

1 SLIDR and SLOPPR: Flexible identification of spliced leader
2 *trans*-splicing and prediction of eukaryotic operons from
3 RNA-Seq data

4 Marius A. Wenzel^{1*}, Berndt Mueller² and Jonathan Pettitt²

5 December 18, 2020

6 ¹ School of Biological Sciences, University of Aberdeen, Zoology Building, Tillydrone Ave, Aberdeen
7 AB24 2TZ, UK;

8 ² School of Medicine, Medical Sciences and Nutrition, University of Aberdeen, Institute of Medical
9 Sciences, Foresterhill, Aberdeen, AB25 2ZD, UK

10 *corresponding author

11 **Abstract**

12 Background: Spliced leader (SL) *trans*-splicing replaces the 5' ends of pre-mRNAs with the spliced
13 leader, an exon derived from a specialised non-coding RNA originating from a different genomic
14 location. This process is essential for resolving polycistronic pre-mRNAs produced by eukaryotic
15 operons into monocistronic transcripts. SL *trans*-splicing and operons have independently evolved
16 multiple times throughout Eukarya, but our understanding of these phenomena is limited to only a
17 few well-characterised organisms, most notably *C. elegans* and trypanosomes. The primary barrier to
18 systematic discovery and characterisation of SL *trans*-splicing and operons is the lack of computational
19 tools for exploiting the surge of transcriptomic and genomic resources for a wide range of eukaryotes.

20 Results: Here we present two novel pipelines that automate the discovery of SLs and the prediction
21 of operons in eukaryotic genomes from RNA-Seq data. SLIDR assembles putative SLs from 5' read
22 tails present after read alignment to a reference genome or transcriptome, which are then verified by
23 interrogation of sequence motifs expected in *bona fide* SL RNA molecules. SLOPPR identifies RNA-
24 Seq reads that contain a given 5' SL sequence, quantifies genome-wide SL *trans*-splicing events and
25 predicts operons via distinct patterns of SL *trans*-splicing events across adjacent genes. We tested
26 both pipelines with organisms known to carry out SL *trans*-splicing and organise their genes into
27 operons, and demonstrate that 1) SLIDR correctly identifies known SLs and often discovers novel
28 SL variants; 2) SLOPPR correctly identifies functionally specialised SLs, correctly predicts known
29 operons and detects plausible novel operons.

30 Conclusions: SLIDR and SLOPPR are flexible tools that will accelerate research into the evolu-
31 tionary dynamics of SL *trans*-splicing and operons throughout Eukarya, and improve gene discovery
32 and annotation for a wide-range of eukaryotic genomes. Both pipelines are implemented in Bash and
33 R and are built upon readily available software commonly installed on most bioinformatics servers.
34 Biological insight can be gleaned even from sparse, low-coverage datasets, implying that an untapped
35 wealth of information can be derived from existing RNA-Seq datasets as well as from novel full-isoform
36 sequencing protocols as they become more widely available.

37 Keywords: spliced-leader *trans*-splicing, eukaryotic operons, polycistronic RNA processing, RNA-Seq,
38 genome annotation, chimeric reads, 5' UTR

39 Background

40 Spliced leader (SL) *trans*-splicing is a eukaryotic post-transcriptional RNA modification whereby the 5'
41 end of a pre-mRNA receives a short “leader” exon from a non-coding RNA molecule that originates
42 from elsewhere in the genome [1, 2]. This mechanism was first discovered in trypanosomes [3] and
43 has received much attention as a potential target for diagnosis and control of a range of medically
44 and agriculturally important pathogens [1, 4, 5]. SL *trans*-splicing is broadly distributed among many
45 eukaryotic groups, for example euglenozoans, dinoflagellates, cnidarians, ctenophores, platyhelminthes,
46 tunicates and nematodes, but is absent from vertebrates, insects, plants and fungi [2]. Its phylogenetic
47 distribution and rich molecular diversity suggest that it has evolved independently many times throughout
48 eukaryote evolution [6–9].

49 One clear biological function of SL *trans*-splicing is the processing of polycistronic pre-mRNAs generated
50 by eukaryotic operons [2]. In contrast to prokaryotes, where such transcripts can be translated imme-
51 diately as they are transcribed, a key complication for eukaryotic operons is that nuclear polycistronic
52 transcripts must be resolved into independent, 5'-capped monocistronic transcripts for translation in the
53 cytoplasm [10]. The *trans*-splicing machinery coordinates cleavage of polycistronic pre-mRNA and pro-
54 vides the essential cap to initially un-capped pre-mRNAs [11, 12]. This process is best characterised in
55 the nematodes, largely, but not exclusively due to work on *C. elegans*, which possesses two types of SL
56 [13]: SL1, which is added to mRNAs derived from the first gene in operons and monocistronic genes;
57 and SL2, which is added to mRNAs arising from genes downstream in operons and thus specialises in
58 resolving polycistronic pre-mRNAs [11–13].

59 The same SL2-type specialisation of some SLs for resolving downstream genes in operons has been re-
60 ported in other nematodes [14–19], but is not seen in other eukaryotic groups. For example, platyhelminth
61 *Schistosoma mansoni* and the tunicates *Ciona intestinalis* and *Oikopleura dioica* each possess only a sin-
62 gle SL, which is used to resolve polycistronic RNAs but is also added to monocistronic transcripts [20–22].
63 Similarly, the chaetognath *Spadella cephaloptera* and the cnidarian *Hydra vulgaris* splice a diverse set of
64 SLs to both monocistronic and polycistronic transcripts [23, 24]. Remarkably, all protein-coding genes
65 in trypanosomes are transcribed as polycistronic RNAs and resolved using a single SL, making SL *trans*-
66 splicing an obligatory process for all mRNAs [25]. In contrast, dinoflagellates use SL *trans*-splicing for
67 all nuclear mRNAs, but only a subset of genes are organised as polycistrons [26, 27]. Although SL
68 *trans*-splicing also occurs in many other organisms including rotifers, copepods, amphipods, ctenophores,
69 cryptomonads and hexactinellid sponges, operons and polycistronic RNAs have not been reported in
70 these groups [7, 8, 28, 29].

71 All these examples illustrate a rich diversity in the SL *trans*-splicing machinery and its role in facilitating
72 polycistronic gene expression and broader RNA processing. A major barrier in dissecting the evolutionary
73 history of these phenomena is the difficulty in systematically quantifying SL *trans*-splicing events. Ident-
74 ifying the full SL repertoire would traditionally require laborious low-throughput cloning-based Sanger
75 sequencing of the 5' ends of mRNAs [e.g., 16, 30]. High-throughput RNA-Seq data is an attractive al-
76 ternative resource that may often already exist for the focal organism. Some studies have demonstrated
77 that SLs can, in principle, be identified from overrepresented 5' tails extracted directly from RNA-Seq
78 reads [31, 32]. The recent SLFinder pipeline uses overrepresented *k*-mers at transcript ends as guides
79 (“hooks”) for annotating potential SL genes in genome assemblies [33]. SLFinder can detect known SLs

80 in several model organisms, but since it does not take into account the known functionally important
81 sequence features of SL RNAs, its outputs may be noisy, incomplete and swamped by pseudogenes [33].

82 Once an SL sequence repertoire has been established, the next steps are to quantify SL *trans*-splicing
83 events genome-wide and to establish functional links between these events and operonic gene organisation.
84 Several studies have demonstrated that 5' information from RNA-Seq reads can be exploited to quantify
85 SL *trans*-splicing events [28, 34], and the SL-QUANT pipeline has automated this task for *C. elegans*
86 and other nematodes [35]. Similarly, it has been demonstrated in the nematodes *Pristionchus pacificus*
87 and *Trichinella spiralis* that genome-wide patterns of SL *trans*-splicing events can be exploited to predict
88 novel operons from SL splicing ratios [18, 19]. However, no software exists to implement these prediction
89 strategies and render them universally applicable beyond the Nematoda.

90 Here we present two fully-automated pipelines that address all these shortcomings and present a unified
91 and universal approach to examining SL *trans*-splicing and operonic gene organisation from RNA-Seq
92 data in any eukaryotic organism. First, SLIDR is a more efficient, sensitive and specific alternative to
93 SLFinder, implementing fully customisable and scalable *de novo* discovery of SLs and associated SL RNA
94 genes. Second, SLOPPR implements a generalised and more flexible solution to quantifying genome-wide
95 SL *trans*-splicing events than SL-QUANT. Uniquely, it provides algorithms for inference of SL sub-
96 functionalisation and customisable prediction of operonic gene organisation. Both pipelines can process
97 single-end or paired-end data from multiple libraries that may differ in strandedness and read config-
98 uration, thus allowing for flexible high-throughput processing of large RNA-Seq or EST datasets from
99 multiple sources. These pipelines present a complete one-stop solution for systematically investigating
100 SL *trans*-splicing and operon organisation in all eukaryotes.

101 Implementation

102 SLIDR: Spliced leader identification from RNA-Seq

103 SLIDR extracts evidence of SLs directly from RNA-Seq reads that contain unmapped 5' tails after
104 alignment to a genome or transcriptome reference (broadly similar to 32). Unlike other methods, SLIDR
105 then implements several optional plausibility checks based on functional nucleotide motifs in the SL RNA
106 molecule, i.e., splice donor and acceptor sites, *Sm* binding motifs and a number of stem loops. These
107 features are expected to be present due to shared evolutionary ancestry of SL RNAs with the snRNAs
108 involved in intron removal by *cis*-splicing [6, 36]. For each plausible SL sequence, expressed SL RNA
109 genes and SL *trans*-spliced genes are annotated in the reference, providing a means of manual inspection
110 of the SL *trans*-splicing landscape if desired.

