

Operons Are a Conserved Feature of Nematode Genomes

Jonathan Pettitt,¹ Lucas Philippe, Debjani Sarkar, Christopher Johnston, Henrike Johanna Gothe, Diane Massie, Bernadette Connolly, and Berndt Müller

School of Medical Sciences, University of Aberdeen, Institute of Medical Sciences, Foresterhill, Aberdeen AB25 2ZD, United Kingdom

ORCID IDs: 0000-0001-7988-6833 (J.P.); 0000-0002-0672-4371 (C.J.); 0000-0003-1279-5421 (B.M.)

ABSTRACT The organization of genes into operons, clusters of genes that are co-transcribed to produce polycistronic pre-mRNAs, is a trait found in a wide range of eukaryotic groups, including multiple animal phyla. Operons are present in the class Chromadorea, one of the two main nematode classes, but their distribution in the other class, the Enoplea, is not known. We have surveyed the genomes of *Trichinella spiralis*, *Trichuris muris*, and *Romanomermis culicivorax* and identified the first putative operons in members of the Enoplea. Consistent with the mechanism of polycistronic RNA resolution in other nematodes, the mRNAs produced by genes downstream of the first gene in the *T. spiralis* and *T. muris* operons are *trans*-spliced to spliced leader RNAs, and we are able to detect polycistronic RNAs derived from these operons. Importantly, a putative intercistronic region from one of these potential enoplean operons confers polycistronic processing activity when expressed as part of a chimeric operon in *Caenorhabditis elegans*. We find that *T. spiralis* genes located in operons have an increased likelihood of having operonic *C. elegans* homologs. However, operon structure in terms of synteny and gene content is not tightly conserved between the two taxa, consistent with models of operon evolution. We have nevertheless identified putative operons conserved between Enoplea and Chromadorea. Our data suggest that operons and “spliced leader” (SL) *trans*-splicing predate the radiation of the nematode phylum, an inference which is supported by the phylogenetic profile of proteins known to be involved in nematode SL *trans*-splicing.

THE organization of open-reading frames into operons, such that multiple, distinct gene products are produced from a single, polycistronic transcript, is commonplace in prokaryote genomes (Jacob *et al.* 1960). Operons are also found in eukaryotes, although their distribution is sporadic and it does not seem likely that they represent an ancestral eukaryotic trait (Lawrence 1999; Hastings 2005). In prokaryotes, translation of multiple open reading frames in a polycistronic RNA occurs through multiple independent translation initiations. In eukaryotes, the polycistronic RNAs must first be processed into individual mRNAs before being translated. This

creates a problem in that the processed, downstream mRNAs would lack a cap structure necessary for RNA stability and translation. A number of eukaryotes are able to circumvent this problem through the *trans*-splicing of a short “spliced leader” (SL) RNA onto the 5′ end of the mRNA. Because the precursor SL RNAs that donate the SL are trimethylguanosine-capped, the *trans*-splicing event provides the cap structure for the mRNA. Thus by providing a mechanism that allows the formation of monocistronic, capped mRNAs from polycistronic RNA, SL *trans*-splicing enables the organization of eukaryotic genes into operons. It is striking that, at least to date, all eukaryotes in which operon usage is widespread also undergo SL *trans*-splicing (Johnson *et al.* 1987; Spieth *et al.* 1993; Davis and Hodgson 1997; Blumenthal *et al.* 2002; Ganot *et al.* 2004; Guiliano and Blaxter 2006; Satou *et al.* 2008; Marlétaz *et al.* 2008; Dana *et al.* 2012; Protasio *et al.* 2012; Tsai *et al.* 2013), suggesting that the resolution of polycistronic RNA is dependent upon SL *trans*-splicing.

Although operon organization is widespread in numerous eukaryotic taxa, the evolutionary mechanisms that have resulted in this form of gene organization are not well

Copyright © 2014 by the Genetics Society of America
doi: 10.1534/genetics.114.162875

Manuscript received February 12, 2014; accepted for publication June 6, 2014; published Early Online June 13, 2014.

Available freely online through the author-supported open access option.

Supporting information is available online at <http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.114.162875/-/DC1>.

The *T. spiralis* and *T. muris* mRNA 5′ ends with/without SL sequences are deposited in GenBank (accession nos. KF442418–KF442435, KF511776, KF511777, and KF768019).

¹Corresponding author: School of Medical Sciences, Institute of Medical Sciences, University of Aberdeen, Aberdeen AB25 2ZD, UK. E-mail: j.pettitt@abdn.ac.uk

understood. The most detailed analysis of the problem has come from studies in *Caenorhabditis elegans*, which led to the hypothesis that operon organization allows the marshaling of multiple genes under the control of a single promoter. This makes cells better able to cope with situations when transcription factors are present in limiting concentrations, such as recovery from growth arrest (Zaslaver *et al.* 2011). However, it is far from clear whether this is the only mechanism responsible for the evolution of operon organization, and the general applicability of this hypothesis to other members of the nematode phylum is not known.

To better understand the relationship between operon evolution and SL *trans*-splicing, it is necessary to determine the distribution of operon organization across the nematode phylum. Nematodes can be divided into two major classes: Enoplea and Chromadorea (Holterman *et al.* 2006; Meldal *et al.* 2007), with the latter class being much better characterized in terms of gene expression mechanisms, largely because it contains *C. elegans*. Both SL *trans*-splicing and operons have been identified in multiple nematodes within the Chromadorea (Evans *et al.* 1997; Lee and Sommer 2003; Guiliano and Blaxter 2006). However, the presence of operons has not been reported in nematodes from the other taxon.

We have previously identified SL *trans*-splicing in the enoplean nematodes *Trichinella spiralis* and *Prionchulus punctatus* (Pettitt *et al.* 2008; Harrison *et al.* 2010), suggesting that they may also possess operons. The draft genome of *T. spiralis* should be a useful resource for identifying operons in this nematode (Mitreva *et al.* 2011). However, identification of operons is not straightforward. The original discovery of operons in *C. elegans* was dependent upon the discovery of a specific spliced leader, SL2, which is *trans*-spliced to most mRNAs derived from genes downstream of the first gene in operons (Spieth *et al.* 1993), but not all nematodes use a specialized SL RNA to resolve polycistronic RNAs (Guiliano and Blaxter 2006). Thus, this feature cannot be considered diagnostic for mRNAs derived from nematode operons. The other feature common to operonic genes is that, at least in *C. elegans*, the distance between genes in an operon (the intergenic region, ICR) is unusually short, with a mean ICR size of 126 bp (Blumenthal *et al.* 2002). Again, this trait is not definitive: the ICR size can be considerably larger in the operons of other nematodes (Guiliano and Blaxter 2006; Ghedin *et al.* 2007) and even in *C. elegans* operons exist with large ICR distances (Morton and Blumenthal 2011).

Previous approaches to identify operons in *T. spiralis* (Mitreva *et al.* 2011) looked for pairs of *T. spiralis* genes whose homologs were in the same operon in *C. elegans*. This resulted in a limited set of 16 neighboring pairs of genes that potentially correspond to *T. spiralis* operons; however, further characterization of these candidate operons was not undertaken. We have used conserved synteny, coupled with the fact that mRNAs derived from downstream genes in operons are dependent on SL *trans*-splicing to elucidate a set of putative *T. spiralis* operons. Detailed analysis of two of these putative operons indicates that they display all the molecular charac-

teristics expected of loci that generate polycistronic RNA. Taken together our data indicate that the organization of genes into operons was present in the last common ancestor of the Chromadorea and Enoplea.

Materials and Methods

Bioinformatic identification of *T. muris* SL RNA genes

Trichuris muris SL RNA genes *Tmu*-SL1, *Tmu*-SL2, *Tmu*-SL3, and *Tmu*-SL9 were identified by searching the *T. muris* genome dataset with *T. spiralis* SL sequences using the BLASTN tool. Hits were considered if two of the three following criteria were met: a candidate Sm protein binding site was detected (AATTTTGTG), the 5' splice site sequence was conserved (AGGT), and a run of at least three Ts was found located ~100 bp from the end of the putative SL sequence. *T. muris* SL RNA genes *Tmu*-SL4, *Tmu*-SL5, *Tmu*-SL6, *Tmu*-SL7, *Tmu*-SL8, and *Tmu*-SL10 were identified by searching the *T. muris* genome dataset with the *Tmu*-SL1, *Tmu*-SL2, *Tmu*-SL3, or *Tmu*-SL9 sequences using the BLASTN tool and fulfilling the same criteria as above. Genes for *Tmu*-SL1, *Tmu*-SL2, *Tmu*-SL6, *Tmu*-SL8, *Tmu*-SL9, *Tmu*-SL10, and *Tmu*-SL11 were also identified with a PERL script (Pettitt *et al.* 2008) used previously to identify *T. spiralis* SL genes, except that the parameters for the Sm binding site were changed to AATTTTGTG/TG.

Analysis of *T. spiralis* SL containing ESTs and Identification of Putative Conserved Operons

The ESTs were identified earlier (Pettitt *et al.* 2008). To identify the corresponding gene from which each EST was derived, the EST sequences were mapped onto the *T. spiralis* draft genome sequencing using BLASTN. The corresponding gene was annotated as being in a putative operon if its upstream or downstream neighbor genes were on the same DNA strand with an intergenic distance of ≤ 1 kb. If the neighbors were on the same strand, but between 1 kb and 5 kb away, they were recorded as ambiguous. Otherwise the genes were annotated as nonoperonic. A minority of the ESTs matched to more than one predicted *T. spiralis* gene.

To identify the operonic status of the *C. elegans* homologs of each SL *trans*-spliced *T. spiralis* EST, BLASTX searches were carried out. We used an *E*-value cutoff of 10^{-5} to determine homology. In addition, if we obtained similar *E*-values for multiple *C. elegans* genes, we excluded that EST from the analysis.

T. spiralis and *Romanomermis culicivorax* homologs of operonic gene pairs conserved between *C. elegans* and *Brugia malayi* (Ghedin *et al.* 2007) were identified using BLASTP searches with the *C. elegans* upstream homologs from each pair as a query in searches against the respective gene predictions. The predicted coding region of the *T. spiralis*/*R. culicivorax* gene immediately downstream of the gene identified by this search was then used as a query sequence in a "reciprocal" BLASTP search against the *C. elegans* gene predictions. Since it was apparent that *T. spiralis* genes, which are separated by unusually short intergenic distance (such as might be expected in genes organized into operons), are prone to misannotation

and conflation into a single gene prediction, we also carried out an additional step using both *C. elegans* gene pairs as query sequences in BLASTP searches against the *T. spiralis* gene predictions, looking for cases where both *C. elegans* genes return matches to the same *T. spiralis* gene. Manual examination of such putative gene prediction errors, guided by the results of the sequence similarity searches, was then used to identify the intercistronic regions in each case. In all BLASTP searches, we used an *E*-value cutoff of 10^{-5} for the establishment of homology.

Phylogenetic profiling of SL trans-splicing snRNPs

Homologs of *C. elegans* *sna-1*, *sut-1*, and *sna-2* were identified by carrying out BLAST searches against the National Center for Biotechnology Information (NCBI) nonredundant database, except in the case of *R. culicivora* and *T. muris*, where BLAST searches were carried out against datasets downloaded from http://www.nematodes.org/genomes/romanomermis_culicivora/ and <http://www.sanger.ac.uk/resources/downloads/helminths/trichuris-muris.html>, respectively. Phylogenetic tree construction was carried out with the online implementation of PhyML (Dereeper *et al.* 2008) using default settings.

Nematode isolation and RNA preparation

T. spiralis RNA was produced as described (Pettitt *et al.* 2008). *T. muris* RNA was a generous gift from Allison Bancroft and Richard Grenic (University of Manchester).

Analysis of RNA 5' ends

The 5' ends of cDNAs were obtained through 5' RACE using the GeneRacer kit (Invitrogen), according to the manufacturer's instructions. Gene-specific primers used are given in Supporting Information, Table S3, and the cDNAs, amplified by PCR, using either GoTaq polymerase (Promega) or Expand High Fidelity polymerase (Roche), were cloned into pGEM T-Easy (Promega). The resulting plasmid inserts were sequenced by the University of Dundee Sequencing Service.

Detection of processing intermediates of polycistronic transcripts

RNA was reverse transcribed using SuperScript III Reverse Transcriptase (Invitrogen) and random primers according to the instructions of the manufacturer. In control reactions ("–RT") all reagents were included except the reverse transcriptase. Processing intermediates were normally amplified by two rounds of PCR with nested primer pairs (Table S4) and either GoTaq polymerase (Promega) or Expand High Fidelity polymerase (Roche) and visualized by agarose gel electrophoresis. The identity of the PCR products was determined by cloning into pGEM-T Easy (Promega) and sequencing of plasmid inserts.

Identification of SL RNA 3' ends

SL RNA 3' ends were determined essentially as described previously (Pettitt *et al.* 2008). *T. muris* total RNA (~5 μ g) was poly(A) tailed using yeast poly(A) polymerase, reverse tran-

scribed using an oligo-dT-anchor primer (GCGAGCTCCGCGG CCGCGTTTTTTTTTTTTTTT) and then PCR amplified using an SL-specific primers (GGTTAATTACCCAATTTAAAAG) and an anchor primer (GCGAGCTCCGCGGCCGCG). PCR fragments were inserted into pGEM-T Easy (Promega), and inserts were sequenced at the University of Dundee DNA Sequencing Facility.

SL RNA secondary structure prediction: Secondary structure prediction of *T. muris* SL RNA was performed using MFOLD Version 2.3 (Zuker 2003) using the default folding conditions (1 M NaCl, 37°) and with the constraint that the Sm-binding site (5'-AAUUUUUG-3') was required to be single stranded.

Generation and analysis of synthetic operon constructs

The GFP coding region was amplified from pTG96 using the primers 5'-CAATACAGACTTCCCAGGATTGGCCAAAGGACC CAAA-3' and 5'-GCTCACCATGCTAGCCTATTTGTATAGTTC ATCCATGC-3'. The mCherry coding region, coupled to the *unc-54* 3'-UTR, was amplified from pPD95.75Cherry (a derivative of pPD95.75 in which the GFP coding region was replaced by mCherry) using the primers 5'-ACAAATAGGCT AGCATGGTGAGCAAGGGCGAG-3' and 5'-CGCGCGAGACG AAAGGGCCCAGGAAACAGTTATGTTTGGTAT-3'. The primers were designed so that they had overlapping complementary 5' extensions that introduced an *NheI* restriction site. The two amplicons were purified and fused using a PCR fusion strategy (Hobert 2002). The resulting amplicon consisting of the GFP and mCherry coding regions flanking an *NheI* site was cloned into *SmaI*–*ApaI* cut pTG96 using In-Fusion HD (Clontech Laboratories) to generate pTG96-Op. The ICRs were cloned from PCR products amplified from genomic DNA. The *Tsp-cpt-2~nuaf-3* ICR was amplified from *T. spiralis* genomic DNA using primers 5'-ATACAAATAGGCTAG CACGAATTATCACTTTTATAAC-3' and 5'-TGCTCACCATGC TAGCTTACGCCAACTAGGAAATTATTGA-3', and the *Cel-cpt-2~prx-14* ICR was amplified from *C. elegans* genomic DNA using primers 5'-ATACAAATAGGCTAGCTTGTGTTGAT GACATTTATGTATTTAT-3' and 5'-TGCTCACCATGCTAGCTTT CAACCTGAAGCTTTAAAAT-3'. The resulting PCR products were cloned into *NheI* cut pTG96-Op using In-Fusion HD (Clontech Laboratories) and the resultant plasmids, pPE#LP1 (*Tsp-cpt-2~nuaf-3* ICR clone) and pPE#LP2 (*Cel-cpt-2~prx-14* ICR clone) were sequenced to confirm the integrity of the cloning process. To generate transgenic *C. elegans* strains, the plasmids were co-injected (100 ng/ μ l) with *P_{myo-2}::dTomato* (10 ng/ μ l) into Bristol (N2) wild-type hermaphrodites. For each construct, several lines were obtained, each of which gave identical expression patterns. Single lines for each construct were selected for the experiments reported here: PE612, *feEx304* [*sur-5::gfp::ICR^{Tsp-cpt-2~nuaf-3}::mCherry P_{myo-2}::dTomato*] and PE613, *feEx305* [*sur-5::gfp::ICR^{Cel-cpt-2~prx-14}::mCherry P_{myo-2}::dTomato*]. *Trans*-splicing of reporter gene transcripts was analyzed as described previously (Harrison *et al.* 2010). Briefly, total RNA was reverse transcribed and *trans*-spliced

transcripts were PCR amplified using *C. elegans* SL2-specific (5'-GGTTTAAACCCAGTTACTCAAG-3') and mCherry-specific (5'-CCGTCCTCGAAGTTCATCAC-3') primers. Primers derived from *gpd-1* (5'-CCAACTGTCTGGCACCAC-3' and 5'-GTCTTCTGGGTTGCGGTTAC-3') were used to normalize the reactions. cDNA fragments were cloned into pGEM-T Easy (Promega), and inserts were sequenced at the University of Dundee DNA Sequencing Facility.

