

Supporting Information

Beauchemin et al. 10.1073/pnas.1206683109

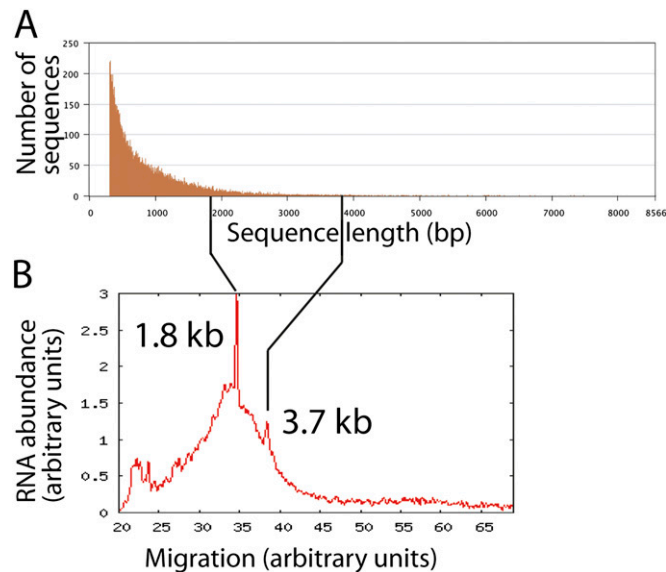


Fig. S1. Size distribution of sequences in the transcriptome and in the mRNA. (A) The size distribution plotted as a histogram of number of sequences for each contig length in the 74,655-sequence transcriptome. (B) The size distribution of the RNA sample used for sequencing as determined using a Bioanalyzer. The sizes of the two peaks of rRNA still visible in the analysis are shown.

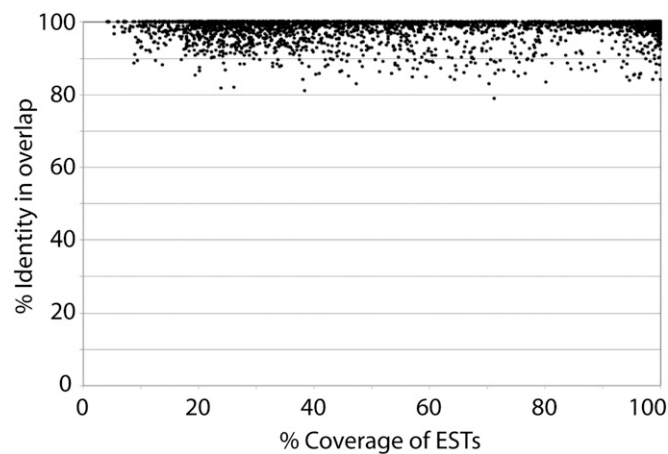


Fig. S2. Degree of sequence identity of *Lingulodinium* ESTs with the transcriptome. The degree of sequence identity as a function of the proportion of the EST sequence covered is shown by a comparison of the transcriptome sequences with 2,111 GC-rich *Lingulodinium* ESTs in GenBank. Each point represents a Sanger EST with a corresponding sequence in the transcriptome. Because of the short average length of the transcriptome sequences, there are many ESTs that are incompletely covered by the transcriptome contigs, and several ESTs have matches with more than one contig.