111 RNA-Seq reads are aligned to the genome or transcriptome reference using HISAT2 [37] or BOWTIE2 [38]
112 in local alignment mode. Since soft-clipped read tails must be long enough to capture full-length SLs
113 (typically about 22 bp in nematodes), a relaxed alignment scoring function is implemented that allows
114 for up to 25 bp tails in a 75 bp read and can easily be customised by the user to accommodate more
115 extreme SL lengths, for example 16 bp in *Ciona intestinalis* [39] or 46 bp in *Hydra vulgaris* [24]. Tails
116 from the read end corresponding to the 5' end of the transcript (inferred from library strandedness) are
117 extracted using SAMTOOLS [40] and dereplicated, 3'-aligned and clustered at 100% sequence identity using
118 VSEARCH [41]. Each cluster thus represents a single putative SL, comprising a collection of 3'-identical
119 read tails of varying length that only differ in their 5' extent (Figure 1).

120 The cluster centroids (longest sequence in each cluster) are then subjected to a number of functional
121 plausibility checks. The centroids are aligned to the genome or transcriptome reference using BLASTN

122 [42] with 100% sequence identity and a relaxed customisable E-value (e.g., 1) to accommodate short
123 queries. Matches that contain the full 3' end of the centroid are retained and the putative full SL RNA
124 sequence (of customisable length) is extracted from the reference using BEDTOOLS [43]. The SL RNA
125 sequence is then inspected for customisable splice-donor (e.g., GT) and *Sm* binding (e.g., AT{4,6}G) sites
126 [44, 45], and secondary structure stem loops are predicted using RNAFOLD [46]. Default criteria expect
127 the *Sm* motif c. 50 bp downstream of the splice donor site, and one stem loop on each side of the *Sm*
128 binding motif [30, 45, 47]. In the reference sequence immediately upstream of the aligned portion of each
129 RNA-Seq read, a splice acceptor site (e.g., AG) is required, corresponding to the *trans*-splice acceptor
130 site of the gene (Figure 1).

131 The locations of splice donor and splice acceptor sites may not be as expected if the 3' end of the SL
132 and the 3' end of the *trans*-splice acceptor site happen to be identical. In these cases, the RNA-Seq
133 read alignment overextends in 5' direction into the *trans*-splice acceptor site and thus 3' truncates the
134 soft-clipped SL read tail (Figure 1). These missing 3' nucleotides can be reconstructed from surplus
135 nucleotides located between the 3' end of the tail BLASTN match and the splice donor site, and must
136 be identical to those surplus nucleotides located between the 5' read alignment location and the splice
137 acceptor site (Figure 1). Following plausibility checks and 3' reconstruction where necessary, all tail
138 cluster centroids are subjected to another round of 3' alignment and clustering at 100% sequence identity
139 in VSEARCH before final SL consensus construction is carried out in R [48]. Final SLs must be supported
140 by at least two reads and must be spliced to at least two genes that are not located in the immediate
141 vicinity (1 kbp distance) of the SL RNA gene [31].

142 SLOPPR: Spliced leader-informed operon prediction from RNA-Seq

143 SLOPPR is designed as a genome-annotation tool that predicts operons from genome-wide distributions
144 of SL *trans*-splicing events at pre-annotated genes. RNA-Seq reads that contain evidence of 5' SLs are
145 identified using a sequence-matching approach equivalent to the “sensitive” mode of SL-QUANT [35]. The
146 operon prediction algorithm is built upon the SL1/SL2-type functional specialisation of SLs observed in
147 many nematodes, but is fully customisable to accommodate other relationships between SLs and operonic
148 genes, even when SL specialisation is absent. Unlike previous approaches that have defined operons in
149 various organisms primarily via short intercistronic distances [17, 18, 21, 49], SLOPPR defines operons
150 principally via SL *trans*-splicing patterns and only optionally takes intercistronic distance into account.
151 SLOPPR can also identify and correct gene annotations where operonic genes are incorrectly annotated
152 as a fused single gene (cf., 19), paving the way for *trans*-splicing-aware genome (re-)annotation.

153 RNA-Seq reads containing SLs are identified using a three-tier strategy. Since such reads cannot align
154 end-to-end to the genome because of the *trans*-spliced 5' SL tail, all reads are first aligned end-to-end to
155 the genome reference using HISAT2 [37] and unmapped reads are retained as candidates. If paired-end
156 reads are used, the read corresponding to the 3' end of the transcript (inferred from library strandedness)
157 must be aligned, and the read corresponding to the 5' end of the transcript must be unaligned [35]. The
158 5' ends of the unaligned candidate reads are then screened for overlap with the 3' portion of any number
159 of supplied SL sequences using CUTADAPT [50]. Finally, those reads that align to the genome end-to-
160 end after the SL tail has been trimmed are then quantified against exons and summarised at the gene
161 level using FEATURECOUNTS [51]. Likewise, background expression levels of all genes are obtained from
162 the original end-to-end read alignments and from candidate reads without SL evidence. This screening
163 strategy is carried out for each RNA-Seq library independently, thus allowing for comparisons among
164 biological replicates during analysis (Figure 2A).

165 The nature of the SL *trans*-splicing process means that SLs must only be present at the first exon of

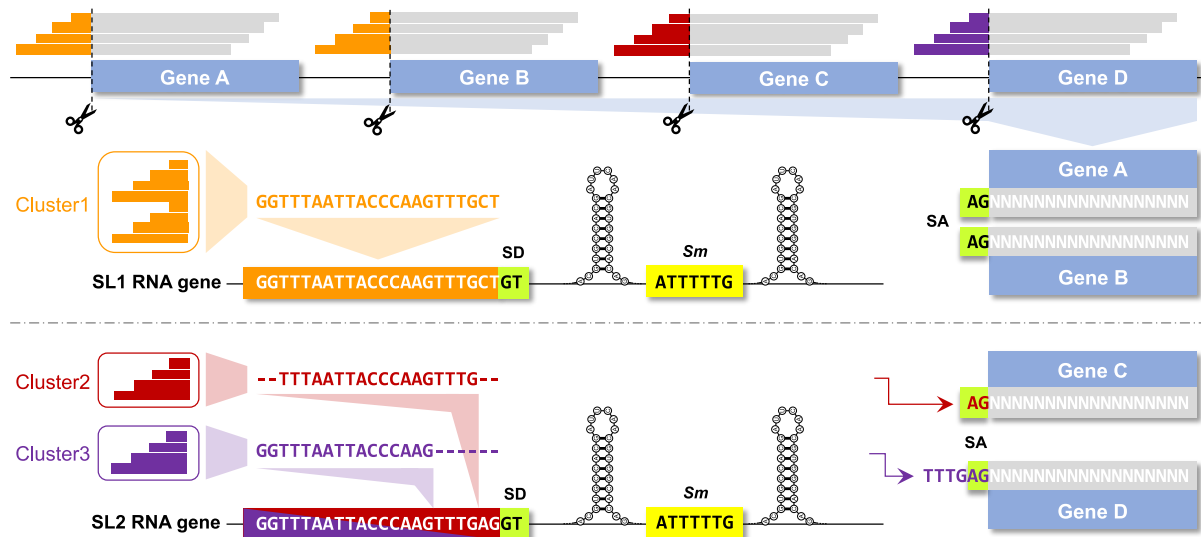


Figure 1: Schematic representation of the SLIDR pipeline (Spliced leader identification from RNA-Seq). Local alignments of reads (grey) to a genomic reference (illustrated by four genes A-D) allow for 5' SL tails to be soft-clipped and extracted (coloured read portions). Similarity clustering of 3' aligned read tails from all genes produces unique consensus SL candidates (cluster centroids), which are required to align to the genomic reference to identify candidate SL RNA genes (illustrated by SL1 and SL2 genes). In SL RNA genes, a splice donor site (SD; for example GT) is expected immediately downstream of the genomic alignment, followed by an *Sm* binding site (for example 5'-ATTTTG-3') bookended by inverted repeats capable of forming stem loops in the RNA transcript. Conversely, the spliced gene requires a splice acceptor site (SA; for example AG) immediately upstream of the 5' read alignment location in the genomic reference. In this illustration, the example SL1 is fully reconstructed from a single read-tail cluster (cluster 1) with GT and AG splice sites in the expected locations (genes A and B). In contrast, the example SL2 highlights how read tails may be 3'-truncated due to overlap with the splice acceptor site (genes C and D) and the upstream *trans*-splice acceptor site sequence at some genes (gene D). These missing nucleotides can be filled in from the *trans*-splice acceptor site region guided by the distance between the 3' tail alignment location and the splice donor site (GT). Note that although cluster 2 is also 5' truncated due to insufficient coverage at gene C, consensus calling with cluster 3 allowed for reconstructing the full SL2 RNA gene.

