

The complex intron landscape and massive intron invasion in a picoeukaryote provides insights into intron evolution

Bram Verhelst^{1,2}, Yves Van de Peer^{1,2,3,*}, Pierre Rouzé^{1,2,*}

¹ Department of Plant Biotechnology and Bioinformatics, Ghent University, Technologiepark 927, B-9052 Ghent, Belgium

² Department of Plant Systems Biology, VIB, Technologiepark 927, B-9052 Ghent, Belgium

³ Department of Genetics, Genomics Research Institute, University of Pretoria, Pretoria, South Africa

* to whom correspondence should be addressed

Supplementary Information

Supplementary Methods

Introner Element (Remnant) Prediction

Starting from handpicked example Introner Elements (IEs), we delineated common motifs (pattern blocks) and assembled them into class-specific pattern files (IEA-1 example listed below). We used PatScan (Dsouza, et al. 1997) to scan the *Micromonas* genomes. For each class, multiple pattern files were constructed, ranging from strict to degenerate. When overlapping matches were detected, only the match belonging to the strictest pattern file was kept. EST and protein alignments were generated using GenomeThreader (Gremme, et al. 2005) (v1.4.6; -minalignmentscore 0.95 –mincoverage 0.89) and the splicing information was used in the automated curation of the final set of predicted IEs i.e. adjusting IE start and stop coordinates to match the exon/intron boundaries.

```
p1=GTGCGT
0...15
p2=ACTGGTYCCCRTACGACC[5,0,0]
0...80
p3=STTTCAAT[2,0,0]
0...40
p4=GCCTTCAACTC[3,0,0]
0...100
p5=AG
```

IE remnants were detected using BLASTN (v2.2.17; -e 1e-05) (Altschul, et al. 1990) using the previously built set of (complete) IEs. For each class, we also built a multiple sequence alignment (MSA), constructed a profile HMM (HMMer v2.3.2), and used it to detect additional instances of degenerated IEs. This HMM approach was also used for members of the IE-B / IE-D class.

Micromonas Re-annotation

IEs were added as an extra evidence track when performing the annotation (EuGene v3.6 (Schiex, et al. 2001)). EST and protein libraries of all Mamiellophyceae were used, as well as extensive sets of manually curated gene models. RNA genes were predicted using tRNAscan-SE (Lowe and Eddy 1997) and Infernal (Griffiths-Jones, et al. 2003). Genes encoding selenoproteins were manually corrected. When compared to the old annotation, the new annotation features fewer but larger gene models. This is mainly due to the IEs that help to span stop codons, allowing adjacent gene models to be merged into one continuous model.

Gene Ontology analysis of IE genes

GO terms for all *Micromonas* proteins were derived using InterPro2GO (Mulder and Apweiler 2007), and GO term over/under-representation of genes carrying IEs, using the GO terms of the entire *Micromonas* proteome as a background, was analysed using the Cytoscape plugin BiNGO (Maere, et al. 2005) (hypergeometric test + FDR correction; significance level 0.05). This GO analysis was only performed on CCMP1545, as the low number of IEs in RCC299 makes the analysis insignificant.

Spliceosomal Components

Spliceosomal components were detected through homology with *A. thaliana* proteins in the Splicing Related Gene Database (<http://www.plantgdb.org/SRGD>) and through the detection of splicing-related GO labels.

Metagenomic Sequence Analysis

When aligning metagenomic sequences (MSs) to the *Micromonas* genomes, the presence or absence of IEs in both the query (the MS) and the genomic sequence can present too big a gap in the alignment for 'regular' alignment programs to cope with. As such, the environmental sequences were aligned to the genomes using a seed-and-align procedure, initiated by a regular BLASTN. Starting from the best-hit, we expanded the genomic space with neighbouring hits. In the end, we used the outer coordinates to extract the corresponding genomic region, and re-aligned it to the environmental sequence using a SMITH-WATERMAN alignment (EMBOSS (Rice, et al. 2000): water).

To be able to draw accurate conclusions on IE presence/absence polymorphisms (PAPs), we performed a quality filtering step. We only continued with alignments that have more than 100 nucleotides labelled as 'non-IE', an identity percentage of more than 50%, and a coverage of more than 60%. After careful consideration, we also decided to leave aside all metagenomic sequences labelled as 'JCVI', as they were assemblies of smaller metagenomic sequences.

RNA secondary structure analysis

For each IE class, a consensus sequence was constructed (EMBOSS (Rice, et al. 2000): cons). Secondary structures were predicted using RNAfold (Vienna RNA (Hofacker 2003)). Secondary structures were compared between classes, but no general model could be obtained.

Protein evidence for IE-D

The protein sequences of genes containing IE-D sequences and their selected orthologs were aligned using Clustal W (default settings: GONNET weight matrix) (Thompson, et al. 1994).

Visualisation tools

For visualisation purposes, we also employed the following tools: seqlogo (Crooks, et al. 2004) and R (v2.13.0) (R Development Core Team 2008).

Supplementary References

- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ 1990. Basic local alignment search tool. *J Mol Biol* 215: 403-410. doi: 10.1006/jmbi.1990.9999
- Crooks GE, Hon G, Chandonia JM, Brenner SE 2004. WebLogo: a sequence logo generator. *Genome Res* 14: 1188-1190. doi: 10.1101/gr.849004
- Dsouza M, Larsen N, Overbeek R 1997. Searching for patterns in genomic data. *Trends in genetics : TIG* 13: 497-498.
- Gremme G, Brendel V, Sparks ME, Kurtz S 2005. Engineering a software tool for gene structure prediction in higher organisms. *Information and Software Technology* 47: 965-978. doi: 10.1016/j.infsof.2005.09.005
- Griffiths-Jones S, Bateman A, Marshall M, Khanna A, Eddy SR 2003. Rfam: an RNA family database. *Nucleic Acids Res* 31: 439-441.
- Hofacker IL 2003. Vienna RNA secondary structure server. *Nucleic Acids Res* 31: 3429-3431.
- Lowe TM, Eddy SR 1997. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res* 25: 955-964.
- Maere S, Heymans K, Kuiper M 2005. BiNGO: a Cytoscape plugin to assess overrepresentation of gene ontology categories in biological networks. *Bioinformatics* 21: 3448-3449. doi: 10.1093/bioinformatics/bti551
- Mulder N, Apweiler R 2007. InterPro and InterProScan: tools for protein sequence classification and comparison. *Methods in molecular biology* 396: 59-70.
- R Development Core Team. 2008. R: A Language and Environment for Statistical Computing. In.
- Rice P, Longden I, Bleasby A 2000. EMBOSS: the European Molecular Biology Open Software Suite. *Trends in genetics : TIG* 16: 276-277.
- Schiex T, Moisan A, Rouzé P. 2001. EuGène: An Eukaryotic Gene Finder That Combines Several Sources of Evidence. *Selected papers from the First International Conference on Computational Biology, Biology, Informatics, and Mathematics: Springer-Verlag.*
- Thompson JD, Higgins DG, Gibson TJ 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22: 4673-4680. doi: 10.1093/nar/22.22.4673

Supplementary Figures

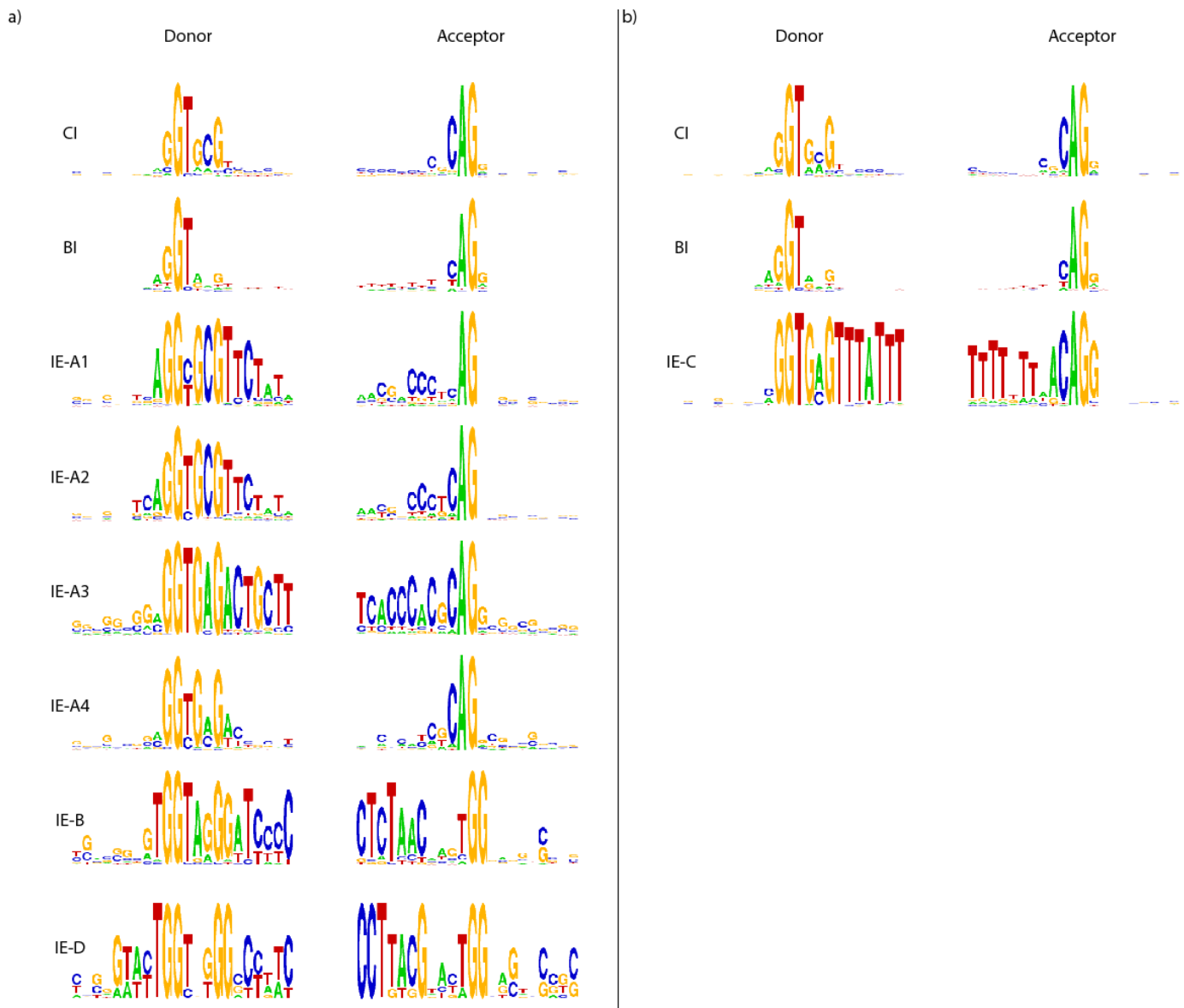


Figure S1. *Micromonas* splice site signals for all intron classes. Shown here are sequence logos for the donor/acceptor site (10 nucleotides upstream and downstream) for CCMP1545 (a) and RCC299 (b).

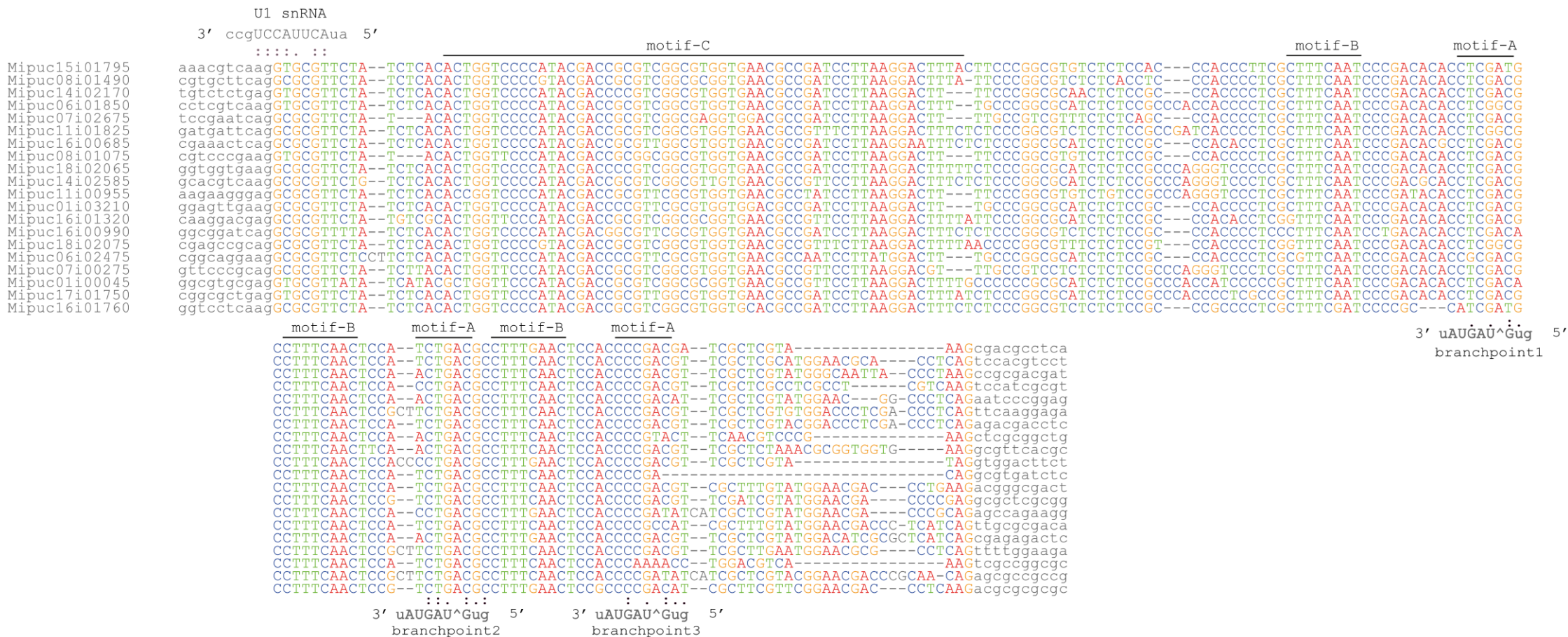


Figure S2. Alignment of 20 random IE-A1 sequences. The motifs (motif-A = branchpoint motif; motif-B: branch-point companion motif; motif-C) are marked, as are the splicing signals (donor site + branch-point) and their base-pairing information from the corresponding spliceosomal RNAs.

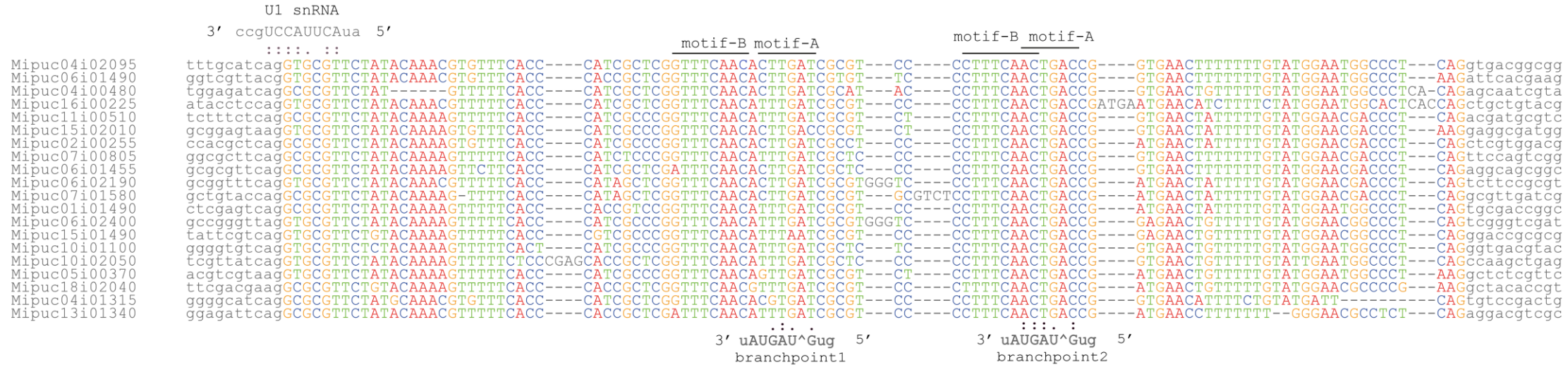


Figure S3. Alignment of 20 random IE-A2 sequences. The motifs (motif-A = branchpoint motif; motif-B: branch-point companion motif; motif-C) are marked, as are the splicing signals (donor site + branch-point) and their base-pairing information from the corresponding spliceosomal RNAs.

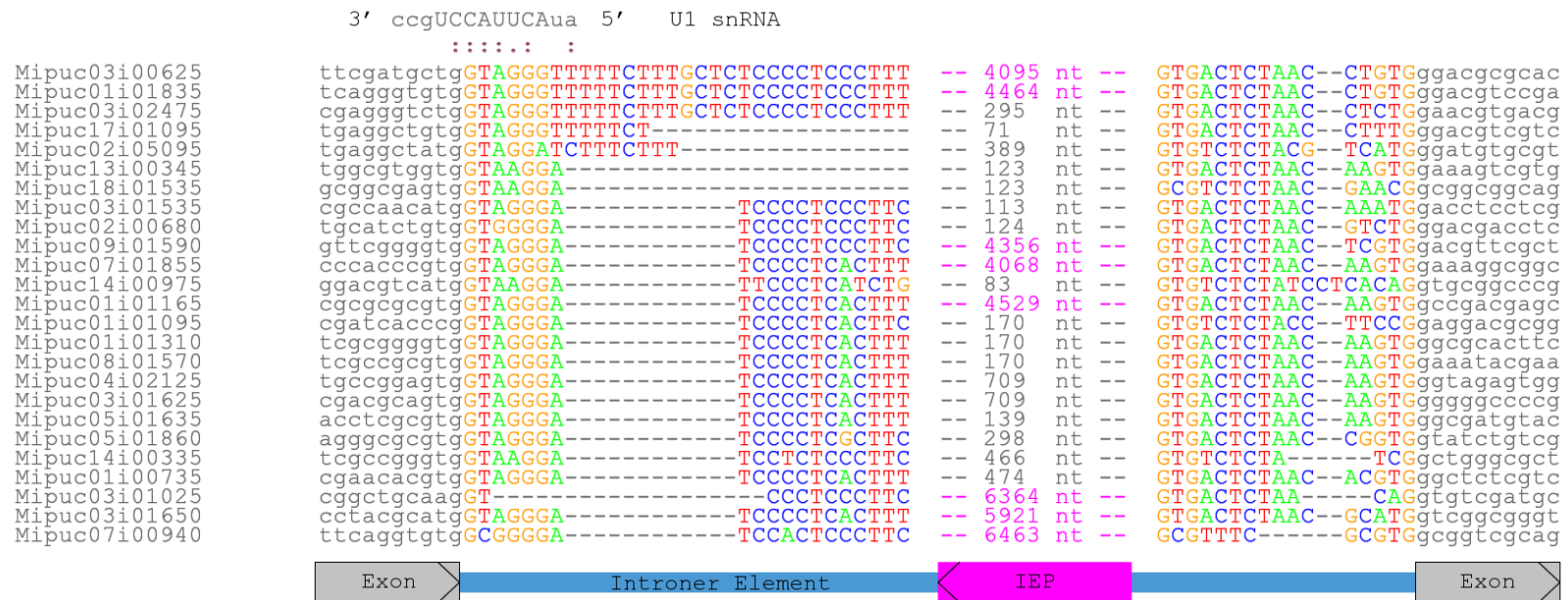
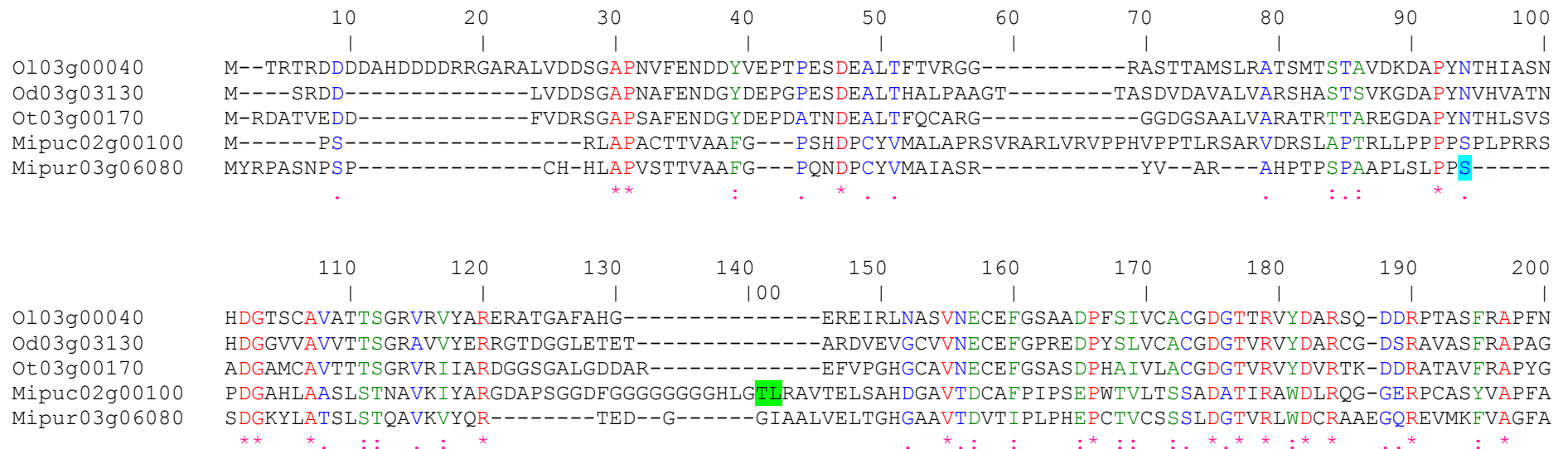


Figure S6. Alignment of all 25 IE-B sequences showing the splice site regions in detail. The structure and orientation of exonic regions (grey) and Intron-Encoded Proteins (purple) is represented schematically beneath the alignment. Base-pairing information regarding the donor site (U1 snRNA) is also provided.

