Comparisons of pairs of orthologous lncRNA loci from *Ectocarpus* and *S.*761 *japonica*. Orthologous loci were detected by comparing FEELnc-predicted lncRNA loci from
762 *Ectocarpus* and *S. japonica* using Blastn with a cut off of 10^{-4} .763 **Table S9** List of 341,426 sequence variants between the genome of the reference male strain
764 Ec32 and the female outcrossing line Ec568.

765

766

767 **Tables**

768

769 **Table 1** Comparison of genome-wide statistics for the v1 and v2 annotations of the770 *Ectocarpus* genome

	v1 annotation	v2 annotation
Genes (including UTRs)		
Number of genes	16,256	17,418
Mean gene length (bp)	6,859	7,542
Longest gene (bp)	122,137	123,931
Shortest gene (bp)	134	150
Exons		
Total number	129,875	134,690
Mean number per gene	7.3	7.96
Max number per gene	171	173
Mean length (bp)	242.2	299.8
Introns		
Total number	113,619	121,264
Mean length (bp)	703.8	739.87
Max length (bp)	25,853	36,147
UTRs		
Genes with only annotated 5' UTR	1,098	918
Genes with only annotated 3' UTR	4,766	3,056
Genes with annotated 5' and 3' UTR	2,484	9,737
Genes without any annotated UTR	7,598	3,715
Mean 5' UTR length (bp)	120.60	139.61
Mean 3' UTR length (bp)	674.74	901.66
Annotation of gene functions		
Genes with predicted functions	5,583	10,688
Genes with associated GO terms	5,989	7,383
miRNA loci	26	64
rRNA loci	n/a	5
snoRNA loci	n/a	656
lncRNA loci	n/a	717

771

772

773 **Table 2** Overview of the modifications to the v1 annotation during the production of the v2774 annotation of the *Ectocarpus* genome

	Number of genes
N° of v1 models with modified CDS region in the v2 annotation	5,336
N° of v1 models with modified CDS and/or UTR in the v2 annotation	11,108
N° of v1 models fused in the v2	784
N° of v1 models split in the v2	19
N° of v1 gene models removed	123

24

N° of new gene models in the v2 annotation

2,030

775

776

777

778

Table 3 Proportions of the different types of alternative splicing events that generate alternative transcripts in *Ectocarpus*

	Mean occurrence per gene	Proportions of alternative splicing events for the genome (%)
Alternative 3' acceptor site	0.481	40.95
Alternative 5' donor site	0.248	21.07
Intron retention	0.139	11.79
Single exon skipping	0.254	21.59
Skipping of multiple exons	0.054	4.58

779

780

For Peer Review

781 **Figures**

782

783 **Fig. 1** Large-scale assembly of the *Ectocarpus* scaffolds into pseudochromosomes based on a
 784 high-density, RAD-seq-based genetic map. Each bar represents one of the 28 chromosomes.
 785 Sequence scaffolds (supercontigs) are drawn to scale and identified with numbers (e.g. 207,
 786 sctg_207). Left or right pointing arrowheads indicate that the scaffolds have been orientated
 787 with respect to the chromosome (i.e. scaffolds with at least two markers separated by at least
 788 one recombination event); unorientated scaffolds are indicated with a spot. Chromosome 13
 789 corresponds to the sex chromosome and the non-recombining sex-determining region is
 790 indicated with a bar.

791

792 **Fig. 2** Representative comparisons of v1 and v2 annotation gene predictions illustrating the
 793 major types of annotation correction carried out during the transition between the two versions.
 794 Protein coding exons are in light or dark green for genome annotation versions v1 and v2,
 795 respectively, UTRs are in grey and introns are indicated by thin black lines. **a** analysis of the
 796 RNA-seq data allowed the identification of UTRs for gene Ec-27_006370. **b** v2 genes Ec-
 797 27_006520 and Ec-05_002440 have been extended and modified compared to their v1
 798 equivalents. **c** v1 genes Esi0002_0099 and Esi0002_0101 were fused to create a single locus,
 799 Ec-01_007860. **d** v1 gene Esi0002_0311 was split to create two loci, Ec-01_006420 and Ec-
 800 01_006425. Arrows indicate gene features that were not identified or misidentified by the v1
 801 annotation.

802

803 **Fig. 3** Comparison of the degree of completeness of gene annotations in the v1 and v2
 804 versions of the *Ectocarpus* genome annotation.