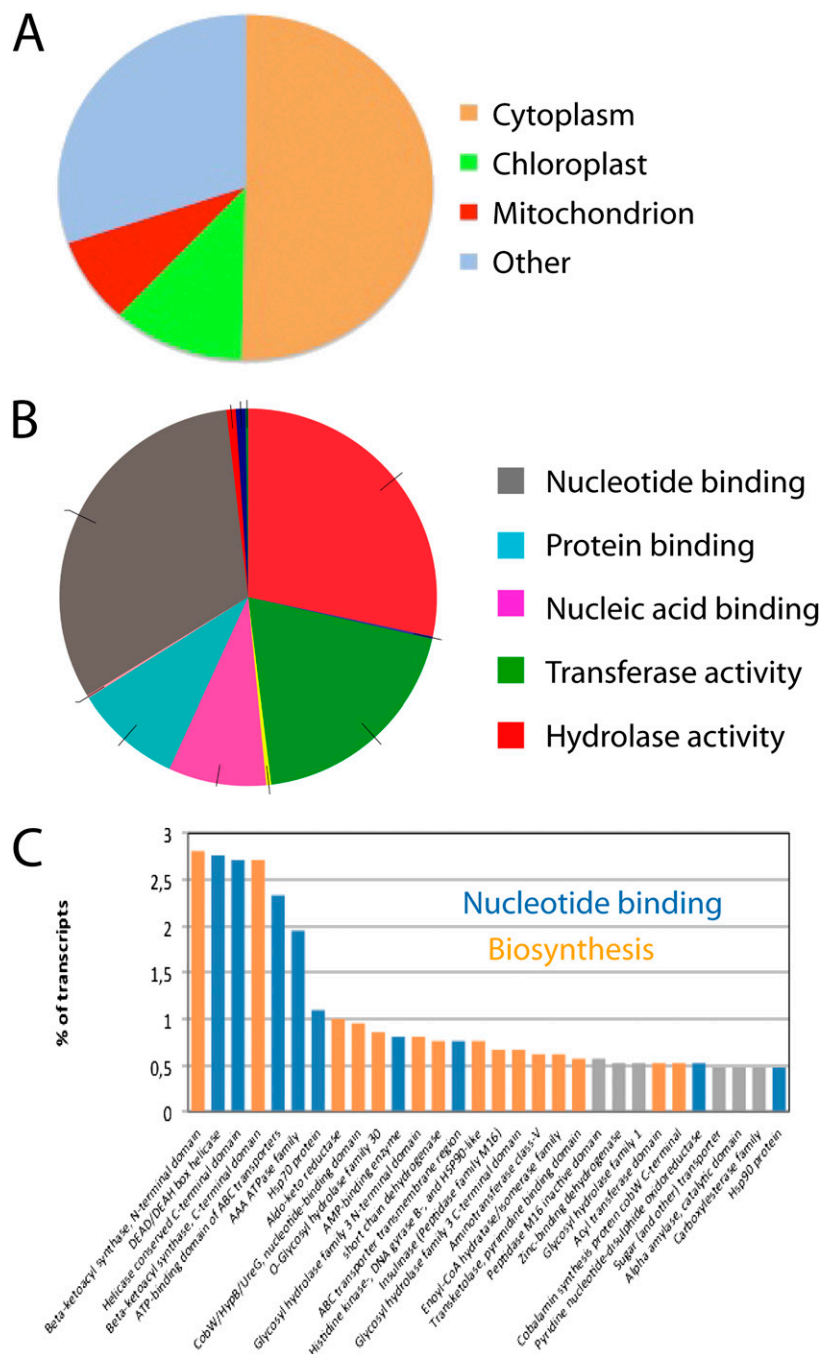


Fig. S5. Characteristics of the bacterial-like sequence in the transcriptome. (A) Of 2,354 sequences identified as putative bacterial sequences, 414 had an annotated match to GenBank. The annotated sequences were classified into the gene ontology compartment categories cytoplasm, mitochondrion, chloroplast, or all other membrane-bound compartments. (B) Functional classification shows enrichment of sequence in nucleotide-binding and enzymatic gene ontology categories. (C) The 30 most abundant protein family domains include principally nucleotide-binding and biosynthetic functions.

Table S2. Number of KEGG pathway sequences found in mammals, plants, alveolata, and diatoms for replication, transcription, splicing, and translation

Process	Subunit	Mammals/ plants		Alveolata*			Diatoms Thaps
		Homsa	Arath	Plafa	Linpo	Tetth	
DNA replication	DNA polymerase						
	α Complex	4	4	4	4	4	4
	δ Complex	4	4	2	2	2	2
	ε Complex	4	2	2	2	2	3
	MCM complex	6	6	6	6	6	6
	RPA	3	2	1	1	1	2
	Clamp/loader	4	4	4	4	4	4
	Other						
	Helicase	1	1	0	1	1	0
	RNaseH1	3	3	1	0	3	1
	Fen1	1	1	1	1	1	1
	DNA ligase	1	1	1	1	1	1
Transcription	RNA polymerase I, II, and III						
	Core	10	9	10	10	9	10
	Specific	13	12	6	6	6	10
	Common	5	5	4	5	4	4
	Basal transcription factors						
	TFIIA	2	2	0	0	0	0
	TFIIB	1	1	1	0	0	0
	TFIID	15	10	1	1	3	4
	TBP	1	1	1	0	1	1
	TFIIE	2	2	0	0	0	1
	TFIIF	3	2	0	0	0	0
	TFIIH (NER)	10	10	9	3	5	8
Translation	Ribosome						
	B	0	7	0	0	1	1
	B/A/E	47	87	40	51	34	48
	A/E	25	25	20	24	21	23
	E	12	12	10	10	8	12
	aa-tRNA synthesis						
	Enzymes	23	23	22	22	22	22
	Basic translation factors						
	Initiation	42	57	29	29	33	37
	Elongation	10	16	8	9	8	9
	Release	11	11	5	3	7	4
	Splicing	Spliceosome					
General		9	8	8	9	9	7
U1		8	7	5	5	5	4
U2		12	10	7	8	9	10
U4/U6		7	7	6	7	7	7
U5		8	8	6	7	7	7
U5/U4/U6		5	5	4	2	4	5
Prp19 complex		9	8	7	5	7	7
Prp19 related		9	8	8	7	5	8
EJC/TREX		6	5	4	3	4	5
Common		3	3	1	1	2	1
Translation related		Ribosome biogenesis					
	90S preribosome	18	18	6	12	15	14
	Nucleus	14	12	9	10	12	12
	Nucleolus	17	15	13	14	14	14
	Cytoplasm	7	6	5	6	6	6
	mRNA transport						
	Nucleus	11	9	9	7	7	8
	Cytoplasm	6	6	2	3	2	3
	NPC	34	22	3	5	8	13
	SMNC	9	2	0	0	0	1
	eIFs	14	10	9	9	7	11
	EJC	16	12	6	9	6	6
TREX	6	6	0	0	2	3	