Figure S8. Protein evidence supporting IE-D introner elements (in lack of any EST evidence). Canonical introns (light blue), IE-A (green), IE-D (dark blue), and chimeric IEs (red) have been highlighted accordingly. Intron phases are indicated on top of the intron positions. Gene IDs correspond to the following organisms: *Ostreococcus lucimarinus* (Ol*), *Ostreococcus sp.* RCC809 (Od*), *Ostreococcus tauri* (Ot*), *M. pusilla* CCMP1545 (Mipuc*), *M.sp.* RCC299 (Mipur*), *Bathycoccus prasinus* (Bathy*), *Physcomitrella patens* (*PHYPHA), *Arabidopsis thaliana* (*ARATH) and *Dictyostelium discoideum* (*DICDI). a) Alignment of WD89 proteins. b) Alignment of Jumonji-1 proteins (Jmj-N and Jmj-C Jumongi domains are highlighted in grey, while the Zn-finger domain has been highlighted in yellow), c) Alignment of sepA-like Serine/Threonine protein kinase, d) Alignment of U5 snRNP helicase, e) Alignment of UAA1 proteins (UDP-N-acetylglucosamine; UMP antiporters).

a) WD89



```

                210      220      230      240      250      260      270      280      290 s      300
                |      |      |      |      |      |      |      |      |      |
O103g00040    ER---DVASASLGGSSASDFLVATAVGPHVAFYDRRQQG----NASLFQDAHSEAVTRVRFHPERRSELYTSSVDGLVCAFDCA--HTPLNDEESLISVMS
Od03g03130    ER---DVCSATLGGGATGSMIAA AVGPHVVFDRRRHGADASHADMFRDAHSEAVTRARFHPERRRELYTSSVDGLMCAFDCS--RAPLNDEDALIAIIS
Ot03g00170    ER---DVPSASLGGDG-DTLVAAAAGAHVVFYDRRTQGGD-ACVGLFQDAHSEVTRVRFHPERRSELYTSSVDSLVCADFCS--RAPLNDEDALITIMS
Mipuc02g00100 TRATGGFATATAGGFG-DHLVVAGWNTQIVFDRRTREG----LRAFEDAHSEDVTRVRFQPTRRNRLFTASVDGLACVFDVGGCPADVDEGLLCVMT
Mipur03g06080 K----EFASATLGGGN-DHLVVGAQNEQVVFWDVRRVGTG----LEVFDESHSEDVTRVRFQPGRRNRLFTAGVDGLACAFDVGGCPADINDEDGLLTMVH
                . : * : * *      . : : . .      : : : * : * * *      * . * : * * * * * * * * * * * * * * * * : : * * : * : :

```

```

                310      320      330      340      350      360      370      380      390      400
                |      |      |      |      |      |      |      |      |      |
O103g00040    AEA AVNNIGFCRS---GSSGDARDALWCVTGVEDAHIFVASS-DRRRVGVHLVHLKNARELVRTAATS-SPEVSASAPEFAAQVDYVLGIHGDGVA-----
Od03g03130    ADA AVNTIGFCRS---GASGDERDAVWCATGVEEAHV FVASS-DRRRVGVHLAHLKNAREVRAATTN-SPRASSTAPEFLAGVDYIVGVHGDGVA-----
Ot03g00170    ADA AVNNIGFCRT---GASGNERDAVWCTTGIEEAHIFLTSSDKRRVGVQLVHLKNARALAQTAHSQS FMD-----FSTQVDYVLGIHGDVD-----
Mipuc02g00100 TDCSIVELGFCGGRFAAGSNDADSIAWVL TGNEDAWAFDAGA-DLETLGSTLAHVPNTRAAAFSAAAQFASSDATASSSFSRQVDYLLGCFRVRPGE GEG
Mipur03g06080 TGCAIVELGFVSG---AGD--DD-VLWLLTGNEDAWFYDAGD-DADTIGNTLAYIPDTRGAAQRSS----FAADIHAGGLSQQVDYLVGCF SRK----EP
                : : : : : * *      . . .      * * * * * : : *      : * * : : : * * . : :      : * * * : * *

```

```

                410      420      430      440      450      460      470      480      490      500
                |      |      |      |      |      |      |      |      |      |
O103g00040    --PGEFLSAGTQSGVVGIFPLVQGPSMRGTSAYVGLGAPAAVLRRHGHRDIVRS-IVWDANATASE SARVPATVGEDALVCAWTPDVGNEERP-----P
Od03g03130    --PGEFLSAGTQAGVVGVFPLVQGPSHARGTADFITLAPPVAVLRRHGHRDIVRC-VAWDANARASEPSRVPLTCGEDSLVCAWTPGTND DAAPP---P
Ot03g00170    --PGEYLSAGTQSGSVGVFPLVQGP EPARGTAGYVILGPPVAVLRNHGHRDIVRA-MSWDGNRHASEAARVPLTCGEDSLVCAWTPGAANAADPP---S
Mipuc02g00100 DGASSVVVAAGTQDGA VGLFPVPPRDGSGD-VARCELAAPTRVLDGGHADI RACVLWDWDADAAA----PATGGEDSRVVLWGAAGGAG-GGDG--SS
Mipur03g06080 GGGSTVMVAAGTQGGAVGVYPV VASTDPAS---SPAVLGAPIAVMDGGHVDIVRA-MAWKPDNAEP----PVTGAEDSRMCAWGEKKAVDPGGDGGYIA
                . : : : * * * * * * * * : : * * * * * : * : *      * * * * * : : *

```

```

                510      520
                |      |
O103g00040    VDRARPPG-----RRHSPY
Od03g03130    VDRTRPPG-----RRHSPY
Ot03g00170    VDRTRPLG-----RRHSPY
Mipuc02g00100 GGGKAARSESAETEYHDRGRRRHSPY
Mipur03g06080 AGGKIGRG-----DEGGRRHSPY
                .      .      * * * * *

```

b) Jumonji-1

	10	20	30	40	50	60	70	80	90	100
Ot08g01040	MSAPT--TRNANHER-----PSTSQSAARRT-----RTVDAAIEDARVFTPTLEEFADPIVY									
Ol08g01020	MSARR--DVDARHAS-----TSEKDSERKPS-----RTTTADIASAPTFRPTLEEFADPIAY									
Bathy04g01120	MLSLN--DKDEDQEGLPLKPTQNDVFNEEEEEERI PRARSKRARRAPKAVYDNDFEAHEIQNAIKASKVITNRCRTDGSDIPLAPIFYPTEEEFADPIAY									
Mipur03g03040	MVAADGGAIASEPPTGRRAKRTIRAPTKFVDEDDGFNDADLQ RALKAS-----MKHVRVVSSTTVPECPVFRPTAEQFADPFAY									
Mipuc02g01760	MAARS--KRLTRPPAKRRD-----ADDAMFDEAEVQRALRAS-----VKLQKRSSSTCMPSCTPHPTPEQFRDPFAY									
	* :			.	:			* . : :	* ** * . * * . *	

	110	120	130	140	150	160	170	180	190	200
	Jumonji-N									
Ot08g01040	LTKIEPLVRRTGICKVIPPARGAKPTWNEDVVRKDVSTFETKLQNVHKLSEGRLFQFGKSYT-KSGYKAMAMAFEKEWAEGRAD--FDACDVN-----									
Ol08g01020	LSSIEARAREAGICKVIPPARGAAPRWNGEAWRRDDARFETKLQNVHLSSEGRTFQFGKEYA-KGEYEAMAKAYEERWAKERPD--VDANDAN-----									
Bathy04g01120	ICSIQSKAEAFGICKIVPPDGYAPNFN RACCFGEKSLVETKHQNVNRLQQGESFPPGKTYVGLEKYKEMADTFEENYKEAHPETFKDIKDEDDL-----									
Mipur03g03040	IKSITPEAMPYGIAKIIPPEGWKPPFN E EAG-GDGI PFDTKLQTVNRLQEGLFEDGERYT-RDSYRDMADAFKRKYLETHRRVADETERLRRENRGWSD									
Mipuc02g01760	IKSITSEGIAAGIVKIVPPEGWRPPFNANAG-GGEI PFDTKLQAVHRLQEGVDFEDGKRYT-RESFRAMADSF AAAMWHAHLSLDAEVKRLMIE-RGFEE									
	: . * .	** * . : . * * *	*	: *		. : . * * *	* * : *		. : . * * :	: . : :

	210	220	230	240	250	260	270	280	290	300
Ot08g01040	-----SVERAFWNMVETQE EKAA-----VEYGNLDLTKEFGTGFVDAHG-----									
Ol08g01020	-----ALERAFWDMVETRSEQAR-----VEYGNLDLTKEFGTGFVDENG-----									
Bathy04g01120	----LKRIDEDEYWRIVETNPNEAKAECGSLIQTKNVNKKGEVLVEYGSVDARRFQSGFAAGISGD-----									
Mipur03g03040	DACEARALEEEFWRIVETDVEKIR-----VEYGSDDLADVYGS GFAKVPLG-----SASAAAGATPDS DSD-----									
Mipuc02g01760	ADARARAVEEEFRIVETNAERVS-----VEYGSDDLADVYGS GFLYPRANYDGNDEISDDVMSRRADADGDGDGDGDGDGDGDG									
	: * :	* : * * * :			**** . * : * :	: : * * .				

	310	320	330	340	350	360	370	380	390	400
	Jumonji-C									
Ot08g01040	-----ERHPWDFEHLYSHPLNLLRVIEHDIPGLTKPWLYLGM LFATFCWHVEDHFLCSVNYLHTGASKTWYGVPGSDAEAFENCARATVPRLF									
Ol08g01020	-----EKHPWDFEHLYSHPLNLLRVVEHDIPGLTKPWLYLGM LFATFCWHVEDHFLCSLNYLHRGA AKTWYGVPGSDAEAFENCARATVPRLF									
Bathy04g01120	-----PEDTEKHPWDMFELSKHPDNLLRVVDDIPGLTTPWVYCGMLFATFCWHVEDHYLASVNYAHKGS AKTWYGI PGSDAEKFEAIAKTAVPSLF									
Mipur03g03040	-----EDGGVPHAWDFGELIRHPSNLLRVVGGDIPGLTRPWLYFGMLFSAFCWHVEDHYLGSVNYMHAGAPKTWY GAPTAAADAFERAVRDI VPGIF									
Mipuc02g01760	DGD CDS DSE DGARRHAWDFSELVNHPSNLLRVVGGDIPGLTRPWIYFGMLFSAFCWHVEDHYLGSVNYLHDGAPKTWYS IPPASASAFERAVRTIVPTRV									
	* . * . :	* * * . * * . * . * . *	: * * * * *	* * * * * *	* * * . * * . * . * . *	* * * * * *	* * * * * *	* * * * * *	* * * * * *	* * * * * *

	410	420	430	440	450	460	470	480	490	500
Ot08g01040	QQAPDILHQIVTMVPPGILID-HGVKVVHTVQHPGEFIVTFPRAYHAGFESHGFNVAEAVNFGHANWLDHGRRRIDVYSTGSEFKRNAVFAHHRLLARA	ET								
O108g01020	EQAPDILHQIVTIVPPGLVD-HGVKVVHTVQQPGEFVVTFFPRAYHAGFESHGFNVAEAVNFGHVNWLDHGRRRIDVYSTGSEFKRNAVFAHHRLLVSR	AAET								
Bathy04g01120	KENPDKLHHITMLVPPGQLIE-NKIKIVKLVQKPGDFVVTFFPRAYHAGFESHGFNVGEAVNFAPVDWIEMGRVACRNYVKNGKRNAVFAHDRVVVTAAKS									
Mipur03g03040	KDAPDLLHRLVTLVPPAVLGEHGVPCQTLQRAGEFVVTWPRAYHAGFESHGWNVGEAVNFGTADWVPMGRAAVNDYQHGVGKRDSIFSHKMIILDTAKA									
Mipuc02g01760	HDTPDLLHRLVTLVPPGVLDAHGVPVFQTLQKPGTFIVTWPRAYHAGFESHGYNVGEAVNFGTAEWVPMGRAAVEAYVTSSEFKRNAVFSHERVLLLETGR									
	.:	**	***:	***	*	:	:	:	:	:

	510	520	530	540	550	560	570	580	590	600
Ot08g01040	FAEVLNAKGL-----LLKSKVMGTVIATLCKELESIVSDEEIIYRSSLVR--RG-----LKMEVVALPNEDD--DAC	CIRCKAIPFLSVVRCK-CLP								
O108g01020	FVEVLGKNAR-----LVKSKAMGAIIVSTLRKELETILSDEEIIYRASLVR--RG-----LNIEIVQAPNEDD--DAC	CIRCKAMPFLSVVRCK-CLP								
Bathy04g01120	LKKIFETTK-----SRGKWMAMHSRVLRTDLETLADELENWQSILNGKQRGDGFIKGDPLRFYKQNIPEMDGPEDC	CVVCKAMPFAAVRCE-CEF								
Mipur03g03040	FVRRYGYGDGSSREDQSLRAPWIIARMADALRAELQIIIEKEQRAGRAVVT--KG-----VKEVAGKENE-ASKHEDEDEN	CALCKAMPHLAVVHCARCFE								
Mipuc02g01760	HARSFASPGGVSDEA---RAPWIIASVARMIRDDLFTIAREQRTGRDAALT--RG-----VRVSADDDCLGGRVTHDHEVV	CAECKSMPYLAVARCETCWK								
	.	:	:	:	:	:	:	:	:	:

	610	620	630	640	650	660	670	680	690	700
Ot08g01040	--TAVRCLRHAMDGCDCASERCLEVRVMGSYIRGL	LRSLLLGEPVEKAGGMAEKDLKLFVASVKTVAVIRSPS-----TTAAPPASKKSKLSD								
O108g01020	--TAVRCLRHAMDACDCAGERTLEIRVDSRLREL	LIKALFFGDGIQTKNDAAKAR--VDFSANVNRVAVNRAPPPKPKVVLKPKTKVKKPPTRAVLASP								
Bathy04g01120	GRSFARCLQHWNRGCDCKQRHRMVEMRMEVDELRAL	AKSLEL-----								
Mipur03g03040	IAEKIEFEQAVARARAEGRSAAAGAAWALQGGSSSER	LRRMGRYH-----RRRRRNRP								
Mipuc02g01760	LKDHIELSQAIIARAQGKSAAGAAALGNDRAKGA	PYGGTTRE-----GRRDLNRP								
	.	:	:	:	:	:	:	:	:	:

	710	720	730	740	750	760	770	780	790
Ot08g01040	LDQTVEG--LQERIEER--QALLSRTEGGKVVWTSERAKAFDKAAEQLGGMFKATSGDIGKILRDEYDIHASRDQIGSRLQKCRNKLRR	EQSTDGDGV							
O108g01020	PPTRIVASKADDAFTARGLPRKRAKCETRRRWTAEMVADFEVAVERLGGVDAATGKKLAEALS---AHDVTRDQCASRLQKHREKIKSNADARATL--								
Bathy04g01120	-----								
Mipur03g03040	VFCLG--HALDGNCGHG-----VTELVMTTNASMDELERVLGALQK-----								
Mipuc02g01760	VFCLK--HALGGACGHP-----ARERVNTRVDVREIEDLARALGEFVGPS-----								

c) SepA

	2	10	20	30	40	50	60	70	80	90	100
SEPA_PHYPA	M-SRHTSGSAFHKS	KTLDNKYMLGDE	IGK	GAYGRVYKGLDLENGDFVAIKQVSL	ENIPPEDLASIMSEIDLLKLNHNRNIVKYQGSFKTKTHLYIILEFV						
MAP3K_ARATH	M-ARQMTSSQFHKS	KTLDNKYMLGDE	IGK	GAYGRVYIGLDLENGDFVAIKQVSL	ENIGQEDLNTIMQEIDLLKLNHKNIVKYLGS�KTTHLHIILEYV						
SEPA_DICDI	M-SKKEP-EEIKKNVTVGN	-YNLGVVIGKGGFGTVYQGLDIEDGDFVAIKQINLTKIPKDQLQGIMNEIDLLKLNHANIIVKYIKYVKTKDNLYIVLEFV									
Mipur05g06740	MELPRSTRRS-----	TLVGSYILGDE	IGK	GAYGVYKAIIDKRDGRVVAIKEIPLAGIDEASLAGVRL	EIDLLGSLSHPNVVGQLGTIRTPSYFYIVLEFV						
Mipuc16g03100	MELPRSTRRS-----	TVVGNVYILGDE	IGK	GAHQVYRAIDKRDGRVAVKEIPLRATSAYDVDAIESECALLRSLSHRNVTRFLGTVRAPEHLYIILELA							
Ot17g01940	M--PQTDQPQ--	HPPRVVGHILGEL	IGV	GATSSVYKAVDQRDGRVSAVKEISLIGVDSQDMERITAELELLSNLEHANVVKYEGARTIGESLYVELEFA							
O116g02130	M--PASPPASPHHAVKRVGNHILGEL	IGSGATSRVHKAVDQRTGSICAVKEIPLRGVPVEQLERITSEVELLSRLEHANIVKYEGAVRVEECLYIMLEYA									
	110	120	130	140	150	160	170	180	190	200	
SEPA_PHYPA	ENGLSANNIKPNKFGALPEN	VVGRYIAQVLEGLVYLHE--	QGV	IHRDIKGANILTTKEGEVKLADFGVATKLTEADI-----							
MAP3K_ARATH	ENGLSANNIKPNKFGPF	PESLVTVYIAQVLEGLVYLHE--	QGV	IHRDIKGANILTTKEGLVKLADFGVATKLNEADF-----							
SEPA_DICDI	ENGLSLSGIK--KFGKF	PETLVCVYIRQVLEGLVYLHE--	QGV	VHRDIKGANILTTKEGKIKLADFGVATKFDDT-----							
Mipur05g06740	EAGSLAASIKANKFGPA	PEALCKVYVAQVLDALAYLHSPRNGIVHRDVKGANLLATKDGCVKLADFGSAARMGEDGR-----	GGAQ	RPGSN							
Mipuc16g03100	ENGLSAGVIKPSRFGPT	PEPLAACVYVAQLLDGLIYLHA--	NGV	THRDIKGANVLLATKDGCVKIDFGVAVRVHGGAA-----	NPPV	NAHH					
Ot17g01940	ENGLSARTVHPSRFGGF	PESLCAVYVAQILRGLAYLHG--	QGV	VHRDIKGANILTTKEGVVKLADFGVATKGSRVGGDGLGRRFGLLN--	EASR	KIASM					
O116g02130	ENGLSARTVHPSRFGAF	PESLCAVYVAQVLRGLAYLHS--	QGV	VHRDIKGANILTTKEGVVKLADFGVATKGGRASG-DGLSGVFGAEGRSEGSDGDASD							
	210	220	230	240	250	260	270	280	290	300	
SEPA_PHYPA	-----NTHSVV	GTPYWM	MAPEVIEMSG--	VSAAS	DIWSVGCTVIELLTCVPPYDLQMPALFRIVQ-DDHPPLPEHVSEVIIDFLRQCFQK						
MAP3K_ARATH	-----NTHSVV	GTPYWM	MAPEVIELSG--	VCAAS	DIWSVGCTIIELLTCVPPYDLQMPALYRIVQ-DDTPPIPDSLSPDITDFLRLCFQK						
SEPA_DICDI	-----SAAAVV	GTPYWM	MAPEIIEIENG--	ATTKS	DIWSVGCTVIELLTCVPPYDLQMPALFRIVQ-DDCPPLPEGISPPPLKDWLMQCFQK						
Mipur05g06740	KP---TKKTTGPVDGEGDV	GTPYWM	MAPEVIEMSG--	GSDPKS	DVWSVACVVVELITGSPPYFDLQMPALFAIVR-DESPPLPPGISPELRGFLSACFRK						
Mipuc16g03100	PA---RSVDDDDGDPGASPV	GTPYWM	MAPEVIEMRPGDDPKS	DVWSVACVVIELLTCVPPYFEMQMPALFAIARGDARERIPPGV-DALRDFLRCFEK							
Ot17g01940	EN--SEEANADGEAKPNDAL	GTPYWM	MAPEVIEMRN--	VTAAD	VWSVGCTIIELLTSNPPYFDLQMPALFRIVR-DKHPPLPAGISDALRDFLMLCFKK						
O116g02130	GRGIGDGTAEAGEGEGDKAL	GTPYWM	MAPEVIEMRS--	VTAAD	IWSVGCTIIELLTSNPPYFDLQMPALFRIVR-DEHPPLPTGISSEALRDFLLLCFRK						
	310	320	330	340	350	360	370	380	390	400	
SEPA_PHYPA	DAKR	RPDAQTL	LLGHAWIRKSRREKRNGVSHGIAH-----	FPRL							
MAP3K_ARATH	DSRQ	RPDAKTLLSHPWIRNSRRALRSSLRHSGTIR-----	YMKE								
SEPA_DICDI	DPNLRISAQKLLKHKWI	QASIKKKPVENGAGGVNGTDSL	GAPANIDDI	AKNITDYNERINKKPSHQKPSIHPKSPKGKVFLLPPEEEEEDEWGDDFSNT							
Mipur05g06740	DPAQRPTASELRSHEWLKGV	TAGAAATGSSSSG-----									
Mipuc16g03100	APER	RPSAAELRRHAWLRDVATPGTTS	DAASRAID-----								
Ot17g01940	DPKDRPAEALLSHTWLTDE-----										
O116g02130	DPKDRPSAEELINHTWLMDEHKVLAETWTKR-----										

```

          410      420      430      440      450      460      470      480      490      500
          |        |        |        |        |        |        |        |        |        |
SEPA_PHYPA  PGS----HDQDLLETYMSTTAIRVPPTVTTSLTRPPAGSRVSESEPEP-LVHNTVLRRTSGGPEEHAQCEAN-----VMRSTSGKFS-----HLLP
MAP3K_ARATH TDSSSEKDAEGSQEVVESVSAEKVEVTKTNSKSKLPVIGGASFRSEKDQSSPSDLGEEGTDSEDDINS DQGP---TLSMHDKSSRQSGT-----CSIS
SEPA_DICDI   PKS IKLPDKKSP LKL TNNK P STPLK QO QPTNNT PVQQQQQQQQQ P PLK L VPK Q PVI ENDD D W G D F N T V S D L S K A V G S L N F N N K K N E T P K P N I K K P T F S
Mipur05g06740 -----PARRTEHLTLDDPRPDSAASSLAGSGRSSPAHSRPTSRGAQRRGSRGDGGDGAWEPAS IAD-----P
Mipuc16g03100 -----DVVVEETVVASVSRVNSGGSTSGMSGGGGGGGGGGKTKNQNNQNNQNHREGDAPSKPTPRRPS-----P
Ot17g01940   -----RSLTEQPRE-----VGALIDRVANIGVEAST-----
O116g02130   -----DASGGRTQNDDAEEG-----VIVKVIDQMAKVG VNGD VTAADDA AAAAN-----S