166 a gene, i.e. the 5' end. Incorrect gene annotations thus become obvious when internal exons receive
167 SL reads. SLOPPR implements an optional gene-correction algorithm that splits gene annotations at
168 exons with distinct SL peaks compared to neighbouring exons (Figure 2B). To obtain exon-based SL
169 counts, genome annotations are converted to GTF using GFFREAD from CUFFLINKS [52], unique exons
170 are extracted using BEDTOOLS [43] and SL reads are quantified with FEATURECOUNTS at the exon level
171 instead of gene level. The peak-finding algorithm is designed to correctly handle reads that may span
172 multiple exons (Figure 2B).

173 The SL read counts obtained from FEATURECOUNTS are normalised against library size using CPM
174 (counts-per-million) against the background gene counts [53]. The normalised SL read-count matrix is
175 then subjected to generalized principal component analysis (PCA) and hierarchical clustering designed
176 for sparse count matrices [54], treating SL read sets as samples and genes as variables. This summary
177 of genome-wide distributions of SL *trans*-splicing events allows for identifying the distinct *trans*-splicing
178 patterns of SL2-type SLs expected from their specialisation to resolve downstream operonic genes. If
179 SL2-type SLs are not known, K-means clustering and linear discriminant analysis are used to assign SLs
180 to one of two synthetic clusters assumed to correspond to SL1-type and SL2-type SLs (Figure 2C). Visual
181 inspection of the clustering results allows the user to determine consistency across biological replicates
182 (if available) and to ascertain functional groups of SLs beyond the SL1/SL2-type groups.

183 Based on the SL clustering results and pre-defined SL1/SL2-type groups (if known), the SL2:SL1 CPM
184 ratio is computed and summarised across all genes that receive both SL types. The operon prediction
185 algorithm is based on finding uninterrupted runs of adjacent genes with SL2-bias, which are designated
186 as downstream operonic genes (Figure 2D). By default, no SL1-type reads at all are allowed, but a more
187 relaxed SL2:SL1 ratio cutoff can be provided. The optimal cutoff is species-specific and could be identified
188 empirically from inspecting the distribution of SL2:SL1 read ratios or from observed read ratios at known
189 operonic genes [19]. After tracts of SL2-biased downstream operonic genes have been designated, each
190 tract can, optionally, receive an additional upstream operonic gene that shows SL1-type bias or absence
191 of SL *trans*-splicing (Figure 2D).

192 Finally, intergenic distances among the predicted operonic genes are computed and compared to
193 genome-wide intergenic distances to diagnose tight physical clustering of operonic genes (Figure 2E).
194 These distances are obtained from the boundaries of consecutive “gene” GFF annotation entries, so
195 their accuracy depends entirely on the provided genome annotations, which should ideally define gene
196 boundaries by poly(A) and *trans*-splice acceptor sites. If desired, operon prediction can take intergenic
197 distances into account, either via a user-supplied distance cutoff or via an automatic K-means clustering
198 method that splits the genome-wide distribution of intergenic distances into two groups, corresponding
199 to tight gene clusters (potential operons) and non-operonic genes. As such, by manually specifying
200 SL1/SL2-type SLs, SL2:SL1 ratio cutoff, upstream gene designation and intergenic distance cutoff, a
201 large gamut of relationships between SLs and operonic genes can be explored, even in situations where
202 no subfunctionalisation of SLs for operon resolution exists, for example in kinetoplastids or tunicates
203 [20–22].

204 Results and Discussion

205 Validation of SLIDR in nematodes

206 In order to assess the performance of SLIDR in identifying SL RNAs, we validated and benchmarked the
207 pipeline in several nematodes where reference genome assemblies are available and the SL repertoire is

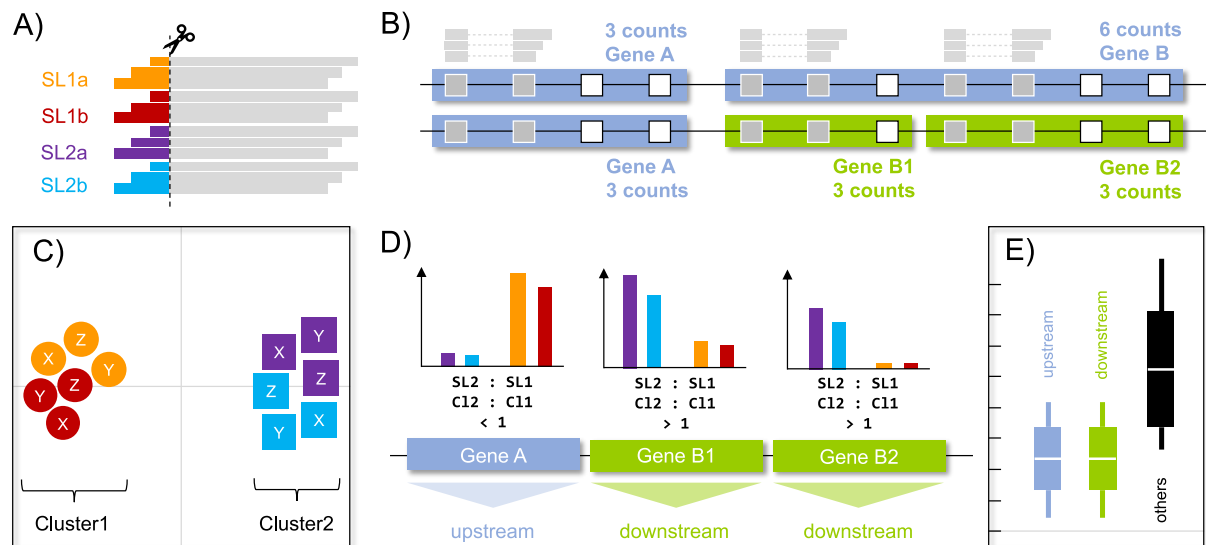


Figure 2: Schematic representation of the SLOPPR pipeline (Spliced leader-informed operon prediction from RNA-Seq). A) Spliced leader tails (example: SL1a, SL1b, SL2a and SL2b) are identified and trimmed from the 5' end of reads that correspond to the 5' end of transcripts. B) Trimmed reads are aligned to the genome, quantified against exons (squares; grey: covered; white: not covered) and counts are summarised by gene (example: two genes A and B). Incorrect gene annotations (fused operonic genes) can optionally be identified and corrected via SL reads at internal exons (example: Gene B is split into B1 and B2). C) SL read sets from multiple libraries (example: X, Y and Z) are ordinated via PCA on genome-wide read counts and grouped into two clusters (K-means clustering) expected to correspond to SL1 (circles) and SL2-type (squares) subfunctionalisation. D) SL2:SL1 read ratios are computed between pre-defined SL groups (SL1, SL2) or inferred clusters (C11, C12). Operons are predicted via tracts of genes receiving SL2 bias (downstream operonic genes) plus an optional upstream gene receiving either an SL1 bias or no SLs at all. E) Intercistronic distances among predicted operons are expected to be reduced compared to intergenic distances among non-operonic genes (others). Operon predictions can optionally be filtered by intercistronic distance using a user-supplied or inferred optimal cutoff.

208 well characterised: *Caenorhabditis elegans*, *Caenorhabditis briggsae*, *Pristionchus pacificus*, *Meloidogyne*
209 *hapla*, *Trichinella spiralis* and *Trichuris muris*. In all cases, we demonstrate that SLIDR detects all
210 known SLs and often discovers novel SL variants (Table 1, Supplemental Table S1). We also provide
211 a proof-of-principle of the transcriptome-mode of SLIDR in *C. elegans*, where the transcriptome is well
212 resolved, and *Prionchulus punctatus*, where no curated reference transcriptome exists. All RNA-Seq data
213 were retrieved from public sources (NCBI SRA or ENA). Illumina adapters (5'-AGATCGGAAGAGC-3')
214 and poor-quality bases (phred 20) were trimmed from all datasets with TRIM_GALORE 0.6.4 [55]. All
215 SLIDR runs used the default splice donor (GT) and splice acceptor (AG) sites and default parameters
216 unless otherwise specified.

217 *Caenorhabditis elegans*

218 *C. elegans* possesses the best understood repertoire of nematode SLs, comprising two types, SL1 and
219 SL2, both of which are encoded by multi-copy gene families with well-described SL RNA structures and
220 *Sm* binding motifs [11–13, 53]. SL *trans*-splicing affects up to 84 % of genes [34, 53], which makes *C.*
221 *elegans* SLs an appropriate benchmark repertoire that any SL detection pipeline should fully resolve.

222 We downloaded the genome assembly GCF_000224145.3 and six out of 24 unstranded 2x100 bp datasets
223 from bioproject PRJNA270896 (SRR1727796–SRR1727801). These are the same datasets used by Calvelo
224 et al. [33] to benchmark their SLFinder pipeline. We ran SLIDR with parameter -S '{40,55}AC?T{4,6}G'
225 and compared the identified SL RNA genes with the reference gene annotations using BEDTOOLS INTER-
226 SECT 2.28.0 [43]. SLIDR identified 5,447 reads that were assembled into a single full-length SL1 sequence
227 and complete 'donatrons' (the intron-like portions of the SL RNAs) were detected for all 10 functional
228 *sls-1* genes (*sls-1.1* and *sls-1.5* are incorrectly annotated 5'-truncated pseudogenes). SLIDR also detected
229 all 11 SL2 sequence variants corresponding to 18 out of 19 annotated *sls-2* genes from as few as eleven
230 reads (Supplemental Table S1). Two reads aligned to the pseudogene *sls-2.19* but dropped out because
231 no splice acceptor sites were present at their gene targets. Thus, SLIDR detected all 18 functional *C.*
232 *elegans* SL2 copies and correctly omitted pseudogenes from its output (Supplemental Table S1).