Results

A putative enoplean operon

As part of the analysis of the transcriptome of the free-living enoplean, *Prionchulus punctatus*, we identified an EST corresponding to an SL *trans*-spliced mRNA. Sequence similarity searches using this sequence identified a single predicted *T. spiralis* gene, Tsp_06075. However, subsequent sequence analysis of Tsp_06075 showed that it corresponds to an erroneous gene prediction, which conflates three genes that are the orthologs of the *C. elegans* genes *zgpa-1* (C33H5.17), *dif-1*, and *aph-1*, respectively. That Tsp_06075 is actually three separate genes was confirmed by sequence analysis of 5' RACE products. It seems likely that the unusually short intergenic regions that exist between these three *T. spiralis* genes caused the gene annotation error (Figure 1). Such short intergenic distances are characteristic of nematode genes that are arranged into operons (Blumenthal *et al.* 2002; Guiliano and Blaxter 2006; Ghedin *et al.* 2007), and we thus decided to investigate the possibility that *Tsp-zgpa-1*, *Tsp-dif-1*, and *Tsp-aph-1* constitute an operon. In parallel, we also analyzed the homologs of these three genes in the closely related enoplean, *T. muris*, which show the same syntenic arrangement, although the intergenic distance between *Tmu-dif-1* and *Tmu-aph-1* is much larger than expected for an ICR. The three *C. elegans* homologs although organized into operons, are not found in the same operon. However, in a close relative of *C. elegans*, *Pristionchus pacificus*, *zgpa-1* and *dif-1* could potentially constitute a single operon, but again the intergenic space between the two genes is also relatively large compared to the average size of ICRs in *C. elegans* operons.

We determined the overall pattern of SL *trans*-splicing of the mRNAs derived from the putative operons in both *T. spiralis* and *T. muris* by analyzing the 5' ends of *zgpa-1*, *dif-1*, and *aph-1* mRNAs using 5' RACE (Figure 1A; Table S2). The analysis of *Tsp-zgpa-1* and *Tmu-zgpa-1* transcripts mapped the mRNA 5' ends to a region 200–250 bp upstream of the start codon, and we failed to detect any SL *trans*-spliced transcripts derived from this gene in either nematode. In contrast, all *dif-1* transcripts analyzed were subject to SL *trans*-splicing in both organisms (Table S2). *Tsp-dif-1* transcripts were *trans*-spliced to *Tsp*-SL10 [note this SL was previously given the designation TSL-10 (Pettitt *et al.* 2008), but we have renamed it to conform to accepted nematode gene nomenclature rules, which employ a species-specific prefix (Beech *et al.* 2010)] and *Tmu-dif-1* transcripts were *trans*-spliced to the newly identified *Tmu*-SL1, *Tmu*-SL4, and *Tmu*-SL12 (Table S2).

Analysis of *aph-1* transcripts showed that in some cases the transcripts are SL *trans*-spliced, but we were also able to detect transcripts that initiated ~200–300 bp upstream of the start codon, indicating that they were not subject to SL *trans*-splicing (Figure 1A; Table S2). It is notable that the distance between *Tmu-aph-1* and *Tmu-dif-1* is relatively large, suggesting the possibility that there are promoter elements immediately upstream of *Tmu-aph-1* that would allow the production of transcripts without the need for SL *trans*-splicing. Such “hybrid operons” have been described in *C. elegans* (Huang *et al.* 2007).

As part of this analysis, we identified spliced leaders in *T. muris*, leading to the discovery of 13 *Tmu*-SLs (Figure 2; Table S2). Previous studies have shown that the primary sequences of spliced leaders in *T. spiralis* are much more variable than those found in the Chromadorea (Pettitt *et al.* 2008), and many lack the conserved motifs that characterize spliced leaders from these latter nematodes. In contrast, those of *P. punctatus* do not show the same diversity and display a greater degree of sequence similarity to the Chromadorid spliced leaders (Harrison *et al.* 2010). Analysis of the 13 distinct *T. muris* spliced leaders, designated *Tmu*-SL1–13, support this view, since the *T. muris* spliced leaders possess the same 5' GGUWW and central CCC motifs that are highly conserved in the *P. punctatus* spliced leaders and Chromadorid SL1 and SL2 families, but missing in most of the *T. spiralis* spliced leaders. The presence of canonical nematode spliced leaders in *T. muris* and *P. punctatus*, despite the fact that the former nematode is more closely related to *T. spiralis*, supports the inference that the *T. spiralis* spliced leaders are derived features.

If *zgpa-1*, *dif-1*, and *aph-1* are components of a *bona fide* operon in the two enoplean nematodes, we would expect to be able to detect the polycistronic RNA from which their mRNAs are derived. Although not a definitive criterion, the presence of polycistronic, partially processed pre-mRNAs is a predicted property of operon usage. We tested for the presence of such RNA molecules in both *T. spiralis* and *T. muris* (Figure 1B; Figure S1) by reverse transcription of total RNA followed by PCR with gene-specific primers. PCR products were then analyzed by agarose gel electrophoresis. In *T. spiralis*, we detected RNA species connecting the open reading frames of *Tsp-zgpa-1* with *Tsp-dif-1* (*Tsp-zgpa-1~dif-1*) and *Tsp-dif-1* with *Tsp-aph-1* (*Tsp-dif-1~aph-1*) (Figure 1B). As we failed to amplify any products in control reactions performed in parallel with RNA subjected to mock reactions without reverse transcriptase (Figure S1), these products represent processing intermediates of polycistronic transcripts.

The *Tsp-zgpa-1~dif-1* intermediates contained the intercistronic region, and two of the intermediates lacked introns. The *Tsp-dif-1~aph-1* processing intermediates detected were all subject to *cis*-splicing of *dif-1* introns, but we failed to detect an intermediate containing the complete ICR. Instead, the ICR was removed by *cis*-splicing of a cryptic splice donor site located in exon 7 of the *dif-1* gene to the SL splice acceptor site of *aph-1* (Figure 1B). Such *cis*-splicing events have

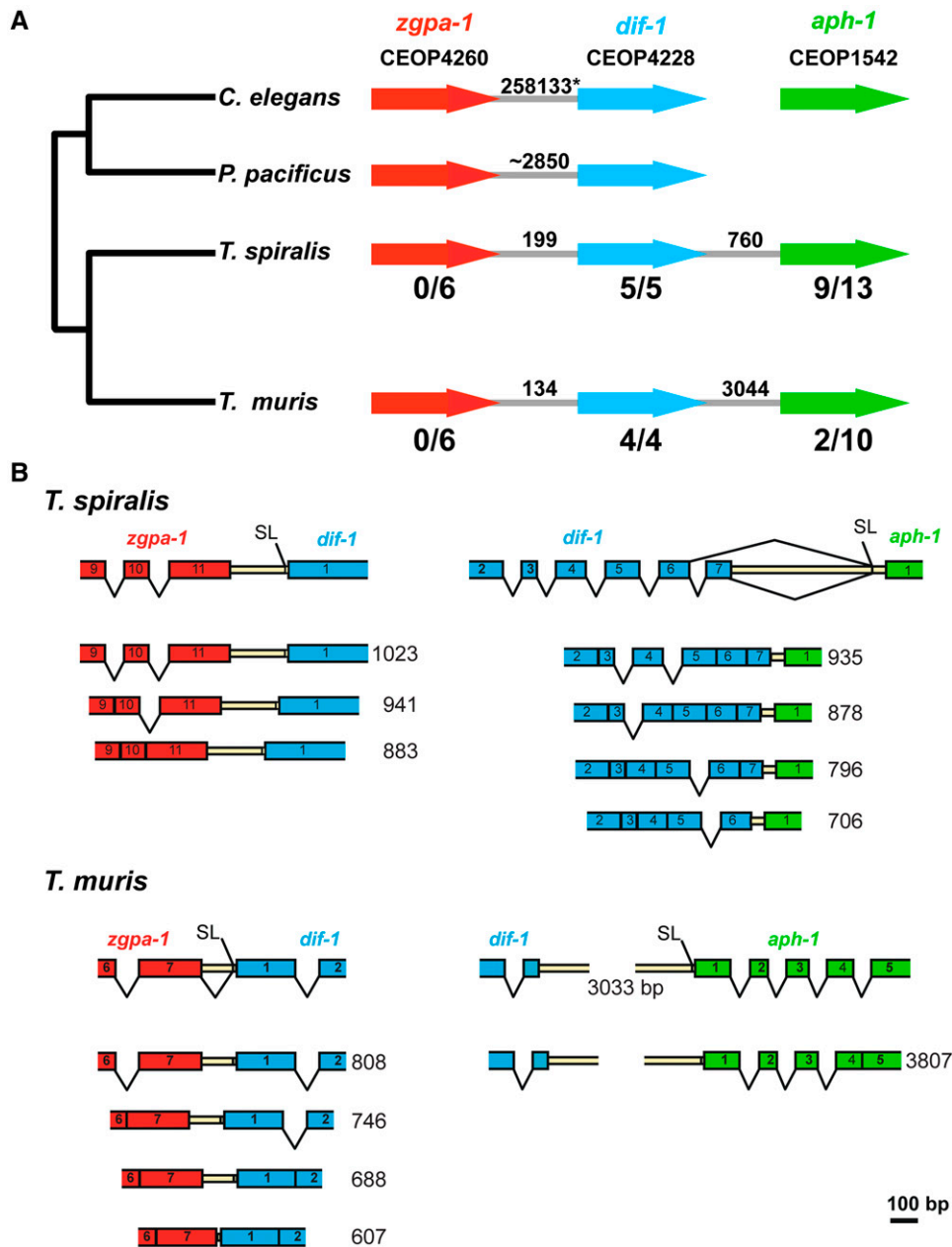


Figure 1 Evidence for the existence of an enoplean operon. (A) Schematic showing the genomic organization of *zgpa-1*, *dif-1*, and *aph-1* in selected nematodes mapped onto their phylogenetic relationships. Arrows represent genes, and the gray lines represent the intergenic regions (ICRs). Numbers above the ICRs represent the distances, in base pairs, between the stop and start codons of the upstream and downstream genes, respectively. The *C. elegans* operon numbers are given where appropriate. Fractions below the *T. spiralis* and *T. muris* genes represent the proportion of cDNAs derived from those genes that begins with a spliced leader sequence (see also Table S2). In *C. elegans*, the three genes are part of different operons. * indicates the distance between genes on chromosome IV. (B) Detecting polycistronic RNAs derived from the *zgpa-1~dif-1~aph-1* operon in enoplean nematodes. The exon-intron structures of the amplicons used to identify polycistronic RNAs are shown, with exons represented by boxes (shaded to identify the genes from which they are derived using the same color coding that was used in A). The intergenic regions are represented by cream-colored boxes. The positions of the SL *trans*-spliced 3' splice sites are indicated. The length of each cDNA is indicated.

also been detected in putative polycistronic RNAs discovered in tapeworm genomes (Tsai *et al.* 2013). Moreover, this demonstrates that *Tsp-dif-1* and *aph-1* are transcribed as a single transcript. In *T. muris* we also detected processing intermediates corresponding to *Tmu-zgpa-1~dif-1* and *Tmu-dif-1~aph-1* transcripts (Figure 1B; Table S2). The latter observation is significant, since the ICR between *Tmu-dif-1* and *Tmu-aph-1* is predicted to be 3033 nt long, a distance substantially longer than the length of an average ICR in *C. elegans*, although ICRs of similar length are also present in some *C. elegans* operons (Morton and Blumenthal 2011).

Identification of additional putative enoplean operons

To more systematically identify enoplean operons, we adopted two approaches. First, we used a set of EST sequences derived

from SL *trans*-spliced *T. spiralis* mRNAs (Pettitt *et al.* 2008) to identify their corresponding genes via sequence similarity searches (we also identified two *T. muris* mRNAs via the same approach: FF145866 and CB277782). For each gene, we then looked for neighboring genes predicted to be transcribed in the same orientation and that lay within 1 kb. Using this approach, we were able to identify multiple potential operons in the *T. spiralis* genome (Table 1; Table S1). We further analyzed this set of genes by identifying the *C. elegans* orthologs of each *T. spiralis* gene and determining whether these correspond to genes within operons. Our analysis revealed that at least 75 of the *C. elegans* orthologs are arranged in operons. This represents 44% of the *C. elegans* genes identified as orthologs of our *T. spiralis* SL *trans*-spliced EST set. Since only 15% of *C. elegans* genes are organized into operons

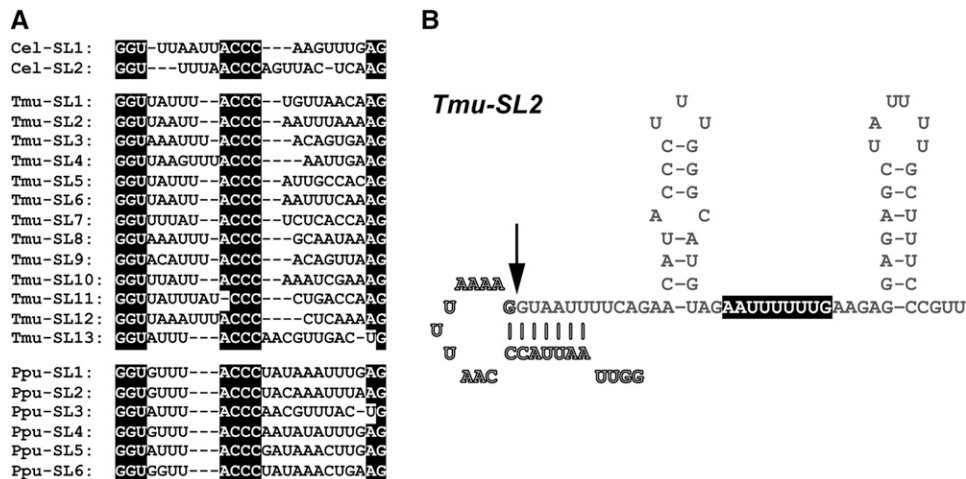


Figure 2 *T. muris* SL sequences and SL RNA structure. (A) *Tmu-SL1–13* genes were identified using a combination of cDNA sequencing and bioinformatics tools as described in *Materials and Methods*. *Tmu-SL12* was found by 5' RACE *trans*-spliced to *nuaf-3* mRNA, and *Tmu-SL13* was found *trans*-spliced to *aph-1* mRNA. In the alignment, only the SL sequences are shown. *T. muris* SL sequences were manually aligned and conserved groups are countershaded. *C. elegans* SL1 and SL2 and the previously identified *P. punctatus* SL sequences were included for comparison. (B) The intron of *Tmu-SL2* was experimentally identified and also found in the genome sequence. The proposed secondary structure was produced using M-fold (Zuker 2003). The SL sequence is shown in outline font and the putative Sm sequence motif is countershaded.