805

806 **Fig. 4** Protein variants predicted to be encoded by alternative transcripts of four genes. **a**
 807 alternative products of the ROCO LRR GTPase gene Ec-06_001640 with different LRR repeat
 808 structures, **b** alternative products of the nucleotide-binding adaptor shared by the NB-ARC
 809 TPR domain containing gene Ec-25_000110 with different TPR domain contents, **c** alternative
 810 products of the Notch domain gene Ec-19_004380 with different Notch repeat structures, **d**
 811 alternative products of the Ankyrin repeat gene Ec-09_000460 with different Ankyrin repeat
 812 structures. Grey lines indicate domains shared between proteins. Roc, Ras of complex proteins

813 domain; DUF4782, domain of unknown function 4782; VPS9, Vacuolar Protein Sorting-
814 associated 9 domain. The LocusID of each isoform is indicated.

815

816 **Fig. 5** *Ectocarpus* lncRNA transcript abundance. On average, lncRNA transcripts are about
817 eight-fold less abundant than those of protein coding genes.

818

819 **Fig. 6** Examples of lncRNA loci conserved between *Ectocarpus* and *Saccharina japonica*.

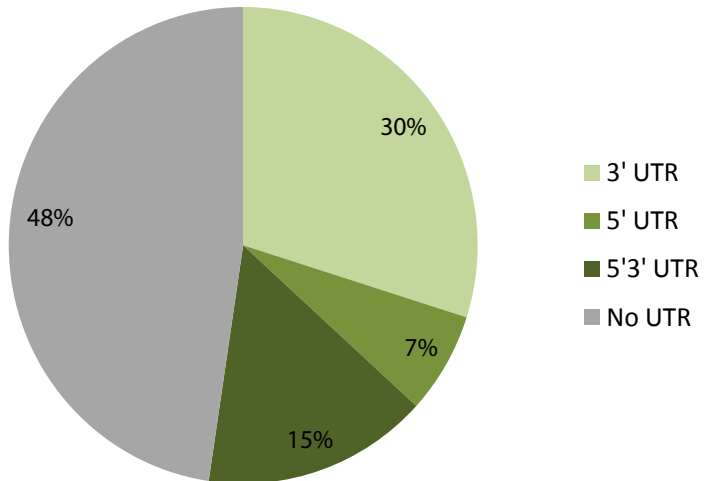
820 lncRNA loci (in blue) are shown for each species, along with the nearest protein-coding locus
821 on the chromosome (in red). Genes above the line, which represents the chromosome, are
822 transcribed to the right, genes below the line to the left. Percent identities over the aligned
823 regions of *Ectocarpus* and *S. japonica* lncRNA transcripts are indicated. Ec, *Ectocarpus*, Sj, *S.*
824 *japonica*.