233 These results are clearer than those reported by SLfinder using the same data [33]. Although SLFinder
234 detected the full SL1 sequence in the transcriptome data and the same *sls-1* genes in the genome as
235 SLIDR, the transcript tails ("hooks") were noisy, did not allow for identifying splice donor sites and
236 required some manual curation [33]. Most strikingly, SLFinder only detected 5 SL2 sequence variants
237 and only 8 out of 19 *sls-2* genes [33]. SLFinder did detect splice donor sites in most of these cases, though
238 the sites were often overlapped by the hook sequences [33]. Thus, SLIDR not only detected all *C. elegans*
239 SL sequences and functional genes with superior sensitivity and specificity from very few reads, but also
240 obviated the requirement by SLFinder to assemble *de novo* transcriptomes at high computational cost
241 (six assemblies were required for the SLFinder analysis).

242 Finally, we tested SLIDR with a transcriptome reference instead of a genome reference. In this situa-
243 tion, SLIDR cannot confirm splice acceptor sites because non-coding sequence regions are not expected
244 to be present in a transcriptome reference; however, if the reference happens to contain SL RNAs it
245 will still be possible to find splice donor and *Sm* binding sites. We used the curated transcriptome
246 GCF_000002985.6_WBcel235 (contains *sls-1* and *sls-2* RNAs) and and three stranded 2x150 bp RNA-
247 Seq libraries from bioproject PRJEB28364 (ERR2756688, ERR2756689, ERR2756736) from NCBI to
248 provide the best possible input data. SLIDR was run with the same *Sm* motif regular expression as
249 above and was able to detect one SL1 variant (*sls-1.10* RNA) and nine SL2 variants from all 18 func-
250 tional *sls-2* RNAs. These results are comparable to those obtained with a genome reference but the
251 inability to confirm splice-acceptor sites means that a large number of false positive candidate SLs were

Table 1: Spliced leader identification (SLIDR) results in seven nematodes and three other eukaryotes. Identifiers for reference genomes/transcriptomes and RNA-Seq libraries are presented alongside numbers of quality-trimmed reads (QC), the *Sm* motif regular expression notation used to filter SLs, numbers of expected SLs detected, numbers of novel SLs identified and numbers of expected SL RNA genes detected. Question marks represent unknown or poorly characterised SL RNA gene numbers.

Species	Reference	Bioproject	QC reads	<i>Sm</i> motif regex	SLs		SL RNA genes	
					(detected expected)	novel SLs	(detected expected)	
<i>Caenorhabditis elegans</i>	GCF_000224145.3	PRJNA270896	150,041,952	.{40,55}AC?T{4,6}G	12 12	-	28 28	
<i>Caenorhabditis elegans</i>	GCF_000002985.6 ¹	PRJEB28364	45,443,137	.{40,55}AC?T{4,6}G	10 12	-	19 28	
<i>Caenorhabditis briggsae</i>	PRJNA10731.WBPS5	PRJNA104933	21,265,790	.{40,55}AT{4,6}G	2 7	1	28 18+?	
<i>Caenorhabditis briggsae</i>	PRJNA10731.WBPS5	PRJNA489172	78,497,095	.{40,55}AT{4,6}G	6 7	2	28 18+?	
<i>Caenorhabditis briggsae</i>	PRJNA10731.WBPS5	PRJNA231838	59,271,467	.{40,55}AT{4,6}G	7 7	-	28 18+?	
<i>Pristionchus pacificus</i>	Hybrid1	SRP039388	81,387,873	.{40,55}[AG]T{4,6}[AG]	6 11	6	219 203	
<i>Pristionchus pacificus</i>	GCA_000180635.3	PRJNA338247	330,856,071	.{40,55}[AG]T{4,6}[AG]	6 11	6	640 203	
<i>Meloidogyne hapla</i>	PRJNA29083.WBPS14	PRJNA229407	169,404,394	.{30,80}AT{4,6}G	5 5	9	15 ?	
<i>Meloidogyne hapla</i>	PRJNA29083.WBPS14	PRJEB14142	62,113,573	.{30,80}AT{4,6}G	5 5	9	15 ?	
<i>Trichinella spiralis</i>	PRJNA12603.WBPS10	PRJNA510020	201,164,867	.{20,50}AT{4,6}G	15 15	-	21 48	
<i>Trichuris muris</i>	PRJEB126.WBPS15	PRJEB1054	115,460,947	.{25,50}AT{4,6}G	13 13	6	20 13	
<i>Priomchulus punctatus</i>	de novo Trinity ¹	PRJEB7585	71,087,651	.{25,60}AC?T{4,6}G	2 6	-	2 ?	
<i>Ciona intestinalis</i>	GCF_000224145.3_KH	PRJNA396771	108,760,969	.{2,25}AGCTTTGG	1 1	-	11 15	
<i>Ciona intestinalis</i>	GCF_000224145.3_KH	PRJNA396771	108,760,969	.{20,50}AT{4,6}G	1 1	3	2 15	
<i>Ciona intestinalis</i>	GCF_000224145.3_KH	PRJNA297221	49,713,049	.{2,25}AGCTTTGG	1 1	-	11 15	
<i>Ciona intestinalis</i>	GCF_000224145.3_KH	PRJNA433724	19,072,215	.{2,25}AGCTTTGG	1 1	-	11 15	
<i>Ciona intestinalis</i>	GCF_000224145.3_KH	PRJNA376667	194,089,029	.{2,25}AGCTTTGG	1 1	2	14 15	
<i>Ciona intestinalis</i>	GCF_000224145.3_KH	PRJNA376667	194,089,029	.{20,50}AT{4,6}G	1 1	1	3 15	
<i>Hydra vulgaris</i>	Hm105	PRJNA497966	125,523,578	.{10,35}[AG]ATTTT[CG][AG]	2 12	1	4 ?	
<i>Hydra vulgaris</i>	Hm105	PRJNA641135	35,234,240	.{10,35}[AG]ATTTT[CG][AG]	3 12	3	5 ?	
<i>Schistosoma mansoni</i>	PRJEA36577.WBPS14	PRJNA225599	38,015,650	.{10,30}AGTTTCTTTGG	1 1	-	112 ?	

¹transcriptome reference

252 reported (103 candidate SLs in total; Supplemental Table S1). Nevertheless, the sensitivity of SLIDR
253 even in transcriptome mode remains higher to that of SLFinder.

254 *Caenorhabditis briggsae*

255 *C. briggsae* is a close relative of *C. elegans* that possesses a similar SL repertoire and shows considerable
256 synteny of operons [15, 17]. *Cbr*-SL1 is encoded by a repetitive gene cluster (about 65 copies; 56) linked
257 to 5S rRNA genes [57], whereas SL2-type SLs are encoded by 18 genes and represent six distinct SL
258 variants (*Cbr*-SL2, *Cbr*-SL3, *Cbr*-SL4, *Cbr*-SL10, *Cbr*-SL13, *Cbr*-SL14), of which four are shared with
259 *C. elegans* [15]. Only 37 % of genes are SL *trans*-spliced [17], compared to 70–84 % in *C. elegans* [34, 53],
260 though this is likely simply a reflection of differential transcriptome read depth.

261 We downloaded genome assembly PRJNA10731.WBPS5 from WormBase and the same two unstranded
262 2x42 bp libraries (SRR440557, SRR440441) from bioproject PRJNA104933 that Uyar et al. [17] used to
263 identify genome-wide SL *trans*-splicing events. These data are particularly difficult to analyse because
264 the short read length is likely to impede identification of the full-length SL. To maximise SL detection,
265 SLIDR was run with parameter -x 2, which would allow a tail of at most 28 bp and leave at least 14 bp
266 for read alignment. Irrespective, SLIDR only detected 5' truncated versions of *Cbr*-SL1 (45 reads spliced
267 to 38 genes) and *Cbr*-SL3 (58 reads spliced to 55 genes).

268 To corroborate these findings with longer and higher-coverage read data, we ran SLIDR on two addi-
269 tional sets of libraries. First, using the high-coverage stranded 2x50 bp library SRR7781208 (bioproject
270 PRJNA489172), SLIDR detected full-length *Cbr*-SL1 (251 reads spliced to 193 genes) and five out of six
271 full-length SL2-type SLs (*Cbr*-SL13 was absent) supported by 8–908 reads and spliced to 8–475 genes
272 (Supplemental Table S1). Second, we used an extensive set of five unstranded 2x75 bp and four stranded
273 2x125 libraries (bioprojects PRJNA231838 and PRJNA306868) in the same SLIDR run. These data
274 supported full-length *Cbr*-SL1 (779 reads spliced to 310 genes) and all six full-length SL2-type SLs (6–72
275 reads spliced to 6–42 genes).

276 Overall, SLIDR detected known SLs with high sensitivity even from relatively unsuitable data in a species
277 with relatively low SL *trans*-splicing frequency. We also note that SLIDR consistently detected five SL
278 RNA gene loci for *Cbr*-SL1 and 23 instead of 18 loci for SL2-type SLs (two additional loci for *Cbr*-SL2
279 and three additional loci for *Cbr*-SL3; 17, 57).