```

```

          510      520      530      540      550      560      570      580      590      600
          |        |        |        |        |        |        |        |        |        |
SEPA_PHYPA  GQKDAELKDSLSGTEICTWKNDIHGMKTEVN-G-----NHVSTQVLSLSYLHPFV FV-----ELSLWCAVKPSTVHS----SKFTDGPLDDN
MAP3K_ARATH SDAKGT SQD VLEN HEKY DRDE I PGNLE TEASE GRRN--TLATKLVGKEYSIQSSHSFSQKG-----EDGLRKA VKTPSSFGGNE LTRFSDPPGDAS
SEPA_DICDI   EDEDEDDDDDGFGSGGDEDDDFGDIPTSIKLNPKFGSNIKGNSSGSANTTNSSTTVVQQPKLTVSNNNNNNK LPLSPRQPSSGNVKEGINHGSTGSK
Mipur05g06740 QP--RERRRRRPP-----PSSPFQDVKVAMN-----GDADADAEVADE-----
Mipuc16g03100 SSPFQDVGKKQPVRIVRGVAANLDDVVVALRPEERS-----ASLAASSSSSSSVASDG-----
Ot17g01940   -----MKRASSSSDALNLLG-----ISTSNGMTKE-----
O116g02130   SP-----LIAA---SERSSSQKDL ENFMKSD-----LGNFMKSDSGVTKA-----

```

```

          610      620      630      640      650      660      670      680      690      700
          |        |        |        |        |        |        |        |        |        |
SEPA_PHYPA  TNGLFFG-----SHEASPSVP-AVPNHGLPPRGGPSLEIHATAAKMKAKTAQSQA EAESLLTAKTN-GYESSIDLDDD DSVRTGLVSYVDKILMYKC
MAP3K_ARATH LHDLFHPLDKVPEGKTNEASTSTPTANVNQGDSPVADG--GKNDLATKLRARIAQKQMEGETGHSQDGGDLFRLMMGV LKDDVLNIDDLV FDEKVP PENL
SEPA_DICDI   SGGVIIDQWG---EDGEEDNDWGDVATVNFDPKVI RKGTVNKPDLSTR LKNRIALSETALSNSFNNGNDEDEDIFADDFDEDDDED FDL DKNLMK DNY
Mipur05g06740 -----VAAPAAVSVRL--EYPNTLNTLNTGFDGDVA---RATRDARWSD LAGLVRTAVADGV TAAAV
Mipuc16g03100 -----L-----KTAKAPTAAPTEEEVRAGDRLASTVKRIAATLRASGS---GSIDDTAGKKGARTPDALAGGFARACQ
Ot17g01940   -----QIDLNETIVSVM DRLRS AF-----ADTSSSENVEMAG-----
O116g02130   -----QIDFNKSLISIVDRLKDAFGTS-----T--SDVAAIDDAEQAG-----

```

```

          710      720      730      740      750      760      770      780      790      800
          |        |        |        |        |        |        |        |        |        |
SEPA_PHYPA  I-VQALEVSRMLMAMKLD ETEEVILAACQKLSSIFREFPKQKQDFMKPHGLI PLMDMLDMN--NNRVIHAVLQV LNLITDDNVELLESACRAGLIP-IMM
MAP3K_ARATH FPLQAVEFSRLVSSLRPDESEDAIVTSS LKLVAMFRQRP GQKAVFVTQNGFL PLMDLLDIP--KSRVICAVLQLINEIVKDN TDFLENACL VGLIP-LVM
SEPA_DICDI   A-RMSSEILKLMNLLTPEQPEEVISSACTQLITMFKENSEQKTL LIRRHGVIPI MEMLEVSNIQSHV LCSI LKVVNQI IDN NMEIQENLCLVGGIP- AIM
Mipur05g06740 D-AALTSADTAGALAAALHAAASKDGDDATAVRE RGN NID DANAAA SL DAAA AVLAATRRK-----RADGGDAPGETLAA FVTFGAASSTLS
Mipuc16g03100 D-LVDALARERGATRVSV DACVAANGVMGALVGALEASS PAAAATAAATAAAAVIDRANRG-----GGGSGSGS-EALATFELLGGIPLIME
Ot17g01940   --AD FVQLMDSSTSFNV TNETLFTTRANC TFMNVINDDAAAETTRTVA FECVAAASAASR-----DVLTSFVTMGLVPRALE
O116g02130   --AEFVRAMETATSVLVANETLFTTRACGVFVNVLENESASARMRTVALECVTAASMLSR-----DILTSFVTLAVLPRALE

```

```

      810      820      830      840      850      860      870      880      890      900
      |      |      |      |      |      |      |      |      |      |
SEPA_PHYPA SFASP--ECSKEMRMQAAYFVQKICHKSSLTLMQMFIA CRGLPVLVGLLEKDYARHRE-MVHMAIDGIWQVLELPGLSLKND FCHIFARSSILARLVDTLH
MAP3K_ARATH SFAGFERDRSREIRKEAAYFLQQLCQSSPLTLMQMFISCRGIPVLVGFLEADYAKHRE-MVHLAIDGMWQVFKLKKSTSRNDFCRIA AKNGILLRLVNTLY
SEPA_DICDI KFSGP--EYPASVRLETASFISKMCSTSTLTLMQMFIA CKGLPILVDFLLSPYAESKR-LVWMAVDIVNVFELQSPTPKNDFCRLF SKCGLLK-----
Mipur05g06740 LLRSDDASVGVGVKIAAAKVVRELSRGGVVALRCLASCDAPRALAGVLAQSAMSPS--RRELARVCVDVARFLTTLHTRAVTGSIPPSEASRMAS-----
Mipuc16g03100 RLARQTGPETRESRVACAKAARSVAAGGAVARRCLAACGAPKALAA TLSSPGFYSGEGNRELTRACVETCAEMMRAHEDEERGRGGGGEGVDFSGRR---
Ot17g01940 VLRRENASAR--AKIAALRLCRVTAKAGPVF SRCLVSCGALPILVGVLDDEFGGSGRELTKYALSAIYIAAEIERGDQMPRQLGQAI SLSLAQAG-----
O116g02130 VLKKKDSSMR--IKLAALRLCRATAKAGPVFTRCLVACGALPILVDVLDYGYSGNGRELTKYALGAIYIAAEIERGDQLPKQLAQTVSLALASAG-----

      910      920      930      940      950      960      970      980      990      1000
      |      |      |      |      |      |      |      |      |      |
SEPA_PHYPA TLNEASRVP-----NGSMTTEQSRPLSSEK-AFVKSQSGP--IDQSRIVN-----KEVFRSRSGQLD-HTPVRVTHDGN-----
MAP3K_ARATH SLSEATRLASISGDALILDGQTPRARSGQLDPNPNPIFSQRETSPSVIDHPDGLKTRNGGGEESHALTSNSQSSDVHQPDALHPDGRPR LSSVVADATE
SEPA_DICDI -----T-LP-----
Mipur05g06740 -----AASSGGISSGDVDFDEDDGETR-----GGGVVFD AFARAGLIPALV-----
Mipuc16g03100 -----KKPDDAAAAATTTTTDDDAIEDVPAGRT-----AGWGM RVAIARAGLLPPLV-----
Ot17g01940 -----LFPKLVALLGSTHTASLEDEANP-----EVMDDTRSEG-----
O116g02130 -----LFPKLVTL LGATHASLEDAAMIT-----DDIDDARSEG-----

      1010      1020      1030      1040      1050      1060      1070      1080      1090      1100
      |      |      |      |      |      |      |      |      |      |
SEPA_PHYPA -----QRNLLGLNNGDSSRMLPWQAHRACGL----P-----EFSWQNLPGQLEYTGQHS G---HLTHMKLAISSE--RLSDPFQPGYGESTWV
MAP3K_ARATH DVIQQRHRI SLSANRTSTDKLQKLAEGASNGFPVTQPDQVRPLLSLLEKEPPSRKISGQLDYVKHIAGIERHESR LPLLYASDEKKTNGDLEFIMAEFAEV
SEPA_DICDI -----
Mipur05g06740 -----NAIRDLNAAANE-----ERGVGPG-----
Mipuc16g03100 -----RALSSLHAASNE-----EAGVGPGGGG-----
Ot17g01940 -----GSS-----MTGSIG-----
O116g02130 -----GSS-----MTGSIG-----

      1110      1120      1130      1140      1150      1160      1170      1180      1190      1200
      |      |      |      |      |      |      |      |      |      |
SEPA_PHYPA LGDGDGCRRRERDDEARSD-----ANTSSSTSQTTS GQLS-----R-----LQETHNPEEAQEYLGKVANL LLEFSRAD-----TAV
MAP3K_ARATH SGRGKENG NLDTAPRYSSKTMTKKVM AIERVASTCGIASQTASGVLSGSVLNARPGSTTSSGLLAHALSADV SMDYLEKVADL LLEFARAE-----TTV
SEPA_DICDI -----IVLR-----DSIADGEAAATYPDR IINLFIMFSAAD-----SVV
Mipur05g06740 -----RKATNDTEKSESEKS-----ESPSSVYRERVAD ILLDATRSRNRDRGCAAS
Mipuc16g03100 -----GFGGLEKNSAAAASKTKTKSSPLAGRNANPD-----LPPHVAATVSGTCRDLVAD ILLSATSPG--HGARDA
Ot17g01940 -----GTSTHSVSRGLYS-----SGGRYYRELVAETLYR VAKRG--DGCPEV
O116g02130 -----GASTISRSRSGGS-----SGGKYYRELVAETLYR VAKRG--EGCAEV

```

	1210	1220	1230	1240	1250	1260	1270	1280	1290	1300	
SEPA_PHYPA			00					00			
MAP3K_ARATH	KKHMCS	TSLQLRLLQML-NTLPPPI	LKILECIKQLSQDPFTLEYLQLAEAMKHLIPFLEAR	DGPYVGR	HNQV	LNALHNLCKIN	-----	-----	-----	-----	
SEPA_DICDI	RKTMSAVEVIRPILDTL-SQLMPEQLAKV	LKS	IKQLSMDHNTLANLQ	NAGAIRFMV	PFLGR	RTGFAV	AEIHNHVLN	TMFHLC	RID	-----	
Mipur05g06740	RFALCELQAMHGLLALAGSPLPKSTS	AKLLRVVGH	LARDNCKDAMQ	RAGAVPKLV	RFLQW	EDPA	----	TRE	ALRALY	NLCRG	
Mipuc16g03100	IEALCETSTLHGLLALVGAPLPAST	SKKILALVRR	LARERNAHEP	MQRAGAI	PKLAR	FLOW	EDDD	----	HRE	TAMVALFHL	CGGGGGGGGGEDAGAASSA
Ot17g01940	SKSLVDVRIIHGCLAQL-SV	VPRSTATKLLGLVHLLSKQ	PDSFDV	LQHS	GAIPKLVK	CV	VEGNVKS	GH	AQSWEL	LKALHNL	CAMN
O116g02130	SKALVDVRIIHGCLAQL-SV	VPRSTATKILGLVHLLSQ	TNAFHVLQ	NAGAI	PKLVK	CV	VEGNFRS	GH	AQSWEM	LKALHNL	CAVN
	1310	1320	1330	1340	1350	1360	1370	1380	1390	1400	
SEPA_PHYPA	-----KRRQE	QAAECGII	PHLMHFIKI	-----	-----	-----	-----	NSPLKQFALPLL	CDMAHASR	TRELLRTN	
MAP3K_ARATH	-----KRRQE	QAAENGII	PHLMLFVMS	-----	-----	-----	-----	DSPLKQYALPLL	CDMAHASR	NSREQLRAH	
SEPA_DICDI	-----PERQY	QAAIDGII	PHLQYFITS	-----	-----	-----	-----	HSPLNQFALP	IICDLAHS	-KKARSELWKN	
Mipur05g06740	-----DASALE	QAAVAGVT	PHLIAVAAPDVFN	-----	KLGLADELGAATGT	NGGGTHWTRGG	PGAEAQIERLAP	LAA	PFLCD	MASSRRTRGELARH	
Mipuc16g03100	AEDPGVAARRE	QAAIAGAVPHL	VAVATPEIAHG	DATVDVSS	SAASAGRGV	VNGLHGDGDH	WRR	----	HPLSAAKLRAL	ASSMLCAFAYSSR	KTRLELSKH
Ot17g01940	-----KERQE	QAAVAGLI	PILVKIVA	-----	ERN	-----	-----	VAHTVQQS	SAMKLA	VPLLCD	MASASRKTREILLEKH
O116g02130	-----KERQE	QAAVAGLI	PILVQII	IASARDEKNG	-----	DGSMSTATNGDK	VNTAA	----	VATTARSSAMT	LAVPLLCD	MASTRKTREILLEKH
	1410	1420	1430	1440	1450	1460	1470	1480	1490	1500	
SEPA_PHYPA	RGLDFYLSLLDDEVWAVT	-----	ALDSLAVCLAHDNEQRKVE	QALLQKDALQRLVA	FFQSCG	----	APSFVHILEPFLKII	ISKSVRL	NNTALAV	SGLTTP	
MAP3K_ARATH	GGLDVYLSLLDDEYWSVI	-----	ALDSIAVCLAQD	VDQ-KVEQAFLK	DAIQKLVNFFQ	NCP	----	ERHFVHILEPFLKII	TKSSS	INKTLALNGLTP	
SEPA_DICDI	NGVAFYLSLLEERYQVN	-----	ALDSLAVWITDETHK	--VENI	IATNENIKKLIQLFT	NAE	----	SQSFAGILEP	LLKIIQIS	SIPVNILLGTSNFIT	
Mipur05g06740	DALDAYLSIARNKPSASSPPGLQLA	AVRAVAGW	MRDEPWK--	VEARLAE	PDAA	AASAI	DPRR	----	IPPIDPE	VLDTLREIARE	SPRLCAALASGGAMA
Mipuc16g03100	RALDAYLSLVASGEGGAG	----	SATALRAVSAWQRDE	PWA--	AEARLME	PD	AVASIAAALQ	PRPGDVLLPCD	-AFL	TLLGLLKR	SPRVCAGVAQGGGMA
Ot17g01940	GALDAYVELISVESGWTL	-----	ALNAVGSWLAIEPWK	--	AEARFLE	PD	AINSILDV	LDEST	----	SQD-NIMQ	ALLNLISTSPRLCQALANEHFIP
O116g02130	GALDITYVQLIAVNSGWTL	-----	ALNAVGSWLAIEPWK	--	VEARLLETDA	IDSILEV	LDEST	----	SQA-NVLE	ALLDLISK	SPRLCQALANEDFIP
	1510	1520	1530	1540	1550	1560	1570	1580	1590	1600	
SEPA_PHYPA	LLVARLENQDA	-----	IAR	-----	-----	-----	-----	LTKL	KIRAVY	EHHP-RPKQLIVEHDL	
MAP3K_ARATH	LLIARLDHQDA	-----	IAR	-----	-----	-----	-----	LTKL	KIKAVY	EKHP-KPKQLIVENDLP	
SEPA_DICDI	KIIDKLGHTNP	-----	QVR	-----	-----	-----	-----	LTKL	KIITS	LYEHP-NAKKMIQEFLI	
Mipur05g06740	PLVEALGAPSSATAAVSSRATHNAGTVDW	NNGAVLRR	SNSRRLELRPGRAGAGS	RDDPSRPHDARKDPRRT	LAYVRL	LLAQLADHRRDAV	DESGR	GRHDLV			
Mipuc16g03100	PLVEALGPPSSSAAAASGLSRSS	-----	GGAFR	--G-RLLE	THPG	-----	RALLFLRLLAVVY	EHHP	-RRE	EEAAAGYDVR	
Ot17g01940	PLMDALTTAAT	-----	KPTTR	-----	-----	-----	ITLLKMLG	VVHEHAQ	-RP	KELIIRYDVV	
O116g02130	SLMDAITSPTT	-----	KPTIR	-----	-----	-----	ITLLKTLG	VVHEHAQ	-RP	KELIIRHDLV	

	1610	1620	1630
	00		
SEPA_PHYPA	TKLQRLIEDRRD	GGERSGGQVLVKQMAYSLLRALHINTVL	
MAP3K_ARATH	QKLQNLIEERRD	QSRGGQVLVKQMATSLLKALHINTIL	
SEPA_DICDI	PIIQKIADT--	DKS-----VLVQKMASKLLEAFNANTVI	
Mipur05g06740	GRLRGVDDDESAD	SERAAEDVR--TEAERLLAELRR----	
Mipuc16g03100	ARLRGVLERGGGG	GREAAAV--DVAEELLRKMR-----	
Ot17g01940	ARLKRLIEGVTD	ER-SSGAVA--QLVDKILRSMRLTRVA	
O116g02130	ARLKSFLT	DHGDERQSHTVIA--QLVDKILRSMRLSRVV	

d) U5 Helicase

	10	20	30	40	50	60	70	80	90	100
Mipuc05g04870	MG--GGAEGKRGAAAP-----	PRLSTFLRFTFLAT	TND-----	ARDAVDVDAEFVARAAI	LARAGDGGAGGHG-AGRYHATASLAARVLGGSDASP					
Mipur08g03240	MG--HGRDGVSGGSQSGSTRPA	PRLSALLRSLATADANE-----	DVDALYVARRGALDAAIASGSATGPAPGRGGKGRYHATSSLAR-----	L						
O108g01930	M-----ATTHAVA-----	PRLTTFLLRAARG-----	LDDVDVDAERARRNETLR	TTRRDATTD-----	DTPSVRLAR-----	R				
Ot08g01860	M-----SSIVVR-----	PRLTTFLLRASRAS---S-----	VDDIVDVEAERARRTRGFEGVTRGGGERER----	EGPSARLAS-----	V					
Bathy01g04990	MAPLSTSSSSSSSSSSSSRRPP	PRLTTFLLRARC	GGKSGDGSDDS	IDRYWYVSDVDEAFVARAKQIGSATTKTSSPSLSFRQTAKKATKNRDLALLKLSDV						
Q9SYP1_ARATH	M-----ANLGGGAEAH-----	ARFKQYEYRANSS-----	LVLTTDNRPRD	THEPTGEPETLWGKIDPRS----	FGDRVAKGRP-----	QEL				
A9RTW1_PHYPA	M-----AHLGGGAEAH-----	ARFKQYEYRANSS-----	LVLTTDTRPRD	THEPTGEPESLYGRIDPRS----	FGDRVYHGRA-----	NDL				
	110	120	130	140	150	160	170	180	190	200
Mipuc05g04870	PSPSSPSSPSSSSRDPP	KDARDALDAFLPVVVKVT--GGDTLLPDQLAD	ASLVAFDAIAGLCEESIAAAERWDEETTYGRDP--DLAKKKAAAE----	E						
Mipur08g03240	ICPASPDSP-----	PREVRAALERFLRAVVRVT--GGDTLLPDALAD	ASLLAFRTVASLARDANAAAAARWAEERYGRDA--DLAKRRGQAE----	R						
O108g01930	VLRCE	DGRE-----AKDVVKTCRAFVSVVV	RML--GSDTMTETESGDAASAWDAASGFVEVIANERRAVE	SERGYVAD--HETARKIFEVR----	E					
Ot08g01860	VFGNDGSDR-----	DKEVASACRAFMSAVLRAL--GGDTMTEDALGDAVVAWEVLRDNVETLGRERRDLE	ERQPYVDDAVFETERKLF	FEVR----	E					
Bathy01g04990	ALKTSSSSLED-----	FDEETRLALDLVADASIQLRLLENDLISPDELLE	APADLYAAVATNCDVDWKLEADLE	DVVEYSNAP--ELERKKVQKRRGNE						
Q9SYP1_ARATH	EDKLLKSKKK-----	ERDVVDDMVNIRQSKR-----	RRLREESVLTDTDDAVYQPKTKETRAAYEAMLGLIQKQLGGQPP-----	SIVSGAAD-----						
A9RTW1_PHYPA	EELTKHRRKR----	EVKEKEKGSNAEGLKKARKRL---	RGMQEESVLSIVDDGMYRPKTKETRAAYEALLSTIQQFGDQPQ-----	DILRGAAD-----						

