



HAL
open science

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*.

Bernard Billoud, Zofia Nehr, Aude Le Bail, Bénédicte Charrier

► To cite this version:

Bernard Billoud, Zofia Nehr, Aude Le Bail, Bénédicte Charrier. Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*.. *Nucleic Acids Research*, 2014, 42 (1), pp.417-29. <10.1093/nar/gkt856>. <hal-01002390>

HAL Id: hal-01002390

<https://hal.sorbonne-universite.fr/hal-01002390v1>

Submitted on 20 Jun 2014

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.



HAL Authorization

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Bernard Billoud^{1,2}, Zofia Nehr^{1,2}, Aude Le Bail^{1,2,3}, Bénédicte Charrier^{1,2}

1. Université Pierre et Marie Curie (UPMC), UMR 7139 Végétaux marins et Biomolécules, Station Biologique, CS 90074, F29688, Roscoff, France.

2. Centre National de la Recherche Scientifique (CNRS), UMR 7139 Végétaux marins et Biomolécules, Station Biologique, CS 90074, F29688, Roscoff, France.

3. Current address: Department of Biology, University of Erlangen-Nuremberg, Staudtstrasse 5, 91058 Erlangen, Germany.

Corresponding author: Bernard Billoud, Morphogenesis of Macro-Algae, UMR7139, Station Biologique de Roscoff, CS 90074, F29688, Roscoff, France. Tel: +33 2 98 29 56 53; Fax: +33 2 98 29 23 24; E-mail: Bernard.Billoud@sb-roscoff.fr

Abstract

We used an *in silico* approach to predict microRNAs genome-wide in the brown alga *Ectocarpus siliculosus*. As brown algae are phylogenetically distant from both animals and land plants, our approach relied on features shared by all known organisms, excluding sequence conservation, genome localisation and pattern of base-pairing with the target. We predicted between 500 and 1500 microRNAs candidates, depending on the values of the energetic parameters used to filter the potential precursors. Using quantitative PCR assays, we confirmed the existence of 22 microRNAs among 72 candidates tested, and of 8 predicted precursors. In addition, we compared the expression of microRNAs and their precursors in two life cycle states (sporophyte, gametophyte) and under salt stress. Several microRNA precursors, Argonaute and DICER mRNAs were differentially expressed in these conditions. Finally, we analysed the gene organisation and the target functions of the predicted candidates. This showed that *E. siliculosus* miRNA genes are, like plant miRNA genes, rarely clustered and, like animal miRNA genes, often located in introns. Among the predicted targets, several widely conserved functional domains are significantly over-represented, like kinesin, NB-ARC and tetra-tricopeptide repeats. The combination of computational and experimental approaches thus emphasises the originality of molecular and cellular processes in brown algae.

Introduction

MicroRNAs (miRNAs) are short, single-stranded RNA molecules, which are able to regulate gene expression by interfering with messenger RNAs (mRNAs). Since their discovery in the nematode *Caenorhabditis elegans* (1, 2), their taxonomic coverage has been extended to other animals, plants, green algae and viruses (reviewed in 3). The miRNAs from different lineages belong to a common class of functional molecules, in which several ubiquitous families share extensive sequence similarity. However, their biogenesis, primary and secondary structures, and mode of action are not exactly the same between plants and animals (4). Distinct from both the opisthokonts (animals, fungi) and the archaeplastida (plants, green and red algae), the heterokonts form a large eukaryotic phylum comprising unicellular (*e.g.* diatoms), syncytial (*e.g.* oomycetes) and multicellular organisms (*e.g.* brown algae). Phylogenetic analyses showed that heterokonts diverged from the ancestors of the opisthokont and the archaeplastida phyla more than one billion years ago (5), enabling a large field of alternative molecular strategies to adapt, develop and evolve. In this perspective, identifying and studying miRNAs in heterokonts is likely to uncover new features and mechanisms. Advanced genomic studies were carried out in several organisms of the heterokont phylum: oomycetes *Phytophthora sp.* (6), *Hyaloperonospora sp.* (7) and *Pythium sp.* (8); diatoms *Thalassiosira sp.* (9) and *Phaeodactylum tricornutum* (10). Yet, no miRNAs were reported in these organisms; only 13 precursors were found in the latter (11). The brown alga *Ectocarpus siliculosus* (Ectocarpales, Phaeophyceae; see 12) genome has recently been published (13) and the identification of several specific proteins involved in the miRNA biogenesis (Argonaute: *AGO1* and DICER: *DCLI*) provided a good presumption of the presence of miRNAs in this alga. Among the results obtained by the annotation consortium, a deep-sequencing approach followed by a computational filtering allowed the detection of nine different miRNA candidates for which targets could be predicted, and fourteen other miRNA candidates without a predicted target. However, organisms with a miRNA machinery are expected to be able to produce and use hundreds of miRNAs (3).

Computational methods aiming at the *ab initio* identification of miRNAs in a newly sequenced genome (reviewed in 14, 15, 16) are either based on sequence conservation, or inspired by the knowledge about their biogenesis and function. Sequence conservation can only be used when an extensive repertoire of miRNAs is available in a closely related species, which is not the case for *E. siliculosus*. Similarly, the phylogenetic position of heterokonts, apart from both animals and plants, does not permit to refer to animal- or plant-specific structural features or miRNA biogenesis processes. Thus, our identification and filtering criteria must rely on the common features of miRNAs shared by the organisms investigated to date. These features are: (i) the primary transcript

1 is processed into a precursor called the pre-miR. Alternatively, pre-miRs can derive from introns of
2 protein-coding genes (17, 18). In any case, the pre-miR is folded as a stable hairpin, which exhibits
3 specific structural features (19); (ii) this structured RNA is recognised and processed by an enzyme
4 called DICER (20), which excises a short (possibly imperfect) duplex made of the miRNA and its
5 complementary strand: the miRNA*. The cutting points may vary by one or two nucleotides, thus
6 allowing the precursor to generate alternative duplexes (21); (iii) the miRNA strand is specifically
7 incorporated into a ribonucleoproteic complex called RISC (RNA-induced silencing complex,
8 including the protein Argonaute), and serves as a guide to bind the complex to the target mRNA,
9 using an imperfect sequence complementarity. It is possible that one single miRNA interacts with
10 several mRNAs, and reciprocally, one mRNA can be targeted by more than one miRNA (22, 23,
11 24). The subsequent inhibition of target translation can be caused by the interference of the RISC
12 complex with regulatory elements, the destabilisation of the mRNA and/or its cleavage at the RISC
13 binding site. In much less documented cases, the target expression is enhanced by an unknown
14 process (25). The repression by cleavage requires a better complementarity between the miRNA and
15 the target mRNA, and is mostly (but not exclusively) found in plants, while the other processes are
16 typical of (but not exclusive to) animals (26).

17 Here, we apply a genome-wide approach to extend the set of miRNA candidates in *E. siliculosus*
18 using a computational identification and filtering, followed by an experimental search of selected
19 candidates under a variety of conditions. Relying on common features, we identified a
20 comprehensive list of 568 miRNA candidates in the genome of *E. siliculosus*, from which 22 were
21 experimentally validated. We also analysed their specific features in terms of sequence, genomic
22 organisation and putative biological functions.

23 **Material & Methods**

24 **Data collection and preparation**

25 The genome sequences (super-contigs with length > 2 kbp), primary annotation and protein
26 sequences (version June 2010) were retrieved from the *E. siliculosus* annotation website
27 (<https://bioinformatics.psb.ugent.be/gdb/ectocarpus/>). These data were processed to obtain a
28 suitable partition of the sequences (shown as “re-assignment” on figure 1). We identified all the
29 ribosomal RNAs, including those which were not annotated, using BLAST (27) to search for
30 sequences similar to the EMBL entries EF990201 (partial rRNA gene of *E. siliculosus* [28]) and
31 D16558 (complete rRNA gene of the closely related species *Scytosiphon lomentaria* [29]). We also
32 identified the tRNAs using the tRNAscan-SE 1.23 software (30) downloaded from the author's web
33 site (<http://lowelab.ucsc.edu/tRNAscan-SE/>). The rRNA and tRNA sequences were stored as

1
2 separate datasets, and masked in the genomic sequence. In the remaining genome sequence, mRNA
3 annotation was found to miss an important information: in a total of 16,254 mRNAs, no 3'UTR was
4 annotated for 8,677, no 5'UTR was annotated for 12,598 mRNAs, and 7,602 mRNAs had no
5 annotation for UTR at all. In all these cases, we added putative UTRs to the annotated mRNA
6 sequence, by extending the first exon in 5' and/or the last exon in 3'. The length of the putative
7 5'UTRs was computed as the 95th percentile of the length distribution of the annotated 5'UTRs
8 upstream of the first coding exon, *ie* 374 nts. Similarly, 3'UTR of 1,734 nts were added
9 downstream the last exons of genes without an annotated 3'UTR. As a summary of this preparation
10 step, the genome was sorted into five sets: rRNAs, tRNAs, mRNAs, introns and intergenic
11 sequences.
12
13
14
15
16
17
18

19 **Constitution of a large set of potential miRNAs**

20
21 The program findMiRNA (31) was downloaded from its author's website
22 (<http://sundarlab.ucdavis.edu/mirna/>). We introduced a small modification to the code, in order to
23 predict the best secondary structure at a temperature of 13°C (the usual temperature for
24 *E. siliculosus* cultures, see 32) instead of the default 37°C. This adapted version of findMiRNA was
25 used to process (i) the “intergenic” sequences (on both strands) and (ii) the introns (on the transcript
26 strand only), both using the mRNAs as a reference (shown as “findMiRNA” on figure 1).
27
28
29
30
31

32 **Reference RNA folding parameters**

33
34 In order to build a set of reference values for the structural and topological properties of folded
35 RNAs in *E. siliculosus*, we isolated local sub-structures from rRNAs and mRNAs. To do so, we
36 predicted their complete optimal folding at 13°C (see above) using the RNAfold 1.8 software (33)
37 downloaded from the Vienna RNA Package website (<http://www.tbi.univie.ac.at/~ivo/RNA/>). Sub-
38 structures made of 40 to 200 nucleotides were extracted out of these complete structures. As tRNA
39 sizes lie within this range, we used their full-length sequences.
40
41
42
43

44 We determined the structural and topological properties of the complete tRNAs and of the rRNA
45 and mRNA sub-structures using the genRNAsStats and RNAspectral programs (19) downloaded
46 from their author's website ([http://web.bii.a-star.edu.sg/~stanley/Publications/Supp_materials/06-
47 004-supp.html](http://web.bii.a-star.edu.sg/~stanley/Publications/Supp_materials/06-004-supp.html)), and tuned to use a folding temperature of 13°C (see above). This step is shown as
48 “RNAfold” and “analyse” on figure 1. Its output is a list of five parameter values computed for each
49 analysed structured RNA sequence: normalised minimum free energy (Nmfe), normalised Shannon
50 entropy (NQ), normalised base-pair distance (ND), degree of compactness (NF) and normalised
51 base-pairing propensity (Nbp). Valid pre-miRs have a lower Nmfe, NQ, ND, NF and a higher Nbp
52 than other structured RNAs (19). Thus, we computed threshold values for each of these parameters
53 as follows. As an example, for Nmfe, we computed three values: the Nmfe values above which
54
55
56
57
58
59
60

1 were found 90% of the rRNA, mRNA and tRNA. The lowest among these three values was retained
2 as Nmfe90. We computed NQ90, ND90 and NF90 by applying the same procedure. Nbp90 was
3 computed the same way, only substituting the minimum among the three values having 90% of the
4 distribution above them by the maximum among the three values having 90% of the distribution
5 below them. For a more discriminant filter, we also computed Nmfe95, NQ95, ND95, NF95 and
6 Nbp95 by the same method, using 95% of the corresponding distributions instead of 90%. This step
7 is shown as “distribution analysis” on figure 1, and the resulting values are shown in table 1. These
8 statistical treatments were performed with R (34), obtained from the Comprehensive R Archive
9 Network (<http://cran.r-project.org/>).