280 *Pristionchus pacificus*

281 *P. pacificus* possesses seven SL1-type (*Ppa*-SL1) and four SL2-type (*Ppa*-SL2) SLs, which are encoded by
282 187 and 16 gene loci respectively [18]. Using *Ppa*-SL1a and *Ppa*-SL2a-enriched RNA-Seq, Sinha et al. [18]
283 showed that 90 % of genes are SL *trans*-spliced. Since it is unknown to what extent non-enriched RNA-Seq
284 data may underestimate this proportion, we took the opportunity to compare SLIDR between enriched
285 and non-enriched RNA-Seq data. We downloaded the Hybrid1 genome assembly [58] and SNAP anno-
286 tations from <http://www.pristionchus.org>, and the three enriched and non-enriched unstranded 2x75 bp
287 libraries from bioproject SRP039388 [18]. SLIDR was run with parameter -S '{40,55}[AG]T{4,6}[AG]'
288 to capture both SL1 and SL2 *Sm* binding motifs [14].

289 Using the SL1-enriched library alone, SLIDR detected *Ppa*-SL1a (1,184 reads spliced to 143 genes), *Ppa*-
290 SL1b (42 reads spliced to 24 genes) and three novel SL1 variants (Supplemental Table S1). Despite the
291 SL1 enrichment, there was also evidence of *Ppa*-SL2b/m and one novel SL2 variant, together representing
292 111 reads spliced to 46 genes. Contamination was also obvious for the SL2a-enriched library, where the

293 most frequent SL was *Ppa*-SL1a (336 reads spliced to 69 genes) followed by *Ppa*-SL2a (112 reads spliced
294 to 44 genes). SLIDR also detected four novel SL1 variants (three of which were also detected in the
295 SL1-enriched library) and *Ppa*-SL2b/c. The non-enriched library produced very similar results at similar
296 read and gene depths, comprising *Ppa*-SL1a, *Ppa*-SL2a/c and two novel SL1 and SL2 variants each
297 (Supplemental Table S1). Overall, 190 SL1 RNA genes and at least 29 SL2 RNA genes were detected,
298 which is a slight increase compared to those reported by Sinha et al. [18].

299 These results highlight that the SL enrichment via biotin pulldown may not have worked quite as ef-
300 fectively as suggested by qPCR control experiments [18]. Although it is impossible to quantify the
301 degree of contamination with non-*trans*-spliced transcripts, the SLIDR results suggest that the 90%
302 SL *trans*-splicing rate may be an overestimate. Using these three low-coverage libraries, SLIDR es-
303 timates the SL *trans*-splicing rate to only about 1%. To explore how SLIDR performs with higher
304 coverage data, we then ran SLIDR with six unstranded 2x150 bp libraries from bioproject PRJNA338247
305 (SRR4017216–SRR4017221) and a better resolved genome assembly plus annotations (GCA_000180635.3).
306 SLIDR detected the same breadth of known and novel SL1 and SL2 variants as above, but recovered
307 considerably more reads and SL *trans*-spliced genes (e.g., 16,085 *Ppa*-SL1a reads spliced to 5,148 genes).
308 Up to 6,498 genes are SL *trans*-spliced in these libraries, resulting in an estimated rate of c. 25 %, which
309 is still substantially below the postulated rate of 90 % [18]. Due to the superior genome assembly, SLIDR
310 detected at least 611 SL1 RNA genes and at least 29 SL2 RNA genes (Supplemental Table S1).

311 Overall, SLIDR detected a rich diversity of known and previously unreported SL1-type and SL2-type SLs
312 beyond the canonical *Ppa*-SL1a and *Ppa*-SL2a variants [15, 18]. SLIDR suggested that five out of the
313 seven *Ppa*-SL1 loci predicted by Sinha et al. [18] are not expressed. While likely not exhaustive, these
314 results highlight the sensitivity of SLIDR in detecting functional SL variants even from low coverage data.

315 *Meloidogyne hapla*

316 The plant-root knot nematode *M. hapla* possesses the canonical *C. elegans* SL1 and four additional
317 variants, all of which are *trans*-spliced to a minority of only 10 % of genes [31]. We used genome
318 PRJNA29083.WBPS14 and the same 32 SRA runs from bioproject PRJNA229407 that were used to
319 discover these SLs [31]. These data are particularly difficult to analyse since they are 75 bp single-ended,
320 unstranded and originate from mixed-culture RNA samples containing primarily material from the host
321 plant *Medicago truncatula*. Since reads from unstranded single-end libraries originate from the 5' end of
322 the transcript only 50% of the time, usable coverage is effectively halved. SLIDR was run with parameters
323 -S '{30,80}AT{4,6}G' -R 90 to allow for larger variation in *Sm* binding motif location and longer SL
324 RNA.

325 SLIDR detected all five known SLs and discovered at least nine novel SLs, suggesting that the SL
326 repertoire in this organism is much larger than previously identified (Supplemental Table S1). However,
327 only at most 176 reads were detected per SL, and at most 143 genes were SL *trans*-spliced. This is
328 consistent with low incidence of *trans*-splicing in this organism [31] and is not due to the RNA-Seq
329 data. We confirmed this with longer reads (100 bp single-end) from a different bioproject (PRJEB14142;
330 biosamples SAMEA4003664 and SAMEA4003666): SLIDR detected the same known and novel SLs
331 with even fewer reads (at most 74) and fewer *trans*-spliced genes (at most 68), suggesting that very
332 high coverage datasets would be necessary to exhaustively characterise SL *trans*-splicing events in this
333 organism (Supplemental Table S1). Nevertheless, SLIDR detected known and novel SLs even at this low
334 coverage, illustrating high sensitivity even with poor data.

335 *Trichinella spiralis*

336 The parasite *T. spiralis* possesses a diverse and unusual set of 15 SLs that are encoded by up to 48 genes
337 [30] and are spliced to c. 30 % of all genes [19]. Three out of these 15 SLs (*Tsp*-SL2, *Tsp*-SL10 and
338 *Tsp*-SL12) are SL2-type SLs specialised for resolving downstream genes in operons [19]. We downloaded
339 genome PRJNA12603.WBPS10 and three RNA-Seq libraries (SRR8327925-SRR8327927) from bioproject
340 PRJNA510020 [19]. SLIDR was run with parameter -S '{20,50}AT{4,6}G' to accommodate for the
341 smaller distance of the *Sm* binding motif to the splice donor site [30]. SLIDR detected all 15 known
342 SLs and a total of 21 SL RNA genes (Supplemental Table S1), which is an increase over the original 19
343 SL RNA genes identified from cDNA evidence [30]. These numbers also suggest that many of the 29
344 additional genomic loci predicted by Pettitt et al. [30] may not be functional. We note relatively low
345 numbers of SL reads (12-516) and SL *trans*-spliced genes (14-387), which is consistent with the notion
346 that SL *trans*-splicing affects at most about 30 % of genes in this organism (Supplemental Table S1).

347 *Trichuris muris*

348 *T. muris* is a gastrointestinal parasite closely related to *Trichinella spiralis* and possesses 13 SLs that,
349 unlike those of *T. spiralis*, resemble *C. elegans* SLs and are encoded by 13 genes [59, 60]. Three of
350 these SLs (*Tmu*-SL1, *Tmu*-SL4 and *Tmu*-SL12) are SL2-type SLs [59]. The genome-wide extent of SL
351 *trans*-splicing in this organism is unknown. We downloaded genome assembly PRJEB126.WBPS15 from
352 WormBase and five unstranded 2x100 bp libraries from bioproject PRJEB1054. SLIDR was run with
353 -S '{25,50}AT{4,6}G' to account for a shorter distance of the *Sm* binding motif to the splice donor site
354 [59]. SLIDR detected all 13 known SLs from 9–301 reads spliced to 10–249 genes (Supplemental Table
355 S1). Additionally, at least six novel SLs were identified from 4–1,117 reads spliced to 12–805 genes. The
356 numbers of SL RNA genes ranged between 1 and 3, suggesting that some of the SLs are encoded by
357 multi-copy genes. Overall, more than 2,000 genes received SLs, which would suggest an SL *trans*-splicing
358 rate of about 15 % (Supplemental Table S1).

359 *Prionchulus punctatus*

360 A limited SL repertoire of *P. punctatus* has been determined using 5-RACE of cDNA and comprises
361 six SLs that show structural similarity with *C. elegans* SL2 [16]. However, since no genome assembly
362 exists, the genomic organisation of SL genes and the extent of SL *trans*-splicing are unknown [16]. Only
363 two RNA-Seq libraries are available (SRA accessions ERR660626, ERR660627, bioproject PRJEB7585)
364 and no reference transcriptome assembly exists. We tested the performance of SLIDR with these two
365 libraries (2x100 bp) using a *de novo* transcriptome assembly obtained from the same libraries. Illumina
366 adapters and poor-quality bases (phred 30) were trimmed using TRIM_GALORE 0.6.4 [55], transcripts
367 were assembled using TRINITY 2.8.5 [61] and clustered at 100 % sequence similarity using CDHIT 4.8.1
368 [62]. The final assembly comprised 141,825 transcripts with an N50 of 786 bp (184-16,745 bp) and total
369 transcriptome size of 74.31 Mbp.