Table S2. Cont.

Process	Subunit	Mammals/ plants		Alveolata*			Diatoms Thaps
		Homsa	Arath	Plafa	Linpo	Tetth	
	mRNA surveillance pathway						
	Nucleus	34	26	9	15	10	17
	Cytoplasm	15	12	8	10	7	9
Total		540	545	323	349	342	396

aa, aminoacyl, Arath, *Arabidopsis thaliana*; eIFs, eukaryotic initiation factors; EJC, exon junction complex; Homsa, *Homo sapiens*; Plafa, *Plasmodium falciparum*; Linpo, *Lingulodinium polyedra*; MCM, mini chromosome maintenance; NER, nucleotide excision repair; NPC, nuclear pore complex; RPA, replication protein A; SMNC, survival motor neuron complex; TBP, TATA-binding protein; Tetth, *Tetrahymena thermophila*; TFI, transcription factor II; Thaps, *Thalassiosira pseudonana*; TREX, transcription/export.

*Apicomplexans, dinoflagellates, and ciliates.

Table S3. Nuclear- and plastid-encoded reference sequences from GenBank used for comparison of synonymous (dS) and nonsynonymous (dN) mutations

Gene name	Accession no.	Length (bp)	dS	dN
p43	AY423581	1,429	14 (0.05)	33 (0.04)
Phosphoribulokinase	AY772247	1,461	64 (0.22)	32 (0.03)
Histone-like protein	AF482694	511	13 (0.17)	18 (0.08)
Luciferase	AF085332	4,000	25 (0.03)	45 (0.02)
GAPDH (plastid isoform)	AF028560	1,433	150 (0.47)	35 (0.04)
Actin	AY423582	1,407	59 (0.23)	15 (0.02)
RuBisCo	GONR15B	1,912	136 (0.36)	33 (0.03)
Carbonic anhydrase	EU044834	1,636	18 (0.06)	23 (0.02)
Cyclin	AY618995	1,825	6 (0.02)	30 (0.03)
Cellulase	GQ258705	1,425	40 (0.12)	30 (0.03)
Glucose phosphate isomerase	DQ812892	1,875	10 (0.03)	18 (0.01)
Fructose-1,6-bisphosphatase	DQ508159	1,235	12 (0.05)	32 (0.04)
Sedoheptulose-1,7-bisphosphatase	DQ508153	1,492	58 (0.18)	44 (0.05)
Superoxide dismutase	AF289824	744	21 (0.13)	5 (0.01)
Peridinin-chlorophyll a-protein	JO692699	1,127	7 (0.025)	11 (0.013)
Luciferin-binding protein	GONLBPA	2,217	37 (0.08)	62 (0.04)
psaA	DQ264850	2,506	72 (0.14)	273 (0.16)
psaB	DQ264852	2,174	41 (0.09)	146 (0.09)
psbA	DQ264844	1,074	10 (0.04)	32 (0.04)
psbB	DQ264845	1,559	29 (0.08)	103 (0.09)
psbC	DQ264846	1,418	43 (0.17)	112 (0.14)
psbD	DQ264847	1,236	11 (0.04)	55 (0.06)
atpA	DQ264853	1,609	20 (0.06)	69 (0.07)
atpB	DQ264857	771	17 (0.10)	82 (0.14)
petB	DQ264849	842	21 (0.13)	29 (0.06)
petD	DQ264848	545	13 (0.11)	48 (0.13)

Column dS shows the total number of synonymous changes (ratio of synonymous substitutions per synonymous site) while column dN shows the total number of nonsynonymous changes (ratio nonsynonymous substitutions per nonsynonymous site).