17 Selection of miRNA candidates

18 We computed the values of the five parameters described above on the predicted secondary
19 structures of the 697,657 potential pre-miRs resulting from the search with findMiRNA. Each
20 candidate was retained in the pre-miR90 set only if its values of Nmfe, NQ, ND and NF were
21 respectively lower than Nmfe90, NQ90, ND90 and NF90, and its value of Nbp was higher than
22 Nbp90. These threshold are displayed as “filter 90” in table 1. Similarly, we used the values of
23 Nmfe95, NQ95, ND95, NF95 and Nbp95 to design the “filter 95” (table 1) and to obtain the
24 pre-miR95 subset of candidates. This step is shown as “filter” on figure 1.

25 A Support Vector Machine (SVM) approach was implemented using the R package e1071
26 (<http://cran.r-project.org/web/packages/e1071>) to filter the potential pre-miRs, using a discriminant
27 kernel function adjusted for maximal separation between known pre-miRNAs (1000 animal and
28 plant precursors from miRBase) and non-miRNA (1000 tRNAs, rRNA and mRNA hairpins). We
29 used the same five parameters as those on which the filters were applied (see above). We did not
30 succeed in obtaining a robust set of pre-miRNA candidates as an output: when several runs of the
31 adjustment procedure were performed using random “positive” and “negative” sequences, the
32 results of the prediction were different. This approach was not carried on further.

45 Expression level analysis

46 We sorted the sequences (tiles) of the *E. siliculosus* high-resolution transcription map (tiling array,
47 Gene Expression Omnibus GSE19912) into seven sets: Exon, tRNA, Intergenic not candidate,
48 Intergenic miR candidate, Intronic not candidate, Intronic miR candidate, Others (discarded from
49 further analysis). The expression level for each tile was computed as the logarithm of its RNA
50 expression normalised to its DNA expression signal. The statistical significance of the difference
51 between “miR candidate” and “not candidate” sets was tested using the Student's t-test implemented
52 in R (34).

Culture conditions and treatments

E. siliculosus uni-algal strain 32 (CCAP accession 1310/4, origin san Juna de Marcona, Peru) was cultivated in 10 L plastic flasks in a culture room at 13°C using filtered and autoclaved natural seawater enriched in Provasoli nutrients (32). Light was provided by daylight fluorescence tubes with a photon flux density of 40 $\mu\text{mol}\cdot\text{m}^{-2}\cdot\text{s}^{-1}$ for 14 hours per day. Cultures were bubbled with filtered (0.22 μm) compressed air to avoid CO₂ depletion. To conduct the chemical treatment experiments, sporophyte materials were transferred into Petri dishes containing artificial seawater enriched with Provasoli (ASW; see 32) for at least 18 hours before treatments in order to acclimatise the cultures to the change of growth conditions. They were then treated with different chemicals for 6 hours. To perform saline stresses, sporophytes were transferred for acclimatization to ASW for one week before applying the salt stresses. Hypo-saline stress corresponded to 56 mM and hyper-saline stress to 1470 mM NaCl in ASW (ASW contains 450 mM NaCl). Treatments were applied for 6 hours before harvesting the tissues in liquid nitrogen for RNA extraction. Gametophytes were cultured in the same conditions as the sporophytes and collected before they reached maturity. For each treatment or condition, total RNAs were extracted from three independent biological replicates.

RNA extraction and cDNA synthesis

RNAs were extracted as described in Le Bail *et al.* (35). They were RNase-free DNase I-treated (Turbo DNase, Ambion), cleaned up, diluted in RNase-free water and quantified using a NanoDrop ND-1000 spectrophotometer. RNA integrity was verified on 1.5% agarose gel stained with ethidium bromide. From each RNA sample, 2 μg was polyadenylated by the poly(A) polymerase (PAP) of the Poly(A) Tailing Kit (Ambion) according to the manufacturer instructions, and reverse transcribed to cDNA using oligo(dT)₁₂₋₁₈. Two types of reverse transcriptases of the “First Strand synthesis for RT-PCR” kits (Invitrogen) were used: (i) the Superscript™ to detect and quantify the expression level of the 72 miRNAs in the sporophyte tissues; (ii) the Thermoscript™, which is particularly relevant to polymerise cDNAs from RNAs with stable secondary structures such as the pre-miRs. In order to allow comparison between pre-miR expression levels and miRNA, miRNA target genes, RNase genes (*AGO1*: D7FQK3 and *DCL1*: D7FZW2), and reference genes *EEF1A2* (D7FZS6), *TUA* (D8LPR8) and *UBCE* (D7G3Z7), the Thermoscript-amplified cDNAs were used to test the transcript level of these molecules in salt stress conditions and in gametophyte tissues.

Real-time PCR

Oligonucleotide sequences were designed using Perl Primer (<http://perlprimer.sourceforge.net>) in the 3' UTR of the *AGO1* and *DCL1* mRNAs (Supplementary table 1A), in the coding sequence of

1 the predicted target mRNAs (Supplemental table 1B), and for the miRNAs (Suppl Table 1C) and
2 the pre-miRNAs (Suppl Table 1D). The RT-qPCR reactions were performed in a 96-well
3 thermocycler (Chromo4 System thermocycler; BioRad Laboratories) with SYBRgreen reaction mix
4 from ABgene (AB-1162/B; ABgene France, Courtabœuf), for 15 min at 95°C, followed by 41 runs
5 of 15 s. at 95°C, 30 s. at 60°C, and 30 s. at 72°C. Each sample was technically duplicated. The
6 amplification efficiency was tested using a dilution series of either genomic DNA (for the *AGO1*,
7 *DCLI*, miRNA target genes, and the pre-miRs; see below) or of poly-adenylated cDNAs (for the
8 miRNAs). The specificity of amplification was checked with a dissociation curve obtained by
9 heating the samples from 65°C to 95°C (measurement every 0.3°C). Experiments were carried out
10 on three independent biological replicates. Pre-miRs were amplified from equivalent amount of
11 cDNAs and RNAs, as well as on genomic DNA. As negative controls, non-reverse-transcribed
12 RNAs were used in addition to water. The PCR products were loaded on an 4% agarose gel, stained
13 with ethidium bromide, and was sequenced using a Sanger-based method on a ABI 3130XL Genetic
14 Analyser (Applied Biosystem, Life Technologies Corporation, USA).

25 Analysis of RT-qPCR data

26 The normalisation of the PCR signals corresponding to mRNAs (*AGO1*, *DCLI* and target genes)
27 and pre-miRs was conducted following Hellemans *et al.* (36), except that instead of working on Cq
28 values averaged over replicates, we normalised each measured amount to the amount of reference
29 genes (*EEF1A2*, *TUA* and *UBCE*) in the same replicate. MiRNAs were normalised using the whole
30 set of miRNAs as a reference set (as recommended in 37). The comparison between samples
31 (gametophyte, hypo-, hyper-saline) and the control (sporophytes in normal culture conditions) was
32 performed as a bidirectional t-test on the log-transformed normalised expression levels of the three
33 replicates, using the Welch correction for inequality of variances. Samples with a p-value < 0.05
34 were retained as significantly different from the control.

43 Genomic DNA extraction

44 *E. siliculosus* genomic DNA was prepared as described in Le Bail *et al.* (35) and was used as a
45 quantification reference for the RT-qPCR experiment and the calculation of PCR efficiency. A
46 dilution series ranging from 47 to 60730 copies (6 times dilution on 5 points) of the *E. siliculosus*
47 genomic DNA was prepared and tested for the *E. siliculosus* *AGO1*, *DCLI*, *EEF1A2*, *TUA* and
48 *UBCE* genes as well as for the pre-miR sequences.