370 We ran SLIDR with a relaxed *Sm* location range (-S '{25,60}AC?T{4,6}G') but only discovered two
371 (*Ppu*-SL1 and *Ppu*-SL3) out of six known SLs, supported by only 96/16 reads and spliced to 29/7
372 genes respectively (Supplemental Table S1). While these results are little more than initial proof-of-
373 concept, it must be noted that the success of this *de novo* strategy depends critically on the presence
374 of SL RNA sequences in the transcriptome data. Since SL RNAs are not polyadenylated, RNA-Seq
375 library preparation protocols that rely on poly(A) selection will not capture SL RNAs, which limits
376 the use of publically available datasets that were not generated with ribosomal depletion protocols [63,

377 64] or poly(A)-tailing prior to library preparation [16]. Thus, we expect SLIDR to underperform in
378 transcriptome mode unless a high-quality transcriptome and high-coverage RNA-Seq data are available.

379 Validation of SLIDR in other eukaryotes

380 Although we designed the algorithms in SLIDR on the basis of known SL RNA structure in nematodes,
381 most of the filters based on sequence motifs can be fully customised or even disabled to relax stringency
382 if required. Here we demonstrate that SLIDR performs equally as well in other eukaryotes with SL
383 repertoires and SL RNA structures that are divergent from the nematode consensus. We used the same
384 datasets for *Ciona intestinalis*, *Hydra vulgaris* and *Schistosoma mansoni* that Calvelo et al. [33] used
385 to benchmark their SLFinder pipeline in order to carry out a detailed comparison of the two pipelines.
386 For all three species, SLIDR produced much clearer results with more sensitivity and specificity than
387 SLFinder, similar to what we observed above in *C. elegans*.

388 *Ciona intestinalis*

389 The tunicate *C. intestinalis* possesses a single 16 bp spliced leader 5'-ATTCTATTTGAATAAG-3' that is
390 spliced to at least 58% of expressed genes [21, 39, 65]. The SL RNA is very short (46 bp) and contains the
391 *Sm*-binding motif 5'-AGCUUUGG-3' [66]. The SL RNA has been suggested to be encoded by a highly
392 repetitive gene family comprising at least 670 copies, though the reference genome contains at most 15
393 of them due to assembly constraints [67]. SLFinder detected this single SL after extensive parameter
394 tweaking and found two distinct gene variants comprising 14 loci with a splice donor site [33].

395 We downloaded the same genome assembly (GCF_000224145.3) and the same three out of six 100 bp
396 paired-end datasets from bioproject PRJNA396771 (SRR5888437, SRR5888438 and SRR5888439) from
397 NCBI as Calvelo et al. [33]. We noticed after preliminary read alignments to the genome that the libraries
398 are not reverse stranded as described in the SRA entries, but are, in fact, unstranded. SLIDR was run
399 with the parameters `-x 0.6 -e 5 -O 5 -R 30 -S '{2,25}AGCTTTGG'` to enforce shorter soft-clipping
400 (maximum 24 bp given 100 bp reads), a BLAST e-value cut-off of 5 (to allow short matches of c. 11 bp),
401 maximum 5 bp outtron overlap, 30 bp RNA length excluding the SL, and the *Sm*-like motif located up to
402 25 bp downstream of the GT splice donor site.

403 SLIDR identified a single SL from only 15 reads (spliced to 15 genes) despite very high genome alignment
404 rates of 93-95 %; this SL represents the known SL sequence with some evidence of extra 5' nucleotides
405 (5'-taaggcATTCTATTTGAATAAG-3'). All but one of the eleven loci identified by SLIDR were on
406 chromosome NC_020175.2 (one was on NC_020166.2), and all loci were part of the 264 bp repeat unit
407 that contains functional SL copies in the genome [67]. In contrast, the 14 loci detected by SLFinder were
408 not located within the 264 bp repeat (none of the loci were on the correct chromosomes) and are therefore
409 probably pseudogenes, which are rife in the genome [67]. Since the poor SL detection rate is at odds with
410 the expected 58 % SL *trans*-splicing rate [65], we re-ran SLIDR with the remaining three libraries of the
411 same bioproject (SRR5888437, SRR5888438 and SRR5888439), but found no improvement. Similarly, two
412 libraries from two different bioprojects (SRR6706554, SRR2532443) yielded only slightly better results,
413 detecting the same SL and SL RNA genes from 59 and 227 reads spliced to 54 and 132 genes respectively
414 (Supplemental Table S1). Removing the filter for the *Sm* motif yielded the same SL sequence and detected
415 >100 pseudogenes, but did not increase the number of SL *trans*-spliced genes (Supplemental Table S1).
416 This difficulty in detecting SL *trans*-splicing would be compatible with the observation that expression
417 levels of SL *trans*-spliced genes are 2-3x lower than those of genes that are not SL *trans*-spliced [39].

418 However, when we tried 13 libraries from bioproject PRJNA376667 we obtained substantially superior
419 results: the known SL was detected from 38,467 reads spliced to 5,824 genes and originating from the same
420 eleven loci as above. Additionally, two novel variants of this SL were detected from 262 reads spliced to 209
421 genes and originating from three novel loci, totalling 14 out of 15 SL RNA loci (Supplemental Table S1).
422 Assuming 15,254 genes in the genome [21], these results indicate a SL *trans*-splicing rate of 40 %, which
423 is much closer to the expected 58 % [65]. This substantial variability of SL *trans*-splicing rates between
424 biosamples may be due to variability among life stages, tissues or RNA-Seq library preparation methods.
425 We also noted that the 13 libraries had substantially lower genome alignment rates (30-70 %) than the
426 other libraries, despite much greater evidence of SL *trans*-splicing, which could be due to contamination
427 with material from organisms other than *C. intestinalis*.

428 When we omitted the *Sm* motif filter during exploration of the initial libraries (PRJNA396771), we made
429 a curious discovery of potentially novel SLs that resemble nematode SLs instead of the canonical *C.*
430 *intestinalis* SL. After increasing the SL RNA length (-R 60) and filtering for a nematode-like *Sm* motif (-S
431 '.{20,50}AT{4,6}G') SLIDR detected one major 21 bp candidate (5'-CCGTTAAGTGTCTTGCCCAAG-
432 3') defined by 2,068 reads and spliced to 6 genes, alongside two additional plausible candidates at much
433 lower read depth (3–39 reads spliced to 4–6 genes) (Supplemental Table S1). That same major SL was
434 also detected among the 13 libraries of bioproject PRJNA376667, defined by 897 reads and spliced to 20
435 genes (Supplemental Table S1). It was beyond the scope of this study to fully resolve and describe these
436 novel SLs, but these preliminary results do highlight that SLIDR is much more sensitive than SLFinder,
437 which found no evidence of these SLs in the same libraries. In summary, although both SLFinder and
438 SLIDR detected the correct published SL sequence, SLFinder was unable to detect functional gene loci,
439 whereas all loci detected by SLIDR are consistent with the known SL RNA properties for this species
440 [66].

441 *Hydra vulgaris*

442 The cnidarian *H. vulgaris* possesses two types of spliced leaders that are added to at least one third of
443 all genes: the first type (SL-A) is 24 bp long and is part of an 80 bp SL RNA [68], whereas the second
444 type is much longer (46 bp SL, 107 bp SL RNA) and comprises a total of eleven SL variants across six
445 SLs (SL-B to SL-G) [24]. The *Sm* binding sites differ between SL-A (5'-GAUUUUCGG-3') and all other
446 SLs (5'-AAUUUUGA-3' or 5'-AAUUUUCG-3') [68]. SLFinder detected the full sequence of SL-B1 and
447 at least 21 loci, all of which were 5' truncated and are thus probably pseudogenes. SLFinder found no
448 evidence of SL-A or any of the other ten SLs [33].

449 We downloaded the genome assembly Hm105 and the same five stranded 2x100 bp datasets from biopro-
450 ject PRJNA497966 (SRR8089745–SRR8089749) [33]. We ran SLIDR with the parameters -x 1.5 -R 60 -S
451 '.{10,35}[AG]ATTTT[CG][AG]', which cover both *Sm* binding site motifs and should allow for detecting
452 both the short and long SLs. SLIDR detected the full SL-B1 sequence from 799,327 reads spliced to
453 18,418 genes and identified two gene loci, both of which were also identified by SLFinder [33]. SLIDR
454 also detected two 5' truncated versions of SL-D (10,696/4,494 reads, 2,727/1,677 spliced genes and two
455 gene loci, none of which were identified by SLFinder) and a novel variant of SL-B1 (864 reads spliced to
456 224 genes; coded by a single locus that was also identified by SLFinder). These results highlight that
457 SLIDR detected more SL variants than SLFinder and only reported functional SL RNA genes, whereas
458 SLFinder identified a large number of truncated gene loci, which are likely to be pseudogenes.

459 Since the SL-B-type SLs are exceptionally long, one would require longer reads than 100 bp to detect
460 full-length SLs with confidence. We tested another RNA-Seq library of 2x150 bp reads (SRR12070443)
461 and were able to detect full-length SL-B1, SL-D and SL-E based on 5,446–261,507 reads spliced to

462 1,598–12,688 genes (Supplemental Table S1). We also detected a novel SL-A-type variant (1,916 reads
463 spliced to 730 genes) and at least one novel SL-B-type variant (392 reads spliced to 202 genes). Similar
464 to the libraries above, we observed strikingly large numbers of SL *trans*-spliced genes – by far the highest
465 among any species detailed in this study. Contrary to previous estimates that only c. 33 % of c. 20,000
466 protein-coding genes are SL *trans*-spliced [24], our SLIDR results suggest that at least 63-92 % (12,688-
467 18,418) genes may be SL *trans*-spliced. While these results indicate the need for further study, they
468 demonstrate the level of inference possible with SLIDR.