	210	220	230	240	250	260	270	280	290	300
Mipuc05g04870	RERARRMKTALAKALGPVDE--KHVPDLIDVAKRLVHIQKKRKGTSDDRDRVDASAAEFGAAAVPARREFGADVVFNAP-RSRAVGSRASTSASTLIPRPA									
Mipur08g03240	RDFMRRKRAALVAELGPVSD--QHLDEFDAAVNELIALQDAHGGGPGGDGANDGDGDD-G-----YEFGSDFPFAAPGRGRAPASSSSR-TSS-----									
O108g01930	REVRKDMMEALRGALGPVEANAPTVEFEFESLVEDLLRMRGE-FDEFHTLQSTLDAEPTTSG-----RIFGQNIKFKRA--GDAQPSDRIR-----									
Ot08g01860	RELRRDFGEKLRGALGPIEASAATLRELEQHVETLLKFRERITDLDYMIASSAAPSTSG-----RSFGHNIKFKCR--GDVQPSDRVR-----									
Bathy01g04990	REAFGFVKARLGGVLESETGEAAQRWKDFWQSANALRRVGSSTSTGGDLNGTSMGVGGKKE-----TYYGEDMFFVPPRFYGANDFDSSN-----									
Q9SYP1_ARATH	-EILAVLKNDAFRNPEKKMEIEKLLNKIENHEFDQLVSIKGLITDFQEGGDSGGGRANDD-----EGLDDDLGVAVEFEENEEDDEESD-----									
A9RTW1_PHYPA	-EVLGVLKNDRFRDLDDKKKEIEKLLNSMSNERFAQLVAIGKLIISDYSEGGDAGAEGAG-----EALDDDI GVAVEFEFEENEEDDES D-----									
	310	320	330	340	350	360	370	380	390	400
Mipuc05g04870	PAASASAAASASAAAKLAATMKAGLRAAGEAEAEKRRGENAGTSP-LDASASRATSFGLPPAEQSANEQLRWLNEKCVVEHVGG---V-AGAGWEDVAQG									
Mipur08g03240	---SAQAPASTSSAARMAEQMRAGLLAATVSDADGKKESELYGAS---DANSASEGEFIS-IPAAEASATAQLAWLKEKRCVAISG-----ANGWEETAQQ									
O108g01930	-----DLKRAMKTNLSEASEAERLRIESEMHPDE-----RMSGFGEPEHVESDATTQTLTWLKRQCEGFVANSAHA-LDQTWEAVAGT									
Ot08g01860	-----DLKRAMKANLDELDEVARARIKSEKISHPN-----ELFGEPEQREESDVTTLTWLKRQCEGFVANSAHA-LDQSWAVASA									
Bathy01g04990	---QHCDQISSGLESTILSNMTRGLDSANAFASLELERSRTPSFHELDALG---TESEVSTSLQAGSGSILSWLREQCEIFRENNQNSDAFAGWTDICGN									
Q9SYP1_ARATH	-----PDMVEEDDDEEDD-EPTRTGGMQVDAGINDE-----DAGDANEGTNLNVQDIDAYWLQRKISQAYEQ-----QIDPQQCQV									
A9RTW1_PHYPA	-----DEVQEESDGEEDGQDTRQASAMQMGQDDE-----DMEEADEG--LNVQDIDAYWLQRKISQAHG-----DIDPQQSQK									
	410	420	430	440	450	460	470	480	490	500
Mipuc05g04870	VARATLSDASR-SDDDVAAEFLFDLGDGGVELIMGAIERRVAMCAALRKRIITTLRETLGGGKDDDGEDRDRAAAGGPGRVTVTISSTTDKQIEKLRKKEER									
Mipur08g03240	LARALLAAETR-SDDEVAAEFLFDLGDGGVETIAGAIERRAAINAALRRRIGTLRETLGG-AGGDGDGGGGARKDAPMAQVTIQSTSDIKMEKLRKKEER									
O108g01930	VGRAIMNASS--SDDECAAELYEYLG DYGFELIAGVVERRTVLTSAIKNRAQLLRDALNAQSGADRA-----GPSITRVVTTITSTLDKQIEKARRKKEER									
Ot08g01860	VGRSIMNSTI--SDDACAAELYEYLG DYGFELIAGVVERTELASAIKKRAQMLRETLAQSGVD-----GPSVARVVTINSTLDKQIEKMRKKEER									
Bathy01g04990	ICRVCLNTSS--SDDQVASEFLFDLGDGGFDLVASVCERRGQLADAIIRRRLQALKEAFDPEESAQDD-----YAKSKSAVSVRSTTDVAMEKIRKKEER									
Q9SYP1_ARATH	LAEELLKILAEGDDRVVEDKLLMHLQYEEKFSLVKFLLRNRLKVVWCTRLARAEDQEERNRIEEMRGL--GPELTAIVEQLHATRATAKEREENLQKSIN									
A9RTW1_PHYPA	LAEDVLSKLAEGDDREVENRLVILLDYDKFDLIIKLLLRNRLKVVWCTRLARAEDDARKKIEEEMSNG--GPVLAGILEQLHATRATAKERQKNLERSIR									
	510	520	530	540	550	560	570	580	590	600
Mipuc05g04870	KVGRRIAQQGGEPLLEWLAA-SGVGFGVLCCEGDWEAAAASGGGR-GEDDVVFAGLR--VGGGSGRKALPAGTTRKVVH-KGYEEVAVPAAKVAPVGDAEERF									
Mipur08g03240	KVGRRIAQQGGEPLLEWLAN-SGVGFAALCEGDWEAAANQPS---TEDDIWAGLYGLGGGGGGGKALPAGTTRKVVH-KGYEEVHVPAGERAPVGEHERF									
O108g01930	KANRKLASGSGASIMEWLQA-VGVGFDAALCEGDWENQQAPSSSS-NPDDILAGLRGLGRGMDGGRKALPAGTTRIVHPEGYEESVPAEPDPVAAGERS									
Ot08g01860	KANRKLASG--ENVMEWLQA-VGVGFDAALCEGDWENQQTPSSSS--PDDILASLRGLGSLGGGRKALPAGTTRMVEHPEGYEESVPAEPDPVAAGERS									
Bathy01g04990	RNKRRRAAGGHGEYLLWFNTDSGLGYGAFCDLPLPNRNEAAPGSVDDIILNSLRGLGLGTDGGKALPAGTTRKVL-EGYEEIYVPAIPDAVADGELQ									
Q9SYP1_ARATH	EEARRLKDETGGDGGRR-----DVADRSESGWVKQRQMLDLESIAFDQG-GLLMANKKCDLPPGSYRSHG-KGYDEVHVPWVS-KKVDRENEKL									
A9RTW1_PHYPA	EEAKLRDDGGEAADRGRK---D--REVGVGGGESGWLKGQRQLLDLEQLTFHQG-GLLMANKKCELPPLSYRTPK-KGYEEVHVPVHLKPKPFAEGEEL									

	610	620	630	640	650	660	670	680	690	700																																																																																										
Mipuc05g04870	V	A	I	E	E	L	D	D	W	A	Q	L	A	F	A	G	M	T	S	L	N	R	I	Q	S	K	I	Y	P	A	A	F	R	S	N	E	N	L	L	V	C	A	P	T	G	A	G	K	T	N	I	A	M	L	T	V	L	H	E	I	G	A	H	F	D	D	D	G	E	W	N	G	D	--	D	F	K	I	V	Y	V	A	P	M	K	A	L	A	A	E	V	T	N	A	F	S	R	R	L	
Mipur08g03240	V	P	I	E	E	L	D	D	W	A	Q	P	A	F	A	G	M	K	S	L	N	R	I	Q	S	R	I	Y	E	A	A	Y	H	S	N	E	N	L	L	V	C	A	P	T	G	A	G	K	T	N	I	A	M	M	T	V	L	H	E	I	G	H	I	E	Y	G	E	L	A	Y	G	A	--	D	F	K	I	V	Y	V	A	P	M	K	A	L	A	A	E	V	T	G	A	F	S	R	R	L		
O108g01930	V	A	I	E	E	L	D	E	W	A	Q	P	A	F	Q	I	R	M	L	N	R	I	Q	S	K	I	F	P	Q	A	Y	H	T	N	E	N	L	L	V	C	A	P	T	G	A	G	K	T	N	I	A	M	L	T	V	L	H	E	I	G	L	H	I	D	E	N	G	D	Y	L	P	E	--	D	F	K	I	V	Y	V	A	P	M	K	A	L	A	A	E	V	T	D	A	F	S	R	R	L		
Ot08g01860	V	A	I	E	E	L	D	E	W	A	Q	P	A	F	K	G	I	K	L	N	R	I	Q	S	R	I	F	P	T	A	Y	H	T	N	E	N	L	L	V	C	A	P	T	G	A	G	K	T	N	I	A	M	L	S	I	L	H	E	I	G	L	H	I	D	E	N	G	D	Y	L	P	E	--	D	F	K	I	V	Y	V	A	P	M	K	A	L	A	A	E	V	T	E	T	F	G	R	R	L		
Bathy01g04990	V	S	V	S	Y	L	P	A	W	A	Q	T	A	F	K	G	I	Q	T	F	N	R	I	Q	S	K	I	F	E	C	A	Y	T	S	N	E	N	L	V	C	A	P	T	G	A	G	K	T	N	I	A	M	L	C	A	M	Q	E	I	A	K	H	F	D	E	E	N	N	C	L	H	E	H	D	D	F	K	I	V	Y	V	A	P	M	K	A	L	A	A	E	V	T	R	T	F	Q	K	R	L	
Q9SYP1_ARATH	V	K	I	T	E	M	P	D	W	A	Q	P	A	F	K	G	M	Q	L	N	R	V	Q	S	K	V	Y	D	T	A	L	F	K	A	E	N	I	L	L	C	A	P	T	G	A	G	K	T	N	V	A	M	L	T	I	L	Q	Q	L	E	M	N	R	N	T	D	G	T	Y	N	H	G	--	D	Y	K	I	V	Y	V	A	P	M	K	A	L	V	A	E	V	V	G	N	L	S	N	R	L		
A9RTW1_PHYPA	V	K	I	S	D	M	P	D	W	A	Q	P	A	F	K	G	M	K	S	L	N	R	V	Q	S	K	V	Y	E	T	A	L	F	T	S	E	N	L	L	L	C	A	P	T	G	A	G	K	T	N	V	A	M	L	T	I	L	H	E	L	G	L	R	K	Q	L	D	G	T	F	D	L	S	--	S	F	K	I	V	Y	V	A	P	M	K	A	L	V	A	E	M	V	G	N	F	S	E	R	L	
	710	720	730	740	750	760	770	780	790	800																																																																																										
Mipuc05g04870	A	P	L	G	I	T	V	R	E	L	T	G	D	T	Q	L	T	K	K	E	L	E	E	T	M	I	V	T	T	P	E	K	W	D	V	I	T	R	K	G	G	E	V	S	V	A	S	T	L	G	L	L	I	D	E	V	H	L	L	N	D	E	R	G	P	V	I	E	T	L	V	A	R	T	H	R	Q	V	E	T	T	Q	S	M	I	R	I	V	G	L	S	A	T	L	P	N	P	M		
Mipur08g03240	E	P	L	G	I	Q	V	R	E	L	T	G	D	T	Q	L	T	K	K	E	M	E	E	T	H	M	I	V	T	T	P	E	K	W	D	V	I	T	R	K	G	G	E	V	S	V	A	S	S	L	R	L	L	I	D	E	V	H	L	L	N	D	E	R	G	P	V	I	E	T	L	V	A	R	T	H	R	Q	V	E	T	T	Q	S	M	I	R	I	V	G	L	S	A	T	L	P	N	P	A	
O108g01930	A	P	L	D	I	V	V	A	E	L	T	G	D	T	Q	M	S	K	R	E	L	E	T	Q	M	I	V	T	T	P	E	K	W	D	V	I	T	R	K	G	G	E	V	S	V	A	S	T	L	R	L	L	I	D	E	V	H	L	L	N	D	E	R	G	P	V	I	E	T	L	V	A	R	T	L	R	Q	V	E	Q	T	Q	S	M	I	R	I	V	G	L	S	A	T	L	P	N	P	V		
Ot08g01860	A	P	L	D	I	V	V	A	E	L	T	G	D	T	Q	M	S	K	R	E	L	E	T	Q	M	I	V	T	T	P	E	K	W	D	V	I	T	R	K	G	G	E	V	S	V	A	S	T	L	R	L	L	I	D	E	V	H	L	L	N	D	E	R	G	P	V	I	E	T	L	V	A	R	T	L	R	Q	V	E	Q	T	Q	S	M	I	R	I	V	G	L	S	A	T	L	P	N	P	L		
Bathy01g04990	D	E	L	G	M	V	C	R	E	L	T	G	D	T	Q	L	S	K	R	E	L	E	T	H	V	I	V	T	T	P	E	K	W	D	V	I	T	R	K	G	G	E	V	S	V	A	S	T	L	R	L	L	I	D	E	V	H	L	L	N	D	E	R	G	P	V	I	E	T	L	V	A	R	T	R	R	Q	V	E	Q	T	Q	S	M	I	R	I	V	G	L	S	A	T	L	P	N	P	R		
Q9SYP1_ARATH	K	D	Y	G	V	I	V	R	E	L	S	G	D	Q	S	L	T	G	R	E	I	E	T	Q	I	I	V	T	T	P	E	K	W	D	I	I	T	R	K	S	G	D	R	Y	T	Q	L	V	R	L	L	I	D	E	I	H	L	L	H	D	N	R	G	P	V	L	E	S	I	V	A	R	T	L	R	Q	I	E	T	T	K	E	N	I	R	L	V	G	L	S	A	T	L	P	N	Y	E			
A9RTW1_PHYPA	E	P	Y	G	V	T	V	R	E	L	T	G	D	A	T	L	S	R	G	Q	I	E	E	T	Q	I	I	V	T	T	P	E	K	W	D	I	I	T	R	K	S	G	D	R	Y	T	Q	M	V	K	L	L	I	D	E	I	H	L	L	H	D	N	R	G	P	V	L	E	S	I	V	A	R	T	V	R	Q	I	E	T	T	Q	E	M	I	R	L	V	G	L	S	A	T	L	P	N	Y	E		
	810	820	830	840	850	860	870	880	890	900																																																																																										
Mipuc05g04870	D	V	A	K	F	L	G	V	S	--	D	A	G	L	F	V	F	D	Q	S	Y	R	P	I	P	L	T	Q	V	F	V	G	V	T	E	G	N	A	M	K	R	L	N	L	M	A	E	I	A	Y	D	K	C	A	G	A	L	K	S	G	K	Q	A	M	V	F	V	H	S	R	K	D	T	V	K	T	A	R	Q	L	A	E	L	A	N	--	A	E	G	G	V	E	L	F	G	C	A	E		
Mipur08g03240	D	V	A	K	F	L	G	V	S	--	D	A	G	L	F	V	F	D	Q	S	F	R	P	I	P	L	T	Q	M	F	V	G	V	T	E	G	N	A	M	K	R	Q	M	L	M	A	Q	I	A	Y	D	K	C	T	A	A	L	R	S	G	K	Q	A	M	V	F	V	H	S	R	K	D	T	V	K	T	A	K	Q	L	G	E	I	A	A	N	D	Q	T	Q	G	G	L	E	L	F	--	A	P	E
O108g01930	D	V	A	R	F	L	G	V	N	N	D	A	G	L	F	V	F	D	Q	S	Y	R	P	I	P	L	T	Q	K	F	I	G	V	T	E	K	N	S	M	K	R	Q	T	L	M	A	Q	I	A	Y	N	K	A	C	E	A	L	R	N	G	K	Q	A	M	V	F	V	H	S	R	K	D	T	V	K	T	A	R	Q	L	A	E	F	A	A	---	Q	D	G	M	E	L	F	--	S	N	N			
Ot08g01860	D	V	A	R	F	L	G	V	N	N	D	A	G	L	F	V	F	D	Q	S	Y	R	P	I	P	L	T	Q	K	F	I	G	V	T	E	K	N	S	M	K	R	Q	T	L	M	T	Q	I	A	Y	N	K	A	C	E	A	L	K	N	G	K	Q	A	M	V	F	V	H	S	R	K	D	T	V	K	T	A	R	Q	L	A	E	F	A	A	---	Q	G	G	L	E	L	F	--	S	N	E			
Bathy01g04990	D	V	A	R	F	L	G	V	T	E	G	K	L	F	V	F	D	Q	S	Y	R	P	I	P	L	T	Q	V	F	I	G	V	S	E	T	N	A	M	K	R	Q	N	V	T	I	R	I	A	F	K	K	A	C	E	A	L	R	K	G	Q	A	M	V	F	V	H	S	R	K	D	T	V	K	T	A	R	Q	L	A	E	I	A	G	E	---	E	G	E	L	E	L	F	--	E	N	D				
Q9SYP1_ARATH	D	V	A	L	F	L	R	V	D	L	K	K	G	L	F	K	F	D	R	S	Y	R	P	V	P	L	H	Q	Y	I	G	I	S	V	K	K	P	L	Q	R	F	Q	L	M	N	D	L	C	Y	Q	K	V	L	A	G	A	G	--	K	H	Q	V	L	I	F	V	H	S	R	K	E	T	S	K	T	A	R	A	I	R	D	T	A	M	A	---	N	D	T	L	S	R	F	--	L	K	E			
A9RTW1_PHYPA	D	V	A	L	F	L	K	V	D	E	K	K	G	L	F	Y	F	D	N	S	Y	R	P	C	P	L	A	Q	Y	I	G	V	T	V	R	K	P	L	Q	R	F	Q	L	M	N	D	I	C	Y	E	K	V	M	E	V	A	G	--	K	H	Q	V	L	I	F	V	H	S	R	K	E	T	A	K	T	A	R	A	I	R	D	A	A	L	A	---	N	D	T	L	G	R	F	--	L	K	E			
	910	920	930	940	950	960	970	980	990	1000																																																																																										
Mipuc05g04870	D	D	E	G	K	K	R	F	K	T	E	I	D	R	S	R	N	N	E	L	K	E	L	V	G	K	G	F	G	C	H	N	A	G	M	L	R	S	D	R	T	L	V	E	K	L	F	A	G	V	V	K	V	L	V	C	T	A	T	L	A	W	G	V	N	L	P	A	H	T	V	V	I	K	G	T	Q	L	Y	D	P	Q	K	G	G	F	R	D	L	G	V	L	D	V	Q	Q	I	F	G	
Mipur08g03240	N	H	P	D	F	T	T	W	K	K	E	V	E	R	S	R	N	N	E	L	K	E	L	F	H	R	G	F	G	C	H	N	A	G	M	L	R	S	D	R	T	L	V	E	R	L	F	S	A	G	V	V	K	V	L	C	C	T	A	T	L	A	W	G	V	N	L	P	A	H	T	V	V	I	K	G	T	T	L	Y	D	P	S	K	G	G	F	R	D	L	G	V	L	D	V	Q	Q			