54 Sequence conservation in the miRNA candidates

55 The known mature miRNAs sequences were obtained from miRBase
56 (<http://microrna.sanger.ac.uk>). We excluded the miRNA* sequences and extracted two subsets of
57
58
59
60

1 miRNA sequences, one corresponding to *Metazoa* and the other to *Viridiplantae*. We computed the
2 Levenshtein distances using our own Java implementation of the classical dynamic programming
3 algorithm. These computations were performed for the predicted miRNAs from sets Mir95 and
4 Mir90 vs the miRBase subsets, and for the miRBase subsets between them. We selected the lowest
5 score for each predicted *E. siliculosus* miRNA against *Metazoa* or *Viridiplantae*, the lowest score
6 for each plant miRNA against all *Metazoa*, and the lowest score for each metazoan miRNA against
7 all *Viridiplantae*. In order to compute Hausdorff distances, we also selected the lowest scores for
8 each *Metazoa* or *Viridiplantae* miRNA against the predicted *E. siliculosus* miRNA.
9

16 Protein domain analysis

17 The protein domains were searched using Interproscan (38), in the whole proteome of
18 *E. siliculosus*. For each domain, we compared the number of proteins containing at least one
19 instance of this domain among the whole proteome of *E. siliculosus*, and among the candidate
20 targets. The p-value retained to estimate the over-representation of a motif occurring in m proteins
21 in the whole proteome and in k proteins in the candidate targets was computed as $P_{N,n,m}(X \geq k)$ where
22 $P_{N,n,m}$ is the distribution function of the hypergeometric law, *ie* the probability to obtain at least k
23 positive instances within a sample of n individuals drawn out of a population of N individuals
24 containing a total of m positives instances. Protein domains were considered over-represented if the
25 p-value was lower than 0.05.
26
27
28
29
30
31
32

34 Results

37 *In silico* identification of miRNAs, pre-miRs and target candidates

38 The global strategy for the *in silico* analysis is shown in figure 1. We searched for miRNAs in two
39 sets of non-protein coding RNA sequences: (i) 12,798 “intergenic” sequences, *ie* sequences which
40 are neither in genes nor in rRNAs or tRNAs and (ii) 112,513 intron sequences. Each of these two
41 sets of sequences was analysed together with the 16,254 mRNAs of *E. siliculosus*, using the
42 findMiRNA software (31). The intergenic sequences were searched on both strands ($2 \times 68,917,369$
43 nucleotides), while for introns, only the transcript strand was used (78,816,594 nucleotides). As an
44 output, we obtained 864,679 results from the intergenic sequences, and 403,761 from the introns.
45 Each of these ~1.27 million potential miRNA candidates is associated to one folded precursor and
46 one target sequence located within a mRNA. The relationship between these three types of
47 molecules is however not univocal: the results contain 516,147 unique miRNAs, 697,657 unique
48 pre-miRs and 16,250 unique target mRNAs.
49
50
51
52
53
54
55
56

57 According to Ng Kwang Loong *et al.* (19), the pre-miRs exhibit values for several structural and
58 topological parameters which differ from those computed on other structured RNAs. We made use
59
60

1 of this feature to filter the potential candidates. To do so, we determined on tRNAs, rRNAs and
2 mRNAs, the distribution of values for five discriminant parameters: normalised minimum free
3 energy (Nmfe), normalised Shannon entropy (NQ), normalised base-pair distance (ND), degree of
4 compactness (NF) and normalised base-pairing propensity (Nbp). As we had no prior idea about the
5 number of microRNA genes in the *E. siliculosus* genome, we used the classical cut-off values of
6 95th and 90th percentiles to decide whether a potential pre-miR was sufficiently different from the
7 ncRNAs. Likewise, the reference energetic values we used were not extracted from literature, but
8 were computed on *E. siliculosus* RNAs. Thus, for each parameter, we computed two threshold
9 values, one corresponding to the 95th percentile, the other the 90th percentile of its distribution (see
10 table 1). Noteworthy, only the combination of all of the filters allowed a drastic reduction in the
11 number of candidates. Yet, the Nmfe appeared to be the most discriminant filter, as it allowed the
12 smallest number of pre-miRs to pass through. Two sets of pre-miRs were derived from these values:
13 (i) a pre-miR was retained in the set Pre95 when the values for the five parameters were beyond
14 their respective threshold, corresponding to the 95th percentile; (ii) a second set, named Pre90, was
15 similarly defined by reference to the 90th percentiles, but excluding those sequences already
16 contained in Pre95. Pre95 contained 597 pre-miRs, which were able to generate 568 different
17 miRNAs (set Mir95, see supplemental tables 2 and 3A,B), which in turn were predicted to interact
18 with 498 target mRNAs (set Tg95, see supplemental table 4). Similarly, Pre90 contained 943
19 pre-miRs, Mir90 is made of 922 miRNAs (supplemental tables 2 and 3A,B), predicted to interact
20 with 1153 target mRNAs (set Tg90, see supplemental table 4).

21 As an alternative procedure, we built a Support Vector Machine (SVM), similar to an approach
22 which has proven efficiency, for instance to predict human pre-miRs (39). When we applied this
23 method in *E. siliculosus*, the selected pre-miR candidates were not the same among repeats of the
24 procedure adjusted on different sets of “positive” and “negative” sequences. A thorough analysis of
25 the data, function parameters and results showed that this was probably due to the distribution of the
26 discriminant factors, which contained a significant number of extreme values. Therefore, the
27 approach based on cut-off filters set on quantiles (for a similar technique, see for instance 40)
28 appeared to be more suitable to these data.

29 **Experimental validation of the predictions**

30 The main drawback of the *in silico* approach is the large number of false positive instances it
31 produces (14). For this reason, we performed an experimental check, to evaluate the ratio of correct
32 predictions. It is expected that miRNA genes located in so-called “intergenic” regions are expressed
33 at a detectable level, while the rest of these regions should only be detected as “noise” in
34 quantitative detection experiments. Similarly, miRNAs issued from introns should be retained, in
35

1 contrast to regular introns which are destroyed after excision. In both cases, the detectable amount
2 of pre-miRs should be statistically distinguishable from the expression level of the regions from
3 which they are issued. We used the high-resolution transcriptome map to isolate, among the
4 intergenic or intronic sequences, those corresponding to the predicted pre-miR (both MiR90 and
5 MiR95). As a comparison, we also included data from the exonic regions and tRNAs. The
6 distribution of the expression levels of these various sets are shown on figure 2. In both intergenic
7 and intronic sequences, the expression level of the predicted pre-miRs is significantly higher than
8 the expression level of the other sequences (Student t-test, $\alpha=10^{-2}$). The intronic predicted miRNAs
9 have an expression level similar to that of exons. In addition, we observed that in both the
10 intergenic and the intronic sequences, the predicted pre-miRs in the set MiR95 are expressed at an
11 even higher level than those in the set MiR90 (not shown). These results showed that the sets of
12 predicted miRNAs were strongly biased towards highly expressed and stable RNA sequences,
13 which confirmed the statistical enrichment of these sets in actual miRNA sequences.
14
15

16
17
18
19
20
21
22
23
24 Among the 1488 miRNAs retained in Mir95 or Mir90, we extracted at random a subset of 36
25 sequences from each set, and quantified their expression by RT-qPCR (41) in *E. siliculosus*
26 sporophyte filaments. We could detect a specific expression characterised by a unique dissociation
27 curve with the expected half-dissociation temperature (T_m) for a total of 22 different miRNAs
28 (Table 2). Figure 3 shows the variable relative expression level (compared to tRNA-Leu) of the
29 detected miRNAs. As a control, we attempted to amplify 13 randomly chosen non-retained potential
30 miRNAs (*i.e.* instances filtered out from the structural filtering step). In agreement with the
31 predictions, we could not detect any of them (not shown). Among the 22 validated miRNAs, 16
32 were from the set Mir95, and 6 from the set Mir90. Thus, these data show that increasing the
33 threshold from the 90th to the 95th percentile enhances the ratio of experimental validation from
34 ~30% (22/72) to ~44% (16/36). Extrapolating this result to the whole prediction of 568 miRNAs in
35 Mir95 suggests that our most stringently filtered set contains ~252 valid miRNAs.
36
37
38
39
40
41
42
43

44 In order to reinforce the biological relevance of this experimental validation, we attempted to
45 detect the precursors of the 22 detected miRNAs, using a PCR-based approach (figure 4A; see
46 supplemental table 1D for pre-miR sequences and position of the oligonucleotides). We could detect
47 amplification of a PCR product for eight of them, the five most expressed ones being displayed in
48 figure 4B. Sequencing these PCR products showed that their primary structure was confirmed in all
49 cases.
50
51
52
53
54

55 **Genomic organisation and sequence conservation of the predicted miRNAs**

56 The experimental validation of 44% of predicted miRNAs in the set Mir95 allowed the
57 assumption that this pool of miRNAs and the corresponding pre-miRs was a relevant population to
58
59
60

1 investigate the genomic organisation and sequence conservation. Among the 597 predicted
2 pre-miRs in the set pre95, 407 (68.2%) came from intergenic regions and 190 (31.8%) from introns.
3 Similarly, in the set pre90, 313 precursors were intronic (33.2% of the 943). We also searched for
4 miRNA gene clusters, which we defined as three or more pre-miRs in a row (not necessarily on the
5 same strand), separated by no more than 5 kilo-bases. We identified three such clusters, each made
6 of three genes (figure 5A). In one case, two pre-miRs of the same cluster shared extensive similarity
7 (figure 5B). Altogether, these observations suggest that pre-miRs of *E. siliculosus* are not
8 predominantly organised as clusters.
9

10 In order to assess a putative sequence conservation between *E. siliculosus* miRNAs and those
11 found in other organisms, we compared the miRNA sequences in Mir95 and Mir90 to the whole
12 content of miRBase (17341 miRNAs, with 10099 different sequences). No identical sequence was
13 found. In order to estimate the proximity between these sequences, we computed the smallest
14 Levenshtein (edition) distance between each sequence in Mir95 or Mir90 and the miRBase entries
15 issued from *Metazoa* and *Viridiplantae*. The results in figure 6 show that the predicted miRNAs of
16 *E. siliculosus* were found to be different from both animal (figure 6A) and land plant (figure 6B)
17 miRNAs. The mean lowest distance from miRNAs in Mir95+90 to *Metazoa* (6.92) appeared to be
18 lower than to *Viridiplantae* (7.74), but this difference could not be interpreted as a higher sequence
19 similarity with *Metazoa* than with *Viridiplantae*. Instead, it was due to the fact that the number of
20 sequences was higher in the former (11,411) than in the latter (3,246), thus increasing the likelihood
21 that any sequence finds a more similar closest relative in *Metazoa* than in *Viridiplantae*. To enable
22 comparisons, we computed the mean lowest distance for *Viridiplantae vs Metazoa* (7.11) and
23 *Metazoa vs Viridiplantae* (7.72). As expected, the same bias was observed, while the pairwise
24 distances themselves were obviously the same. An other measure of the divergence between
25 sequence sets is the maximum distance between each element in one set and the closest element in
26 the other set, known as the “Hausdorff distance”. We computed that the Hausdorff distance between
27 Mir95+90 and the miRBase entries issued from *Metazoa* and *Viridiplantae* were 12 and 11,
28 respectively. The Hausdorff distance between *Metazoa* and *Viridiplantae* was 15. Again, the
29 differences between these values corresponded to an expected bias, as bigger sets have a higher
30 probability to contain at least one highly divergent sequence. Noticeably, the edition distances were
31 high compared to the length of the sequences considered, showing that in each of these sets, there
32 was at least one sequence displaying a high level of divergence with any miRNA in the other sets
33 with which the comparison was being performed. To summarise, these data showed that the
34 distance between the predicted *E. siliculosus* miRNAs and known miRNAs in *Metazoa* or
35 *Viridiplantae* were similar to the distance between the miRNAs of *Viridiplantae vs Metazoa*.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Biological function of miRNAs in *E. siliculosus*

In order to propose biological functions for *E. siliculosus* miRNAs, we first examined their target mRNAs. We searched for protein motifs which were over-represented in the set of predicted targets. The most significant results are shown in table 3 (a complete set of results is provided as supplemental table 5). We grouped the 15 over-represented patterns into 7 classes, according to their cellular function. The most represented classes were related to kinesin molecular motors and to tetratricopeptide repeats involved in nuclear protein import and mitotic spindle, suggesting altogether a role in nucleus organisation and dynamics (42, 43). Interestingly, proteins displaying an LNR domain (44) were also over-represented, suggesting that cell differentiation processes could be subject to a control by miRNAs.