469 *Schistosoma mansoni*

470 The platyhelminth *S. mansoni* possesses a single, relatively long (36bp) SL with an unusually long *Sm*
471 binding site (5'-AGUUUUCUUUGG-3') and a total RNA length of 90 bp [69]. The transcripts from at
472 least 46 % of genes undergo *trans*-splicing by this SL [22]. SLFinder detected this SL but missed the first
473 two A nucleotides; similarly, all detected loci were considerably 5' truncated [33].

474 We downloaded the genome assembly PRJEA36577.WBPS14 and two unstranded 2x100 bp datasets
475 from bioproject PRJNA225599 (SRR1020297–SRR1020298) [33]. SLIDR was run with parameters -x
476 1.25 -R 55 -S '{10,30}AGTTTTCTTTGG' to allow for detecting this large SL and the unusual *Sm*
477 binding site. SLIDR detected 18,568 reads that were assembled to the full length SL with two extra 5'
478 nucleotides (Supplemental Table S1). This SL was encoded by 112 genes, all of which correspond to the
479 *Smansoni_pSL-1* gene cluster on chromosome SM_V7_6 detected by SLFinder [33]. Contrary to the 5'
480 truncation that the SLFinder output suggested, these loci do contain the full length SL sequence, which
481 was correctly identified by SLIDR. SLIDR also detected two sequence variants (from only 2-3 reads) of
482 the SL corresponding to two loci, none of which were detected by SLFinder. SLFinder also reported 9
483 additional loci without clear splice donor sites, suggesting that these may be pseudogenes [33]. A total of
484 2,745 genes were SL *trans*-spliced, which is slightly higher than 2,459 genes previously identified from a
485 large-scale RNA-Seq data (250 million reads; 22), but represents only 21 % instead of the expected 46 %
486 SL *trans*-splicing rate [22].

487 Interestingly, SLFinder detected a genomic locus where the terminal ATG nucleotides of the SL sequence
488 were replaced by ACG, though this was not informed by evidence from the RNA-Seq data [33]. This
489 illustrates a key difference between the two software pipelines: SLFinder is primarily a genome annotation
490 pipeline for SL loci that uses RNA-Seq evidence as initial anchors (“hooks”) to search for all possible
491 gene loci [33]. In contrast, SLIDR aims to extract SL evidence directly from RNA-Seq data and uses the
492 genome only to extract additional evidence for functional sequence components of putative SLs; SLIDR
493 therefore only annotates gene loci that are expressed in the RNA-Seq libraries under consideration and
494 ignores alternative but unexpressed loci. Both approaches are clearly complementary, though our analyses
495 suggest that SLIDR is more robust in detecting functional SLs and SL RNA genes.

496 **Validation of SLOPPR in nematodes**

497 We have previously used SLOPPR to comprehensively discover operons in the genome of a nematode,
498 *T. spiralis*, for which there was only limited evidence for operon organisation [19]. Here we validate
499 and benchmark SLOPPR in the nematodes *C. elegans*, *C. briggsae* and *P. pacificus*, all of which have
500 well-characterised operon repertoires and use SL2-type *trans*-splicing to resolve mRNAs transcribed from
501 these operons. SLOPPR correctly classified SL1- and SL2-type SLs in all three species and identified
502 large proportions of known operons alongside several novel candidate operons (Table 2, Supplemental
503 Table S2). We further validated SLOPPR by confirming the presence of SL2-type SLs in the nematode

504 *T. muris*, for which the genome-wide landscape of SL *trans*-splicing and operon organisation is unresolved
505 [59].

506 *Caenorhabditis elegans*

507 *C. elegans* is the benchmark model organism for eukaryotic operons: up to 20 % of genes are situated
508 in operons [53]. Downstream operonic genes are readily diagnosable by an 80%–95% bias toward SL2,
509 though there are exceptions where downstream genes receive much lower proportions of SL2 [53].

510 We designed the operon prediction algorithm in SLOPPR on the basis of SL2:SL1 ratios at genes and
511 benchmarked its performance with curated *C. elegans* operons. We downloaded the genome assembly and
512 genome annotations from WormBase (PRJNA13758.WS276). These annotations contain 19,999 coding
513 genes and 1,542 operons, which are absent from the NCBI version of the genome. We used a large dataset
514 of 24 unstranded 2x100 bp RNA-Seq runs from bioproject PRJNA270896 (SRR1727796–SRR1727819).
515 We provided SLOPPR with the canonical SL1 sequence and 11 SL2 variants as supplied by SL-QUANT
516 [35], and with GFF annotations for the 1,542 reference operons. SLOPPR was run with the default
517 SL2:SL1 ratio of infinity, thus enforcing absence of SL1 at downstream operonic genes.

518 SLOPPR classified 36 % of genes as strictly SL1 *trans*-spliced, only 1 % as strictly SL2 *trans*-spliced
519 and 15 % as *trans*-spliced by both SL1 and SL2. The clustering algorithm correctly identified SL1-
520 and SL2-type subfunctionalisation of SLs (Supplemental Table S2). From these classifications a total of
521 434 operonic genes were identified in 213 operons, with median intergenic distance (distance between
522 “gene” GFF annotations) of 105 bp. Of these operons, 166 (77 %) matched reference operons, but these
523 represented only 11 % of the 1,524 total operons. This may be because the overall SL *trans*-splicing rate
524 of 52 % was below the expectation of 70–84 %, and the proportion of genes receiving a mixture of SLs
525 was much higher than the expected 6 % [53]. We thus relaxed the SL2:SL1 ratio threshold to 2, which
526 predicted 721 operonic genes in 345 operons (99 bp median intergenic distance), of which 295 (85 %)
527 matched reference operons (19 % of 1,542 operons).

528 While these numbers illustrate that SLOPPR predicts bona fide operons and also finds novel candidate
529 operons, they also demonstrate that this dataset is not nearly large enough to provide exhaustive insight
530 into the SL *trans*-splicing landscape. Tourasse et al. [34] carried out a meta-analysis of SL *trans*-splicing
531 in *C. elegans* using 1,682 RNA-Seq datasets comprising more than 50 billion reads, of which 287 million
532 reads contained evidence of SLs. Even at this huge coverage 97.4% of SL *trans*-splicing events were
533 supported by fewer than 100 reads and a vast number of events with very low read counts could not
534 be distinguished from biological noise in the splicing process. This highlights the inherent limitations of
535 standard RNA-Seq protocols and indicates that it is unrealistic to expect that all SL *trans*-splicing events
536 and operonic genes be detected using limited amounts of RNA-Seq data.

537 *Caenorhabditis briggsae*

538 *C. briggsae* is an important comparative model to *C. elegans*, but its gene and operon repertoires are less
539 resolved than those of its relative. The current genome annotations (PRJNA10731.WBPS14) contain
540 only 48 confirmed dicistronic operons. In contrast, Uyar et al. [17] used tight gene clusters that receive
541 SL2 to predict 1,034 operons, of which 51 % were syntenic with *C. elegans*. We decided to examine
542 SLOPPR with both annotation sets. The current CB4 assembly and annotations were downloaded from
543 WormBase (PRJNA10731.WBPS14); the CB3 assembly was downloaded vom UCSC to ensure that the
544 annotations by Uyar et al. [17] were compatible.

Table 2: Operon prediction (SLOPPR) results in four nematode and two tunicate species. Identifiers for reference genomes and paired-end RNA-Seq libraries are presented alongside numbers of quality-trimmed reads (QC), numbers of reads with a spliced leader (SL), percentage of genes receiving an SL, numbers of available reference operons, numbers of predicted operons, specificity (predicted operons matching reference operons), sensitivity (fraction of reference operons detected), numbers of operonic genes, and median intergenic distance among operonic genes. All runs did not require upstream operonic genes to be SL *trans*-spliced. Since the two tunicates *C. intestinalis* and *O. dioica* do not use SL2-type *trans*-splicing, operonic genes were automatically filtered by intergenic distance (maximum 84 bp for *C. intestinalis* and 413 bp for *O. dioica*; see main text).