(Allen *et al.* 2011), we would expect only 15% of *T. spiralis* genes in our dataset to match *C. elegans* operonic homologs if they were selected at random. It is difficult to determine the reason for the increased likelihood of matches to *C. elegans* operonic genes among the *T. spiralis* SL *trans*-spliced ESTs; it may be that the corresponding *T. spiralis* gene set is biased for highly expressed genes, for instance. However, it is consistent with the possibility that this dataset is enriched for transcripts derived from operonic genes.

As an alternative approach, which would potentially identify operons that have been conserved since the separation of the Enoplea and Chromadorea, we looked for *T. spiralis* homologs of a set of putative operons conserved between *C. elegans* and *B. malayi* (Ghedini *et al.* 2007). Of the 107 operonic gene pairs screened, we identified 12 *T. spiralis* gene pairs that displayed conserved synteny, and whose component genes were separated by an average intergenic distance of 607 bp (Table 2). To determine whether any of these operons are also conserved in another enoplean species, we examined the organization of the corresponding homologous genes in the genome of *R. culicivora* (Schiffer *et al.* 2013). This analysis revealed that 4 of the 12 gene pairs were arranged in putative operons (assuming a maximum ICR distance of up to 1 kb) in this nematode (Table 2).

Taken together, our analysis indicates that there are multiple *T. spiralis* gene pairs whose genomic arrangement is consistent with their corresponding to operons. Moreover, it is possible to identify gene pairs conserved between *T. spiralis*, *R. culicivora*, *B. malayi*, and *C. elegans*, suggesting that these represent operons that were present in the last common ancestor of the three species.

Characterization of a conserved nematode operon

The analysis of one of the two SL *trans*-spliced *T. muris* mRNAs (GenBank accession no. FF145866) led to the identification of a putative operon conserved between multiple

nematodes (Figure 3A). The two genes contained in these putative operons have *C. elegans* homologs, *cpt-2* and *nuaf-3*, respectively, which are in the same operon (CEOP4424) (Figure 3A), although there is an additional gene, *prx-14*, located between these genes in CEOP4424 that is not present in the putative homologous *T. spiralis* operon. Examination of the genomic organization of the homologous genes in a selection of nematode species confirmed the evolutionary conservation of the synteny of the *cpt-2* and *nuaf-3* homologs (Figure 3A). This analysis also showed that insertion of *prx-14* into the *cpt-2~nuaf-3* operon was a relatively recent event, since it is present only in *C. elegans* and other closely related *Caenorhabditis* species. We also find, based on the head-to-tail organization and spacing between coding regions, that there is variation in the composition of both operons in different species, and these genes in *R. culicivora* are not in operons.

We further focused on the *cpt-2~nuaf-3* operon in *T. spiralis* and *T. muris* to determine the pattern of SL *trans*-splicing exhibited by the mRNAs derived from this operon and to verify that we were able to detect cDNAs consistent with the production of polycistronic RNAs.

Determination of the 5' ends of *cpt-2* and *nuaf-3* transcripts by 5' RACE revealed that *nuaf-3* mRNA is subject to

Table 1 Operonic status of genes that match ESTs derived from *T. spiralis* SL *trans*-spliced transcripts

Species	Location in operon			
	Upstream	Downstream	Ambiguous	Nonoperonic
<i>T. spiralis</i>	30	35	40	54
<i>C. elegans</i>	29	46	5 ^a	89

The status of each *T. spiralis* gene was determined using criteria given in *Materials and Methods*. The status of *C. elegans* genes was obtained from WormBase (Release WS237). Eleven EST matches were absent from the *T. spiralis* data relative to the *C. elegans* data as the corresponding *T. spiralis* gene could not be identified (see *Materials and Methods*). Full details of the individual EST sequence matches are given in Table S1.

^a Genes not annotated as operons, but having intergenic spacing with respect to their neighbors, that suggests they may be organized in an operon.

Table 2 Putative conserved nematode operons

<i>T. spiralis</i>		<i>R. culicivora</i>		<i>C. elegans</i>			<i>B. malayi</i>	
Upstream gene	Downstream gene	Upstream gene	Downstream gene	Upstream gene	Downstream gene	Operon no.	Upstream gene	Downstream gene
Tsp_00685	^a	Nonoperonic		<i>mrps-17</i>	<i>C05D11.9</i>	3372	Bm1_13520	Bm1_13525
Tsp_03140	Tsp_03139	t32947	t32944	<i>T26E3.4</i>	<i>par-6</i>	1672	Bm1_48785	Bm1_48780
Tsp_05540	Tsp_05541	Nonoperonic		<i>K11B4.1</i>	<i>K11B4.2</i>	1764	Bm1_55805	Bm1_55810
Tsp_06077	Tsp_06076	Nonoperonic		<i>Y62E10A.2</i>	<i>Y62E10A.6</i>	4540	Bm1_54855	Bm1_54850
Tsp_06996	^a	t05598/9	t05596	<i>sel-1</i>	<i>mrps-5</i>	5365	Bm1_45745	Bm1_45750
Tsp_09103	Tsp_09102	t35569	t35568	<i>snu-23</i>	<i>ZK686.3</i>	3452	Bm1_13735	Bm1_13740
Tsp_09506	^a	Nonoperonic		<i>H20J04.6</i>	<i>mog-2</i>	2124	Bm1_15855	Bm1_15860
Tsp_09539	^a	Nonoperonic		<i>E02H1.5</i>	<i>E02H1.6</i>	2436	Bm1_24720	Bm1_24715
Tsp_10673	Tsp_10674	t34344.1	t34344.2	<i>B0491.1</i>	<i>B0491.7</i>	2532	Bm1_10780	Bm1_10785
Tsp_10698	Tsp_10702	Nonoperonic		<i>trpp-8</i>	<i>vha-10</i>	1264	Bm1_12140	Bm1_12135
Tsp_10959	^a	Nonoperonic		<i>ubxn-2</i>	<i>Y94H6A.8</i>	4665	Bm1_36515	Bm1_36520
Tsp_11898	^a	Nonoperonic		<i>lst-6</i>	<i>sqv-7</i>	2276	Bm1_24075	Bm1_24080

^a Single gene annotation matches to both *C. elegans* genes in the operon pair consistent with annotation error caused by short intergenic spacing.

SL *trans*-splicing (Table S2) in both nematodes. However, we failed to detect any SL *trans*-spliced *cpt-2* transcripts, similar to the situation with *zgpa-1* transcripts in the *zgpa-1~dif-1~aph-1* operon.

We were also able to detect processing intermediates derived from the putative polycistronic *T. spiralis* *cpt-2~nuaf-3* transcripts. As for the *zgpa-1~dif-1~aph-1* operon, in addition to unprocessed, polycistronic transcripts, we detected *cis*-spliced intermediates lacking *cpt-2* and *nuaf-3* introns and a transcript in which the ICR was removed by splicing from a cryptic 3' splice site in *cpt-2* to the *nuaf-3* SL splice acceptor site (Figure 3B).

The *T. spiralis* *cpt-2~nuaf-3* ICR can mediate polycistronic RNA processing in *C. elegans*

Analysis of the intercistronic region between *Tsp-cpt-2* and *nuaf-3* downstream of the polyadenylation signal of *Tsp-cpt-2* revealed a clear Ur element and there are several U-rich regions, characteristics of the ICRs in *C. elegans* operons (Graber *et al.* 2007). To investigate the possibility that this region is able to function in polycistronic RNA processing, we determined whether the ICR from it could be recognized and processed if heterologously expressed in *C. elegans*. We generated an artificial operon consisting of *sur-5::gfp* (Gu *et al.* 1998) and mCherry genes flanking the ICR from *Tsp-cpt-2~nuaf-3*. Transgenic animals carrying this construct expressed nuclear GFP and cytoplasmic mCherry, consistent with the processing of the two coding regions under the direction of the *Tsp-cpt-2~nuaf-3* ICR. We confirmed that this involved *trans*-splicing to SL2, as expected for polycistronic RNA processing in *C. elegans*, by showing that we could detect SL2 *trans*-splicing to the mCherry mRNA in RNA derived from transgenic animals (Figure 4). Thus, the predicted ICR between *Tsp-cpt-2* and *nuaf-3* is recognized and used as a substrate for polycistronic RNA processing in *C. elegans*.

Conservation of SL *trans*-splicing snRNPs between Enoplea and Chromadorea

Our data show that both SL *trans*-splicing and operons were likely present in the last common ancestor of the nematode

phylum. This suggests that the processing machinery necessary for the coordination of these processes was already in place prior to radiation of the nematode phylum. To address this, we sought to determine the conservation of known protein components that are specifically involved in SL *trans*-splicing. Previous studies have shown the existence of two interacting proteins, conserved between *Ascaris suum* and *C. elegans*, that are components of the SL small nuclear ribonucleic particle (snRNP) (Denker *et al.* 2002; MacMorris *et al.* 2007). The two proteins, termed *SNA-1* and *SNA-2* in *C. elegans*, form a complex with SL RNA. In addition, a paralog of *SNA-1*, *SUT-1*, forms novel snRNPs containing *SNA-2* and a family of nematode-specific RNAs, designated Sm Y (MacMorris *et al.* 2007; Jones *et al.* 2009). The function of Sm Y RNAs is not known, but they are associated with SL *trans*-splicing (Maroney *et al.* 1996), and a role in the recycling of Sm proteins following SL *trans*-splicing has been proposed (MacMorris *et al.* 2007).

We were able to identify credible homologs of *sna-2* in the genomes of all nematodes in which searches were carried out, including *T. spiralis* and *R. culicivora*, indicating that this gene encodes a nematode-wide SL *trans*-splicing component (Figure S2). Searches of the same datasets identified clear *sna-1* and *sut-1* orthologs in the chromadorean nematodes (Figure 5), extending the previously reported phylogenetic distribution of *sut-1* (MacMorris *et al.* 2007). In contrast, the only enoplean genome in which we were able to identify a *sut-1* homolog was *R. culicivora* (Figure 5). We could not detect *sna-1* homologs in any of the enoplean genomes that we assayed. These data suggest that the gene duplication event that gave rise to *sna-1* and *sut-1* occurred after the separation of the two major nematode taxa; however, we cannot rule out that the possibility that our failure to detect *sna-1* homologs is due to the incomplete state of the enoplean genome drafts.

Discussion

The incidence of operons as a means to coordinate gene expression has been investigated in only one of the two main nematode taxa, leaving unanswered the question about when

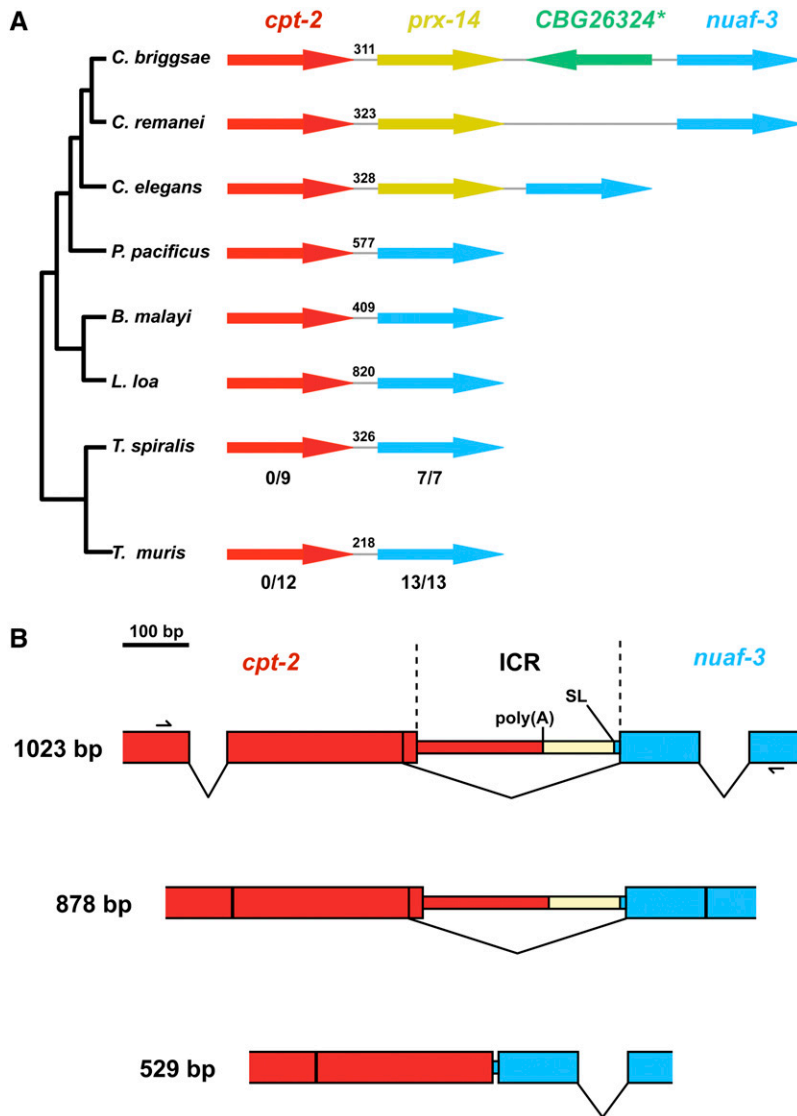


Figure 3 Evidence for an evolutionarily conserved nematode operon. (A) The structure of the *cpt-2~nuaf-3* genomic regions from a range of nematode species mapped onto the nematode phylogeny. Genes are represented by arrows, and the gray lines represent the intergenic regions (ICRs). Numbers above the ICRs represent the distances, in base pairs, between the stop and start codons of the upstream and downstream genes, respectively. The *C. elegans* operon numbers are given where appropriate. Fractions below the *T. spiralis* and *T. muris* genes represent the proportion of cDNAs derived from those genes that begins with a spliced leader sequence. (B) Detecting polycistronic RNAs derived from the *cpt-2~nuaf-3* operon in *T. spiralis*. The exon–intron structures of the amplicons used to identify polycistronic RNAs are shown, with exons represented by boxes (shaded to identify the genes from which they are derived using the same color coding that was used in A). The region removed during operon processing is represented by cream-colored boxes. The positions of the SL *trans*-splice 3' spliced sites are indicated.

this mechanism first occurred during nematode evolution (Guiliano and Blaxter 2006; Ghedin *et al.* 2007; Liu *et al.* 2010). The work presented here provides strong support for the existence of operons in the Enoplea, and that operons were likely present in the ancestor of the Chromadorea and the Enoplea. The unequivocal identification of operons is not straightforward in those nematodes that do not utilize a specialized spliced leader to resolve polycistronic RNAs; the use of SL2 in *C. elegans* and its close relatives has greatly facilitated the identification of operons in these species. In contrast, in enoplean species we can only infer the presence of operons through multiple lines of evidence.