As a second step, we searched for conditions able to induce or repress the expression of the pre-miRs, as well as the corresponding miRNAs and their target genes. In parallel, we studied the expression of *AGO1* and *DCL1*. Because *E. siliculosus* is a marine macro-alga, we tested two salt stress conditions on the sporophyte organism: hyper- and hypo-osmotic stresses. In addition, a morphologically different phase of the *E. siliculosus* life cycle, the gametophyte organism, was tested. The transcript level of *AGO1* and *DCL1* were significantly higher (Student t-test, $\alpha=0.05$) in response to hyper salt stress conditions. In contrast, both the hypo-osmotic stress and the gametophytic stage did not modify their expression (figure 7A). The expression profile of four pre-miRs was affected by growth conditions (figure 7B). A differential expression level was statistically validated for two miRNAs: one for a hyper- and one for a hypo-saline stress. In contrast, the difference of expression between the gametophytes and the sporophytes appeared to be lower. The changes in miRNA expression were not statistically supported (not shown). We also quantified the transcript level of the target genes, by using oligonucleotides downstream from the predicted miRNA recognition site. In these experimental conditions, we noticed that variations in expression were not statistically higher than inter-individual variations between the biological replicates (not shown).

Discussion

This study presents the first genome-wide scale list of candidate miRNAs for an organism of the heterokont phylum. Our *in silico* search for new miRNA candidates in *E. siliculosus* was based on structural considerations, without any *a priori* on the sequence conservation or on the expression level of mature and/or precursor RNAs. Many features which are usually used to identify miRNA precursors in plants or animals had to be discarded, because they were specific to one of these two

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

kingdoms (45). The initial search step, however, was performed using findmiRNA, a software designed to detect a nearly perfect complementarity between the miRNA candidates and their target(s), as it is usual in plants. The rationale behind this choice is that a less constrained search would have allowed to detect one or more target(s) for nearly any oligonucleotide with a length in the range expected for a miRNA (data not shown). Nevertheless, the higher expression level of the candidates, compared to sequences of the same origin (inter-genic or intronic), constituted an emerging property of the predicted set of microRNAs. Our experimental confirmation suggested that the selection based on structural features of the precursors was efficient in reducing the ratio of false positives. This work also illustrated the fact that a large scale prediction of miRNAs requires the combined use of computational and experimental analyses, whatever the order in which they are used (46). In any case, the relevance of the predictions relies on contextual information. For instance, any identification of miRNAs requires a clear distinction between coding (mRNAs) and non-coding (“intergenic” and intronic) sequences. Its accuracy is therefore strongly dependent on the initial assignment of nucleotides to these two sets. In particular, the annotation of UTRs can be critical, as miRNA target sites are expected to be found in UTRs (47, 48). This is, however, one of the most difficult tasks in the primary annotation of a newly sequenced genome, and its output is not fully reliable. We tried to overcome these impediments by re-assigning the regions flanking the first and/or last exon to the mRNA sequences, in the case where the mRNA was devoid of a 5'UTR and/or a 3'UTR. Although this procedure lies on the reasonable hypothesis that the structure of the unknown UTRs is similar to that of the experimentally observed ones, it might nevertheless add some errors (both false positive and false negative) to the analysis. In these conditions, and after experimentally validating by RT-qPCR a sub-set of 72 miRNAs candidates, we could extrapolate our prediction to a conservative number of 252 valid miRNAs. This number is likely underestimated as, in contrast to the *in silico* approach, the experimental detection of miRNAs is highly specific but suffers from a lack of sensitivity. Hence, many undetected candidates might be false negatives (15). In any case, the validation of each candidate and its implication in a given process would require a complex combination of *ad hoc* experiments.

47
48
49
50
51
52
53
54
55
56
57
58
59
60

The miRNAs identified in *E. siliculosus* display several specificities. First, they do not share significant sequence similarities with miRNAs already known in other species. Indeed, many miRNAs are species- or lineage-specific, and *E. siliculosus* is the first heterokont in which miRNAs are known. More precisely, the predicted miRNAs of *E. siliculosus* are as different from their closest animal miRNAs as plant miRNAs are, and as different from their closest plant miRNAs as animal miRNAs are. This is supported by the position of brown algae in the tree of life, distant from both the *opisthokonta* and the *archaeplastida*. Secondly, their position within the genome was peculiar. While in the metazoan species studied to date pre-miR clustering is frequent (49), and

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

examples are also known in plants (see for instance 50, 51), in *E. siliculosus* we found only three miRNA gene clusters, each comprising three genes. Interestingly, this low prevalence of miRNA gene clusters can be related to the low frequency of tandem repeats in the genome of *E. siliculosus* (13). In addition, only one of these clusters contained two pre-miRs sharing extensive similarity, whereas metazoan miRNA clusters are often made up of genes of the same family. From this criterion, the *E. siliculosus* miRNAs seem to be closer to plant miRNAs than to animal miRNAs. Conversely, about one third of the predicted miRNA precursors were located in introns of protein-coding genes, a feature shared with human miRNAs, which are intronic in 25% to 40% of the cases (52), in contrast to plant miRNAs (53). Finally, several pre-miRs were detected, suggesting that these molecules have a sufficient long lifetime, like most animal pre-miRs have, but usually not plant pre-miRs (54). Therefore, the genomic organisation and biogenesis of *E. siliculosus* miRNAs share features with either animals or plants, again illustrating its original evolutionary history.

In previously studied cases, the mechanisms by which miRNAs inhibit their target mRNA can be divided into two main classes: mRNA cleavage or translation repression. The AGO protein in the RISC complex is able to conduct mRNA cleavage if it contains a nucleolytic triad made of three conserved residues: Asp(760), Asp(846), Asp/His(986) (numbering of *A. thaliana* AGO1) (55). The AGO protein of *E. siliculosus* does contain these three residues, namely Asp(703), Asp(775), His(912). It is thus expected to perform the endonucleolytic cleavage of the target mRNA. Despite numerous attempts (RACE-PCR on the target mRNA candidates), cleavage of the predicted target genes by miRNAs could not be demonstrated (data not shown). These negative results do not allow to rule out the possibility that miRNAs in *E. siliculosus* direct the cleavage of their target. However, future work should consider the hypothesis that, although the required residues are present in AGO, the mechanism by which miRNAs regulate their targets in *E. siliculosus* might rely on translation inhibition rather than on mRNA cleavage. A similar situation has been shown to occur in human (56). The actual mechanism of this effect remains to be demonstrated.

Expression studies performed by RT-qPCR revealed a possible involvement of the miRNA machinery in physiological processes. The four operating levels of this machinery corresponding to (i) the RISC RNase Argonaute and the RNase DICER, (ii) the pre-miR, (iii) the miRNA and (iv) the miRNA targets, were investigated using this approach. Both Argonaute (*AGO1*) and DICER (*DCLI*) genes and some of the pre-miRs tested were induced in response to a modification in salt concentration. These changes in *AGO1* and *DCLI* expression upon stress distinguish *E. siliculosus* from other organisms. In animals, various stress conditions result in a decrease in *DCLI* expression (57). In plants, the expression pattern of miRNA related proteins is often complex, as it involves multi-copy genes, with divergent expression patterns within each family (58). For instance, in

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

A. thaliana, all four *DCL* transcript levels are depleted upon a salt stress, but each gene exhibits a distinct time course and intensity for this regulation (59). In *Oryza sativa*, only one *AGO* gene among nineteen is induced in stress conditions, while none of the eight *DCL* genes is affected (60). Although a systematic study in a broad range of animals and plants remains to be conducted, the induction of *DCL* by a salt stress seems to be an exception, and might be a distinctive feature of this marine brown algae, along with the existence of one single instance of *AGO1* and *DCLI* genes. This can be related to the particular environmental conditions macro-algae have to face. Indeed, these organisms are living attached to rocks or to other algae, and hence are subject to salt concentration variations during the day depending on tides and evaporation. We could not validate any differential expression of the corresponding target mRNAs. Therefore, if a miRNA-regulated process is involved in response to salt stress, we propose that its mechanism should rely on translational repression rather than cleavage and degradation of the target. In addition to a role in salt stress, miRNA-mediated gene expression regulation could also be involved in developmental processes. The target prediction by findmiRNA allows for substantial mismatches in miRNA-mRNA base pairing. This is in agreement with the addressing of the RISC complex to its target. However, this loose constraint might generate many false positive targets. For this reason, we cannot analyse individual targets in the whole set of predictions, most of which have not been experimentally validated. However, unless there are reasons to suspect a coincidental specific bias towards a given class of proteins, the over-representation of a process among the predicted target can point to possible roles for miRNAs in cell processes. The predicted targets of *E. siliculosus* miRNAs represent a wide variety of functions and protein families, notably involved in nucleus dynamics, cell polarity and differentiation. These molecular functions, extensively conserved through the tree of life have been shown to be regulated by miRNAs in other organisms: kinesins are regulated by miRNAs in human (61), like the NB-ARC containing protein APAF1 (62). Similarly, expressed proteins containing kinesin and NB-ARC domains have been predicted to be regulated by miRNAs in *A. thaliana* and *O. sativa* (63). MiRNAs also regulate methyl-transferases (64) and Notch-related proteins (65). Interestingly, the recent characterisation of the *E. siliculosus* morphogenetic mutant “*étoile*” supported a role of Notch-related proteins in cell differentiation, with a mechanism which remains to be identified (66). In addition, we predicted that miRNAs could be able to regulate some of the functional families identified as stress-responsive by a transcriptomic study in *E. siliculosus* (67). Finally, in contrast to the already known miRNAs in *E. siliculosus* (13), we did not identify any bias towards the proteins containing Leucine-rich repeats.