	1010	1020	1030	1040	1050	1060	1070	1080	1090	1100
Mipuc05g04870	RAGRPGFDTS	GEGVIVTEHKKLAHYLSLLTHSTPIESQFISCLADNLNAEIVLGTVTNVKEGAQWLGYSYLHTRMEKNPLAYGITWDDVKLDPGLGEHRR								
Mipur08g03240	RAGRPGFDTS	GEGVIVTEHKKLAHYLALLTHSTPIESQFISCLADNLNAELVLTVCVSVKEGAQWLGYSYLHTRMEKNPLAYGLTWDDVNLDPGLVRRHR								
O108g01930	RAGRPGFDTS	GEGVIVTEHKNLAHYVSMMLTHSTPIESQFVSNLADNLNAEVLGTVTNVREGAQWLGYSYLHTRMEKNPLAYGLTWDDIRLDPGLLDHRR								
Ot08g01860	RAGRPGFDTS	GEGVIVTEHKNLAHYIAMLTHSTPIESQFISNLADNLNAEVLGTVTNVREGAQWLGYSYLHTRMEKNPLAYGLTWDDVRLDPGLLDHRR								
Bathy01g04990	RAGRPGFDTS	GEGVIVTEHKKLTKYVAMLTHSTPIESQFIECLADNLNAEIVLGTVTNVREGAQWLSYSYLHTRMEQNPLGYALTWDEVRLDPGLIEHRR								
Q9SYP1_ARATH	RAGRPQYDQH	GEGIIITGYSLELQYYLSLMNEQLPIESQFISKLADQLNAEIVLGTVQNAREACHWLGYTYLYIRMVRNPPLYGLAPDALAKDVVLEERRA								
A9RTW1_PHYPA	RAGRPQFDTY	GEGIIITGHSELQYYLSLMNQQLPIESQYISKLADNLNAEIVLGSVQDAREACDWLGYTYLYIRMLKNPPLYGVSREALEADPSLEERRA								
	1110	1120	1130	1140	1150	1160	1170	1180	1190	1200
Mipuc05g04870	00	KLVKEAARTLDRAKMIRFDE	SGQLYQTEAGRIASHFYIKQTSMEMFDEHLKRHMSVPEVFMVSHAG	EFENISPREDEMP	ELET	LR	DR	DK	KN	ACPIEVKA
Mipur08g03240	KLVTEAARTLHRAKMVRFDE	KSGFIYQTEAGRIASHFYIKQASME	LFDLQRHMSMPEVFMVAQATEFENI	APREDEMP	ELEAL	RR	NR	KG	ACPLEIKA	
O108g01930	KLIKEAARVLDRAKIRFDE	SGQLYQTEAGRTASHFYIRVNSME	FDGLMHRHMTLPDIFHMI	SHSSEFENIVPREDEI	PELET	LR	NR	RR	VVPIDIKA	
Ot08g01860	KLIKEAARTLDRAKMIRFDE	SGQLYQTEGGRTASHFYIRVSSME	FDLSMHRHMTLPEVLHMI	SHSSEFENIVPREDEI	PELET	LR	DR	RR	IIPVEIKA	
Bathy01g04990	NLIKTAARKLHKAKMIRFDE	QSGQLYQTEAGRIASHFYIKVTSME	LFEEMNRHMSLPEVLHVI	SHSSEFENIAPREDEMP	ELEAL	RR	NR	RS	ACPIEIKG	
Q9SYP1_ARATH	DLIHS	AATILDKNNLVKYDRKSGYFQV	TDLGR	IASYYYITHGTIATYNEHLKPT	MGDIDL	YRLF	SLSDEFKYVTVRQDE	KMELAKLL	DRVP----	IPIKE
A9RTW1_PHYPA	DLVHS	AAIVLDRNNLVKYDRKSGYFQV	TDLGR	IASYYYISHGSMATYNEHLKPT	MGDIEL	CRLF	SLSEEFKFVTVREE	KMELAKLL	DRVP----	IPVKE
	1210	1220	1230	1240	1250	1260	1270	1280	1290	1300
Mipuc05g04870	00	TLADKAGKVNLLQVYIS	RRAMEAFSLIADSSYISQNASRICRALYELCLRRGWPSLAE	TLLTLLKTVDLRIWPHQHTLRQFETT	LS	PD	TL	YR	LE	TRDAT
Mipur08g03240	TLADRAGKVNLLMQVYIS	RRAMEAFSLVADSSYISQNASRICRALFELCLRRGWPSLAE	ELLTLSKAVDLRIWPHQHALLRQFE	QTLSPETLYKLEERQAT						
O108g01930	SLTDKVGKVNLLQVYIS	RRASMQFSLIADSMYISQNASRICRALFELCLRRGWPSLAE	QLLTVSKSCDLRIWPHQHELRQFE	KS	LP	EV	LF	KL	EEKKAT	
Ot08g01860	SLTDRVGKVNLLQVYIS	RRANMQFSLIADSMYISQNASRICRALFELCLRRGWPSLAE	QLLTVSKACDLRIWPHQHELRQFE	KT	LP	EV	LY	KL	EEKKAT	
Bathy01g04990	DMSDKIAKVNLLQVYVSR	KRLEFSLVADSSYISQNASRICRALFELMLKRGWPSLAE	TLLTLSKAVDRRLWPHHSP	LRQFENTL	KP	ET	IY	KL	EEKD	AT
Q9SYP1_ARATH	TLEEPSAKINVLLQAYIS	QLKLEGLSLTSDMVYITQS	AGRLVRALYEIVLKR	GWAQLAEKALNLSK	MV	GK	RM	WS	VQ	TP
A9RTW1_PHYPA	SLEEPSAKINVLLQAYIS	QLKLEGLSLTSDMVFITQS	AGRLMRALFEIVLKR	GWAQLAEKALT	LCKM	VS	RR	MWS	SQ	TP
	1310	1320	1330	1340	1350	1360	1370	1380	1390	1400
Mipuc05g04870	VERLW	DMSPEIGSLLRLNTDV	GKKVKGCLEALPHLAMEASVQ	PITRSVLRVSVTLT	PDFIWRDSQHGGIQRWLVW	VEDPVNEHIYHTET	FNL	SK	KQ	HKHE
Mipur08g03240	VERLFDMSAQEIGSMLRLNTAV	GQKVRGCELEPHLTMEATVQ	PITRSVLRVTVALTPEFKWRDAV	HGGLQRWLVW	VEDPVNEHIYHNET	FML	SK	KL	HGE	
O108g01930	LDRRLWMSASEIGSMLRLNTQ	IGGQVKSCMRAMPHLNMTAV	QIPITRSVLRVSVTLTPEFEWRDAV	HGALQRWLIW	VEDPVNEHIYHSET	FNL	SK	KQ	SRD	
Ot08g01860	LDRRLWMSGGEIGSMLRLNA	QIGGQIKSCMRAMPHLNMTAT	VQIPITRTVLRVSVTLIPEFEWRD	QLHGALQRWLIW	VEDPVNEHIYHSET	FNL	SK	KQ	CRD	
Bathy01g04990	VDR	LIDVSAKEVGDLLRLNAV	VGAQVKRCVEQLPHVNLEAV	VRPITRSVLRVSATLTPEFMWR	DEVHGQAQWLIW	VEDPVNEHIYHTET	FNL	SK	KQ	YKE
Q9SYP1_ARATH	WERYYDLSAQELGELIRS	-PKMGKPLHKFIHQFPK	VTL	SAHVQIPITRTV	LNVELT	VPDF	WDEK	IHKY	VE	PFWI
A9RTW1_PHYPA	WERYYDLS	SQEI	GELIRY-PKM	GKSIHRYIHQFPK	LELAHVQIPITRS	VLRV	DLTITPDF	QWDE	KYHGY	VE

	1410	1420	1430	1440	1450	1460	1470	1480	1490	1500		
Mipuc05g04870	GKQHMAFTIPIFEPMP	QYFLRATSESWLGCET	FLELRFDGLVLPQKH	PPHTDLLDLTLP	PRSA	LN	-----	EKYESLYAKKFT	HFNAIQ	TQAFHTL		
Mipur08g03240	GKQHLAFTIPIFEPV	PPQYFLRATSESWLGCET	FLELNFNELVLPDRG	PAHTELLDL	PPVPRQAL	YPENPELGRKE	FFDL	YEGKFEFF	NKVQ	TQAFNTL		
O108g01930	GAQYLAFTIPIFEPV	PPQYFLRAMSETWLGCES	FVELNFQHLLI	LPEEHP	PPHTELLDL	DPLPRSA	LN	-----	PVYESMYEGKFT	HFNAIQ		
Ot08g01860	GAQYLAFTIPIFEPV	PPQYFLRAISENWLGCES	FVELNFQHLLI	LPEEHP	PPHTELLDL	DPLPRSA	LN	-----	PVFESMYEKKFT	HFNAIQ		
Bathy01g04990	GRMTLAFTIPIFDPR	PPQYFLRATHLYWLGCES	FLELDLEDIVLPTE	PPNTELLDL	EPLPRSA	LN	-----	PTYESLYEKKFT	HFNAIQ	TQAFHTL		
Q9SYP1_ARATH	EDHTLHF	FTVPIFEPL	PPQYFVRVSDKWLGSET	VLPVSRHLI	LPEKY	PPPTELLDL	QPLPVTALRN	-----	PNYEILYQ	-DFKHFNPVQ	TQVFTVL	
A9RTW1_PHYPA	EDHNLS	FTVPIYEPL	PPQYFVRVSDRWLGSET	VLPVSRHLI	LPEKY	PPPTELLDL	QPLPVSA	LN	-----	PSYEVLYQ	-KFRHFNP	IQ

	1510	1520	1530	1540	1550	1560	1570	1580	1590	1600				
Mipuc05g04870	FHTNVN	VLLGAPT	GS	GKTI	ISAE	LAMMRTFR	DE-P-GGKVVYIAP	LKALVRER	IEDWRKHL	CPVLGKRL	VELTGDYTPDLR	ALLSADIIVAT	TPEKWDG	ISR
Mipur08g03240	FHSESN	VLLGAPT	GS	GKTI	ISAE	LAMMAAFRDH	-P-GGKIIYIAP	LKALVRER	IEDWKGLCKVL	NKKL	VELTGDYTPDIR	ALQGADII	VC	TPEKWDG
O108g01930	YHTDTN	VLLGAPT	GS	GKTI	ISAE	LMMKVF	RDS-P-GSKVVYIAP	LKALVRER	IKDWRKNLCPTL	GLRM	VELTGDYTPDLR	ALLQADII	V	STPEKWDG
Ot08g01860	YHTDTN	VLLGAPT	GS	GKTI	ISAE	LMMKVF	RDY-A-GSKVVYIAP	LKALVRER	IKDWRKNLCPTL	GLRM	VELTGDYTPDLR	ALLQADII	V	STPEKWDG
Bathy01g04990	YHTNHN	VLLGAPT	GS	GKTI	ISSE	LTKMFR	DE-PPGSKVVYIAP	LKALVRER	VDDWKKYFCPTV	NKKM	VELTGDYTPDLR	ALLRADII	VAT	TPEKWDG
Q9SYP1_ARATH	YNTNDNV	LVAAPT	GS	GKTI	ICAE	FAILRNHHE	GPDATMRVVYIAP	LEAI	AKEQFRI	WEGKFGKGL	GLRV	VELTGETAL	DLKLEK	QIIIST
A9RTW1_PHYPA	YNTDDNV	LVAAPT	GS	GKTI	ICAE	FVLRMLQK	-EAGGRCVYIAP	VEALAKER	LRDWE	SKFGRTL	GVRV	VELTGETAT	DMKLEK	QIIIST

	1610	1620	1630	1640	1650	1660	1670	1680	1690	1700																																
Mipuc05g04870	N	Q	RAYVQK	VS	LV	VIDE	I	HLLGADR	GP	I	LEVIVSR	MRYISART	KQPV	RIVGL	STALANARDL	GDWLG	I	DE	-----	GLFNFR	PSV	RPV	PLE	CHI	QGF	PG																
Mipur08g03240	Q	W	QARSY	VT	KV	SLV	VIDE	I	HLLGADR	GP	I	LEVIVSR	MR	FISTR	TERP	V	RIVGL	STALANANDL	ADWLG	I	EKQ	EGPKS	GLFN	FK	PSV	RPV	PLE	CHI	QGY	PG												
O108g01930	N	W	QRAYV	TK	V	SLV	VIDE	I	HLLASDR	GP	I	LEVIVSR	MRYISART	G	SNVRI	GL	STALANARDL	GDWLG	I	D	E	-----	GLFNFR	PSV	RPV	PLE	CHI	QGF	PG													
Ot08g01860	N	W	QRAYV	KK	V	SLV	VIDE	I	HLLASDR	GP	I	LEVIVSR	MRYISART	G	SNVRI	GL	STALANARDL	GDWLG	I	E	E	-----	GLFNFR	PSV	RPV	PLE	CHI	QGF	PG													
Bathy01g04990	N	W	QRSYV	SK	V	SLV	VIDE	I	HLLGADR	GP	I	LEVIVSR	MRYISART	K	SKIR	I	VGL	STALANARDL	GDWLG	I	E	N	D	-----	GLFNFR	PSV	RPV	PLE	CHI	QGF	PG											
Q9SYP1_ARATH	R	W	QRKYV	Q	V	SLF	IVDEL	H	LIGG	Q	HPV	LEVIVSR	MRYISS	Q	VINKIR	I	VAL	STSLANAKDL	G	E	W	I	G	ASSH	-----	GLFNFR	PP	G	V	R	P	V	P	L	E	I	H	I	Q	V	D	I
A9RTW1_PHYPA	R	W	QRKHV	Q	V	SLF	VVDEL	H	LIGG	E	GPV	LEVIVSR	MRYIG	S	Q	TENQIR	I	VAL	STSLANAKDL	G	W	I	G	ASSH	-----	GLFNFR	PP	G	V	R	P	V	P	L	E	I	H	I	Q	V	D	I

	1710	1720	1730	1740	1750	1760	1770	1780	1790	1800																																																																																								
Mipuc05g04870	KFYCPR	MM	T	M	N	K	P	T	Y	A	A	I	R	T	H	S	P	-E	K	P	T	L	V	F	V	S	S	R	R	Q	T	R	L	T	A	M	D	L	I	A	Y	A	A	A	D	-E	R	P	E	G	F	V	H	M	S	A	N	E	L	A	G	V	R	R	A	R	D	P	A	L	K	H	C	L	Q	F	G	I	G	I	H	H	A	G	L	S	P	E	D	R								
Mipur08g03240	KFYCPR	MM	T	M	N	K	P	T	Y	A	A	I	R	T	H	S	P	-E	K	P	A	L	V	F	V	S	S	R	R	Q	T	R	L	T	A	I	D	L	I	A	Y	A	A	A	D	-E	R	P	D	T	F	V	H	M	D	P	Y	E	M	M	H	L	A	K	V	K	S	P	E	L	R	H	T	L	Q	F	G	V	G	L	H	H	A	G	L	A	P	E	D	R								
O108g01930	KFYCPR	MM	T	M	N	K	P	T	Y	A	A	I	R	T	H	S	P	-E	K	P	T	L	V	F	V	S	S	R	R	Q	T	R	L	T	A	L	D	L	I	A	Y	A	A	A	D	-E	R	P	D	G	F	V	H	M	S	D	E	L	T	M	H	L	S	K	V	K	D	P	A	L	K	H	T	L	Q	F	G	I	G	L	H	H	A	G	L	T	P	E	D	R								
Ot08g01860	KFYCPR	M	S	M	N	K	P	T	Y	A	A	I	R	T	H	S	P	-T	K	P	A	L	V	F	V	S	S	R	R	Q	T	R	L	T	A	L	D	L	I	A	Y	A	A	A	D	-E	R	P	D	G	F	V	H	M	S	N	E	E	L	S	I	H	L	S	K	V	K	D	P	A	L	K	H	T	L	Q	F	G	I	G	L	H	H	A	G	L	T	P	E	D	R							
Bathy01g04990	KFYCPR	M	L	S	M	N	K	P	T	Y	A	A	I	R	T	H	S	P	-L	K	P	A	L	V	F	V	S	S	R	R	Q	T	R	L	T	A	L	D	L	I	A	Y	A	A	A	D	-E	N	P	D	A	F	V	H	C	N	S	Q	E	L	E	Q	R	I	A	K	I	Q	D	P	A	L	K	H	T	L	Q	F	G	I	G	L	H	H	A	G	L	S	P	E	D	R						
Q9SYP1_ARATH	S	S	F	E	A	R	M	Q	A	M	T	K	P	T	Y	T	A	I	V	Q	H	A	K	N	K	K	P	A	I	V	F	V	P	T	R	K	H	V	R	L	T	A	V	D	L	M	A	Y	S	H	M	D	N	P	Q	S	P	D	F	L	L	G	K	L	E	L	D	P	F	V	E	Q	I	R	E	T	L	K	E	T	L	C	H	G	I	G	Y	L	H	E	G	L	S	S	L	D	Q	
A9RTW1_PHYPA	A	N	F	E	A	R	M	Q	A	M	T	K	P	T	Y	T	A	I	V	H	V	K	K	Q	E	P	A	L	I	F	V	P	T	R	K	H	A	R	L	T	A	L	D	L	V	T	Y	A	T	V	N	G	N	G	K	S	P	F	L	H	C	A	E	A	D	L	A	P	F	L	S	K	V	K	D	E	A	L	I	H	A	L	L	Q	G	I	G	Y	L	H	E	G	L	S	A	I	E	Q