Our work demonstrates the existence of clusters of genes ordered in a head-to-tail arrangement with short, less than 1 kb, intergenic distances, consistent with their being organized into operons. The fact that we can identify homologous pairs of closely spaced genes conserved between enoplean and chromadorean nematodes suggests that at least in these cases there has been selective pressure to retain short intergenic

distances, consistent with what would be expected if they were part of operons. Analysis of the transcripts produced by these putative operons shows that, as expected, genes predicted to be downstream in operons produce mRNAs that are SL *trans*-spliced. Further supporting evidence comes from the fact that we can detect unprocessed, polycistronic pre-mRNAs derived from these putative operons. Thus these genes exhibit the molecular properties expected of operons. Most significantly, the intergenic region between *Tsp-cpt-2* and *nuaf-3* acts as a substrate for the polycistronic RNA processing machinery of *C. elegans*, providing the strongest evidence that it is part of an operon in *T. spiralis*.

The identification of SLs from *T. muris* has extended our understanding of nematode SL evolution. Previous studies have shown that *T. spiralis* and *T. pseudospiralis* have unusually diverse SLs that do not readily correspond to those SLs found within Chromadorea (Pettitt *et al.* 2008, 2010). Another enoplean, *P. punctatus*, possesses SLs that resemble the specialized SL2 associated with *trans*-splicing and polycistronic RNA

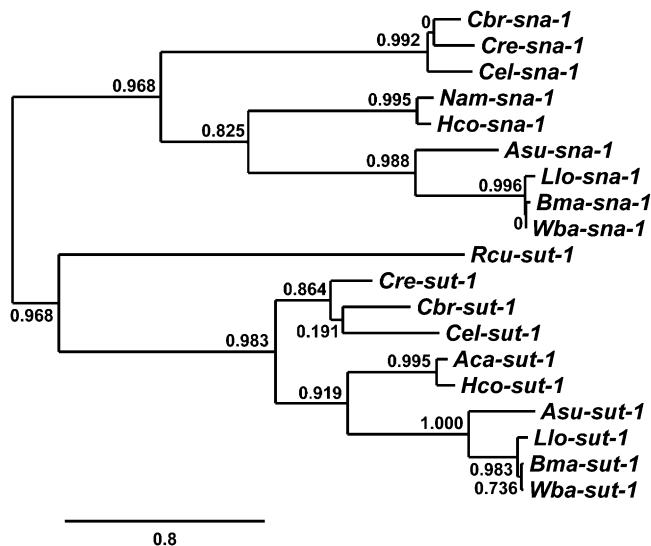


Figure 5 Evolutionary relationship between snRNPs associated with SL *trans*-splicing. Unrooted PhyML tree showing the relationship between *sna-1* and *sut-1* homologs identified in selected nematodes. Genes were named on the basis of their *C. elegans* homologs. The following species-specific prefixes were used: Aca, *Angiostrongylus cantonensis*; Asu, *Ascaris suum*; Bma, *Brugia malayi*; Cel, *C. elegans*; Cbr, *C. briggsae*; Cre, *C. remanei*; Hco, *Haemochus contortus*; Llo, *Loa loa*; Nam, *Necator americanus*; Rcu, *Romanomermis culicivorax*; and Wba, *Wuchereria bancrofti*. The numbers at each node are approximate likelihood ratio test statistics.

genes, or the first gene in an operon. Other nematodes within the Rhabditina also show the same specialization (Evans *et al.* 1997; Lee and Sommer 2003), but other taxa within the Chromadorea appear to use the same set of SLs for all SL *trans*-spliced mRNAs (Guiliano and Blaxter 2006). Within the Enoplea we clearly see multiple different SLs used in the *trans*-splicing of mRNAs derived from downstream genes in operons, but our data are not comprehensive enough to determine whether some SLs are preferentially used to process the transcripts arising from such genes.

We have identified multiple examples of gene pair synteny conserved between *T. spiralis*, *R. culicivorax*, *B. malayi*, and *C. elegans*, suggesting these correspond to conserved operons present in the last common ancestor of the Enoplea and Chromadorea. Nevertheless, the majority of the putative operons that we have identified in *T. spiralis* are not syntenic with *C. elegans* operons, although in many cases we see that the component genes of these putative *T. spiralis* operons have *C. elegans* orthologs that are found in operons. An example of such an operon is that comprising *Tsp-zgpa-1*, *Tsp-dif-1*, and *Tsp-aph-1*, whose *C. elegans* orthologs are located in different operons. We also observed putative *T. spiralis* operon genes whose *C. elegans* orthologs are not part of operons. The changes in the operon complements between the two nematodes could arise from the lineage-specific rearrangements of an ancestral set of nematode operons, but could also be accounted for by differential *de novo* operon generation in the two clades, or a mixture of the two processes. A key question is whether the synteny shown between the putative

enoplean operons and their chromadorean orthologs is significant, *i.e.*, are the genes that comprise these operons more constrained to be located in the same operon than other operonic genes, or is the conservation of synteny merely random chance? The availability of tools to engineer the *C. elegans* genome (Golic 2013) might allow this question to be addressed by assaying the function of selected operons compared with their individual component genes each expressed under their own promoters.

Finally, it is clear from this work and previous studies that SL *trans*-splicing and operon organization arose prior to the divergence of the Enoplea and Chromadorea. An important question is whether “nematode” SL *trans*-splicing and operons predate the foundation of the phylum. It will thus be important to establish whether they are also present in the other, so far uncharacterized, phyla that are closely related to the nematodes.

Acknowledgments

We thank Allison Bancroft and Richard Grecnis for providing *T. muris* RNA. The Bristol (N2) *C. elegans* N2 strain was provided by the *Caenorhabditis* Genetics Center, which is funded by National Institutes of Health Office of Research Infrastructure Programs (P40 OD010440). This work was funded by a Biotechnology and Biological Sciences Research Council project grant BB/J007137/1.

Literature Cited

- Allen, M. A., L. W. Hillier, R. H. Waterston, and T. Blumenthal, 2011 A global analysis of *C. elegans trans*-splicing. *Genome Res.* 21: 255–264.
- Beech, R. N., A. J. Wolstenholme, C. Neveu, and J. A. Dent, 2010 Nematode parasite genes: what’s in a name? *Trends Parasitol.* 26: 334–340.
- Blumenthal, T., D. Evans, C. D. Link, A. Guffanti, D. Lawson *et al.*, 2002 A global analysis of *Caenorhabditis elegans* operons. *Nature* 417: 851–854.
- Dana, C. E., K. M. Glauber, T. A. Chan, D. M. Bridge, and R. E. Steele, 2012 Incorporation of a horizontally transferred gene into an operon during cnidarian evolution. *PLoS ONE* 7: e31643.
- Davis, R. E., and S. Hodgson, 1997 Gene linkage and steady state RNAs suggest *trans*-splicing may be associated with a polycistronic transcript in *Schistosoma mansoni*. *Mol. Biochem. Parasitol.* 89: 25–39.
- Denker, J. A., D. M. Zuckerman, P. A. Maroney, and T. W. Nilsen, 2002 New components of the spliced leader RNP required for nematode *trans*-splicing. *Nature* 417: 667–670.
- Dereeper, A., V. Guignon, G. Blanc, S. Audic, S. Buffet *et al.*, 2008 Phylogeny.fr: robust phylogenetic analysis for the non-specialist. *Nucleic Acids Res.* 36: W465–W469.
- Evans, D., D. Zorio, M. MacMorris, C. E. Winter, K. Lea *et al.*, 1997 Operons and SL2 *trans*-splicing exist in nematodes outside the genus *Caenorhabditis*. *Proc. Natl. Acad. Sci. USA* 94: 9751–9756.
- Ganot, P., T. Kallesøe, R. Reinhardt, D. Chourrout, and E. M. Thompson, 2004 Spliced-leader RNA *trans* splicing in a chordate, *Oikopleura dioica*, with a compact genome. *Mol. Cell. Biol.* 24: 7795–7805.

- Ghedini, E., S. Wang, D. Spiro, E. Caler, Q. Zhao *et al.*, 2007 Draft genome of the filarial nematode parasite *Brugia malayi*. *Science* 317: 1756–1760.
- Golic, K. G., 2013 RNA-guided nucleases: a new era for engineering the genomes of model and nonmodel organisms. *Genetics* 195: 303–308.
- Graber, J. H., J. Salisbury, L. N. Hutchins, and T. Blumenthal, 2007 *C. elegans* sequences that control trans-splicing and operon pre-mRNA processing. *RNA* 13: 1409–1426.
- Gu, T., S. Orita, M. Han, 1998 *Caenorhabditis elegans* SUR-5, a novel but conserved protein, negatively regulates LET-60 Ras activity during vulval induction. *Mol. Cell. Biol.* 18: 4556–4564.
- Guiliano, D. B., and M. L. Blaxter, 2006 Operon conservation and the evolution of trans-splicing in the phylum Nematoda. *PLoS Genet.* 2: e198.
- Harrison, N., A. Kalbfleisch, B. Connolly, J. Pettitt, and B. Müller, 2010 SL2-like spliced leader RNAs in the basal nematode *Prionchulus punctatus*: New insight into the evolution of nematode SL2 RNAs. *RNA* 16: 1500–1507.
- Hastings, K. E. M., 2005 SL trans-splicing: easy come or easy go? *Trends Genet.* 21: 240–247.
- Hobert, O., 2002 PCR fusion-based approach to create reporter gene constructs for expression analysis in transgenic *C. elegans*. *Biotechniques* 32: 728–730.
- Holterman, M., A. van der Wurff, S. van den Elsen, H. van Megen, T. Bongers *et al.*, 2006 Phylum-wide analysis of SSU rDNA reveals deep phylogenetic relationships among nematodes and accelerated evolution toward crown Clades. *Mol. Biol. Evol.* 23: 1792–1800.
- Huang, P., E. D. Pleasance, J. S. Maydan, R. Hunt-Newbury, N. J. O’Neil *et al.*, 2007 Identification and analysis of internal promoters in *Caenorhabditis elegans* operons. *Genome Res.* 17: 1478–1485.
- Jacob, F., D. Perrin, C. Sánchez, J. Monod, and S. Edelman, 1960 The operon: a group of genes with expression coordinated by an operator. *C. R. Acad. Sci. Paris* 250: 1727–1729.
- Johnson, P. J., J. M. Kooter, and P. Borst, 1987 Inactivation of transcription by UV irradiation of *T. brucei* provides evidence for a multicistronic transcription unit including a VSG gene. *Cell* 51: 273–281.
- Jones, T., W. Otto, M. Marz, S. Eddy, and P. Stadler, 2009 A survey of nematode SmY RNAs. *RNA Biol.* 6: 5–8.
- Lawrence, J., 1999 Selfish operons: the evolutionary impact of gene clustering in prokaryotes and eukaryotes. *Curr. Opin. Genet. Dev.* 9: 642–648.
- Lee, K.-Z., and R. J. Sommer, 2003 Operon structure and trans-splicing in the nematode *Pristionchus pacificus*. *Mol. Biol. Evol.* 20: 2097–2103.
- Liu, C., A. Oliveira, C. Chauhan, E. Ghedin, and T. R. Unnasch, 2010 Functional analysis of putative operons in *Brugia malayi*. *Int. J. Parasitol.* 40: 63–71.
- MacMorris, M., M. Kumar, E. Lasda, A. Larsen, B. Kraemer *et al.*, 2007 A novel family of *C. elegans* snRNPs contains proteins associated with trans-splicing. *RNA* 13: 511–520.
- Marlétaz, F., A. Gilles, X. Caubit, Y. Perez, C. Dossat *et al.*, 2008 Chaetognath transcriptome reveals ancestral and unique features among bilaterians. *Genome Biol.* 9: R94.
- Maroney, P. A., Y. T. Yu, M. Jankowska, and T. W. Nilsen, 1996 Direct analysis of nematode cis- and trans-spliceosomes: a functional role for U5 snRNA in spliced leader addition trans-splicing and the identification of novel Sm snRNPs. *RNA* 2: 735–745.
- Meldal, B. H. M., N. J. Debenham, P. De Ley, I. T. De Ley, J. R. Vanfleteren *et al.*, 2007 An improved molecular phylogeny of the Nematoda with special emphasis on marine taxa. *Mol. Phylogenet. Evol.* 42: 622–636.
- Mitrevva, M., D. P. Jasmer, D. S. Zarlenga, Z. Wang, S. Abubucker *et al.*, 2011 The draft genome of the parasitic nematode *Trichinella spiralis*. *Nat. Genet.* 43: 228–235.
- Morton, J. J., and T. Blumenthal, 2011 Identification of transcription start sites of trans-spliced genes: uncovering unusual operon arrangements. *RNA* 17: 327–337.
- Pettitt, J., N. Harrison, I. Stansfield, B. Connolly, and B. Müller, 2010 The evolution of spliced leader trans-splicing in nematodes. *Biochem. Soc. Trans.* 38: 1125–1130.
- Pettitt, J., B. Müller, I. Stansfield, and B. Connolly, 2008 Spliced leader trans-splicing in the nematode *Trichinella spiralis* uses highly polymorphic, noncanonical spliced leaders. *RNA* 14: 760–770.
- Protasio, A. V., I. J. Tsai, A. Babbage, S. Nichol, M. Hunt *et al.*, 2012 A systematically improved high quality genome and transcriptome of the human blood fluke *Schistosoma mansoni*. *PLoS Negl. Trop. Dis.* 6: e1455.
- Satou, Y., K. Mineta, M. Ogasawara, Y. Sasakura, E. Shoguchi *et al.*, 2008 Improved genome assembly and evidence-based global gene model set for the chordate *Ciona intestinalis*: new insight into intron and operon populations. *Genome Biol.* 9: R152.
- Schiffer, P. H., M. Kroiher, C. Kraus, G. D. Koutsovoulos, S. Kumar *et al.*, 2013 The genome of *Romanomermis culicivorax*: revealing fundamental changes in the core developmental genetic toolkit in Nematoda. *BMC Genomics* 14: 923.
- Spieth, J., G. Brooke, S. Kuersten, K. Lea, and T. Blumenthal, 1993 Operons in *C. elegans*: polycistronic mRNA precursors are processed by trans-splicing of SL2 to downstream coding regions. *Cell* 73: 521–532.
- Tsai, I. J., M. Zarowiecki, N. Holroyd, A. Garciarrubio, A. Sanchez-Flores *et al.*, 2013 The genomes of four tapeworm species reveal adaptations to parasitism. *Nature* 496: 57–63.
- Zaslaver, A., L. R. Baugh, and P. W. Sternberg, 2011 Metazoan operons accelerate recovery from growth-arrested states. *Cell* 145: 981–992.
- Zuker, M., 2003 Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic Acids Res.* 31: 3406–3415.