55
56
57
58
59
60

In summary, the list of *E. siliculosus* miRNAs proposed in this study is a solid starting point for further investigations aiming at deciphering in detail their biological roles, as well as the molecular mechanisms by which they operate. In this perspective, brown algae represent a source of novelty

1 because of their extraneous phylogenetic position.
2
3

4 **Accession numbers**

5
6
7 The validated pre-miRNAs were deposited in miRBase under the following accession numbers (see
8 also Supplementary table 1D): pre95_0213a = esi-MIR8618b; pre95_0055a = esi-MIR8619;
9 pre95_0064a = esi-MIR8620; pre95_0207a = esi-MIR8622b; pre90_0829a = esi-MIR8623b;
10 pre90_0257a = esi-MIR8623d; pre95_0365a = esi-MIR8624a; pre95_0400a = esi-MIR8625.
11
12
13

14 **Supplementary data**

15
16
17 Supplementary Data are available at NAR online: Supplementary tables 1-5.
18

19 **Acknowledgements**

20
21
22 The Interproscan analysis was run on the the ABIMS platform (Station Biologique de Roscoff).
23 The authors wish to thank Marie-Hélène Mucchielli-Giorgi for her expertise in SVM analysis,
24 Martine Boccara (ENS, Paris), Martin Crespi, Christine Lelandais-Brière and Florian Frugier (ISV,
25 CNRS Gif-sur-Yvette) for fruitful discussion and critical reading of the manuscript.
26
27
28
29

30 **References**

- 31
32 1. Lau, N.C., Lim, L.P., Weinstein, E.G. and Bartel, D.P. (2001) An abundant class of tiny RNAs
33 with probable regulatory roles in *Caenorhabditis elegans*. *Science*, **294**, 858-862.
34
- 35
36 2. Lee, R.C. and Ambros, V. (2001) An extensive class of small RNAs in *Caenorhabditis elegans*.
37 *Science*, **294**, 862-864.
38
- 39
40 3. Griffiths-Jones, S., Saini, H.K., van Dongen, S. and Enright, A.J. (2008) miRBase: tools for
41 microRNA genomics. *Nucleic Acids Res*, **36**, D154-8.
42
- 43
44 4. Axtell, M.J., Westholm, J.O. and Lai, E.C. (2011) Vive la différence: biogenesis and evolution of
45 microRNAs in plants and animals. *Genome Biol*, **12**, 221.
46
- 47
48 5. Yoon, H.S., Grant, J., Tekle, Y.I., Wu, M., Chaon, B.C., Cole, J.C., Logsdon, J.M.J., Patterson,
49 D.J., Bhattacharya, D. and Katz, L.A. (2008) Broadly sampled multigene trees of eukaryotes.
50 *BMC Evol Biol*, **8**, 14.
51
- 52
53 6. Haas, B.J., Kamoun, S., Zody, M.C., Jiang, R.H.Y., Handsaker, R.E., Cano, L.M., Grabherr, M.,
54 Kodira, C.D., Raffaele, S., Torto-Alalibo, T. et al. (2009) Genome sequence and analysis of the
55 Irish potato famine pathogen *Phytophthora infestans*. *Nature*, **461**, 393-398.
56
- 57
58 7. Baxter, L., Tripathy, S., Ishaque, N., Boot, N., Cabral, A., Kemen, E., Thines, M., Ah-Fong, A.,
59
60

- 1 Anderson, R., Badejoko, W. et al. (2010) Signatures of adaptation to obligate biotrophy in the
2 Hyaloperonospora arabidopsidis genome. *Science*, **330**, 1549-1551.
- 3
4
5 8. Lévesque, C.A., Brouwer, H., Cano, L., Hamilton, J.P., Holt, C., Huitema, E., Raffaele, S.,
6 Robideau, G.P., Thines, M., Win, J. et al. (2010) Genome sequence of the necrotrophic plant
7 pathogen *Pythium ultimum* reveals original pathogenicity mechanisms and effector repertoire.
8 *Genome Biol*, **11**, R73.
- 9
10
11
12 9. Armbrust, E.V., Berges, J.A., Bowler, C., Green, B.R., Martinez, D., Putnam, N.H., Zhou, S.,
13 Allen, A.E., Apt, K.E., Bechner, M. et al. (2004) The genome of the diatom *Thalassiosira*
14 *pseudonana*: ecology, evolution, and metabolism. *Science*, **306**, 79-86.
- 15
16
17
18 10. Bowler, C., Allen, A.E., Badger, J.H., Grimwood, J., Jabbari, K., Kuo, A., Maheswari, U.,
19 Martens, C., Maumus, F., Otiillar, R.P. et al. (2008) The *Phaeodactylum* genome reveals the
20 evolutionary history of diatom genomes. *Nature*, **456**, 239-244.
- 21
22
23
24 11. Huang, A., He, L. and Wang, G. (2011) Identification and characterization of microRNAs from
25 *Phaeodactylum tricornutum* by high-throughput sequencing and bioinformatics analysis. *BMC*
26 *Genomics*, **12**, 337.
- 27
28
29 12. Charrier, B., Coelho, S.M., Le Bail, A., Tonon, T., Michel, G., Potin, P., Kloareg, B., Boyen, C.,
30 Peters, A.F. and Cock, J.M. (2008) Development and physiology of the brown alga *Ectocarpus*
31 *siliculosus*: two centuries of research. *New Phytol*, **177**, 319-332.
- 32
33
34 13. Cock, J.M., Sterck, L., Rouzé, P., Scornet, D., Allen, A.E., Amoutzias, G., Anthouard, V.,
35 Artiguenave, F., Aury, J., Badger, J.H. et al. (2010) The *Ectocarpus* genome and the independent
36 evolution of multicellularity in brown algae. *Nature*, **465**, 617-621.
- 37
38
39 14. Lindow, M. and Gorodkin, J. (2007) Principles and limitations of computational microRNA
40 gene and target finding. *DNA Cell Biol*, **26**, 339-351.
- 41
42
43 15. Mendes, N.D., Freitas, A.T. and Sagot, M. (2009) Current tools for the identification of miRNA
44 genes and their targets. *Nucleic Acids Res*, **37**, 2419-2433.
- 45
46
47 16. Allmer, J. and Yousef, M. (2012) Computational methods for ab initio detection of microRNAs.
48 *Front Genet*, **3**, 209.
- 49
50
51 17. Lim, L.P., Lau, N.C., Weinstein, E.G., Abdelhakim, A., Yekta, S., Rhoades, M.W., Burge, C.B.
52 and Bartel, D.P. (2003) The microRNAs of *Caenorhabditis elegans*. *Genes Dev*, **17**, 991-1008.
- 53
54
55 18. Lai, E.C., Tomancak, P., Williams, R.W. and Rubin, G.M. (2003) Computational identification
56 of *Drosophila* microRNA genes. *Genome Biol*, **4**, R42.
- 57
58
59 19. Ng Kwang Loong, S. and Mishra, S.K. (2007) Unique folding of precursor microRNAs:
60