	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900																																																																																										
Mipuc05g04870	I	C	E	E	L	F	A	E	C	K	I	Q	V	L	V	C	T	S	T	L	A	W	G	V	N	L	P	A	H	L	C	I	K	G	T	E	F	Y	D	G	K	S	R	R	Y	V	D	F	P	I	T	D	V	L	Q	M	M	G	R	A	G	R	P	Q	F	D	K	S	G	C	C	V	I	M	V	H	E	P	K	K	A	F	Y	K	K	F	L	Y	E	P	F	P	V	E	S	S	L	A	D	
Mipur08g03240	L	C	E	E	L	F	L	K	C	K	I	Q	V	L	V	C	T	S	T	L	A	W	G	V	N	L	P	A	H	L	V	V	I	K	G	T	E	F	Y	D	G	K	T	R	R	Y	V	D	F	P	I	T	D	V	L	Q	M	M	G	R	A	G	R	P	Q	F	D	T	S	A	V	A	V	I	M	V	H	E	P	K	K	A	F	Y	K	K	F	L	Y	E	P	F	P	V	E	S	S	L	A	D
O108g01930	L	C	E	E	L	F	A	Q	C	K	I	Q	V	L	V	T	T	S	T	L	A	W	G	V	N	L	P	A	H	L	V	V	I	K	G	T	E	F	Y	D	G	K	T	K	R	Y	D	F	P	I	T	D	V	L	Q	M	M	G	R	A	G	R	P	Q	F	D	K	S	G	C	C	V	I	L	V	H	E	P	K	K	T	F	Y	K	K	F	L	Y	E	P	F	P	V	E	S	S	L	A	E	
Ot08g01860	L	C	E	E	L	F	A	Q	C	K	I	Q	V	L	V	T	T	S	T	L	A	W	G	V	N	L	P	A	H	L	V	V	I	K	G	T	E	F	Y	D	G	K	T	K	R	Y	D	F	P	I	T	D	V	L	Q	M	M	G	R	A	G	R	P	Q	F	D	K	S	G	C	C	V	I	L	V	H	E	P	K	K	T	F	Y	K	K	F	L	Y	E	P	F	P	V	E	S	S	L	A	E	
Bathy01g04990	V	A	E	Q	L	F	A	E	C	K	I	Q	V	L	V	S	T	S	T	L	A	W	G	V	N	L	P	A	H	L	V	V	I	K	G	T	E	F	Y	D	G	K	T	K	R	Y	D	F	P	I	T	D	V	L	Q	M	M	G	R	A	G	R	P	Q	F	D	K	S	G	C	C	V	L	V	H	E	P	K	K	N	F	Y	K	K	F	L	Y	E	P	F	P	V	E	S	S	F	N	E		
Q9SYP1_ARATH	I	V	T	Q	L	F	E	A	G	R	I	Q	V	C	V	M	S	S	L	C	W	G	T	P	L	T	A	H	L	V	V	M	G	T	Q	Y	D	G	R	E	N	S	H	S	D	Y	P	V	P	D	L	L	Q	M	M	G	R	A	S	R	P	L	D	N	A	G	K	C	V	I	F	C	H	A	P	R	K	E	Y	Y	K	K	F	L	Y	E	A	F	P	V	E	S	Q	L	Q	H				
A9RTW1_PHYPA	V	V	T	S	L	L	T	A	E	A	I	Q	V	C	V	A	T	S	S	M	C	W	G	M	T	L	S	A	H	L	V	V	M	G	T	Q	F	Y	D	G	R	E	N	A	H	T	D	Y	P	I	T	D	L	L	Q	M	M	G	R	A	S	R	P	Q	V	D	T	S	G	K	C	V	I	L	C	H	A	P	R	K	E	Y	Y	K	K	F	L	Y	E	P	F	P	V	E	S	H	L	D	H	
	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000																																																																																										
Mipuc05g04870	N	L	P	D	H	F	N	A	E	V	V	A	G	T	I	R	S	K	Q	D	A	V	D	Y	L	T	W	T	Y	F	F	R	L	V	Q	N	P	S	Y	D	C	E	G	V	E	H	A	E	L	N	A	F	L	S	R	L	V	E	N	A	L	V	M	L	E	D	A	R	C	V	E	I	G	E	-	D	D	S	V	A	P	L	L	L	G	R	I	A	S	Y	Y	L	Q	H	P	S	V			
Mipur08g03240	Q	L	P	D	H	F	N	A	E	V	V	A	G	T	I	R	S	K	Q	D	A	V	D	Y	L	T	W	T	Y	F	F	R	L	V	K	N	P	S	Y	D	L	E	S	V	E	H	D	A	L	N	A	F	L	S	R	L	V	E	N	A	L	Q	L	E	D	A	Q	C	L	T	I	G	E	-	D	D	S	L	E	P	A	T	M	G	R	I	A	S	F	Y	L	Q	H	P	S	V				
O108g01930	N	L	C	D	H	F	N	A	E	I	V	S	G	T	I	K	T	K	Q	D	A	V	D	Y	L	T	W	T	Y	F	F	R	L	L	K	N	P	T	Y	N	L	D	T	I	Q	T	D	N	L	N	E	Y	L	S	D	L	V	E	N	A	L	L	L	E	D	A	R	C	I	A	I	D	E	E	D	D	G	L	E	P	L	M	L	G	R	V	A	S	Y	Y	L	Q	Y	P	S	V				
Ot08g01860	N	L	C	D	H	F	N	A	E	I	V	S	G	T	I	K	T	K	Q	D	A	V	D	Y	L	T	W	T	Y	F	F	R	L	L	K	N	P	T	Y	N	L	D	T	I	E	A	D	K	M	N	E	Y	M	S	D	L	V	E	G	A	L	L	L	E	D	A	R	C	I	A	I	D	D	D	D	S	L	E	P	L	M	L	G	R	V	A	S	Y	Y	L	Q	Y	P	S	V					
Bathy01g04990	C	L	E	D	H	F	N	A	E	V	G	G	A	I	K	S	K	Q	D	A	V	D	Y	L	T	W	T	Y	F	F	R	A	M	K	N	P	T	Y	N	L	E	D	T	N	H	E	T	V	N	S	Y	L	S	E	M	V	E	N	T	M	E	T	L	A	S	A	K	C	L	A	I	N	E	D	D	S	I	K	P	L	M	L	G	R	I	A	S	F	Y	L	N	F	K	T	M					
Q9SYP1_ARATH	F	L	H	D	N	F	N	A	E	V	V	A	G	V	I	E	N	K	Q	D	A	V	D	Y	L	T	W	T	F	M	Y	R	R	L	P	Q	N	P	N	Y	N	L	Q	V	S	H	R	H	L	S	D	H	L	S	E	L	V	E	N	T	L	S	D	L	E	A	S	K	C	I	E	V	E	D	-	E	M	E	L	S	P	L	N	L	G	M	I	A	S	Y	Y	I	S	Y	T	T	I			
A9RTW1_PHYPA	Y	L	H	D	H	L	N	A	E	V	V	R	T	I	E	N	K	Q	D	A	V	D	Y	L	T	W	T	F	M	Y	R	R	L	T	Q	N	P	N	Y	N	L	Q	V	S	H	R	H	L	S	D	H	L	S	E	L	V	E	S	T	L	S	D	L	E	S	S	K	C	V	A	I	E	D	-	D	M	D	L	S	P	L	N	L	G	M	I	A	A	Y	Y	I	S	Y	T	T	I				
	2010	2020	2030	2040	2050	2060	2070	2080	2090	2100																																																																																										
Mipuc05g04870	A	L	F	A	S	S	L	S	H	A	N	T	V	E	Q	L	L	K	T	L	C	G	V	A	E	Y	D	E	L	P	V	R	H	N	E	D	K	V	N	A	E	L	A	I	R	V	K	E	A	G	G	F	A	V	D	A	R	L	A	D	D	P	H	T	K	A	N	L	L	F	Q	A	H	F	L	R	V	L	P	M	S	D	Y	V	T	D	T	K	S	V	L	D	Q	A	I	R	I	I	Q	
Mipur08g03240	A	L	F	A	S	S	L	G	P	D	T	S	L	E	Q	L	L	G	I	L	C	G	V	A	E	Y	D	E	L	P	V	R	H	N	E	D	K	V	N	A	E	L	A	R	Q	V	E	D	A	G	G	F	K	V	D	A	R	L	A	D	D	P	H	T	K	A	N	L	L	F	Q	A	H	F	L	R	L	Q	L	P	M	S	D	Y	V	T	D	A	K	G	V	L	D	Q	A	V	R	I	L	Q
O108g01930	A	L	F	A	S	N	I	K	A	N	S	S	L	E	S	L	L	E	T	L	C	G	V	A	E	Y	D	E	L	P	V	R	H	N	E	D	K	L	N	T	E	L	A	E	V	A	D	A	G	G	F	Q	V	D	I	R	L	A	E	D	P	H	V	K	T	S	L	L	F	Q	C	H	F	L	R	L	P	L	P	L	S	D	Y	T	D	T	K	S	V	L	D	Q	A	I	R	I	L	Q		
Ot08g01860	A	L	F	A	S	N	I	K	A	N	S	S	L	E	D	L	L	E	T	L	C	G	V	A	E	Y	D	E	L	P	V	R	H	N	E	D	R	H	N	T	E	L	A	Q	V	A	D	A	G	G	F	Q	V	D	V	R	L	A	E	D	P	H	V	K	T	S	L	L	F	Q	C	H	F	L	R	L	P	L	P	V	S	D	Y	T	D	T	K	S	V	L	D	Q	A	I	R	I	L	Q		
Bathy01g04990	A	V	F	S	K	R	L	K	K	S	N	T	L	E	D	V	L	T	T	L	C	D	V	A	E	Y	D	E	I	P	V	R	H	N	E	D	K	L	N	A	D	L	A	I	N	V	L	K	A	G	G	Y	Q	V	D	R	R	A	Y	D	D	P	H	V	K	A	S	L	L	F	Q	A	H	F	L	R	L	P	L	P	M	S	D	Y	H	T	D	T	K	S	V	L	D	Q	S	Q	R	I	L	Q
Q9SYP1_ARATH	E	R	F	S	S	L	L	S	S	K	T	K	M	K	G	L	L	E	I	L	T	S	A	S	E	Y	D	M	I	P	I	R	P	G	E	-----	D	T	V	R	R	L	I	N	H	Q	R	F	S	F	E	N	P	K	C	T	D	P	H	V	K	A	N	A	L	L	Q	A	H	F	S	R	Q	N	I	G	G	-	N	L	A	M	D	Q	R	D	V	L	L	S	A	T	R	L	L	Q				
A9RTW1_PHYPA	E	L	F	S	S	L	T	A	K	T	K	L	K	G	L	L	E	I	L	S	N	A	S	E	Y	T	R	L	P	M	R	P	G	E	-----	E	L	I	R	K	L	V	M	H	Q	R	F	S	M	D	K	P	K	F	T	D	P	H	V	K	A	N	A	L	L	Q	A	H	F	A	R	H	S	V	S	G	-	N	L	A	L	D	Q	R	D	I	L	I	D	A	S	R	L	I	Q					
	2110	2120	2130	2140	2150	2160	2170	2180	2190	2200																																																																																										
Mipuc05g04870	A	I	I	D	V	A	S	D	A	G	W	L	H	T	A	L	N	A	M	R	L	M	Q	M	V	M	Q	G	R	F	L	T	D	S	P	L	T	T	L	P	H	V	D	A	E	V	A	G	K	L	R	R	G	-	G	V	K	S	L	P	Q	F	V	T	R	A	I	K	D	R	A	G	A	K	K	A	L	C	A	A	G	L	S	G	R	T	A	E	E	T	T	N	V	A	A	R	Y	P	S	V
Mipur08g03240	A	I	I	D	V	C	A	E	S	G	W	L	A	T	C	L	H	A	M	N	L	M	Q	M	V	M	Q	G	R	F	I	T	D	P	S	C	M	S	V	P	G	V	D	E	Q	K	A	A	S	L	S	G	S	-	G	Y	E	A	L	P	Q	L	V	D	A	C	V	N	K	E	A	A	R	K	A	L	T	N	A	G	L	K	Q	R	Q	V	D	E	A	V										

```

                2210      2220      2230      2240      2250      2260      2270      2280      2290      2300
                |        |        |        |        |        |        |        |        |        |
Mipuc05g04870  MMRASSVKTSRASAGGGKADEGVVEVHLKRLHARGGKDDGGNGGGGRSSAPRAVCPLFPKLKEEGWWLVLGD-RISGELLLALRRVGFGGAASAKLTYAAPD
Mipur08g03240  DIDVKLSK-----DGTEVEVNLRRTSKSAG--GGGKGGGRGSAPRAILPRYPKVKEEGWWLLIGD-RNNRELLSLKRVGFGQSARAKLAVDRSA
O108g01930    DAKATTETTKGIN-----GEKTVHVKLRRIGKKCG-----SK--APTSYTPRFPKIKEEGWWIVVGD-TANDELLALRRISFGDAANVKLKCPSGS
Ot08g01860    HVEASIETVKGG-----GDTTVHVQIRRRIGKKCG-----SK--APTSYTPRFPKIKEEGWWIVVGD-TANNELLALRRISFGDRADVKLKCPPSA
Bathy01g04990 DMKATLVEDKVNSSDG--RRNVSVKVSLKRSGKSG-----RKTAPRAYAPRFPKQKDEGWWIVLGEKRRTGELVAMRRAQYADTFDAVLKIDNFP
Q9SYP1_ARATH  DLTYEIVGSEEVNP-G---KEVTLQVMLER--DMEG-----RTEVGPVDSLRYPKTKEEGWWLVVGD-TKTNQLLAIKRVSLQRKVKVLDFTAPS
A9RTW1_PHYPA  DLAHEVLDNDDISP-G---DTVTLQVTLER--EMEG-----RQELSPVDAPRFPKPKEEGWWLVVCE-PKSNQLLAIKRVSLQRRSKVKLDFTAPN

```

```

                2310      2320      2330      2340      2350      2360      2370      2380      2390      2400
                |        |        |        |        |        |        |        |        |        |
Mipuc05g04870  APIGGGRGPELDLVVHLVSDCYVGMDQDLGVSEGLPAS--IDAEEDGDSSDDDGFWLPAAAVAERMKA-ASETARRLASESESESESEDEDEAFWEDETPT
Mipur08g03240  NAV-----FEPDLHVYLLISDCYVGLDQEVEVARGAGAVGAEDAGDTGDTDEQGFWLSPEQVAARLAARATERATEDTSDEDDFWEMPAAPAGERYSA
O108g01930    SSRAR-----PDLVVFLMSDSYIGLDQEVKIDSNTMVD---EDSSDEFAEDDDTFWALP-----PD-STEPFWLGEG-----
Ot08g01860    SPRPRR----QTLAVYVVSDSYIGLDQEILINADDFVE---VSD-DEVDDNADTFWLLP-----PT-QTEPFWLGEAD-----
Bathy01g04990  RGMS-----VTDITVFIMSDTYIGLDQEVLVSNTDDKR---FLSSAGVADRHRFFEERE-----S---DSEAEGNFWQDEDELS-----
Q9SYP1_ARATH  EPGE-----KSYTLYFMCDSYLGCDQEYSFSVDVKG-----SGAGDRMEE-----
A9RTW1_PHYPA  EVGR-----KTYTLLFFMCDAYLGCDQENEFTIDVKEG--VDAEDDGNAMEE-----

```

```

                2410      2420      2430
                |        |        |
Mipuc05g04870  EEVAAAAVDQDAFFWEGEGAYLAAGGDGEKKT
Mipur08g03240  VMAAKMPVEEDPPFFWENEREYLDAK-----
O108g01930    -----ENTLLT-----
Ot08g01860    -----ENSLLT-----
Bathy01g04990  -----DEDIPDF-----
Q9SYP1_ARATH  -----
A9RTW1_PHYPA  -----

```

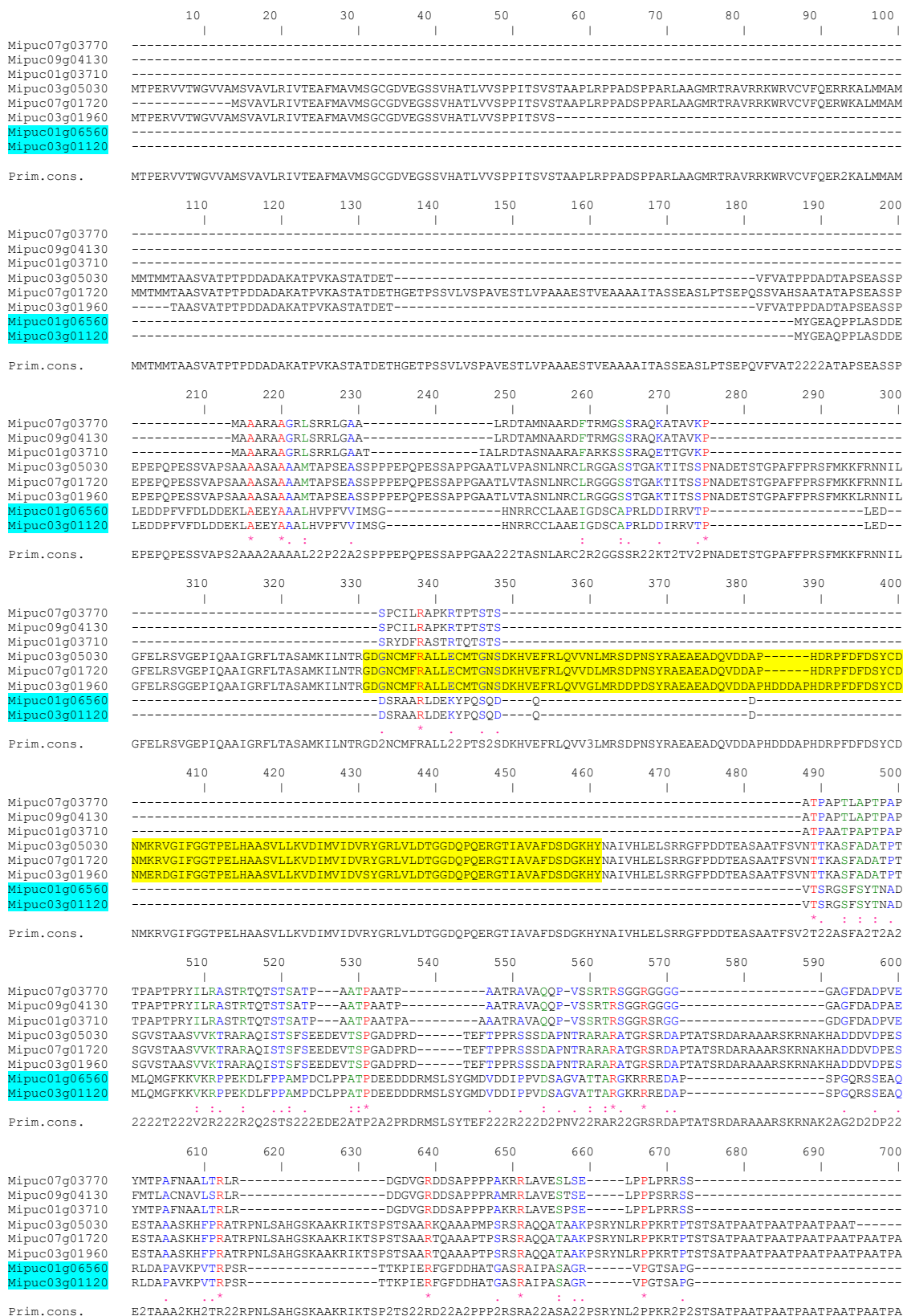
e) UAA1

```

                10      20      30      40      50      60      70      80      90      100
                |        |        |        |        |        |        |        |        |        |
Mipuc16g00610  MKTRARDAGTRGGD-----RDADARGRGAPRGRRGA-----GVFSTR----ARERSGAR----AAEARTDATSTR-----AAV
Mipur10g04930  MSDRPAYHLVPGRED-----PPTNHDPPPMNPPLGHRSSD-----HVGSVR----LLEALGAPSRRDDGEPEKGASILS-----SSG
O116g00570    MGVAASARASRPSDFVDADDDAWRRADAGARARRGRRETSAGS----FSAIPSVSSFGGLMNYERQGARVTRRSPSNEGGDVEQGLADWGRDDASI
Ot17g00460    MDDAASDGRGR-----TDDARWTSASAGARDRARGRETRSGADGQRFETIPSVSSFGGLMNYEVFSRRSANEMSSEP-DVERG----GVEPSPV
Bathy06g01490  MGLFDGFSSSSSKPSSS-----LNASSSSKNDAFPMNSILSDD---VETSLVVNASNNN--KHNNRGEEEETITMMDNNNNLSVN-----EMPM
                :                . . .                * :                : .                .                .