Communicating editor: M. P. Colaiácovo

GENETICS

Supporting Information

<http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.114.162875/-/DC1>

Operons Are a Conserved Feature of Nematode Genomes

**Jonathan Pettitt, Lucas Philippe, Debjani Sarkar, Christopher Johnston, Henrike Johanna Gothe,
Diane Massie, Bernadette Connolly, and Berndt Müller**

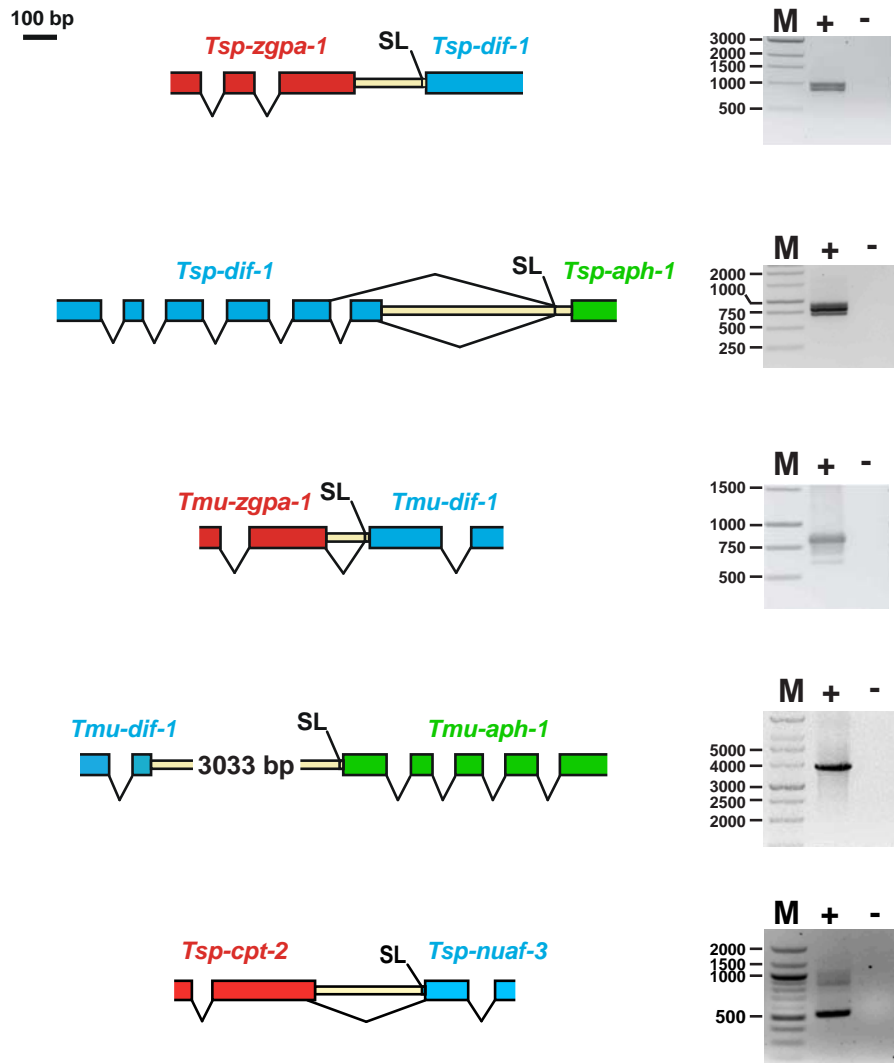


Figure S1 RT-PCR detection of polycistronic RNA. The exon-intron structures of the amplicons used to identify polycistronic RNAs are shown, with exons represented by boxes (shaded to identify the genes from which they are derived using the same colour coding that was used in Figures 1A and 3A). Cream coloured boxes represent the intergenic regions. The positions of the SL trans-splice acceptor sites are indicated. Images of representative RT-PCR products produced from each pair of primers are shown to the right of each amplicon cartoon. Reactions were carried out in the presence or absence of reverse transcriptase, + or -, respectively. M indicates DNA size markers.

```

Tsp-SNA-2      1  -----
Rcu-SNA-2      1  -----
Asu-SNA-2      1  -----MDVIMGEETTS DNPLAE
Llo-SNA-2      1  -----
Bma-SNA-2      1  -----
Wba-SNA-2      1  -----
Hco-SNA-2      1  MGTPQMKGRGLRLCAIACMPTLN GEEVLVGGNYLVGRGRCWRHATVLDKRRGEDGTLEL
Ace-SNA-2      1  -----
Cre-SNA-2      1  -----
Cel-SNA-2      1  -----

Tsp-SNA-2      1  -----
Rcu-SNA-2      1  -----
Asu-SNA-2     18  NMQTDENVLRSNVVESRSRQRSPSAARDDPVDV-----SPTTAHSEIPLNV
Llo-SNA-2      1  -----MDDGTVLEVRI--
Bma-SNA-2      1  -----MDDGTVLEVRI--
Wba-SNA-2      1  -----MDDGTVLEVRI--
Hco-SNA-2     61  YLHFEGDDRRRLDHWLDS SRLKFRNQPSKRITSA YEMREKRIRRAKESPNSASLDYAEIL
Ace-SNA-2      1  -----
Cre-SNA-2      1  -----
Cel-SNA-2      1  -----

Tsp-SNA-2      1  -----
Rcu-SNA-2      1  -----
Asu-SNA-2     64  SETSMTAPPISRGPMTPPGPELVSTPGSSTPPP-----
Llo-SNA-2     12  -----
Bma-SNA-2     12  -----
Wba-SNA-2     12  -----
Hco-SNA-2    121  EEQHKKEVTKVKYIEVVRYGEFEMDTWYVSPYPDEY GKERYLFICDKCFLYMRQERAFKMM
Ace-SNA-2      1  -----
Cre-SNA-2      1  -----
Cel-SNA-2      1  -----

Tsp-SNA-2      1  -----
Rcu-SNA-2      1  -----
Asu-SNA-2     97  -----
Llo-SNA-2     12  -----
Bma-SNA-2     12  -----
Wba-SNA-2     12  -----
Hco-SNA-2    181  LSSCTAKHPPGREIYVDYSENI AVYEV DGEKEKLFQCCLCLFAKLFMDHKTIYFDVTTFL
Ace-SNA-2      1  -----
Cre-SNA-2      1  -----
Cel-SNA-2      1  -----

Tsp-SNA-2      1  -----
Rcu-SNA-2      1  -----
Asu-SNA-2     97  -----CAIVPAESSKSAAYPPQGVDS AAMQNSIFVAPTRTS
Llo-SNA-2     12  -----NGEINGDAACL-----
Bma-SNA-2     12  -----NGEINGDAACL-----
Wba-SNA-2     12  -----NGEINGDAACL-----
Hco-SNA-2    241  FYVVCQLRESGEPRPVGYFSKERTSPDGHNLSCLLVFP AFQRQGFGLIQLSYELSRRE
Ace-SNA-2      1  -----MVSCLLYLLEFRR-----
Cre-SNA-2      1  -----
Cel-SNA-2      1  -----

Tsp-SNA-2      1  -----
Rcu-SNA-2      1  -----
Asu-SNA-2    134  GPATPPEGPGPDASP-----
Llo-SNA-2     23  -----
Bma-SNA-2     23  -----
Wba-SNA-2     23  -----
Hco-SNA-2    301  GIQGSPEKPLSDLGAASFRHYWAYLIVEYLSGFMDTAWIRVSELAKTLGMQAE DVVDTLH
Ace-SNA-2     14  -----
Cre-SNA-2      1  -----
Cel-SNA-2      1  -----

```

```

Tsp-SNA-2      1 -----
Rcu-SNA-2      1 -----
Asu-SNA-2     149 -----PISSSPHSTTIPQSYSASPAASLSTALVSPQ
Llo-SNA-2     23 -----AKLQNTA
Bma-SNA-2     23 -----AKLQNAA
Wba-SNA-2     23 -----AKLQNTA
Hco-SNA-2    361 WLRICDPVMSDAPDDDIELWVVYIKKLDSLRTKAAKPPLMDDKEHEENVQENGSSHDTPE
Ace-SNA-2     14 -----IINKEAMEDQDKEDHVQQNGSHNSTE
Cre-SNA-2     1 -----
Cel-SNA-2     1 -----MI

Tsp-SNA-2      1 -----MVLAKSKDSDTA-----PSVSLARDVSSDDDVDSNVAG
Rcu-SNA-2      1 -----
Asu-SNA-2    180 SSTLSSSTAASNTTQVISARVAEKVTPHRDSECDSNGIELNKTSDPIDTSAKGEKIE
Llo-SNA-2     30 ESISSPSTPCIETSSSLSDTA--ITLIRKELEONFSAKPKINDQSHAENFIKSNEQASD
Bma-SNA-2     30 ESISSPSTPCIETSSSLSDTA--ISLMRKEFEHENCAAKPKSNDQIHVENCVKSSQVL-
Wba-SNA-2     30 ESISSPSTPCIETSSSLSDTA--ISLMRKEFEHENCAAKPKPNDQIHVENCVKSSQVSD
Hco-SNA-2    421 ESITAPLVNCTADSTAIDSG-----SEAVTEGSSAGPASTAETEIVKDKSTPSEVNE
Ace-SNA-2     40 EPAVVSSVNCTAENTAIDSG-----IETAAEGSSREEAAVETGNPATETAEAPCSS
Cre-SNA-2     1 ---MSEVNFVIEQGGEEET-----
Cel-SNA-2     3 LNPVKEYENVDIQSDAFDSV-----ENVAE

Tsp-SNA-2     34 QAQQADNEKEVRVDGLIVPCD-MQLPVLSEAMYACBRFACCEK-KSNALEQEKISVVSL
Rcu-SNA-2     1 -----
Asu-SNA-2    240 SREVHDAKRKFRVVGLTFPGKWIHHEKLMEIFVKSESFHMWFDK--PEKAGVKSGVSV
Llo-SNA-2     88 SRDSSSRRKFRVLGLTFPGKWMTSEKFNDLLEKSESFHIWEDK---TKGGIKHGLSA
Bma-SNA-2     87 --DSNDSKRKFRVLGLTFPGKWMTSEKFNDLLEKSESFHMWFDK--TKGGIKHGLST
Wba-SNA-2     88 SRDLNDSKRKFRVLGLTFPGKWMTSANENDLLEKSESFHMWFDK--TKDCTKCGLST
Hco-SNA-2    474 EPELDN--RTFTVTGLKHPGGWRTHSEINMLARCDSFIREGR--AMKQKEPKPESISF
Ace-SNA-2     93 ESPDVDN--RTFTVTGLKHPGAWRTHSEINMLARCESFNIKFER--GAKQKEPKPESISF
Cre-SNA-2     19 -----VEPAASLVSVNFQNLLNKSIDEADAQESRPENSNERQIKTERIE
Cel-SNA-2     29 DKINDKSLRTIVQNLKHPGGWLRNRQFQMLLNRCVTEDVSESRPEDIDEKQIKYEKVEF

Tsp-SNA-2     92 YDCLLKA--EAFCSAMKESVQRNEIK-----GIVTAKVVNLSAEELPVN
Rcu-SNA-2     1 -----
Asu-SNA-2    298 VEENADEAKAAFASVQKMSLDGASFKLOASRHEVGGGTLGGPSTIGSRSLTAPSLPFRVD
Llo-SNA-2    145 VVETPEEAKLAFAFHKMTFDGNPLKLOASKLFY---D--PVPNSRIVO--PPMPFKIN
Bma-SNA-2    142 VVDTPEEAKLAFAFHKMTFDGNPLKLOASKPFY---D--LTPNSRTIQ--PPMPFKIS
Wba-SNA-2    145 VVDTPEEAKLAFAFHKMTFDGNPLKLOASKPFY---D--LTPNSRTIQ--PPMPFKIN
Hco-SNA-2    530 TFESISSARESETQAQKMRVDQAVKVEASPAFF-----A--PAACRSRPF-----FKIP
Ace-SNA-2    150 TFESISSAREAYTQAQKMRVDQAVKVEACPAFF-----A--PAACRSRPF-----FKIP
Cre-SNA-2     63 LFSSIDAARDAYTSLQKMVIDGFRFEVLDPRFE-----SVETTTQKRQF-----FEIA
Cel-SNA-2     89 VFSTVAGAREAYSTIQKMIDGIRFEALVDPQFE-----SVETFQSSRQF-----FDLS

Tsp-SNA-2    135 VDG-----LLAAANESLIRDIFAEYNVQEIANLSDVDSTTGTK--TATLRKESSE
Rcu-SNA-2     1 -----VKA--ETEISIKKP
Asu-SNA-2    358 VSEADSPRRIYALHLTPSTDCTLLSSILGSETIESI--RENDPIVPSER--QAEIVFRTA
Llo-SNA-2    197 MSDPDASRRIYALHLSATDQTLLSSIFGSECIETI--QLEADPFIASRK--QAEIVFHTV
Bma-SNA-2    194 LSDPDASRRIYALHLSATDQTVLSSIFGSECIETV--QLEVDPFIASRK--QAEIVFHTV
Wba-SNA-2    197 LSDPDVSRRIYALHLSATDQTVLSSIFGSECIETV--HLEVDPFIASRK--QAEIVFHTV
Hco-SNA-2    578 VDEETQKRTIYALDLPTSAERNLLNSIFEADHIE--KLTFLPLRNQHK--QAEVVMHTSE
Ace-SNA-2    198 TEEEVKKRTIYALDLPTSAERNLLNSIFEADHIE--RITFLPLRSQHK--QAEVVMHTSE
Cre-SNA-2    112 DS---RLVFLDAPRSIDEHMIEHFENGEAPI---HFQTLPMASENGLFQVEVLLSNG
Cel-SNA-2    138 NDT---RLVFLDAPRSIDHMISYFFEGKTPI---HFQTLPMPSSNGLFQVEVLLSDK

Tsp-SNA-2    182 QDAENMLEIND--E--LDDGVCSSLKISFLTPENNNNNNVRLSTEFVAQNEKAVCVDD
Rcu-SNA-2     13 ES-----YKNKNEENVRSVDDLTKLEQ-----
Asu-SNA-2    415 EEANDARIELADGFE--IDEGDROSTMML--LTAE-----EYLMHMKS
Llo-SNA-2    254 QEANDARVELADGFE--IDEGDROLTMRL--LTMD-----EYISVMKT
Bma-SNA-2    251 QEANDARMELADGFE--IDEGDROLTMRL--LTAD-----EYISVMKT
Wba-SNA-2    254 QEANDARMELADGFE--IDEGDROLTMRL--LTAD-----EYISVMKT
Hco-SNA-2    633 EQAEAARSE--DGFE--IDDGQQSVLRI--LKEA-----EYAKFVEE
Ace-SNA-2    253 EQADAARSE--DGFE--IDDGQQSVLRI--LTPA-----DYAAFVAE
Cre-SNA-2    165 EAAEKVLGGESK--EK--MHDGENECIVTL--LSPR-----EYALVSKM
Cel-SNA-2    191 ETAEKVLGAGPSK--ES--MSDEDNECIVTL--LSPR-----EYSLMSKM