- quantitative evidence and implications for de novo identification. *RNA*, **13**, 170-187.
20. Hutvagner, G. and Zamore, P.D. (2002) A microRNA in a multiple-turnover RNAi enzyme complex. *Science*, **297**, 2056-2060.
21. Vazquez, F., Blevins, T., Ailhas, J., Boller, T. and Meins, F.J. (2008) Evolution of Arabidopsis MIR genes generates novel microRNA classes. *Nucleic Acids Res*, **36**, 6429-6438.
22. Stark, A., Brennecke, J., Russell, R.B. and Cohen, S.M. (2003) Identification of Drosophila MicroRNA targets. *PLoS Biol*, **1**, E60.
23. Lim, L.P., Lau, N.C., Garrett-Engele, P., Grimson, A., Schelter, J.M., Castle, J., Bartel, D.P., Linsley, P.S. and Johnson, J.M. (2005) Microarray analysis shows that some microRNAs downregulate large numbers of target mRNAs. *Nature*, **433**, 769-773.
24. Watanabe, Y., Yachie, N., Numata, K., Saito, R., Kanai, A. and Tomita, M. (2006) Computational analysis of microRNA targets in *Caenorhabditis elegans*. *Gene*, **365**, 2-10.
25. Vasudevan, S., Tong, Y. and Steitz, J.A. (2007) Switching from repression to activation: microRNAs can up-regulate translation. *Science*, **318**, 1931-1934.
26. Gu, S. and Kay, M.A. (2010) How do miRNAs mediate translational repression?. *Silence*, **1**, 11.
27. Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W. and Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*, **25**, 3389-3402.
28. Phillips, N., Burrowes, R., Rousseau, F., de Reviere, B. and Saunders, G.W. (2008) Resolving evolutionary relationships among the brown algae using chloroplast and nuclear genes. *Journal of Phycology*, **44**, 394-405.
29. Kawai, H., Muto, H., Fujii, T. and Kato, A. (1995) A linked 5S rRNA gene in Scytosiphon lomentaria (scytosiphonales, phaeophyceae). *Journal of Phycology*, **31**, 306-311.
30. Lowe, T.M. and Eddy, S.R. (1997) tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res*, **25**, 955-964.
31. Adai, A., Johnson, C., Mlotshwa, S., Archer-Evans, S., Manocha, V., Vance, V. and Sundaresan, V. (2005) Computational prediction of miRNAs in *Arabidopsis thaliana*. *Genome Research*, **15**, 78-91.
32. Le Bail, A., Billoud, B., Maisonneuve, C., Peters, A., Cock, M. and Charrier, B. (2008) Early development pattern of the brown alga *Ectocarpus siliculosus* (ectocarpales, phaeophyceae) sporophyte. *Journal of Phycology*, **44**, 1269-1281.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
33. Hofacker, I., Fontana, W., Stadler, P., Bonhoeffer, L., Tacker, M. and Schuster, P. (1994) Fast Folding and Comparison of RNA Secondary Structures. *Monatshefte für Chemie*, **125**, 167-188.
34. R Development Core Team (2010) R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria.
35. Le Bail, A., Dittami, S.M., de Franco, P., Rousvoal, S., Cock, M.J., Tonon, T. and Charrier, B. (2008) Normalisation genes for expression analyses in the brown alga model *Ectocarpus siliculosus*. *BMC Mol Biol*, **9**, 75.
36. Hellemans, J., Mortier, G., De Paepe, A., Speleman, F. and Vandesompele, J. (2007) qBase relative quantification framework and software for management and automated analysis of real-time quantitative PCR data. *Genome Biol*, **8**, R19.
37. Mestdagh, P., Van Vlierberghe, P., De Weer, A., Muth, D., Westermann, F., Speleman, F. and Vandesompele, J. (2009) A novel and universal method for microRNA RT-qPCR data normalization. *Genome Biol*, **10**, R64.
38. Zdobnov, E.M. and Apweiler, R. (2001) InterProScan--an integration platform for the signature-recognition methods in InterPro. *Bioinformatics*, **17**, 847-848.
39. Ng Kwang Loong, S. and Mishra, S.K. (2007) De novo SVM classification of precursor microRNAs from genomic pseudo hairpins using global and intrinsic folding measures. *Bioinformatics*, **23**, 1321-1330.
40. Cakir, M. and Allmer, J. (2010) Systematic computational analysis of potential RNAi regulation in *Toxoplasma gondii*. *5th International Symposium on Health Informatics and Bioinformatics (HIBIT)*, **5**, 31-38.
41. Shi, R. and Chiang, V.L. (2005) Facile means for quantifying microRNA expression by real-time PCR. *Biotechniques*, **39**, 519-525.
42. Paddy, M.R. (1998) The Tpr protein: linking structure and function in the nuclear interior?. *Am J Hum Genet*, **63**, 305-310.
43. Tikhonenko, I., Nag, D.K., Robinson, D.N. and Koonce, M.P. (2009) Microtubule-nucleus interactions in *Dictyostelium discoideum* mediated by central motor kinesins. *Eukaryot Cell*, **8**, 723-731.
44. Fiúza, U. and Arias, A.M. (2007) Cell and molecular biology of Notch. *J Endocrinol*, **194**, 459-474.
45. Bologna, N.G., Schapire, A.L. and Palatnik, J.F. (2013) Processing of plant microRNA precursors. *Brief Funct Genomics*, **12**, 37-45.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
46. Mathelier, A. and Carbone, A. (2010) MIRENA: finding microRNAs with high accuracy and no learning at genome scale and from deep sequencing data. *Bioinformatics*, **26**, 2226-2234.
 47. Zhou, X., Duan, X., Qian, J. and Li, F. (2009) Abundant conserved microRNA target sites in the 5'-untranslated region and coding sequence. *Genetica*, **137**, 159-164.
 48. Brodersen, P. and Voinnet, O. (2009) Revisiting the principles of microRNA target recognition and mode of action. *Nat Rev Mol Cell Biol*, **10**, 141-148.
 49. Shomron, N., Golan, D. and Hornstein, E. (2009) An evolutionary perspective of animal microRNAs and their targets. *J Biomed Biotechnol*, **2009**, 594738.
 50. Cui, X., Xu, S.M., Mu, D.S. and Yang, Z.M. (2009) Genomic analysis of rice microRNA promoters and clusters. *Gene*, **431**, 61-66.
 51. Merchan, F., Boualem, A., Crespi, M. and Frugier, F. (2009) Plant polycistronic precursors containing non-homologous microRNAs target transcripts encoding functionally related proteins. *Genome Biol*, **10**, R136.
 52. Rodriguez, A., Griffiths-Jones, S., Ashurst, J.L. and Bradley, A. (2004) Identification of mammalian microRNA host genes and transcription units. *Genome Res*, **14**, 1902-1910.
 53. Zhu, Q., Spriggs, A., Matthew, L., Fan, L., Kennedy, G., Gubler, F. and Helliwell, C. (2008) A diverse set of microRNAs and microRNA-like small RNAs in developing rice grains. *Genome Res*, **18**, 1456-1465.
 54. Jones-Rhoades, M.W., Bartel, D.P. and Bartel, B. (2006) MicroRNAs and their regulatory roles in plants. *Annu Rev Plant Biol*, **57**, 19-53.
 55. Liu, J., Carmell, M.A., Rivas, F.V., Marsden, C.G., Thomson, J.M., Song, J., Hammond, S.M., Joshua-Tor, L. and Hannon, G.J. (2004) Argonaute2 is the catalytic engine of mammalian RNAi. *Science*, **305**, 1437-1441.
 56. Pillai, R.S., Artus, C.G. and Filipowicz, W. (2004) Tethering of human Ago proteins to mRNA mimics the miRNA-mediated repression of protein synthesis. *RNA*, **10**, 1518-1525.
 57. Wiesen, J.L. and Tomasi, T.B. (2009) Dicer is regulated by cellular stresses and interferons. *Mol Immunol*, **46**, 1222-1228.
 58. Khraiwesh, B., Zhu, J. and Zhu, J. (2012) Role of miRNAs and siRNAs in biotic and abiotic stress responses of plants. *Biochim Biophys Acta*, **1819**, 137-148.
 59. Liu, Q., Feng, Y. and Zhu, Z. (2009) Dicer-like (DCL) proteins in plants. *Funct Integr Genomics*, **9**, 277-286.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
60. Kapoor, M., Arora, R., Lama, T., Nijhawan, A., Khurana, J.P., Tyagi, A.K. and Kapoor, S. (2008) Genome-wide identification, organization and phylogenetic analysis of Dicer-like, Argonaute and RNA-dependent RNA Polymerase gene families and their expression analysis during reproductive development and stress in rice. *BMC Genomics*, **9**, 451.
61. Li, G, Luna, C., Qiu, J., Epstein, D.L. and Gonzalez, P. (2010) Targeting of integrin beta1 and kinesin 2alpha by microRNA 183. *J Biol Chem*, **285**, 5461-5471.
62. Frankel, L.B., Christoffersen, N.R., Jacobsen, A., Lindow, M., Krogh, A. and Lund, A.H. (2008) Programmed cell death 4 (PDCD4) is an important functional target of the microRNA miR-21 in breast cancer cells. *J Biol Chem*, **283**, 1026-1033.
63. Lindow, M., Jacobsen, A., Nygaard, S., Mang, Y. and Krogh, A. (2007) Intragenomic matching reveals a huge potential for miRNA-mediated regulation in plants. *PLoS Comput Biol*, **3**, e238.
64. Kisliouk, T., Yosefi, S. and Meiri, N. (2010) MiR-138 inhibits EZH2 methyltransferase expression and methylation of histone H3 at lysine 27, and affects thermotolerance acquisition. *Eur J Neurosci*, **33**, 224-235.
65. Vallejo, D.M., Caparros, E. and Dominguez, M. (2011) Targeting Notch signalling by the conserved miR-8/200 microRNA family in development and cancer cells. *EMBO J*, **30**, 756-769.
66. Le Bail, A., Billoud, B., Le Panse, S., Chenivresse, S. and Charrier, B. (2011) ETOILE regulates developmental patterning in the filamentous brown alga *Ectocarpus siliculosus*. *The Plant Cell*, **23**, 1666-1678.
67. Dittami, S.M., Scornet, D., Petit, J., Ségurens, B., Da Silva, C., Corre, E., Dondrup, M., Glatting, K., König, R., Sterck, L. et al. (2009) Global expression analysis of the brown alga *Ectocarpus siliculosus* (Phaeophyceae) reveals large-scale reprogramming of the transcriptome in response to abiotic stress. *Genome Biol*, **10**, R66.

45 Figures legends

46
47
48
49
50
51
52
53

Figure 1. Flowchart of the *in silico* analysis. Data are represented in boxes, processes in ovals. Colour code: black: initial data; red: mRNAs or parts of mRNAs; green: intergenic / non-protein coding RNA; grey: pre-existing software; purple: our software. The resulting *in silico* predictions were then tested for experimental validation.

54
55
56
57
58
59
60

Figure 2. Expression level of various sequence sets extracted from a whole-genome tiling array experiment. The predicted pre-miRNA (“miR cand.”, lines 2 and 4) are compared to the other sequences (“Not cand.”, lines 1 and 3) of the same genomic origin (“Intergenic” or “Intronic”). Exons and tRNAs are also displayed. Each “miR cand.” line is significantly higher than the

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

corresponding “Not cand.” line.

Figure 3. Detection by reverse-transcriptase quantitative PCR of 22 candidate miRNAs in *E. siliculosus* sporophyte tissues grown in normal culture conditions. The expression level was normalised to the level of tRNA-Leu.

Figure 4. Pre-miR detection by RT-qPCR amplification. (A) Example of primer design for PCR assay. Oligonucleotide used for the RT-qPCR are shown as arrows, and the predicted miRNA sequence is shaded in grey. To ensure a high specificity, one of the oligonucleotide was designed, when possible, in the terminal loop. (B) Visualisation of assay results. For each predicted precursor, the PCR assay was conducted from H₂O, gDNA (2.3 ng), total RNAs (40 ng) and cDNAs (40 ng equivalent of total RNAs) and run on a 4% agarose gel. Their predicted size in base pairs is indicated in the right side of the figure. Left-overs of primers or primer dimers are visible in the bottom part of each photography.

Figure 5. MiRNA clusters in *E. siliculosus* genome. (A) Three clusters were predicted, on three different super-contigs (sctg_XX). (B) Two of the predicted miRNAs composing the second cluster share extensive sequence similarity.

Figure 6. Comparison between miRNA candidates and miRBase. Each chart shows the distribution of the Leveshtein distance for each sequence to the closest entry in a subset of miRBase. (A) Distance for *Viridiplantae* and *E. siliculosus* to *Metazoa*; (B) Distance for *Metazoa* and *E. siliculosus* to *Viridiplantae*.

Figure 7. Transcript levels of miRNA processing proteins, pre-miRs, miRNAs and target genes measured by RT-qPCR in different algal materials. (A) *AGO1* and *DCL1* cDNAs. (B) Pre-miRs detected in figure 4. The data are expressed as fold changes relative to the control. “Control”: sporophytes grown in normal culture conditions, used to set the reference value at 1; “Gametophyte”: gametophytes grown in normal culture conditions; “Hyposaline” and “Hypersaline”: sporophytes subjected to corresponding osmotic stress conditions. Conditions for which the distribution of replicates values significantly differs from the control distribution (Student t-test, $\alpha=0.05$) are denoted by a star.

Table captions

Table 1. Secondary structure filter on pre-miRs. A valid pre-miR is expected to have a lower Nmfe, NQ, ND, NF and a higher Npb than other structured RNAs. For each parameter (rows 1-5), the filtering threshold “Filter n” was defined as the lowest (except Npb: highest) nth percentile of the value distribution among the three RNA reference sets: tRNAs, rRNAs and mRNAs. The number of pre-miRNAs passing each individual filter is indicated on their respective row in the last column. The 6th row shows the number of pre-miRs passing all the filters of the two values of n presented: 95 and 90.