```


Figure S9. Alignment of all intron encoded proteins (IEP) within IE-B sequences. The divergent group of IEPs (blue) and the ‘OTU-like cysteine protease domain’ (yellow) have been highlighted.



```

710      720      730      740      750      760      770      780      790      800
Mipuc07g03770 -----PRLQPSALVVTAPAPMTSAT----PRP---SVARTLPSAAR----KAVTRGAHSVQ
Mipuc09g04130 -----PRLQPSALVVTAPAPMTSAT----PRP---SVARTLPSAAR----KAVTRGAHSVQ
Mipuc01g03710 -----PRLQPSALVVTAPAPMTSAT----PRP---CAASTPSAAR----KPATRGAHSVQ
Mipuc03g05030 -----PLQSVSSRTRSRGRDRGGDGFANPRLQSSALVATSAAPTPSAPRT-KPKSKAKPAARTKSKPKYKGAVTRGAHSVQ
Mipuc07g01720 ATPAATPAATPAATPVATPLQSVSSRTRSRGRDRGGDGFANPRLQSSALVATSAAPTPSAPRT-KPKSKAKPAARTKSKPKYKGAVTRGAHSVQ
Mipuc03g01960 ATPAATPAATPAATPAATPLQSVSSRTRSRGRDRGGDGFANPRLQSSALVATSAAPTPSAPRT-KPKSKAKPAARTKSKPKYKGAHTRGAHVVQ
Mipuc01g06560 -----G---LAASAPRGRATPAPRGRATPARRVPDASAPRGRAAPAPESSESVED----EEAA-RFNAGI
Mipuc03g01120 -----G---LAASAPRGRATPAPRGRATPARRVPDASAPRGRAAPAPESSESVED----EEAA-RFNAGI
          *   . . . . . *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   ATPAATPAATPAATP2ATPLQSVSSRTRSRGRDRGGDGFANPRLQ2SALV2T2AAP22SAPRTSKPR2KAKPAART2S222PYKGAHVTRGAHSVQ

810      820      830      840      850      860      870      880      890      900
Mipuc07g03770 DMHSVETRNYDAGKVARQQRKSPDERRDDAAAEARANAARDAAEERAFKLYETIDELR-AAFEADPLSDDECVRDIKAEHDMVAAAKKRAADEAK
Mipuc09g04130 DMHSVETRNYEAGKVARQQRKSPDERRDDAAAEARANAARDAAEERAFKLYETIDELRAAAVKAVQMSDDECVRDIKAEHDMVAAAKKRAADEAK
Mipuc01g03710 DMHSVETRNYEAGKVARQQRKSPDERRDDAAAEARANAARDAAEERAFKLYETIDELR-AAVKADLMSDDECVRDIKAEHDMVAAAKKRAADEAK
Mipuc03g05030 DMHSVETRNYEAGKVARQQRKSPDERRDDAAAEARANAARDAAEERAFKLYETIDELR-AAFEADPLSDDECVRDIKAEHDMVAAAKKRAADEAK
Mipuc07g01720 DMHSVETRNYEAGKVARQQRKSPDERSRDDAAAEARANAARDAAEERAFKLYETIDELR-AAFEADPLSDDECVRDIKAEHDMVAAAKKRAADEAK
Mipuc03g01960 DMHSVETRNYEAGKVARQQRKSPDERRDDAAAEARANAARDAAEERAFKLYETIDELR-AAFEADPLSDDECVRDIKAEHDMVAAAKKRAADEAK
Mipuc01g06560 ICRAVERRNHDDGQVRKEQRTFVDVKRRAEEALALAAVDEAARVVKENADFEAVRKATLEER-AEFDLNPPEAAYLKRVKQEHKKECDARVRAEAGG-
Mipuc03g01120 ICRAVERRNHDDGQVRKEQRTFVDVKRRAEEALALAAVDEAARVVKENADFEAVRKATLEER-AEFDLNPPEAAYLKRVKQEHKKECDARVRAEAGG-
          :   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   DMHSVETRNYEAGKVARQQRKSPDERRDDAAAEARANAARDAAEERAFKLYETIDELRAAAF2ADP2S2D2CV2DIKAEHDMVAAAKKRAADEAK

910      920      930      940      950      960      970      980      990      1000
Mipuc07g03770 DEASRSVPKPAELKRR-ILTPPEARVVDKIDAFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQY-EIDTMDLGHWFAMCSDASPDFAK-SGR
Mipuc09g04130 DEASRSVPKPAELKRR-ILTPPEARVVDKIDAFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQY-EIDTMDLGHWFAMCSDAYFEAK-SGR
Mipuc01g03710 DRASRSVPKPAELKRR-TLKPALDRVNNKIDKFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQY-DITMDLGHWFAMCSDAYFEAN-TGR
Mipuc03g05030 DEASRSVPKPAELKRR-LKPAERARVDKIDAFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQY-KITDLDLGHWFAMCSDAYFEAN-TGR
Mipuc07g01720 DEASRSVPKPAELKRR-ILTPPEARVVDKIDAFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQY-EIDTMDLGHWFAMCSDAYFETN-TGR
Mipuc03g01960 DEASRSVPKPAELKRR-ILKLAERARVVDKIDKFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQY-KITDMDLGHWFAMCSDAYFEAKS2GR
Mipuc01g06560 ---AVTMPREALQILFHLIKGHLERLQVLAIPWRIAMVSLYKRPRLPCFPATQALALVIPAANLRLDPWIIDLADIGYPTGAMP5DAGMDCS--GN
Mipuc03g01120 ---AVTMPREALQILFHLIKGHLERLQVLAIPWRIAMVSLYKRPRLPCFPATQALALVIPAANLRLDPWIIDLADIGYPTGAMP5DAGMDCS--GN
          *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   DEASRSVPKPAELKRRF1L2PAERARVDKIDAFWRGVAKAYLEHWDPLPWFVSAGALAQNETLELIVEQYPEIDT2DELGHWFAMCSDAYF2A2S2GR

1010     1020     1030     1040     1050     1060     1070     1080     1090     1100
Mipuc07g03770 FTLQTTNRDYADLQRANLHNQWKTVPFRP---GAEMESDFIEKFEATCSNGAKSWKYAITTTPVGDK-----KLQEKFYTNRRGKVVILP
Mipuc09g04130 FTLQTTNRDYADLQRANLHNQWKTVPFRP---GAEMESDFIEKFEATCSNGAKSWKYAITTTPVGDK-----KLQEKFYTNRRGKVVILP
Mipuc01g03710 FALQTTNRDYADLQRANLHNQWKTVPFRP---GAEMESDFIEKFEATCSNGAKSWKYITAAITTPVGDK-----KLQEKFYTNRRGKVVILP
Mipuc03g05030 FALQTTNRDYADLQRANLHNQWKTVPFRP---GAEMESDFIEKFEATCSNGAKSWKYAITTTPVGDK-----KLQEKFYTNRRGKVVILP
Mipuc07g01720 FALQTTNRDYADLQRANLHNQWKTVPST---GKMESEFIERFEATCSNGAKSWKYAITTTPVGDK-----KLQEKFYTNRRGKVVILP
Mipuc03g01960 FALQTTNRDYADLQRANLHNQWKTVPFRP---GAEMESDFIEKFEATCSNGAKSWHYAIRTPGNGK-----KLQEKFYTNRRGKVVILP
Mipuc01g06560 FVLASTNREFPWRLRANLPLGLNWDGDDPDRVSGRWDAEGVYVHTTKCSNGTPEKHRVVSPPQSWV-----HWWDWYPMQYHQMGSKRLPR
Mipuc03g01120 FVLASTNREFPWRLRANLPLGLNWDGDDPDRVSGRWDAEGVYVHTTKCSNGTPEKHRVVSPPQSWVGGATKYSSRGSVHWWDWYPMQYHQMGSKRLPR
          *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   FALQTTNRDYADLQRANLHNQWKTVPFRPDRVGAEMESDFIEKFEATCSNGAKSWKYAITTTPVGDKGATKYSSRGSVHWVWKLQEKFYTNR2GKVVILP

1110     1120     1130     1140     1150     1160     1170     1180     1190     1200
Mipuc07g03770 KDLED---DPLFPLALATATPAMWDATVAGDGLSREGIFSLAFCDFTPRELDWFINILARPRWRDG-TGAVVHGAFVVRKTKRKGQYT-LDMTEKSLKD
Mipuc09g04130 KDLED---DPLFPLALATATPAMWDATVAGDGLSREGIFSLAFCDFTPRELDWFINILARPRWRDG-TDVTVHGAFVVRKTKRKGQYT-LDMTEKSLKD
Mipuc01g03710 KDLED---DPLFPLALATATPAMWDATVAGDGLSREGIFSLAFCDFTPRELDWFINILARPRWRDG-TGAVVHGAFVVRKTKRKGQYT-LDMTEKSLKD
Mipuc03g05030 QNLED---DPLFPLALATATPAMWDATVAGDGLSREGIFSLAFCDFTPRELDWFINILARPRWRDG-TGAVVHGAFVVRKTKRKGQYT-LDMTEKSLKD
Mipuc07g01720 QNLEK---DPLFPLALATATPAMWDATVAGDGLSREGIFSLAFCDFTPRELDWFINILARPRWRDG-TDVTVHGAFVVRKTKRKGQYT-LDMTEKSLKD
Mipuc03g01960 KDLEK---DPLFPLALATATPAMWDATVAGDGLSREGIFSLAFCDFTPRELDWFINILARPRWRDG-TDVTVHGAFVVRKTKRKGQYT-LDMTEKSLKD
Mipuc01g06560 SYVDDPLDRDFMMSALATPAIDGLVVDGTTCCDTRVQHLLNDFTVREHEFISEWLGGRVFRDGTGEIIGHGFPVRR-RHTRNDEHL-LEMYSIDILEE
Mipuc03g01120 SYVDDPLDRDFMMSALATPAIDGLVVDGTTCCDTRVQHLLNDFTVREHEFISEWLGGRVFRDGTGEIIGHGFPVRR-KHTRKQDHW-LEMYSIDILEE
          .   .   .   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   KDLEDPLDDPLFPLALATATP2WMDA2VAGDGLS2EGIFSLAFCDFTPRELDWFINILARPRWRDGTG2VHGAFVVRKTKRKGQYTTLDMTEKSLKD

1210     1220     1230     1240     1250     1260     1270     1280     1290     1300
Mipuc07g03770 AEYAPRFGCGIRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDLCSAAGAWRSRGESFTQIIVDGVIRGFSLEELQAE LLRATGLNFTAADLFGIISK
Mipuc09g04130 AGYAPRFGSGTRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDLCSAAGAWRSRGESFTTHIIVDGVIRGFSLEELQAE LLRATGLNFTDENLFGIITM
Mipuc01g03710 ARYAAPRFGCGIRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDLYIAGAWRSRGESFTQIIVDGVIRGFSLEEMQAE LLRATGLNFTAADLFGIITK
Mipuc03g05030 AGYAPRFGCGIRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDLYSALAWRSRGESFTQIIVDGVIRGFSLEELQAE LLRATGLNFTAADLFGIITK
Mipuc07g01720 AEYAPRFGCGIRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDLCSAAGAWRSRGESFTTHIIVDGVIRGFSLEELQAE LLRATGLNFTAADLFGIISK
Mipuc03g01960 AGYAPRFGCGIRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDLYSALAWRSRGESFTQIIVDGVIRGFSLEELQAE LLRATGLNFTAADLFGIITK
Mipuc01g06560 NGYALHSSVTRTHRNRLVDLLTAPRALVLAHAKVACVEDIDVWKRADEWRASGASDYHIVRELYLIRVFSFEEIQEAVQKHTKMLFDRRLFY----
Mipuc03g01120 NGYALHSSVTRTHRNRLVDLLTAPRALVLAHAKVACVEDIDVWKRADEWRASGASDYHIVRELYLIRVFSFEEIQEAVQKHTKMLFARQTLFY----
          *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   AGYAPRFGCGIRTHGIRVKEFIYTAALPVTRAYPHKGVFPDNDL2SAGAWRSRGESFT2IIVDGVYIRGFSLEELQAE LLRATGLNFTAADLFGIITK

1310     1320     1330     1340     1350     1360     1370     1380     1390     1400
Mipuc07g03770 PLLPPGESSGRKKELGKHHGLGLPFPFRRTCDKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc09g04130 PLLPPGESSGRKKELGKHHGLGLPFPFRRTCDKTNADRVKHTIHDVNAANVIGIVRGGREHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc01g03710 PLLPPGESSGRKKELGKHHGLGLPFPFRRTCDKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc03g05030 PLLPPGESSGRKVVLRQYGLGLPFPFRRTCDKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc07g01720 PLLPPGESSGRKKELGKHHGLGLPFPFRRTCDKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc03g01960 PLLPPGESSGRKVLGEHGLGLPFPFRRTCDKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc01g06560 -----EKLGS---LGLVRAHRTCEKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
Mipuc03g01120 -----EKLGS---LGLVRAHRTCEKTNADRVKHTIHDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ
          *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   PLLPPGESSGRKKELGKHHGLGLPFPFRRTCDKTNADRVKHT2HDVNAANVIGIVRKGEHVAHLHGGLACGGEDSVKHFIKMRGAMYEEGYIFDAKQ

1410     1420     1430     1440     1450     1460     1470     1480     1490     1500
Mipuc07g03770 GLCIFPGGTIDRDTTKLLTQEVLDAAQRYHTGCVENKGFAPPELLPRRTDAQLKQVQFAPLLEKLEVRALIVKLGQWRSRLRRPASSNKARTYFIYTD
Mipuc09g04130 GLCIFPGGTIDRDTTKLLTQEVLDAAQRYHTGCVENKGFAPPELLPHRTDAQLKQVQFAPLLEKLEVRALIVKLGQWRMELRRTPASSNARSFYIYTD
Mipuc01g03710 GLCIFPGGTIDRDTTKLLTQEVLDAAQRYHTGCVENKGFAPPELLPHRTDAQLKQVQFAPLLEKLEVRALIVKLGQWRMELRRTPASSNARSFYIYTD
Mipuc03g05030 GLCIFPGGTIDRDTTKLLTQEVLDAAQRYHTGCVENKGFAPPELLPHRTDAQLKQVQFAPLLEKLRLSALIVKLGQWRMELRRTPASSNARSFYIYTD
Mipuc07g01720 GLCIFPGGTIDRDTTKLLTQEVLDAAQRYHTGCVENKGFAPPELLPRRTDAQLKQVQFAPLLEKLRLSALIVKLGQWRMELRRTPASSNARSFYIYTD
Mipuc03g01960 GLCIFPGGTIDRDTTKLLTQEVLDAAQRYHTGCVENKGFAPPELLPHRTDAQLKQVQFAPLLEKLRLSALIVKLGQWRMELRRTPASSNARSFYIYTD
Mipuc01g06560 YRAIVP5GMIIG--SITDWTBEEIRDQELFVVKCKTFENNMP--LLELPGAEEHREHFKNNWISAPEKSDACKVLSLPYIWRARDNTTTLGFCQCTPFG-
Mipuc03g01120 YRAIVP5GMIIG--SITDWTBEEIRDQELFVVKCKTFENNMP--LLELPGAEEHREHFKNNWISAPEKSDACKVLSLPYIWRARDNTNIBELGFCQCTPFG-
          *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *   *
Prim. cons.   GLCIFPGGTIDRDTTKLLTQEVLDAAQ2RYH2GCVENKGFAPPELLPHRTDAQLKQVQFAPLLEKLR22ALIVKLGQWRMELRRTPASSNARSFYIYTD

```

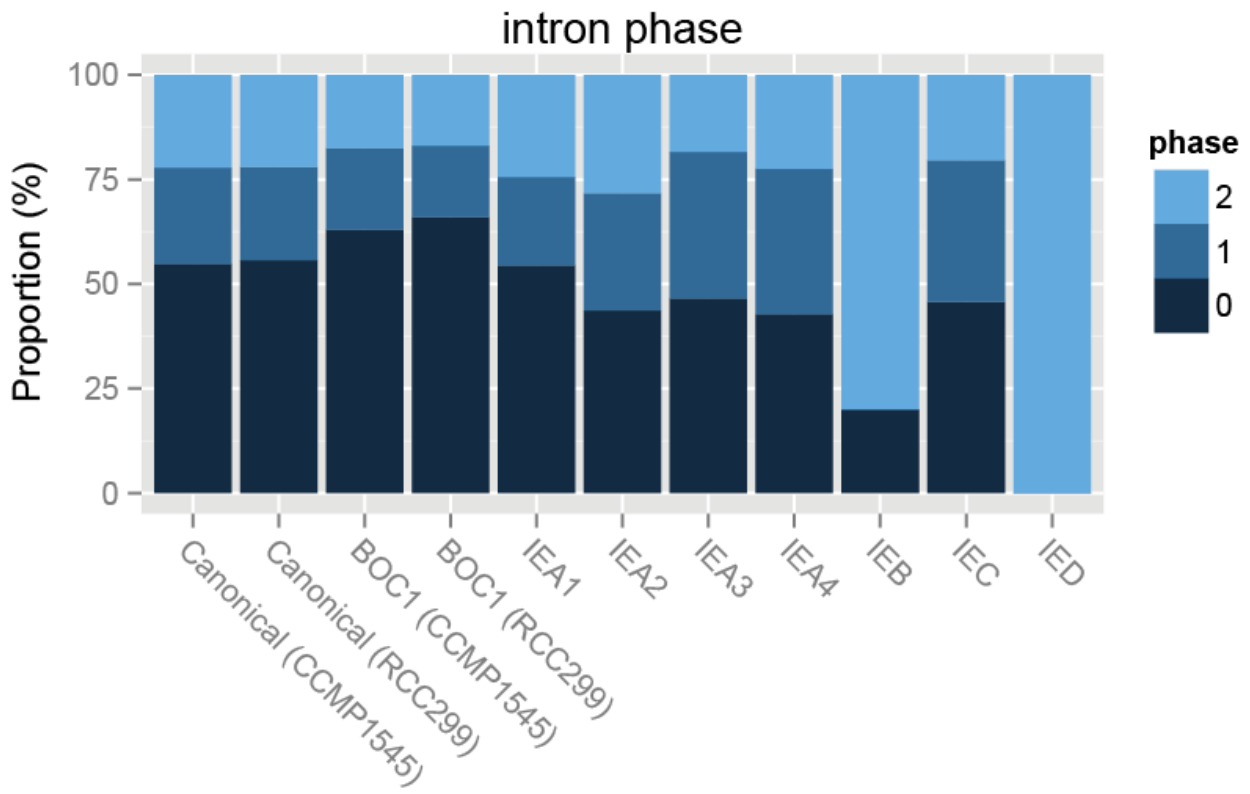



Figure S10. Phase distribution of Introner Elements, BOC1 and canonical introns.

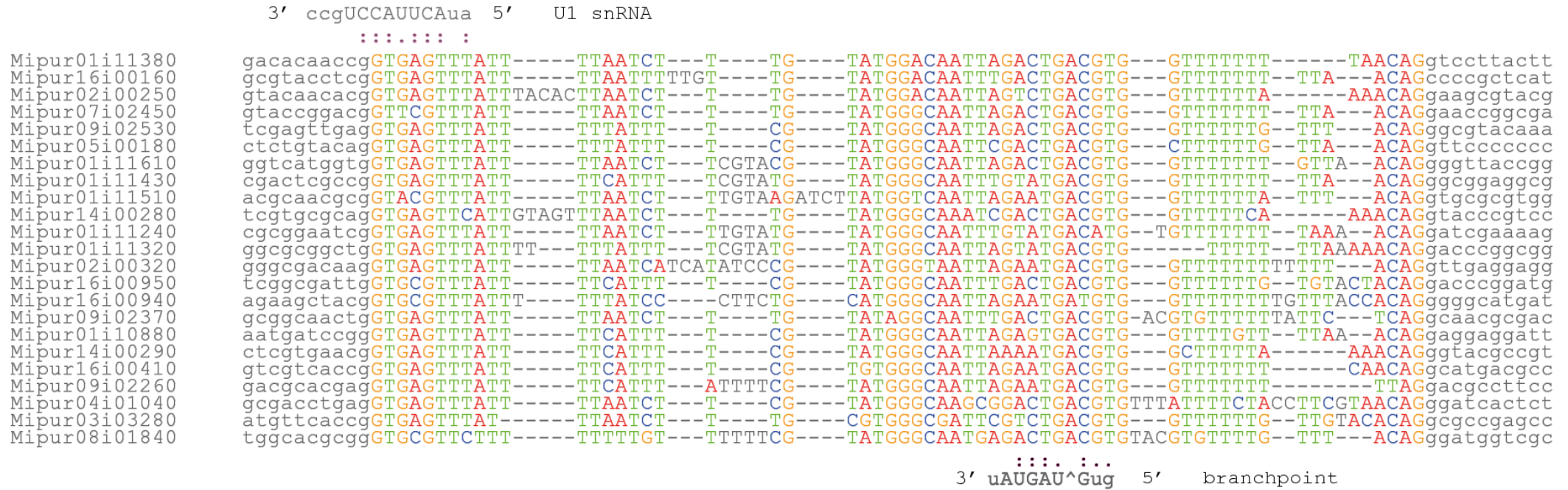


Figure S11. Alignment of typical IE-C sequences, flanked by their exonic regions (grey). Base-pairing information regarding the donor site (U1 snRNA) and the branch-point (U2 snRNA) is also provided.

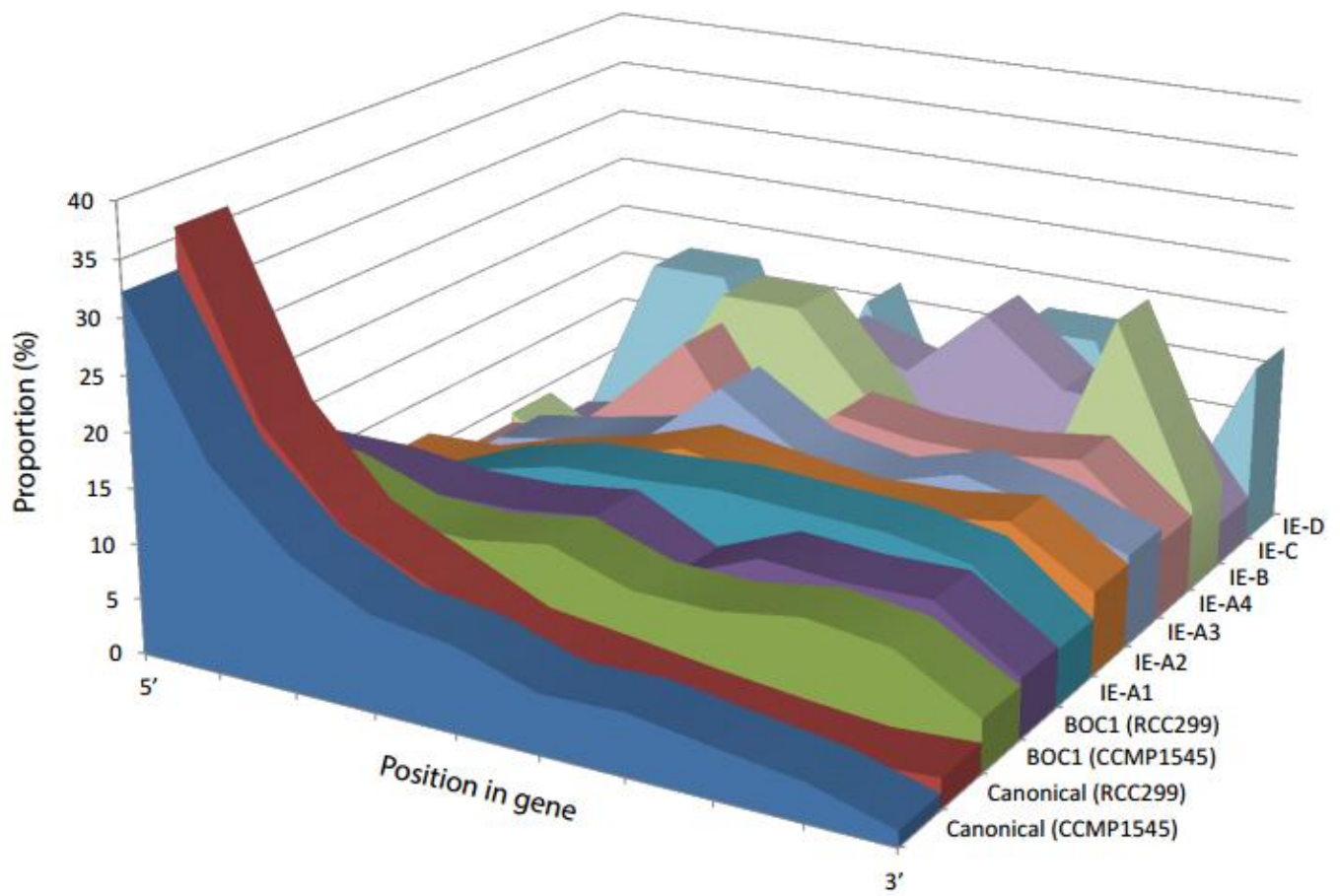


Figure S12. Positioning of Introner Elements, BOC1 and canonical introns inside genes.