825

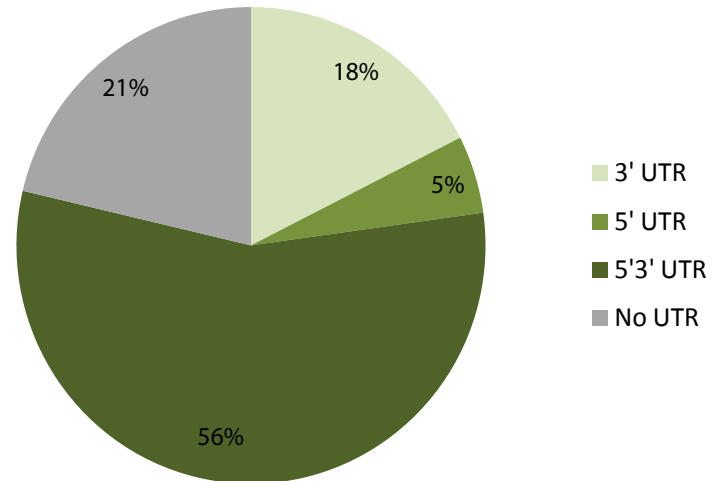
826

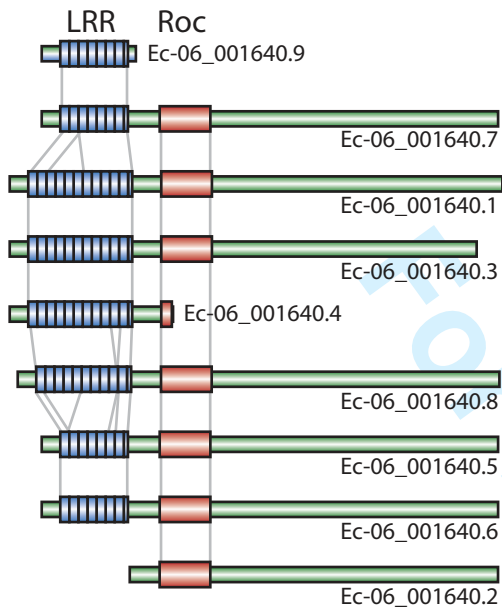
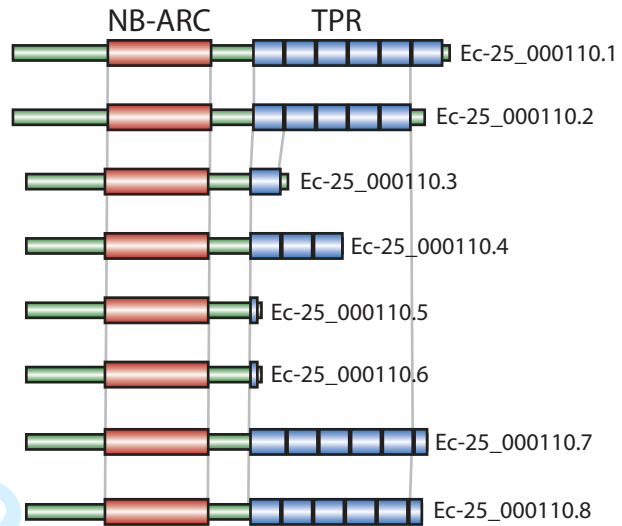
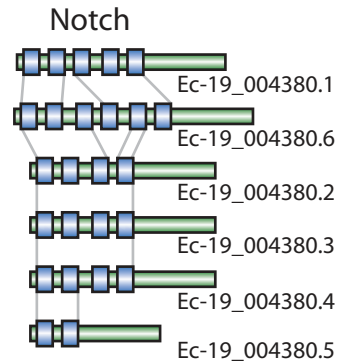
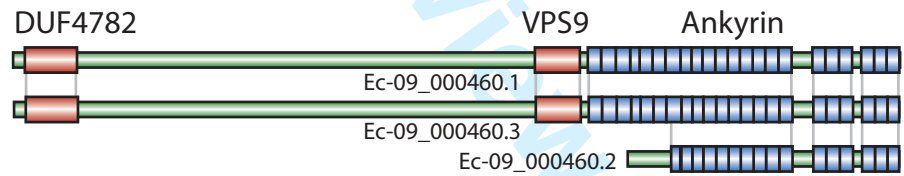


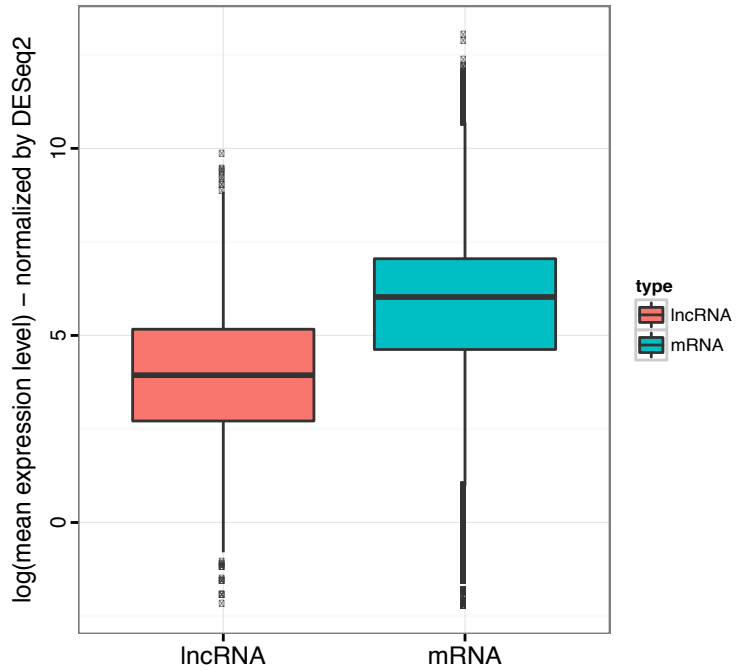
Annotation V1



Annotation V2



a**b****c****d**



Pre-Review