Table 2. Experimentally validated miRNAs. The first part of the name, “miR95” or “miR90”, indicates the set from which the candidate was drawn (see table 1). The mature sequence of each miRNA is shown, together with the list of predicted target proteins, for which the Uniprot accession numbers are indicated.

Table 3. Over-representation of functional motifs in the predicted targets. Motifs were grouped by similar function, and groups were sorted by ascending best p-value. “Proteome” and “Targets” show the number of proteins containing at least one instance of the motif in the whole genome and in the set of targets which we predicted *in silico*, respectively. The “Over-representation” is the ratio of the two previous columns, each normalised to its respective total number of proteins in the set. The p-value is computed using the hypergeometric probability law.

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus***Table 1**

Parameter	Short	Filter 90	# passing 90	Filter 95	# passing 95
Normalised minimum free energy	Nmfe	< -0.750	86485	< -0.826	43885
Normalised Shanon entropy	NQ	< 0.0565	73518	< 0.0441	50731
Normalised base-pair distance	ND	< 0.0244	88088	< 0.0175	52414
Degree of compactness	NF	< 0.382	350532	< 0.325	332241
Normalised base-pairing propensity	Npb	> 0.364	162473	> 0.369	129715
<i>All five filters</i>			1540		597

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Table 2

miRNA name	miRNA sequence	Target(s)				
miR95_0055	GUAGCCAAGAUCGGGUGCACCCGG	D7FJK8				
miR95_0064	UCCUGCAGCCGUGCCGC	D8LF98	D8LGC5			
miR95_0076	ACUUCGCCGCCGCGCAG	D7FT76				
miR95_0168	GUCCCAUGGUGUCCCAUGGGAC	D7FWG6	D8LJ42			
miR95_0195	CCAAGAUCGGGUGCACCCGGUGGAA	D7FWA3				
miR95_0207	CCGCUUGGUUCCGGCGCAGAUUUCG	D8LFK3	D7FWF7	D7FRZ9		
miR95_0213	UCGUACAGUGGCUCAGCCUC	D7FWM6	D7G720			
miR95_0257	UCUCAUGGGCGCCCAUGGG	D7G179				
miR95_0295	GCCGCCGUUCCUGCAGCCGC	D8LS33				
miR95_0324	GUAUGGCCUGUGACCGCUCG	D7G722				
miR95_0365	AGUUGAUAGCUUCCGAGAAUCAGUG	D8LID2				
miR95_0400	AACCAAGAUGGCCUGGAUUUGCG	D8LNB1	D8LNB2	D8LNB6	D7FM32	D8LNC7
miR95_0424	UUGGUUCCGGCGUAGAUC	D7G7S8				
miR95_0439	GUCCGCAGCCGCCGCGC	D8LBP6	D8LXB5	D7FMJ4	D7FGT5	D8LJR3
miR95_0483	AAGUUGAUAGCUUCCGAGAAUCAGU	D8LID2				
miR95_0485	CUGGCGCCCCAGGGCGG	D7FQ65	D8LEB3	D7FML8	D7FI70	
miR90_0079	UAUAGAAGCCGAAAUCAA	D7G3H8	D7FRR2			
miR90_0193	CAUGGAACUCCAUGGAACUCCAUG	D7FM32				
miR90_0201	CCACAGCAGUACCACUAGCUUCA	D7G014	D8LGE5			
miR90_0682	CGCCGACGGUUGCCGUGC	D8LLD2				
miR90_0796	GUCUGGGGACUAUGUUAACC	D7FPI3				
miR90_0805	AGCAAAGUUGAUAGCUUCCGAG	D8LID2				

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus***Table 3**

Bank	Name	Proteome	Targets	Over-rep	P-value	Function
Panther	PTHR19959:SF16	59	22	4.00	5.62E-09	Kinesin
Fprint	PR00381	38	17	4.80	1.20E-08	Kinesin light chain
Panther	PTHR19959	80	23	3.08	6.27E-07	Structural constituent of cytoskeleton
Pfam	PF00931	31	12	4.15	1.09E-05	NB-ARC domain
superfamily	SSF48452	231	41	1.90	4.23E-05	TPR-like superfamily
Pfam	PF07721	21	9	4.59	5.43E-05	Tetratricopeptide repeat
Smart	SM00028	124	26	2.25	6.43E-05	Tetratricopeptide repeats
Prosite	PS50293	166	31	2.00	1.35E-04	TPR repeat region
Pfam	PF00515	61	14	2.46	1.21E-03	Tetratricopeptide repeat
Gene3D	G3DSA:1.25.40.10	326	45	1.48	4.94E-03	TPR-like_helical
Pfam	PF08241	49	12	2.63	1.46E-03	Methyltransferase
Panther	PTHR10108	33	9	2.92	2.55E-03	Methyltransferase
Pfam	PF00066	31	9	3.11	1.58E-03	LNR domain
Prosite	PS00120	13	5	4.12	4.76E-03	Lipases, serine active site
Gene3D	G3DSA:2.160.20.10	98	17	1.86	8.82E-03	Pectin lyase-like

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Figure 1

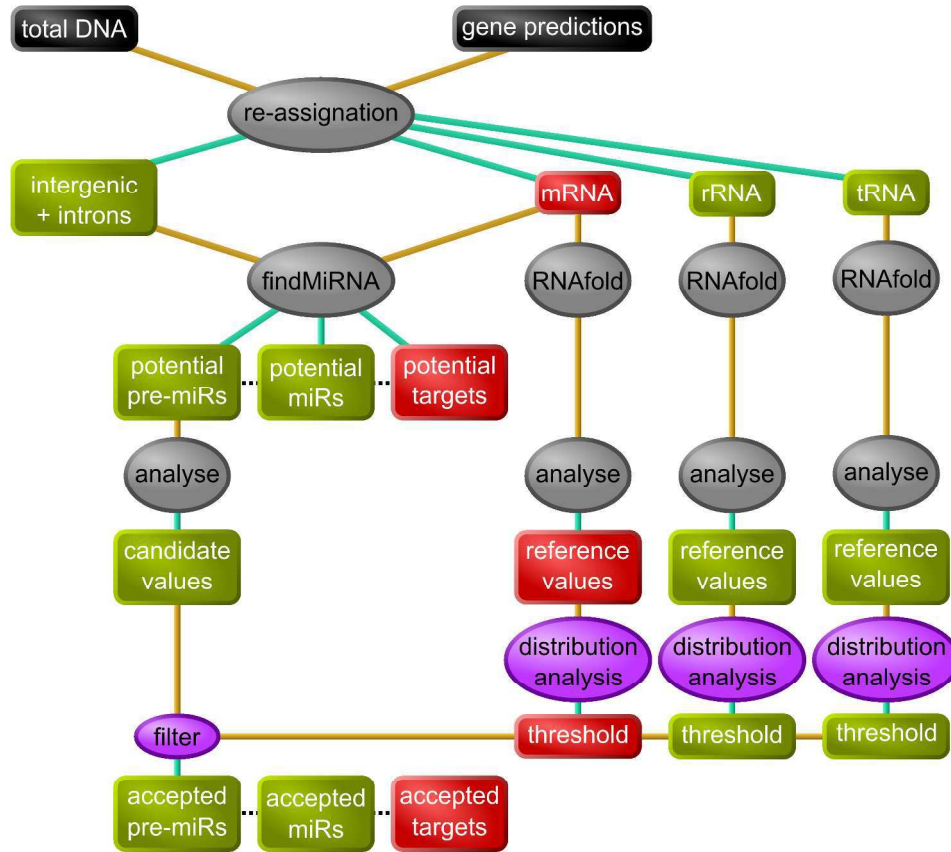


Figure 1. Flowchart of the in silico analysis. Data are represented in boxes, processes in ovals. Colour code: black: initial data; red: mRNAs or parts of mRNAs; green: intergenic / non-protein coding RNA; grey: pre-existing software; purple: our software. The resulting in silico predictions were then tested for experimental validation.

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Figure 2

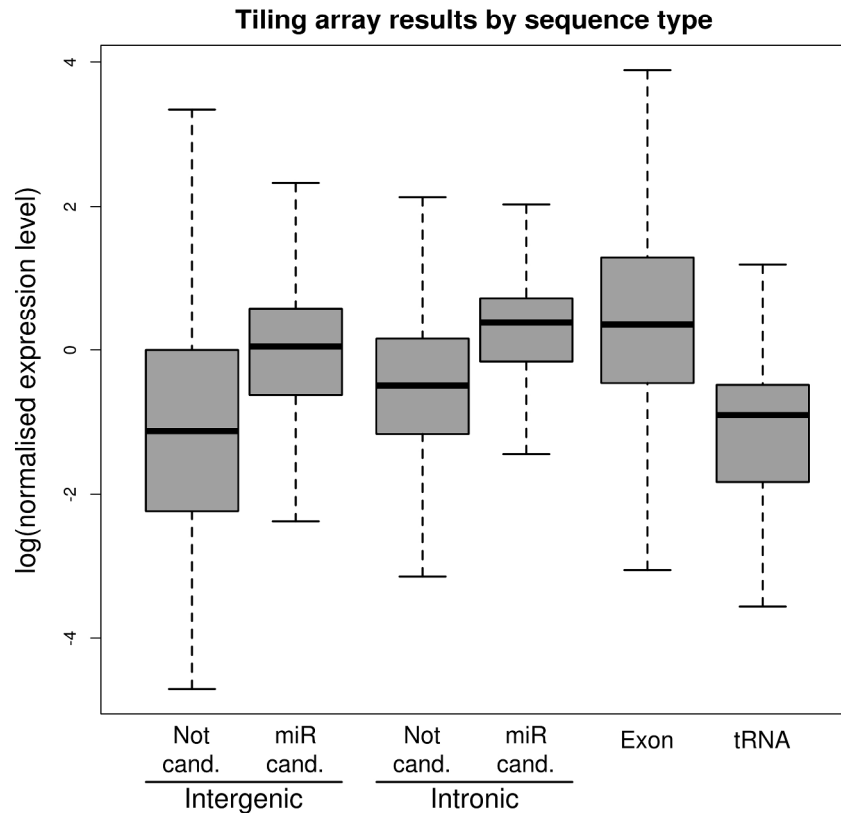


Figure 2. Expression level of various sequence sets extracted from a whole-genome tiling array experiment. The predicted pre-miRNA ("miR cand.", lines 2 and 4) are compared to the other sequences ("Not cand.", lines 1 and 3) of the same genomic origin ("Intergenic" or "Intronic"). Exons and tRNAs are also displayed. Each "miR cand." line is significantly higher than the corresponding "Not cand." line.