```

```

Tsp-SNA-2 240 ETRRLQIEAEFAQSSDTRVDACQLMETNLAADVETADTEIDRTVVKOPPVVHPETVEIDVD
Rcu-SNA-2 36 -----VNHGAQF-----SGESES
Asu-SNA-2 454 GS-----EKAATISQRAVRPSSSPVSTP-----QTITQSPSASFVPAPSEITVE
Llo-SNA-2 293 ENDKILNA-----RVVPVSGSSAPSSSSQPSPSVSSTAIAATSETAVVPPFPAPSEITVE
Bma-SNA-2 290 ENDKILNA-----RVVPVSAASSAPSSSSQPSPSVSSIAVAATSETAVVPPFPAPSEITVE
Wba-SNA-2 293 ENDKILNA-----RVVPVSAASSAPSSSSQPSPSVSSIAVAATSETAVVPPFPAPSEITVE
Hco-SNA-2 670 EQ-----KPIPFVPPPEVDSSSAQEPVATTTQQA-APVPADAKL----APVLDED
Ace-SNA-2 290 EQ-----KPIPFVPPPEVDFPPVKASTPSVSOPTTASPAEVKL----APVLDED
Cre-SNA-2 203 DD-----IRPSTSRKAPVQHSSESPSNTTGGQL----APEIDED
Cel-SNA-2 229 DD-----IQQSTSAKKPLKQSYISENVGQTL----APEIDED

Tsp-SNA-2 300 QISAHREYMRNKRINNAEMNTVKELWVVLNQLT--SEIGWQPEEMLLKGMKAALFYAVQ
Rcu-SNA-2 50 EISRLIDTYEANKRVNFAELHSTEDLWSVHDAAPMDKVGWQPEOTLKDALLMSTLQNALP
Asu-SNA-2 496 DVTALQKHTEETRTNWAELNELGELWKLCDSEVS--RQRKCHPDSMLKPAALLNVLRRHQ
Llo-SNA-2 348 DVTSFVQQYVDDTRTNWAELNELNELWKLCDSEVS--QLHRCHPDSLLKPAALLNVLRRHQ
Bma-SNA-2 345 DVTGFQQYVDDTRTNWAELNELNELWKLCDSEVS--QLHRCHPDSLLKPAALLNVLRRHQ
Wba-SNA-2 348 DVTGFQQYVDDTRTNWAELNELNELWKLCDSEVS--QLHRCHPDSLLKPAALLNVLRRHQ
Hco-SNA-2 714 DVTDMFIKVTDRVNWAEITDVMELYTMCDAVS--AQIGCHPDSVLRPAMLHTLORHLT
Ace-SNA-2 335 DVTDMFIKVTDRVNWAEITDVMELYTMCDAVS--AQVGCHPDSVLRPAMLHTLORHLT
Cre-SNA-2 239 ILLNRLEIIEERKLNFAEINNEKNEELYELCDTVS--AEYNGVPSLTKPAMGTVLRHLN
Cel-SNA-2 262 VISNRLEIIEERKLNFAEINNEKNEELYELCDTVS--SEYNGVPSLTKPAMGVSVLQKHLN

Tsp-SNA-2 358 EFSVQVLAQHLKHLTQLWNDLMLLEKRSETRKELLMAVEEDLSYL-----
Rcu-SNA-2 110 FNDA-HWQONHIEFLIKLWKRELDLECRSKVRKQLLAVEMEMFLAANGQQKLPLOTET
Asu-SNA-2 554 TLPG-HWLRHDVDDLKMKREIQNTKSSHRPTAYINMSAAPYSPR-----
Llo-SNA-2 406 SLPT-HWMRQHVDNLIIRWKEVQTRGVORSYSYMSAAPYIPYD-----
Bma-SNA-2 403 TLPT-HWMRQHVDNLIIRWKEVQNTKGVORSYSYMSAAPYIPYD-----
Wba-SNA-2 406 TLPT-HWMRQHVDNLIIRWKEVQNTKGVORSYSYMSAAPYIPYD-----
Hco-SNA-2 772 EAQT-AWMREHLQDLIKGWKEVKKDAFVDRPFLVOMKAABVFPVA-----
Ace-SNA-2 393 EAQS-AWMREHLQDLIKGWKEVVRNDTFVDRPFLVOMKAADVFPVA-----
Cre-SNA-2 297 RTEH-HWMRDHIEGLLRMWMKMEIMNEQIYERSFVPMETVNYQPMVI-----
Cel-SNA-2 320 RTEH-HWMREHLEGLLRVWVKTEILNEQTYERPSFVPMETVAMQPVV-----

Tsp-SNA-2 404 -----PKASPYCGLKKKKK---YPVV
Rcu-SNA-2 169 KVANAGETAETTTTTTFAEHSSTSLHPPQKSSSTYEPASPSPKPKRKGAGKR--ALIEL
Asu-SNA-2 599 -----RVSPHREPNETKRCQ--ALRMM
Llo-SNA-2 451 -----IPERKNS-NSTRKR--ALATI
Bma-SNA-2 448 -----VPERRNSKISTTRKR--ALATI
Wba-SNA-2 451 -----VSERRNSKISTTRKR--ALATI
Hco-SNA-2 817 -----PSKRRLRGKNSAESR--AGRVL
Ace-SNA-2 438 -----PSKRKFRGKNSAQSR--AGRVM
Cre-SNA-2 342 -----QEEQKETASKGLKRCQKARAQ
Cel-SNA-2 365 -----EEQKTESKGAEKRNQARKQ

Tsp-SNA-2 423 RGVGGLERLRKE-ANDAAQLSTKMKELKQENGRKKQQQQQHLLVFAASAKSAAVKNGETS
Rcu-SNA-2 227 MGVGAFLANVQKEHAVSAPKSS-----DDDDQSGVDDKGTSEFKRKAS--KVNLPV
Asu-SNA-2 619 MGHGVLNLTARLKLATEEGELDV-----DEDEHGNVTIDGQPLSFESWAE--ITKPTSPKA
Llo-SNA-2 470 MGVGAEVLAKVRTMLATEEGELDV-----DEDEHGNVTIDGQPLSFESWAS--LNKPKVPAP
Bma-SNA-2 468 MGVGAEVLAKVRTMLATEEGELDV-----DEDEHGNVTIDGQPLSFESWAN--LNKPKVPAP
Wba-SNA-2 471 MGVGAEVLTKVRTMLATEEGELDV-----DEDEHGNVTIDGQPLSFESWAN--LNKPKVPPT
Hco-SNA-2 837 MGVGAELEAQRQKMTVEEGELDV-----EDEDGNIMLGGEALSFEWAK--LTKTNQEDI
Ace-SNA-2 458 MGVGAELEAQRQKMTVEEGELDV-----EDEDGNIMLGGEALSFEWAK--LTKTNQEDI
Cre-SNA-2 364 MGVGAEFLTANRDRLIVRSQDVEI-----ESDDGDMVVGGEALSFEAWAR--LTKTKASGV
Cel-SNA-2 387 MGVAAFLNANRDRLIVRSQDVEI-----ESDDGDMVVTGGEALSFEWSAR--LTKTKASGV

Tsp-SNA-2 482 STTTTTTTKSARRHASLSKESGEASSSEDEAQEINTVNAKFHNNNNNNNSRPVHNA
Rcu-SNA-2 278 PSSSAPAPKRRPGRSSSSCESGEASSTSSSEDGDDGGRH-----DGSASSDDDRGGT
Asu-SNA-2 673 AVE-----PQPRILPEYRKRIR-----NQRKKEF-----REKMAEA
Llo-SNA-2 524 VYER-----PQLYIDPKMRKLRN-----QWKKETW-----QKKLEEE
Bma-SNA-2 522 VYER-----PQLFIHPKMRKFRN-----QWKKETW-----QKKKEE
Wba-SNA-2 525 VYER-----PQLFIDPKMRKFRN-----QWKKETW-----QKKM-EE
Hco-SNA-2 891 VCARLPGDEPAAKK--SKPEVHDFRKL--KQSQKEW-----KEKRMGAK
Ace-SNA-2 512 VCANMGQEEPPAKKPR--EMHDFKKY--KQSQKEW-----KEKRMGAK
Cre-SNA-2 418 VRANDDGSKKIGNGLSKKQMRQARLVEGMSQEEWKKW-----RNEKSSAR
Cel-SNA-2 441 VRANDDGSKKIGNGLSKKQMRQARLVEGMSQEEWKKW-----RNEKSSNR

```



```

Tsp-SNA-2 542 KNDIPSFHQYINPMVQAADAPEDLAYNHVFGHTFPDTSAHLIPMMTRLMEMFLQFLQQQN
Rcu-SNA-2 329 VHHKSRRRRASAVATASPPPEPPARNRSERRQQQQQRDDLRTHGDD--OCDFRVVQO--
Asu-SNA-2 704 KMQKMREAVEPSKEPEVEEQPEGSKEMEAVDQGSVTESEA-----PEKKRELEEE--
Llo-SNA-2 556 KVKKTQEVVPEIENENEKDEMAP-----DGEVVRATET-----PVIRRELEEE--
Bma-SNA-2 554 KIKKTQEVLEIENEKDEESAP-----DGEVKAETET-----PVIRRELEEE--
Wba-SNA-2 556 KIKKTQEVLEIENEKDEESAP-----DGEVTAETET-----PVIRRELEEE--
Hco-SNA-2 931 KMRIMQOAREELEALDKEKQAPAQTKNPKDVENGKPKKDAEKPS----SKPAKELDE--
Ace-SNA-2 551 KMRIMQOAREELEALDQEKAEQAQV--PKEAENGKDKKENASPEKP--PKPPKELDE--
Cre-SNA-2 464 KADIKQRIIMARGAVIEGEDSDVDGGAEMMDSEATSSSLPTTAMTTAPPV--PPVKELDE--
Cel-SNA-2 487 KADIKQRIIMARGAIDEGEIDNETAVNG-NSQVAIEETVEPSPTL-----AAVKRILDE--

Tsp-SNA-2 602 FQMPLETTTTNNNSNIAQQOQLYQQLQMFQVQVGFGNARPPAPVRQPLNLGFRYPFIPQOQQ
Rcu-SNA-2 385 -----ANDLHSSQNIHHPAEYDQF---IAQVPPPPPAQOQQQTWMAQMMPFMMQO
Asu-SNA-2 753 -----GEMSSD-----SSSPSSSSSSDSDMSQHTSRD
Llo-SNA-2 599 -----GEMSSS-----SSSSSASSSSEDDIDDGPTSRHR
Bma-SNA-2 595 -----GEMSSS-----SSSSSASSSSEDDIDDGATSRHR
Wba-SNA-2 597 -----GEMSSS-----SSSSSASSSSEDDIDDGATSRHR
Hco-SNA-2 984 -----GEIDSEERREEKPKRKRKRVSSSSNSSSSSSSSEDDDG--DARKR
Ace-SNA-2 604 -----GEIDSEERREEKPKRKRK--ASSNSSSSSSSSEDDDG--DARKR
Cre-SNA-2 521 -----GEIDSEEEKKDKTKNKIKRR--ASDSSDSSSSSEDDPDGVPDARKR
Cel-SNA-2 539 -----GEIDSEEEKQASKSKKKK--ATSDSSDSSSSSEDDPDGVPDARKR

Tsp-SNA-2 662 QQQQPQQQANMFRSAAVAANAHIAYGSNRSIIPPRP-----SLTGIHSQAQVHSN
Rcu-SNA-2 431 QFMPMMSQMFQNRMPAASAIN--TPOFQQQFMQLYQVQSQIASGTMMAASSAMSALSSNP
Asu-SNA-2 781 RRRKRTRTRKNQQRSGSFLGSGSPHFESEFKQLYEHRH-----SL--LRSLSA--THKH
Llo-SNA-2 629 RRRKRTRHEKSSQLNSVQAAMGSVNPQFTTLFAQLYEHRH-----TL--CRLLLS--SHKN
Bma-SNA-2 625 RRRKRTRHEKNNQLNPVQAAMGSVNPQFTTLFAQLYEHRH-----TL--CRLLLS--SHKN
Wba-SNA-2 627 RRRKRTRHDKNNQLNPVQAAMGSVNPQFTTLFAQLYEHRH-----TL--CRLLLS--SHKN
Hco-SNA-2 1030 RRRNRKTERLK-----GQVPPPLFORMFNVRH-----EI--INALSA--EHKT
Ace-SNA-2 648 RRRNRKTERMK-----GQVPPPLFORMFNVRH-----VI--INALSA--EHKT
Cre-SNA-2 567 RNRKRKMDRKNPRTTAA-----SAQLDPVFKQMFENRK-----AI--IAQMSP--AHKA
Cel-SNA-2 585 RKLRRKMDRKNPRTTAA-----SSLDPVFKQMFENRK-----TI--LAQMTF--AHKS

Tsp-SNA-2 716 KGVVAVGGGDSATYFRKCSISTAAACVPNDYTNISAVSFLQPTDRSLTVVDVDIRLGRI
Rcu-SNA-2 489 MFAQLAORMGAAS--TSSTATQOMLPQQOHOVVPYGSAGLGPYHYR-----
Asu-SNA-2 832 GFASVILGQILSS----PQGMQSA-----QHDQMAVEMRTFSAHNMHGK-----
Llo-SNA-2 680 GFASVILGQILSS----PQGMQSA-----QONOMNLFIONFSAQNPFK-----
Bma-SNA-2 676 AFASVILGQVLS--PQGMQSA-----QONOMNLFIONFSAQNPFK-----
Wba-SNA-2 678 AFASVILGQVLS--PQGMQSA-----QONOMNLFIONFSAQNPFK-----
Hco-SNA-2 1069 AFASVILHOMKSSG--RTGISHS--DQOQMLTFLSSFSR-----
Ace-SNA-2 687 AFASVILHOMKSSG--RTGISQS--DQOQMLTFLSSFSR-----
Cre-SNA-2 614 AFASVILTOIAQNN--ASSAQQA--KASQVINTMMSGFR-----
Cel-SNA-2 630 AFVSAVLTQIVNNN--PSGVSQA--KASQVINTMMSGFR-----

Tsp-SNA-2 776 GRNGR
Rcu-SNA-2 -----
Asu-SNA-2 -----
Llo-SNA-2 -----
Bma-SNA-2 -----
Wba-SNA-2 -----
Hco-SNA-2 -----
Ace-SNA-2 -----
Cre-SNA-2 -----
Cel-SNA-2 -----

```

Figure S2 Multiple sequence alignment of SNA-2 homologues. Residues conserved in 50% or more of the homologues are shaded black (identities) or grey (similarities). Species prefixes used as follows: Tsp - *T. spiralis*; Rcu - *R. culicivora*; Asu - *A. suum*; Llo - *L. loa*; Bma - *B. malayi*; Wba - *W. bancrofti*; Hco - *H. contortus*; Ace - *A. ceylanicum*; Cre - *C. remanei*; Cel - *C. elegans*.

Table S1 Operonic status of *T. spiralis* genes that produce SL *trans*-spliced transcripts.