Figure S13. Presence/absence polymorphisms in *Micromonas*: IE-A1 PAP (a), IE-A2 PAP (b), IE-A3 PAP (c), and two IE-C PAPs (d).

a) CCMP1545: scaffold_14 (384641..384976)
Metagenomic read: CAM_READ_0212050959 (CAMERA accession)

```

CCMP1545  GCGGCTGCTGCAGTACGCGTTCAGCGCGTTCAGCGCGTGGCGACCGGTTCTGCGCGTTCGCGAGCGAGGACGACAAGGCGACGCTGCTGG
Metagenome GCGGCTGCTGCAGTACGCGTTCAGCGCGTTCAGCGCGTGGCGACCGGTTCTGCGCGTGGCGAGCGAGGACGACAAGGCGACGCTGCTGG

ACGCATTTTTTCGAGCCGATGCGGCGGGCGGCGGCGAGAGCGCCGCCGCTCCGGCGTCGCG-----
|||||
ACGCA-TTTTTTCGAGCCGATGCGGCGGGCGGCGGCGAGAGCGCCGCCGCTCCGGAGGCGAGGTGCGGTTTTATCTCACACTGGTCCCATACGACGGCGTT

-----
IE-A1 PAP
CGCGTGGTGAACGCCGATCCTTAAGGACTTTGCCCGGGCGCATCTCTCCGCCACACCTCGCTTTCAATCCCGACCTCGACGCCTTCAACTCCATCTGA

-----
GACGAAGAGATATGTCCGACGCGCGTCCGACGAACTCGCGACGAAAAGCCCGCTCGGTGGGACG
|||||
CGCCTTCAACTCCACCCGACATCGCTTCGTAGGACGAAGAGATATGTCCGACGCGCGTCCGACGAACTCGCGACGAAAGAGCCGTCGTCGGTGGGACG

AGCACGAGGGGTTCCGGCCCGCAGCGTGGATCGAGGCCACGTCGCGAGCACGACGTCGTCGTAAGGGATTGATCGGGTTGGGTTGGAGGAGCCGTGGAG
|||||
AGCACGAGGGGTTCCGGCCCGCAGCGTGGATCGAGGCCACGTCGCGAGCACGACGTCGTCGTAAGGGATTGATCGGGTTGGGTTGGAGGAGCCGTGGAG

TTCGTGGGA
|||||
TTCGTGGGA

```

b) CCMP1545: scaffold_01 (251138..251395)
Metagenomic read: CAM_READ_0212289655 (CAMERA accession)

```

CCMP1545  CCGCGAGCGGCTGAGCGAGCTCGACCCGGACACGTTTCGAGTTCGAGGCGAGGAAAAATTCACG-----
Metagenome CCGCGAGCGGTTGGGCGAGCTCGACCCGGACACGTTTCGAGTTCGAGGCGAGGAAAAATTCACAGGTGCGGGTTCGACGCGAAACAAACGTGC

-----
IE-A2 PAP
CCGCACCGCTCGGATTCAACATTTGATCGCGTGGGTCCCTTTCAACTGACCGGTGTACTATCCTATTTTTGTACGGAACGCCCTCATCAGCGGCCCTCGG

GTTTTCCAGCGCCACGTTGCGATGGAACGCGCGACGAAGGACATGTCCGCGGGTGGCGCATGCGCGTCTCGCTCGCGAAGGCGCTGTTCCGCCGCGCG
|
G-TTCCAGCGCCACCGTTCGCGATGGAACGCGCGACGAAGGACATGTCCGCGGGTGGCGCATGCGCGTCTCGCTCGCGAAGGCGCTGTTCCGCCGCGCG

ACGCTGCTGCTTCTGAGCAGCCGACGAACCATCTCGACCTCGAGGCGTGGCTGTGGCTCGAAGAACATCTGTCGCGGTATAAAAA-
|||||
ACGCTGCTGCTTCTGAGCAGCCGACGAACCATCTCGACCTCGAGGCGTGGCTGTGGCTCGAAGAACATCTGTCGCGGTATAAAGAG

```

c) CCMP1545: scaffold_14 (643394..644184)
Metagenomic read: AACY020815663 (GenBank Accession)

```

CCMP1545  GTCGAGTCCGAACATGGGCGACCTCGCGGACGTCAGCGAGCTCCTCACGAGCAAAGTTACGGCAGCGACCACAGCGACAGCGAGGGC
Metagenome GTCGAGTCCGAACATGGGCGACCTCGCGGACGTCAGCGAGCTCCTCACGAGCAAAGTTACGGCAGCGACCACAGCGACAGCGAGGGC

GAACGCGCGCGTGGACTTGACGCAAGATTACAACCGTTCGCGCGCAAGGAACGCGGTCGCGGATCATCCTGCAGGAGATTGGCCCGCGGATGGAGCTC
|||||
GAACGCGCGCGTGGACTTGACGCAAGATTACAACCGTTCGCGCGCAAGGAACGCGGTCGCGGATCATCCTGCAGGAGATTGGCCCGCGGATGGAGCTC

GAGCTCGTGAAGGTGAGGAGGGGATGTGCGAGGGCGGGTGTGTATCACGCGTACGTGAAGAAGACGAGGAGGAGGTGTTGGAG-----
|||||
GAGCTCGTGAAGGTGAGGAGGGGATGTGCGAGGGCGGGTGTGTATCACGCGTACGTGAAGAAGACGAGGAGGAGGTGTTGGAGGTGAGACTGCTTC

-----
IE-A3 PAP
CCATACGACCCCGTTCCGCTGGTGCACGCCGTTCTTGAGGACTTTCCCGTCGTCACCTTACCACCGCTTCCCTTTCAACGTTTGACCGGTAAGACG

-----
CTGGAGAATAAGGCGGTGACGAAGGAGGCGCTGCGGAAGCGACGGAGGAGGAACAGGAGCGAACGCTGC
|||||
TTCGACTGACCGATCGCTTCAACCCACGCGAGTTGGAGAATAAGGCGGTGACGAAGGAGGCGCTGCGGAAGCGACGGAGGAGGAACAGGAGCGAACGCTGC

GGAGGAAGGAGAGAGGCGAAGGAGAAGGCGCGGATGGAGGAAGAACTCGCGGGATGAAGAAGAGGAAGCGGATAATAAATCCGCCCGGATAAGCT
|||||
GGAGGAAGGAGAGAGGCGAAGGAGAAGGCGCGGATGGAGGAAGAACTCGCGGGATGAAGAAGAGGAAGCGGATAATAAATCCGCCCGGATAAGCT

```

```

CGACGAGAACGGACGGAAGAAAAAGAAGGGGGCGAGCGGCGCGGAGAAGAGACGCGACGGCGGGCAGCAGCAGGCGCCGTTTAAGCGCGGTGAAAG
|||||
CGACGAGAACGGACGGAAGAAAAAGAAGGGGGCGAGCGGCGCGGAGAAGAGACGCGACGGCGGGCAGCAGCAGGCGCCGTTTAAGCGCGGTGAAAG
|||||
AGTCAGTACAAGGACAGGGAAG-----GCGGCGGGCGGGCGGAAGCGGCGGGCGGGGAAGAAATCGCCGGCGGGGACGAAGGGCA
|||||
AGTCAGTACAAGGACAGGGAAGGCGGGCGGGCGGGCGGAAGCGGCGGGCGGGGAAGAAATCGCCGGCGGGGACGAAGGGCA
|||||
AGGACAAGAACAAGCGCCGGGAAGGGTGC GCGCGCGAAGAGGAAGTAGAGACACACGCTCGGCCACGCGGTACGCGCGACGCTCTCGGCTCTCGCTCGC
|||||
AGGACAAGAACAAGCGCCGGGAAGGGTGC GCGCGCGAAGAGGAAGTAGAGACACACGCTCGGCCACGCGGTACGCGCGACGCTCTCGGCTCTCGCTCGC
|||||
GTCTTCAAAAGTTGTGCGGATTGTACATCCGCATCGCCCCCGTCCGGAGCAA
|||||
GTCTTCAAAAGTTGTGCGGATTGTACATCCGCATCGCCCCCGTCCGGAGCAA

```

d) **RCC299: chrom_06 (160647..160930)**
Metagenomic read: CCMP1764_READ_00758094 (CAMERA accession)

```

RCC299 CTCATCGACGTTTCGGTGGGATCGGTGGCGAACCAGCAGCCGCGGGACGGTACCTGGTCAGCCCGCGCGTCGTGTGCACCTCGCCAACCTG
|||||
Metagenome CTCATCGACGTTTCGGTGGGATCGGTGGCGAACCAGCAGCCGCGGGACGGTACCTGGTCA-CCCACGCGTCGTGTGCACCTCGCCAACCTG
|||||
CCACCCGCGGGCTGAGTTCGGGCAGGTGATTCTTTGCATTTAATCTTTGCATCATGGGCAATCAGATTGACGTGGTTTTTTATTTCGCAGGCTGTGCGGC
|||||
TCACCCGCGGGCTGAGGTCGGCAA-----IE-C PAB-----ACTGTGCGGC
|||||
GCCCGAGCCGTCGCCGTCGTCGACGAGACG-----GCGTACGTCGTG
|||
GCTCGAGCCGTCGCCGTCGTCGACGAGACGTTGAGTTTATTGTTATTTTCGTATGGGCAATTAGTTTGACATCGTGTGTTTGTTCACAGGCGTACGTCGTG
|||||
TCGGAGATGTCAAACCTCCCTCTCCGTCGTCGTCCTCCCGGCG
|||||
TCCGAGATGTCAAACCTCCCTCTCCGTCGTCGTCGTCCTCCCGGCG

```

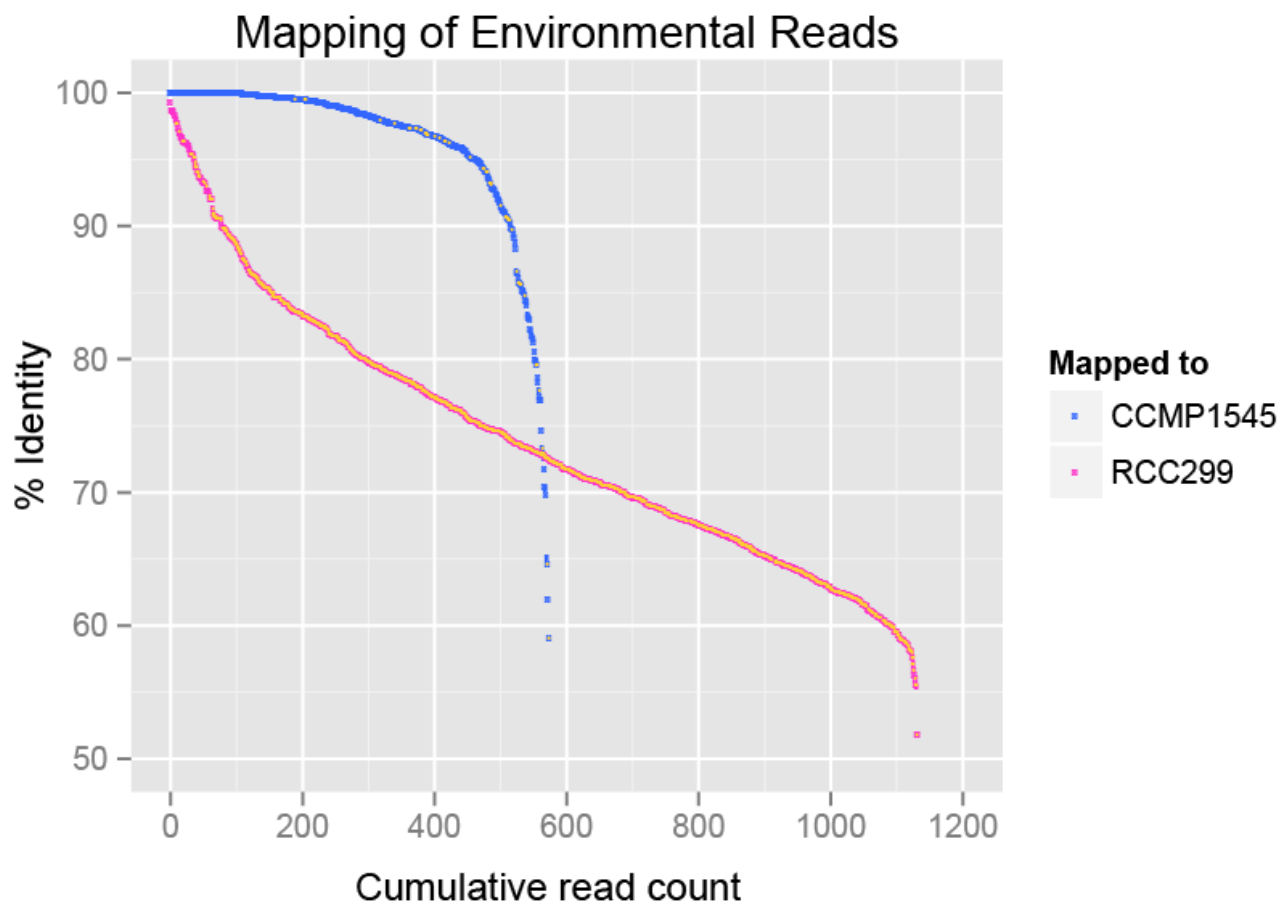


Figure S14. Mapping identities for environmental sequences. Yellow dots indicate alignments showing PAPA.

* or proof-reading ?

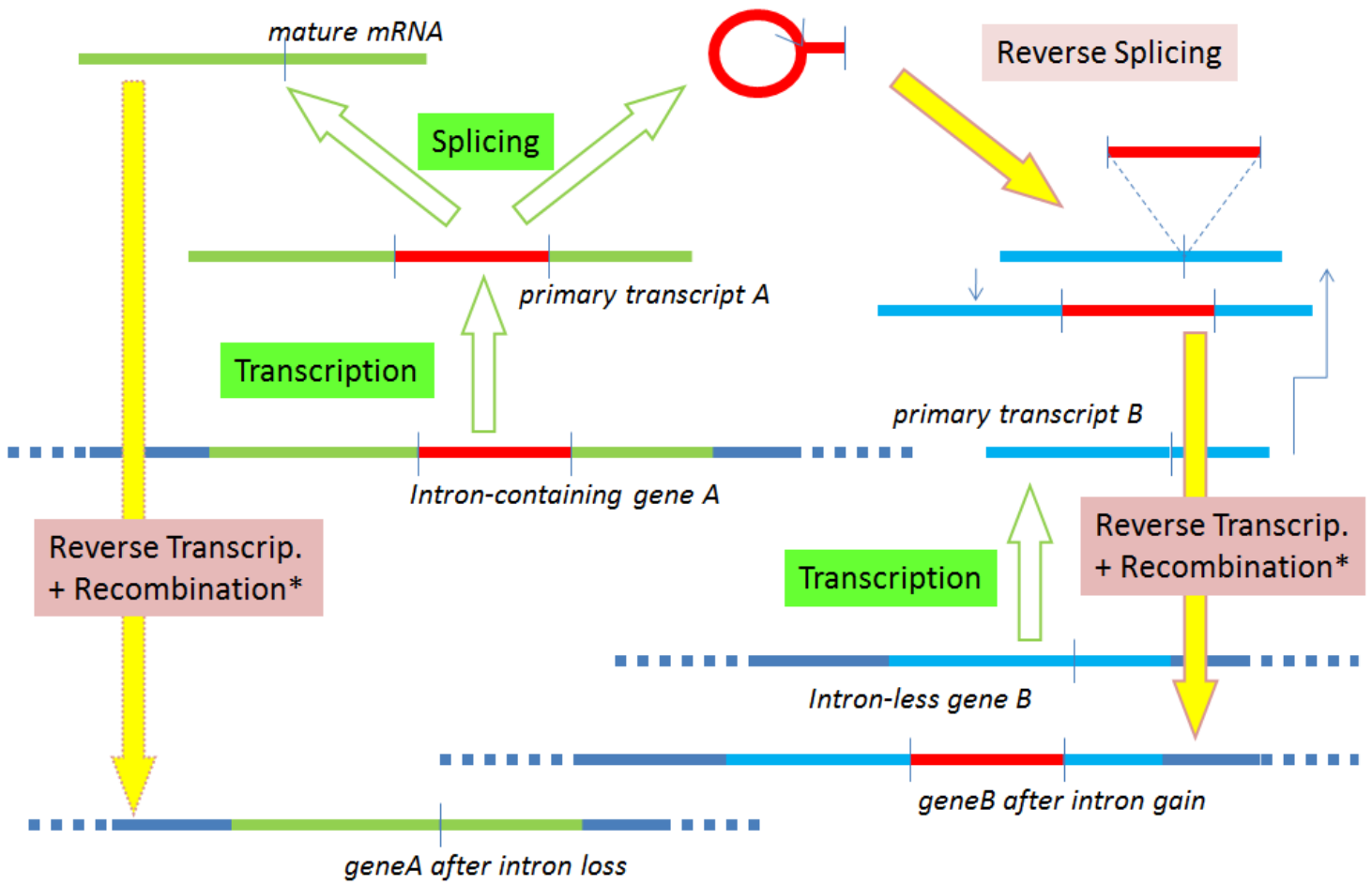


Figure S15. Introner Element replication mechanism.

Mipuc05i01435: IE-A1 + IE-A1

GCGCGTTCTGACTGGTCCCCATACGACCGCGTTTCGCGTGGTGAACGCCGATCCTTAAGGACTTTTGC
CGTCCTCTCTCCCCGCCACCCCTTCCTTTCAATCCCCGCCGCGACGCCTTTCAACTCCTTTCAAC
TCCTCACGCCGGTCCCCGTACGACCCCGTCGGCGCGGTGCACGCCGTTTCCTTAAGGACTTTCTCTCC
CGGCGTGCCTTCGTCTCTCTCCGCCAGGGTTCCTTCGGTTCCAATCCCGACGCGCCTCGACGCCTTT
CAACTCCAACAACGACGCCTTTCAACTTCACCCCGACGCCTTTCAACTTCACCCCGTCAG

Mipuc02i02490: IE-A1 + IE-A2

GCGCGTTCTTCCTCACCGCGCCGGTCCCCGCGCTGACCGGTCCCCGCGCTGACTGGTCCCCGTACGA
CCGCGTTCGCGTCGTGAAC TTCATTCTCAAGGACTTTTCGGCCGATCCTCGCGGACGGACTTCGCCC
GCCCGGCGCGTCTCTCCGCGCACCCCTCGCCGCTTTTCGATCCCCCGACGCGCCTACGCGACGCCTT
TCGACTCCATCTGACGCGCCCTTGAAC TCACCCCGACGCGCTTCGCTTCGCTTCGCTTCGCGCGGA
CCCTCCCTCGACCCTCAGGTGCGTTCTATACAAAAGTTTTTTCACCCACCGCTCGGTATCAACATTTG
ATCGCGTGGGTTCCTTTTCAACTGACCGGTGAAC TATTTTTGTATGGAACGACCCTCAG

Mipuc16i00710: IE-A1 + IE-A3

GCGCGTTCTGTCTCATCACACTGGTCCCCGTACGACCGCGTTGGCGCGGTGAACGCCGTTTCCTTAAG
GACTTTGCCCGCCC GCGTGCGTTTCTCTCCGCCACACCACGGTTTCAATCCCGACAACACACCGCG
ATGCCTTTCAACTTCAACCGACGCCTTTCAACTCCACCCCGACGTTTCATAACTACCCTAAACCCCGT
TCGCGTGGTGAACGCCGTTTCCTTAAGGACTTTTCCCGTCGTCACCTTCACCCGCGCTTCCTTTCA
ACGTTTGACCGGTAAGACGTTTCGACTGACCGATCGCTTCACCCACGCAGGCTCGCGTCGCTCTCCGC
GACGTCGAAG

Mipuc12i00190: IE-A? + IE-A1

GCGCGTTTGGGGCTGACTGGTCCCCATACGACCGCGTTTCGCGTGGTGAACGCCGATCCTTAAGGACT
TTACTTCCC GCGCCCTCCGCTCCGTCTATCATCACGCCGGTCCCCGTGCGACTCAGTCGGCGCG
GTGAACGCCGTTCTCGAGGACTTTGCCCCCGCTCGCGTTTCTCTCCGCCATCGCCCCTCGGTTT
TAATCCCGACGCACCGCGACGCCTCGCGACTTCATCTGACACCCCTGAACGCCACCCACCCCGACG
TCACCTCGTATGGAAAACGACCCTCAG

Figure S16. Examples of merged Introner Elements. Due to loss of internal splice structures when merging, it is hard to exactly delineate borders.