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Figure 3

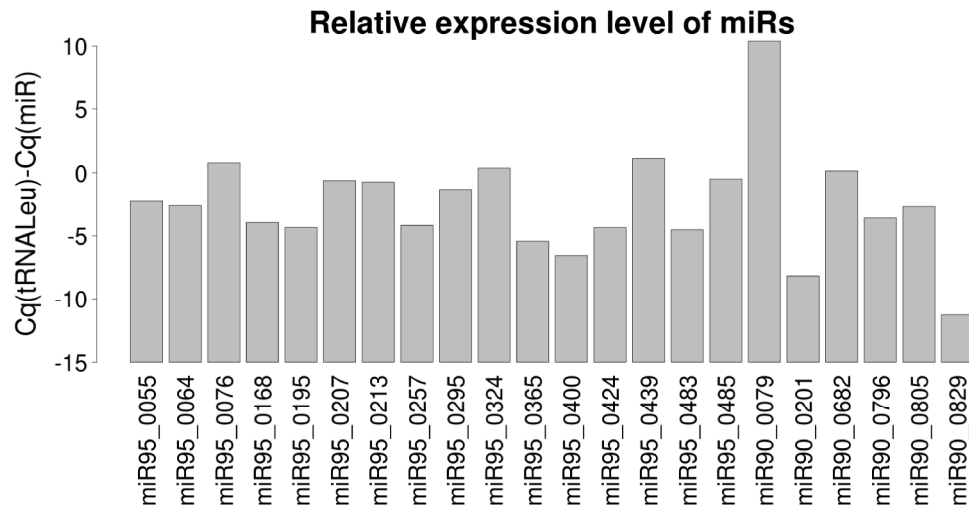


Figure 3. Detection by reverse-transcriptase quantitative PCR of 22 candidate miRNAs in *E. siliculosus* sporophyte tissues grown in normal culture conditions. The expression level was normalised to the level of tRNA-Leu.

Review

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Figure 4

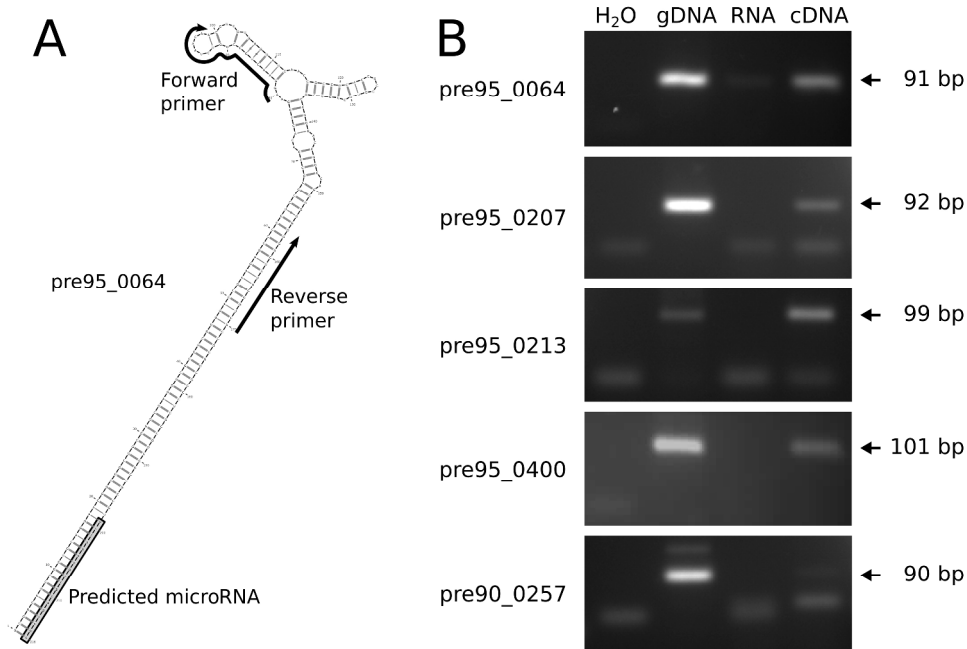


Figure 4. Pre-miR detection by RT-qPCR amplification. (A) Example of primer design for PCR assay. Oligonucleotide used for the RT-qPCR are shown as arrows, and the predicted miRNA sequence is shaded in grey. To ensure a high specificity, one of the oligonucleotide was designed, when possible, in the terminal loop. (B) Visualisation of assay results. For each predicted precursor, the PCR assay was conducted from H₂O, gDNA (2.3 ng), total RNAs (40 ng) and cDNAs (40 ng equivalent of total RNAs) and run on a 4% agarose gel. Their predicted size in base pairs is indicated in the right side of the figure. Left-overs of primers or primer dimers are visible in the bottom part of each photography.

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Figure 6

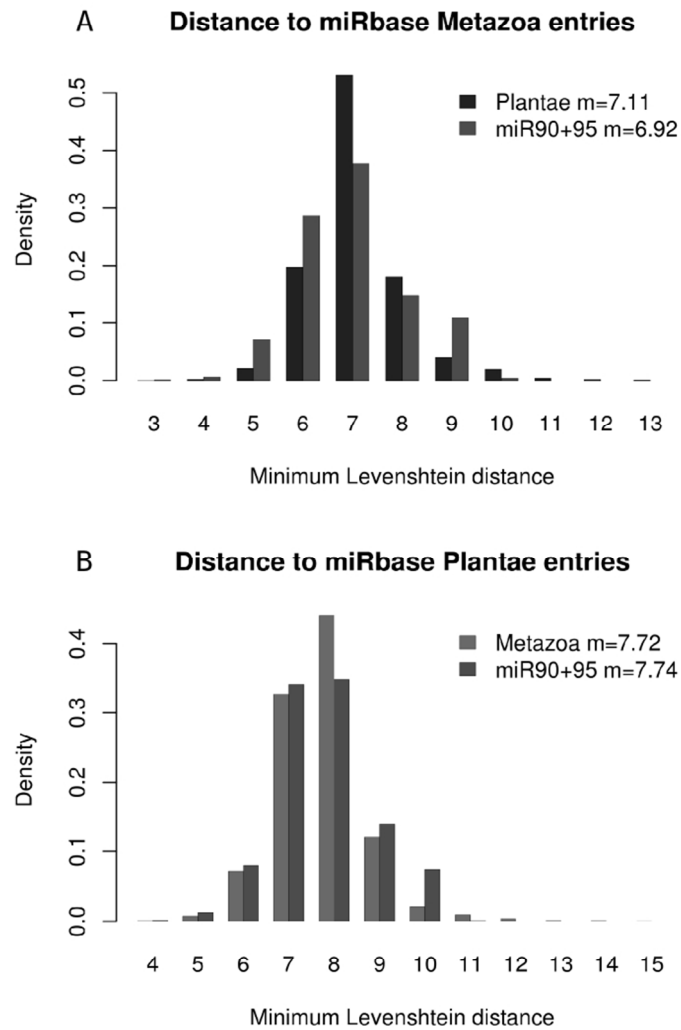


Figure 6. Comparison between miRNA candidates and miRBase. Each chart shows the distribution of the Levenshtein distance for each sequence to the closest entry in a subset of miRBase. (A) Distance for Viridiplantae and *E. siliculosus* to Metazoa; (B) Distance for Metazoa and *E. siliculosus* to Viridiplantae.

Computational prediction and experimental validation of microRNAs in the brown alga *Ectocarpus siliculosus*

Figure 7

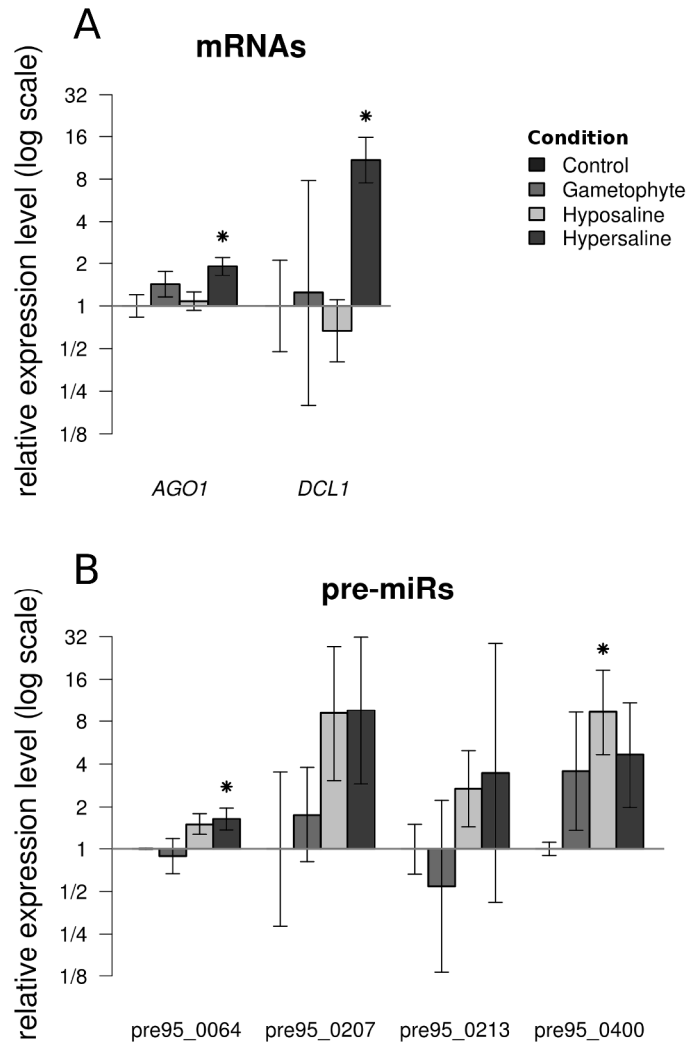


Figure 7. Transcript levels of miRNA processing proteins, pre-miRs, miRNAs and target genes measured by RT-qPCR in different algal materials. (A) AGO1 and DCL1 cDNAs. (B) Pre-miRs detected in figure 4. The data are expressed as fold changes relative to the control. "Control": sporophytes grown in normal culture conditions, used to set the reference value at 1; "Gametophyte": gametophytes grown in normal culture conditions; "Hyposaline" and "Hypersaline": sporophytes subjected to corresponding osmotic stress conditions. Conditions for which the distribution of replicates values significantly differs from the control distribution (Student t-test, $\alpha=0.05$) are denoted by a star.