<i>T. spiralis</i> Gene ¹	<i>Tsp</i> -SL Number	EST	<i>T. spiralis</i> Gene Operon Status ²	<i>C. elegans</i> homologue	<i>C. elegans</i> Gene Operon Status	<i>C. elegans</i> Operon Number
Tsp_04323	3	ES566285	Upstream	ZK829.7	Upstream	4484
Tsp_04801	1	ES565008	Upstream	<i>prmt-1</i>	Upstream	5508
Tsp_07643	3	ES272890	Upstream	Y57G11C.22	Upstream	4661
Tsp_09747	3	ES565435	Upstream	<i>fkf-2</i>	Upstream	1903
Tsp_10562	3	ES566995	Upstream	K10C3.4	Upstream	1566
Tsp_10928	3	ES570931	Upstream	C17E4.6	Upstream	1544
Tsp_08576	2	ES568840	Upstream	<i>ung-1</i>	Upstream	3746
Tsp_00718	5	ES570675	Upstream	<i>ufd-1</i>	Downstream	4512
Tsp_01411	2	ES560816	Upstream	C16A3.4	Downstream	3364
Tsp_02079	3	ES570654	Upstream	<i>idi-1</i>	Downstream	3480
Tsp_06176	12	ES565259	Upstream	T04G9.4	Downstream	X147
Tsp_06763	4	BG353375	Upstream	C18E3.5	Downstream	1152
Tsp_08492	4	ES567181	Upstream	<i>wwp-1</i>	Downstream	1934
Tsp_08574	3	BQ692442	Upstream	<i>vps-28</i>	Downstream	1708
Tsp_08643	3	ES565814	Upstream	<i>pdha-1</i>	Downstream	2336
Tsp_11146	3	ES566810	Upstream	<i>hpo-11</i>	Downstream	1368
Tsp_09859	3	ES570398	Upstream	<i>gpdh-2</i>	Downstream	3795
Tsp_12645	12	ES567036	Upstream	<i>acs-16</i>	Non-Operonic	
Tsp_00181	2	ES569343	Upstream	Y53F4B.42	Non-Operonic	
Tsp_01518	2	ES563585	Upstream	<i>tomm-1</i>	Non-Operonic	
Tsp_02629	2	ES569407	Upstream	<i>tsp-11</i>	Non-Operonic	
Tsp_04773	3	ES565760	Upstream	<i>apd-3</i>	Non-Operonic	
Tsp_05351	1	ES565863	Upstream	<i>smn-1</i>	Non-Operonic	
Tsp_05770	2	ES570548	Upstream	<i>jnk-1</i>	Non-Operonic	
Tsp_09338	1	ES567564	Upstream	<i>rpl-32</i>	Non-Operonic	
Tsp_09639	2	ES568996	Upstream	<i>aldo-1</i>	Non-Operonic	
Tsp_10534	3	ES570566	Upstream	<i>skr-1</i>	Non-Operonic	
Tsp_11514	3	ES568803	Upstream	<i>rps-9</i>	Non-Operonic	
Tsp_12763	3	ES565626	Upstream	<i>cyn-7</i>	Non-Operonic	
Tsp_13841	2	ES273211	Upstream	Y37D8A.17	Non-Operonic	
Tsp_00182	4	ES570621	Downstream	<i>swd-2.2</i>	Upstream	4248
Tsp_00331	8	ES569101	Downstream	<i>ubl-1</i>	Upstream	3088
Tsp_00685	2	ES568435	Downstream	<i>mrps-17</i>	Upstream	3372
Tsp_00837	12	ES568088	Downstream	<i>rrbs-1</i>	Upstream	5416
Tsp_03905	2	ES563877	Downstream	F53F4.10	Upstream	5392
Tsp_04143	10	ES565871	Downstream	C06A8.2	Upstream	2731
Tsp_08528	2	ES569156	Downstream	Y55F3AM.9	Upstream	4642
Tsp_11545	9	ES564614	Downstream	<i>atp-5</i>	Upstream	5272
Tsp_00424	1	ES562658	Downstream	F44E7.4	Downstream	5565
Tsp_00734	2	ES560968	Downstream	F54C8.7	Downstream	3584
Tsp_02509	2	ES569909	Downstream	<i>ard-1</i>	Downstream	4412

Tsp_02873	9	ES569583	Downstream	<i>eif-3F</i>	Downstream	2424
Tsp_03622	12	ES565661	Downstream	<i>clpf-1</i>	Downstream	3108
Tsp_05308	2	ES562011	Downstream	<i>T09B4.9</i>	Downstream	1304
Tsp_06507	2	ES561596	Downstream	<i>raga-1</i>	Downstream	2528
Tsp_06818	12	ES563737	Downstream	<i>Y62E10A.10</i>	Downstream	4544
Tsp_08736	12	ES567003	Downstream	<i>lars-2</i>	Downstream	1396
Tsp_11177	3	ES563835	Downstream	<i>rlbp-1</i>	Downstream	1416
Tsp_08003	4	ES566741	Downstream	<i>ins-1</i>	Non-Operonic	
Tsp_00064	11	ES564912	Downstream	<i>mlp-1</i>	Non-Operonic	
Tsp_00608	1	ES568215	Downstream	<i>ppt-1</i>	Non-Operonic	
Tsp_00609	1	ES564896	Downstream	<i>mdt-19</i>	Non-Operonic	
Tsp_00958	2	ES564918	Downstream	<i>vha-8</i>	Non-Operonic	
Tsp_01274	2	ES564040	Downstream	<i>asg-1</i>	Non-Operonic	
Tsp_03114	3	BG353703	Downstream	<i>pdhb-1</i>	Non-Operonic	
Tsp_03857	3	ES570432	Downstream	<i>exc-7</i>	Non-Operonic	
Tsp_06181	3	ES569208	Downstream	<i>idh-2</i>	Non-Operonic	
Tsp_06812	1	ES564176	Downstream	<i>B0035.1</i>	Non-Operonic	
Tsp_06870	2	ES273592	Downstream	<i>rpl-33</i>	Non-Operonic	
Tsp_07223	4	ES566475	Downstream	<i>amph-1</i>	Non-Operonic	
Tsp_07289	3	ES273487	Downstream	<i>M153.2</i>	Non-Operonic	
Tsp_08373	3	ES565720	Downstream	<i>dyn-1</i>	Non-Operonic	
Tsp_11596	4	ES565698	Downstream	<i>cpl-1</i>	Non-Operonic	
Tsp_12613	2	ES565903	Downstream	<i>lip-1</i>	Non-Operonic	
Tsp_14129	10	ES564415	Downstream	<i>ace-1</i>	Non-Operonic	
Tsp_01343	3	ES273458	Ambiguous	<i>fnta-1</i>	Upstream	4008
Tsp_07410	3	ES569154	Ambiguous	<i>ent-7</i>	Upstream	1542
Tsp_11644	1	BG732210	Ambiguous	<i>rpb-7</i>	Upstream	1924
Tsp_00266	2	ES570691	Ambiguous	<i>C10H11.8</i>	Downstream	1780
Tsp_00680	2	ES561947	Ambiguous	<i>arl-5</i>	Downstream	3636
Tsp_01365	3	ES569529	Ambiguous	<i>hpl-2</i>	Downstream	3701
Tsp_02416	4	BG353715	Ambiguous	<i>B0205.6</i>	Downstream	1644
Tsp_04477	3	ES565704	Ambiguous	<i>kin-10</i>	Downstream	1456
Tsp_05422	3	ES569822	Ambiguous	<i>C50D2.7</i>	Downstream	2012
Tsp_05915	2	ES565057	Ambiguous	<i>D2085.3</i>	Downstream	2380
Tsp_07955	3	ES563974	Ambiguous	<i>prmt-7</i>	Downstream	1528
Tsp_08270	3	ES570232	Ambiguous	<i>Y32H12A.4</i>	Downstream	3802
Tsp_09368	2	ES568561	Ambiguous	<i>lsm-8</i>	Downstream	4176
Tsp_10835	3	ES565721	Ambiguous	<i>C09G9.1</i>	Downstream	4310
Tsp_11447	12	ES569795	Ambiguous	<i>C05C10.3</i>	Downstream	2751
Tsp_10765	2	ES567336	Ambiguous	<i>T25G3.1</i>	Ambiguous	
Tsp_00026	2	ES563422	Ambiguous	<i>F52E4.5</i>	Non-Operonic	
Tsp_00262	2	ES567711	Ambiguous	<i>pitp-1</i>	Non-Operonic	
Tsp_00717	4	ES273452	Ambiguous	<i>T07D1.2</i>	Non-Operonic	
Tsp_01078	2	ES564803	Ambiguous	<i>F21F3.7</i>	Non-Operonic	
Tsp_01228	4	ES570063	Ambiguous	<i>cco-2</i>	Non-Operonic	
Tsp_01287	4	ES562009	Ambiguous	<i>mtrr-1</i>	Non-Operonic	

Tsp_01573	2	ES565382	Ambiguous	<i>F37A8.5</i>	Non-Operonic	
Tsp_01930	1	ES570094	Ambiguous	<i>ifa-1</i>	Non-Operonic	
Tsp_03184	7	BG353788	Ambiguous	<i>F53F4.3</i>	Non-Operonic	
Tsp_03217	13	BQ541428	Ambiguous	<i>cmd-1</i>	Non-Operonic	
Tsp_03724	2	ES569761	Ambiguous	<i>nrf1-1</i>	Non-Operonic	
Tsp_04029	2	ES273134	Ambiguous	<i>BE0003N10.1</i>	Non-Operonic	
Tsp_04084	1	ES568229	Ambiguous	<i>F52D10.2</i>	Non-Operonic	
Tsp_04141	7	ES565838	Ambiguous	<i>T18D3.9</i>	Non-Operonic	
Tsp_04152	3	ES566044	Ambiguous	<i>C15H9.5</i>	Non-Operonic	
Tsp_05204	3	ES567104	Ambiguous	<i>C32E8.5</i>	Non-Operonic	
Tsp_07239	1	ES564125	Ambiguous	<i>swt-1</i>	Non-Operonic	
Tsp_08116	5	BG354896	Ambiguous	<i>F35G12.7</i>	Non-Operonic	
Tsp_09247	3	ES563796	Ambiguous	<i>F31D4.2</i>	Non-Operonic	
Tsp_09527	4	ES566597	Ambiguous	<i>mtch-1</i>	Non-Operonic	
Tsp_10184	2	BG353445	Ambiguous	<i>ck-1</i>	Non-Operonic	
Tsp_10493	10	ES563996	Ambiguous	<i>cpz-1</i>	Non-Operonic	
Tsp_10907	2	ES570466	Ambiguous	<i>kin-19</i>	Non-Operonic	
Tsp_11214	4	ES563030	Ambiguous	<i>C32F10.8</i>	Non-Operonic	
Tsp_01291	3	ES568398	Non-Operonic	<i>rpl-6</i>	Upstream	3836
Tsp_01684	13	BG353107	Non-Operonic	<i>rsp-6</i>	Upstream	4252
Tsp_01730	2	ES563676	Non-Operonic	<i>ubxn-4</i>	Upstream	3520
Tsp_05024	3	ES566087	Non-Operonic	<i>C01G5.6</i>	Upstream	4180
Tsp_05180	2	ES568248	Non-Operonic	<i>lin-53</i>	Upstream	1552
Tsp_06564	4	ES566240	Non-Operonic	<i>npp-1</i>	Upstream	4340
Tsp_06972	2	ES564495	Non-Operonic	<i>rla-0</i>	Upstream	1624
Tsp_07066	2	ES569769	Non-Operonic	<i>ubh-4</i>	Upstream	2328
Tsp_08234	4	ES570299	Non-Operonic	<i>C01A2.5</i>	Upstream	1696
Tsp_08309	4	ES564762	Non-Operonic	<i>mnat-1</i>	Upstream	2070
Tsp_00048	3	ES570469	Non-Operonic	<i>exos-4.1</i>	Downstream	4532
Tsp_00176	7	BG732141	Non-Operonic	<i>ced-9</i>	Downstream	3666
Tsp_00183	3	ES560678	Non-Operonic	<i>cua-1</i>	Downstream	3792
Tsp_01361	4	ES272901	Non-Operonic	<i>mms-19</i>	Downstream	5533
Tsp_02307	3	ES562862	Non-Operonic	<i>Y50D4C.3</i>	Downstream	5559
Tsp_02829	1	BG322177	Non-Operonic	<i>gdi-1</i>	Downstream	4580
Tsp_04483	7	ES273272	Non-Operonic	<i>cts-1</i>	Downstream	3660
Tsp_06831	2	ES566483	Non-Operonic	<i>dpy-11</i>	Downstream	5120
Tsp_09589	2	ES566654	Non-Operonic	<i>F45F2.9</i>	Downstream	5196
Tsp_11375	4	ES561664	Non-Operonic	<i>rba-1</i>	Downstream	1552
Tsp_08545	7	ES273001	Non-Operonic	<i>drh-3</i>	Downstream	1979
Tsp_00870	3	ES569690	Non-Operonic	<i>smr-1</i>	Ambiguous	
Tsp_03751	3	ES569638	Non-Operonic	<i>rps-4</i>	Ambiguous	
Tsp_08048	3	ES563226	Non-Operonic	<i>B0491.5</i>	Ambiguous	
Tsp_09957	2	ES567293	Non-Operonic	<i>rpl-7</i>	Ambiguous	
Tsp_05009	3	ES569586	Non-Operonic	<i>ncs-7</i>	Non-Operonic	
Tsp_10058	3	ES569899	Non-Operonic	<i>cec-1</i>	Non-Operonic	
Tsp_00219	8	ES563383	Non-Operonic	<i>hcf-1</i>	Non-Operonic	

Tsp_00352	2	ES569781	Non-Operonic	<i>C16E9.2</i>	Non-Operonic	
Tsp_01044	1	ES273428	Non-Operonic	<i>kin-2</i>	Non-Operonic	
Tsp_01088	4	BG322014	Non-Operonic	<i>slf-1</i>	Non-Operonic	
Tsp_02150	1	ES273162	Non-Operonic	<i>flp-14</i>	Non-Operonic	
Tsp_03394	2	ES570427	Non-Operonic	<i>T03F1.12</i>	Non-Operonic	
Tsp_04026	1	ES561895	Non-Operonic	<i>pcyt-1</i>	Non-Operonic	
Tsp_05440	4	ES568782	Non-Operonic	<i>sulp-4</i>	Non-Operonic	
Tsp_05646	3	ES569812	Non-Operonic	<i>K09E9.1</i>	Non-Operonic	
Tsp_05819	4	ES570798	Non-Operonic	<i>C25F6.7</i>	Non-Operonic	
Tsp_05837	3	ES561164	Non-Operonic	<i>ZK795.1</i>	Non-Operonic	
Tsp_05990	1	ES273061	Non-Operonic	<i>ppm-1</i>	Non-Operonic	
Tsp_06158	12	ES567235	Non-Operonic	<i>B0272.3</i>	Non-Operonic	
Tsp_06305	11	ES566685	Non-Operonic	<i>ife-3</i>	Non-Operonic	
Tsp_06353	1	ES569025	Non-Operonic	<i>ppk-2</i>	Non-Operonic	
Tsp_07030	3	ES569609	Non-Operonic	<i>mecr-1</i>	Non-Operonic	
Tsp_07118	2	ES563420	Non-Operonic	<i>Y82E9BR.3</i>	Non-Operonic	
Tsp_07593	7	ES567967	Non-Operonic	<i>ZK809.3</i>	Non-Operonic	
Tsp_07622	12	ES563060	Non-Operonic	<i>F17A9.2</i>	Non-Operonic	
Tsp_07825	2	ES569035	Non-Operonic	<i>lipf-6</i>	Non-Operonic	
Tsp_07989	6	BG520790	Non-Operonic	<i>pck-1</i>	Non-Operonic	
Tsp_08023	1	ES563253	Non-Operonic	<i>Y54E10BR.2</i>	Non-Operonic	
Tsp_09739	3	ES273513	Non-Operonic	<i>rab-10</i>	Non-Operonic	
Tsp_10501	12	ES564714	Non-Operonic	<i>F58F12.1</i>	Non-Operonic	
Tsp_10907	2	ES561739	Non-Operonic	<i>lin-19</i>	Non-Operonic	
Tsp_11867	7	BG302133	Non-Operonic	<i>hmg-1.2</i>	Non-Operonic	
Tsp_12507	2	ES568668	Non-Operonic	<i>rps-14</i>	Non-Operonic	
No match	2	BG354854	Unknown	<i>dcn-1</i>	Upstream	3172
Ambiguous	2	BQ693036	Unknown	<i>icln-1</i>	Downstream	4324
Ambiguous	2	ES561812	Unknown	<i>T02C12.3</i>	Downstream	3148
Ambiguous	3	ES563178	Unknown	<i>pkc-1</i>	Downstream	5312
No match	2	ES568194	Unknown	<i>gst-1</i>	Downstream	3572
Ambiguous	3	ES566071	Unknown	<i>F53B1.8</i>	Non-Operonic	
Tsp_00916	1	ES564849	Unknown	<i>clc-3</i>	Non-Operonic	
Tsp_09335	4	ES566864	Unknown	<i>plp-1</i>	Non-Operonic	
No match	12	ES561871	Unknown	<i>glr-1</i>	Non-Operonic	
No match	3	BQ693326	Unknown	<i>K09E9.4</i>	Non-Operonic	
No match	2	ES567398	Unknown	<i>cdc-42</i>	Non-Operonic	

¹ In some cases we found the EST sequence matched to multiple genes (Ambiguous), or we failed to find a match in *T. spiralis* draft genome (No match) and we were thus unable to determine the operonic status.

² A number of genes are present on short contigs, so we were unable to determine its operonic status.

Table S2 The 5' ends of *T. spiralis* and *T. muris* mRNAs identified by 5' RACE. Initiation codons are underlined and identical residues present in all sequences are indicated (*). For the ends of non *trans*-spliced mRNAs the distances from the initiation codon position in the cDNA (with A in ATG set as 0) are indicated. 'n' indicates how many times the sequence was found.

<i>T. spiralis</i> mRNA 5' ends			n	GenBank
<i>zgpa-1</i>	genome	AACTGACTACAGTTTGACAGTAATCCTATTCTGTGATGCAGTCGTGTAATTATTTGCTTCACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG		
	-205	-----GATGCAGTCGTGTAATTATTTGCTTCACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG	1	
	-214	-----CTATTCTGTGATGCAGTCGTGTAATTATTTGCTTCACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG	1	
	-179	-----ACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG	1	
	-192	-----AATTATTTGCTTCACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG	1	
	-229	-----AGTTTGACAGTAATCCTATTCTGTGATGCAGTCGTGTAATTATTTGCTTCACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG	1	
	-228	-----GTTTGACAGTAATCCTATTCTGTGATGCAGTCGTGTAATTATTTGCTTCACGCGTATTTTTGAAGATTAATTGAACTGTGTACATTAATGTG *****	1	KF442420
<i>dif-1</i>	genome	TTGTTTCGGCGAATTGTATACAGATAAATATGGTGGAAAATGATTTCGAAATCGATCGATTGGTGAAACCACGTC AAGATACGGATCCACTGAGAAATTTCTTAGC		
	<i>Tsp</i> -SL10/ <i>dif-1</i>	GGTAATATTTACTGAATTCAGATAAATATGGTGGAAAATGATTTCGAAATCGATCGATTGGTGAAACCACGTC AAGATACGGATCCACTGAGAAATTTCTTAGC * * * *****	5	KF442421
<i>aph-1</i> <i>trans-spliced</i>	genome	ATCAATATGTATATCTTTTAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT		
	<i>Tsp</i> -SL11/ <i>aph-1</i>	--TACCTTTGAACCCACTTCAAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT	1	KF442422
	<i>Tsp</i> -SL1/ <i>aph-1</i>	-AGGTATTTACCAGATCTAAAAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT	2	KF442423
	<i>Tsp</i> -SL7/ <i>aph-1</i>	-AACCTGCACGACTTGTTCGAAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT	2	KF442423
	<i>Tsp</i> -SL12/ <i>aph-1</i>	ACGAATTTACCGTATTTGTCAAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT	2	KF442425
	<i>Tsp</i> -SL14/ <i>aph-1</i>	ATACCGTTCAATTAATTTGAAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT	1	KF442426
	<i>Tsp</i> -SL9/ <i>aph-1</i>	--AGACGTGGTTATTTATTGAAGGATTATATTAATAATCGCCTTGAGAAGCATTCCGTCCGAGTACACGTTTCAACTATGGGATTTCAGAATTCACAGGATAT *****	1	KF442427
<i>non-trans-spliced</i>	genome	ATCTTACTGTAGCTTTATTTGAACGAAAACGTGTAATGTTATTGTTGCTTTCTTTTTAGCTTTTTACTTTTTATGTAAAATCTGTGTTTGTGGTAAAAAGTCA		
	-201	-----GAACGAAAACGTGTAATGTTATTGTTGCTTTCTTTTTAGCTTTTTACTTTTTATGTAAAATCTGTGTTTGTGGTAAAAAGTCA	2	
	-205	-----ATTTGAACGAAAACGTGTAATGTTATTGTTGCTTTCTTTTTAGCTTTTTACTTTTTATGTAAAATCTGTGTTTGTGGTAAAAAGTCA	1	
-211	-----GCTTTATTTGAACGAAAACGTGTAATGTTATTGTTGCTTTCTTTTTAGCTTTTTACTTTTTATGTAAAATCTGTGTTTGTGGTAAAAAGTCA *****	1	KF442428	
<i>cpt-2</i>	genome	TCATTTTACTGTCTGCACATTATGTTTACATAGTGTAAAGAGCTGTTGTATTTTGAGAAACAAATTTGTGTGCTTGAATTGCTTGCATTCTTTTAATTCGAACACATT		
	-204	-----ATTATTGTTTACATAGTGTAAAGAGCTGTTGTATTTTGAGAAACAAATTTGTGTGCTTGAATTGCTTGCATTCTTTTAATTCGAACACATT *****	9	KF442419
<i>nuaf-3</i>	genome	TTATTATATCAATAATTTCTAGTTTGGCGTAAATGTTGGCTATTAGACGAGCTTTACTTAAAAATCGCTGCTTGTTTAATCAAATAATGCAATCAGTGGAAAG		
	<i>Tsp</i> -SL12/ <i>nuaf-3</i>	ACGAATTTACCGTATTTGTCAAGTTTGGCGTAAATGTTGGCTATTAGACGAGCTTTACTTAAAAATCGCTGCTTGTTTAATCAAATAATGCAATCAGTGGAAAG * * * * * *****	7	KF442418

<i>T. muris</i> mRNA 5' ends				
<i>zgap-1</i>	genome	GATCACTTACGAATTACTGTCTGCTGATGCCAGATGCGGGCAGCGCGTATGCTGCGTGAAGTCTTTCGCTCAGTTTTCCATGGATGAAAACGAATATCCCGCGGGCT		
	-232 -254	-----AGATGCGGGCAGCGCGTATGCTGCGTGAAGTCTTTCGCTCAGTTTTCCATGGATGAAAACGAATATCCCGCGGGCT -----ATTACTGTCTGCTGATGCCAGATGCGGGCAGCGCGTATGCCGCGTGAAGTCTTTCGCTCAATTTTTCCATGGATGAAAACGAATATCCCGCGGGCT *****	3 3	KF442429
<i>dif-1</i>	genome	CGAATGCCAATTGATTTGACAGGAGTTCGCAAATGGAGGATGAGGAAGTGGTTCCAAGTGAACATATGCGAACGGATCCCTTGAAGAACTTTGTCTGCTGGTGGC		
	<i>Tmu-SL12/dif-1</i>	-GTTAAATTTACCCCTCAAAGGAGTTCGCAAATGGAGGATGAGGAAGTGGTTCCAAGTGAACATATGCGAACGGATCCCTTGAAGAACTTTGTCTGCTGGTGGC	1	KF442430
	<i>Tmu-SL1/dif-1</i> <i>Tmu-SL4/dif-1</i>	GGTTATTTACCCCTGTTAACAAGGAGTTCGCAAATGGAGGATGAGGAAGTGGTTCCAAGTGAACATATGCGAACGGATCCCTTGAAGAACTTTGTCTGCTGGTGGC GGTTAAGTTTACCCAATTGAAGGAGTTCGCAAATGGAGGATGAGGAAGTGGTTCCAAGTGAACATATGCGAACGGATCCCTTGAAGAACTTTGTCTGCTGGTGGC * *****	1 1 2	KF442431 KF442432
<i>aph-1</i> trans-spliced	genome	GCAATGGGCCCTGTTTGATCGTTTGGCCATAATCTTCAGCGTGGAGAATCATGGGCCCTTTCGGAATTCGTCGGTTGCAGCTTGATAGCCTTTGGACCTTCTTTG		
	<i>Tmu-SL13/aph-1</i> <i>Tmu-SL6/aph-1</i>	-----GGTATTTACCCAACGTTGACTGCGTGGAGAATCATGGGCCCTTTCGGAATTCGTCGGTTGCAGCTTGATAGCCTTTGGACCTTCTTTG -----GGTTAATTTACCCAATTTCAAAGCGTGGAGAATCATGGGCCCTTTCGGAATTCGTCGGTTGCAGCTTGATAGCCTTTGGACCTTCTTTG * ** *	1 1	KF442433 KF442434
non-trans-spliced	genome	ACATAACCTGCGAATCACTTTATTGTCCCTCGAGTGGGCGATTGCTGATCTTTACCTGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC		
	-245	-----GATTGCTGATCTTTACCTGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC	1	KF442435
	-273	-----GAATCACTTTATTGTCCCTCGAGTGGGCGATTGCTGATCTTTACATGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC	1	
	-275	-----GCGAATCACTTTATTGTCCCTCGAGTGGGCGATTGCTGATCTTTACCTGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC	1	
	-269	-----CACTTTATTGTCCCTCGAGTGGGCGATTGCTGATCTTTACCTGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC	1	
	-229	-----ACCTGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC	3	
-228	-----CCTGATTGTCTGGCGTATCACTCACGGGACAGCGTCCCTTCGATGTACATGATC * *****	1		
<i>cpt-2</i>	genome	AGTCAGTTCATCGTGCTATAAATTGTAGCGACGCGACAACAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC		
		AGTCAGTTCATCGTGCTATAAATTGTAGCGACGCGACAACAATGCTCCGTTGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	KF768019
		---AGTTCATCGTGCTATAAATTGTAACGACGCGACAACAATGCTCCGTCCTTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	4	
		---AGTTCATCGTGCTATAAATTGTAGCGACGCGACAACAATGCTCCATCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	
		-----GTGCTATAAATTGTAGCGACGCGACAACAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	
		-----ATATAAATTGTAGCGACGCGACAACAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	
		-----ATAAATTGTAGCGACGCGACAACAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	
		-----ATAAATTGTAGCGACGCGACAACAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	
		-----AAATTGTAGCGACGCGACAACAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC	1	
		-----AACAAATGCTCCGTCGTTGCCCGTCAGCTTGGCTTGGCTGCTGGTTTTTGGGCAGCATGTTCAAGTTGGGAC *****	1	

<i>nuaf-3</i>	genome	GCAACCATCAATTTGGTTTTAGCATGAACGTTAGCAGACGTGCGTCTTTCGTTATTTCGACAACCTTCGTTGCGTCAAGCGTTTCGCTTCGACAGTGCTCTTCCGC		
	<i>Tmu-SL01/nuaf-3</i>	GGTTATTTACCCCTGTTAACAACCATGAACGTTAGCAGACGTGCGTCTTTCGTTATTTCGACAACCTTCGTTGCGTCAAGCGTTTCGCTTCGACAGTGCTCTTCCGC	4	KF511776KF 511777
	<i>Tmu-SL12/nuaf-3</i>	GGTTAAATTTACCCCTCAAAGCATGAACGTTAGCAGACGTGCGTCTTTCGTTATTTCGACAACCTTCGTTGCGTCAAGCGTTTCGCTTCGACAGTGCTCTTCCGC * * *****	9	

Table S3 Primers used for the amplification of 5'RACE products.

mRNA 5' ends	Primers
<i>Tsp-cpt-2</i>	GGCAAGGCAGCTCTGTATCGGAA and Gene Racer 5' primer (CGACTGGAGCACGAGGACTGA).
<i>Tsp-nuaf-3</i>	1 st PCR: CACTTCCATGATAATGCAGCCTGTGG and Gene Racer 5' primer. 2 nd PCR: GCCAGTGATTCATTCCAAGCC and Gene Racer 5' nested primer (GGACTGACATGGACTGAAGGAGTA).
<i>Tsp-zgpa-1</i>	1 st PCR: AGTTCTGTCGCTTCCAACAACGCAT and Gene Racer 5' nested primer. 2 nd PCR: CGCTTCCAACAACGCATCTATACTAC and Gene Racer 5' nested primer.
<i>Tsp-dif-1</i>	1 st PCR: TGCTGCAGCCGAAGAAGTACAACGC and Gene Racer 5' nested primer. 2 nd PCR: CCGAAGAAGTACAACGCAAACAGC and Gene Racer 5' nested primer.
<i>Tsp-aph-1</i>	1 st PCR: ACTTCTTGTAAGAAGACAACAAGAAAACGG and Gene Racer 5' nested primer. 2 nd PCR: GATCGTGCATGATCACGCAGAGA and Gene Racer 5' nested primer.
<i>Tmu-nuaf-3</i>	1 st PCR: GCCATGATGTCAACTCAATC and Gene Racer 5' nested primer. 2 nd PCR: CCGACGACAAATACGCCATTGG and Gene Racer 5' nested primer.
<i>Tmu-zgpa-1</i>	1 st PCR: GGAATAGCTTAACCAACAGACCCTTT and Gene Racer 5' nested primer. 2 nd PCR: AGTAAACAAGATACAAGCCACTCACC and Gene Racer 5' nested primer.
<i>Tmu-dif-1</i>	1 st PCR: TTCGTCTGGGTGCCTTTGTTGGAG and Gene Racer 5' nested primer. 2 nd PCR: GAAAATACAAGGCGAATAGCGGAGC and GR 5' nested primer.
<i>Tmu-aph-1</i>	1 st PCR: ATGAGTGCGGTAGGCCAATAGTTCC and Gene Racer 5' nested primer. 2 nd PCR: AGAAGGTTTAGCAGTGCAACTACGCC and Gene Racer 5' nested primer.

Table S4 Primers used for the detection of processing intermediates of polycistronic transcripts by PCR.

Transcripts	Primers
<i>Tsp-cpt-2~nuaf-3</i>	GAAGCATGTTCCAAGCAGCATTC GCCAGTGATTCATTCCAAGCC
<i>Tsp-zgpa-1~dif-1</i>	AAGCAAAATTGGCGAAAAGGACCGAAG TGCTGCAGCCGAAGAAGTACAACGC
<i>Tsp-dif-1~aph-1</i> (1 st PCR)	TGCTGGAGCTTTGTCAGGTATGATG ACTTCTTGTAAGAAGACAAACAAGAAAACGG
<i>Tsp-dif-1~aph-1</i> (2 nd PCR)	CGGGCGAGAGGATCAAATGC GATCGTGCATGATCACGCAGAGA
<i>Tmu-zgpa-1~dif-1</i> (1 st PCR)	TTGTTGCGGGCAAGGCTTAG TTCGTCTGGGTGCCTTTGTTGGAG
<i>Tmu-zgpa-1~dif-1</i> ((2 nd PCR)	TGATCGAAAGTGTAGAGCGGATACGG GAAAATACAAGGCGAATAGCGGAGC
<i>Tmu-dif-1~aph-1</i> (1 st PCR)	TGCTGGAGCTTTGTCAGGTATGATG ATGAGTGCGGTAGGCCAATAGTTCC
<i>Tmu-dif-1~aph-1</i> (2 nd PCR)	CGGGCGAGAGGATCAAATGC AGAAGGTTTAGCAGTGCAACTACGCC