Species	Genome	Bioproject	QC reads	SL reads	SL genes	Reference		Specificity		Operonic genes	Intergenic distance (bp)
						operons	operons	sensitivity			
<i>Caenorhabditis elegans</i>	GCF_000224145.3	PRJNA270896	621,744,295	937,587	52 %	1,542	345	86 % 19 %	721	99	
<i>Caenorhabditis briggsae</i>	CB3 (UCSC)	PRJNA104933	21,265,790	209,005	27 %	1,035	840	81 % 66 %	1,847	342	
<i>Caenorhabditis briggsae</i>	CB3 (UCSC)	PRJNA231838	59,271,467	286,520	36 %	1,035	752	81 % 59 %	1,631	332	
<i>Caenorhabditis briggsae</i>	CB4	PRJNA104933	21,265,790	235,001	38 %	48	921	5 % 96 %	2,058	111	
<i>Caenorhabditis briggsae</i>	CB4	PRJNA231838	59,271,467	289,834	41 %	48	750	5 % 81 %	1,626	110	
<i>Pristionchus pacificus</i>	Hybrid1	SRP039388	81,387,873	438,557	20 %	2,219	190	61 % 5 %	382	981	
<i>Pristionchus pacificus</i>	Hybrid1	PRJNA338247	330,856,071	312,718	28 %	2,219	205	50 % 5 %	414	872	
<i>Trichouris muris</i>	PRJEB126.WBPS15	PRJEB1054	115,460,947	168,727	35 %	-	718	-	1,492	418	
<i>Ciona intestinalis</i>	KH	PRJNA376667	194,089,029	452,505	51 %	1,328	1,172	94 % 83 %	2,563	1	
		PRJNA269316,									
<i>Oklopleura dioica</i>	GCA_000209535.1	PRJDB5668	250,361,592	5,860,107	9 %	1,765	577	90 % 30 %	1,292	33	

545 We first used SLOPPR with the two unstranded 2x42 bp libraries from Uyar et al. [17]. SLOPPR
546 quantified an overall SL *trans*-splicing rate of c. 27 % using the C3 assembly. SLOPPR predicted 1,346
547 operonic genes in 631 operons with a median intergenic distance of 333 bp, of which 507 operons
548 (80 %) matched the 1,035 reference operons (48 % detected). Relaxing the SL2:SL1 ratio threshold from
549 infinity to two predicted 1,847 genes in 840 operons (342 bp median intergenic distance), of which
550 682 (81 %) matched reference operons (65 % detected). Using the CB4 assembly, SLOPPR quantified
551 an SL *trans*-splicing rate of 38 % and predicted 1,475 operonic genes in 688 operons (112 bp median
552 intergenic distance), of which 37 (5 %) were among the 48 (77 %) reference dicistronic operons.
553 Relaxing the SL2:SL1 ratio to two resulted in 2,058 genes in 921 operons (111 bp median intergenic
554 distance) and recovered 46 out of 48 (95 %) reference dicistrons. These numbers highlight that the CB4
555 gene annotations are superior to those used by Uyar et al. [17] and that SLOPPR predicts *bona fide*
556 operons and novel candidate operons (Supplemental Table S2).

557 However, two concerns were raised during analysis of these 2x42 bp RNA-Seq libraries: First, only 62-69
558 % of reads aligned to the genome and only 50-58 % were properly paired. This is consistent with the
559 short read lengths and poor sequence quality which causes difficulty in aligning these reads with standard
560 aligners [17]. Second, the SL *trans*-splicing patterns varied more between the two libraries (L1 vs. mixed
561 life stages) than they did between SL1 and SL2-type SLs, which caused SLOPPR to cluster the SLs
562 by library instead of SL type (Supplemental Table S2). We thus re-ran the analyses with longer reads
563 from five unstranded 2x75 bp libraries (bioproject PRJNA231838). These libraries supported higher
564 SL *trans*-splicing rates of 36 % and 41 % for the CB3 and CB4 genome assemblies respectively, and
565 SLOPPR correctly identified SL1 and SL2-type clusters among these data using either genome assembly
566 (Supplemental Table S2).

567 Using CB3, SLOPPR predicted 1,202-1,631 operonic genes in 564-752 operons (332-341 bp median inter-
568 cistronic distance), of which 442-606 (78-80 %) were among the 1,035 CB3 reference operons (43-59 %
569 detected). Using CB4, SLOPPR predicted 1,179-1,626 operonic genes in 553-750 operons (110 bp median
570 intergenic distance), of which 24-39 (4-5 %) were among the 48 reference dicistrons (50-81 % detected).
571 These results are somewhat superior to those above, but echo essentially the same patterns. The more
572 recent CB4 assembly clearly has better gene annotations that yield a much lower median intergenic
573 distance of about 110 bp, which is consistent with *C. elegans* [53]. However, this assembly contains only
574 few curated operons and SLOPPR detected most of these. In addition, SLOPPR discovered a large set
575 of novel operons, which warrants further study.

576 *Pristionchus pacificus*

577 *P. pacificus* is another important comparative model to *C. elegans* that resolves operons with SL2-type
578 *trans*-splicing [15]. A comprehensive survey of SL *trans*-splicing events using SL1- and SL2-enriched RNA-
579 Seq data suggested that 90 % of genes are SL *trans*-spliced and a total of 2,219 operons may exist on the
580 basis of tight gene clusters and SL1/SL2 *trans*-splicing ratios [18]. We downloaded the Hybrid1 genome
581 assembly [58] and SNAP genome annotations and operon annotations from <http://www.pristionchus.org>
582 [18].

583 We first ran SLOPPR with those same SL-enriched unstranded 2x75 bp libraries from bioproject SRP039388,
584 supplying the two canonical *Ppa*-SL1a and *Ppa*-SL2a sequences that were used for SL enrichment [18].
585 SLOPPR detected SLs at only 20 % of genes, even when including the non-enriched library (SRR1182510).
586 The SL1-enriched library (SRR1542610) supported SL1 and SL2 *trans*-splicing at 16.17 % and 0.05 %
587 of genes, consistent with SL1-enrichment. However, the SL2-enriched library (SRR1542630) showed no
588 evidence of SL2-enrichment (0.98 % of genes) but comparable SL1 levels to the non-enriched control

589 library (8.2 % of genes). These results echo the SLIDR results using the same libraries (see above) and
590 would suggest a far lower SL *trans*-splicing rate than 90 % [18]. Due to the low SL *trans*-splicing rate,
591 only 234 operonic genes in 117 operons were predicted, of which 99 (84 %) matched the 6,909 operon-like
592 gene clusters and 67 (57 %) matched the 2,219 plausible reference operons [18]. Relaxing the SL2:SL1
593 ratio threshold to two yielded slightly better results with 382 genes in 190 operons, of which 115 (61 %)
594 were among the 2,219 reference operons.

595 To examine whether these poor numbers were due to the specific libraries, we then ran SLOPPR with a
596 more extensive set of six unstranded 2x150 bp libraries from bioproject PRJNA338247 (SRR4017216–SRR4017221).
597 Surprisingly, these libraries yielded similar results, suggesting that 28 % of genes are SL *trans*-spliced and
598 predicting 300 operonic genes in 150 operons, of which 76 (50 %) matched the 2,219 reference operons.
599 Using a relaxed SL2:SL1 ratio of two, 414 genes were identified in 205 operons, of which 102 (50 %) were
600 among the 2,219 reference operons. Both sets of libraries yielded similar median intergenic distances
601 of 785-860 bp, which increased to 872-981 bp when relaxing the SL2:SL1 threshold. These distances are
602 much larger than the 100 bp expected in *C. elegans* [53] but are consistent with the median distance of
603 1,149 bp among all 6,909 gene clusters in *P. pacificus* and very poor synteny of these clusters with *C.*
604 *elegans* (only 37 out of 6,909 clusters are syntenic; 18). SLOPPR also correctly identified SL1- and SL2-
605 type clusters from genome-wide *trans*-splicing patterns in both library sets, confirming that *Ppa*-SL1a
606 and *Ppa*-SL2a are functionally diverged (Supplemental Table S2).

607 All these observations suggest that SLOPPR produces plausible results given the limitations of relying
608 on RNA-Seq reads covering the 5' end of transcripts. The striking discrepancies between SLOPPR and
609 the analyses by Sinha et al. [18] are likely due to their assumption that all reads from the SL-enriched
610 libraries are from SL *trans*-spliced transcripts. Since only a small fraction of RNA-Seq reads originate
611 from the 5' end of transcripts (confirmed by SLIDR and SLOPPR), this assumption cannot be confirmed
612 bioinformatically, and thus it cannot be ruled out that these libraries contained substantial amounts of
613 contaminant non-*trans*-spliced transcripts despite the authors' efforts of confirming their methods with
614 qPCR [18]. If this were the case, their analyses based on SL1/SL2 *trans*-splicing patterns would be flawed,
615 which would explain the poor overlap with the more transparent SLOPPR results. One way of testing
616 this would be to combine SL-enrichment with long-read whole-transcript sequencing on the PacBio or
617 ONT NanoPore platforms. Such work would be instrumental in improving detection of SL *trans*-splicing
618 in any species.

619 *Trichuris muris*

620 *T. muris* is a gastrointestinal parasite of mice and is an important model system for studying mammalian
621 gastrointestinal parasitism. It belongs to the same clade as *T. spiralis* and *P. punctatus* [70]. Comparative
622 work with *T. spiralis* has identified a repertoire of 13 *T. muris* SLs, of which *Tmu*-SL1, *Tmu*-SL4 and
623 *Tmu*-SL12 show structural similarity with *C. elegans* SL2 and are *trans*-spliced to the downstream gene
624 of a bicistronic operon that is conserved among several nematode species [59]. However, the genome-wide
625 landscape of SL *trans*-splicing and operons in *T. muris* is unresolved [59]. Since we thus cannot compare
626 SLOPPR operon predictions against reference operons, we aimed to merely test the hypothesis that the
627 three putative SL2-type SLs show genome-wide SL *trans*-splicing patterns that are distinct from those of
628 the other ten SLs.

629 We downloaded genome assembly PRJEB126.WBPS15 from WormBase and five unstranded 2x100 bp
630 libraries from bioproject PRJEB1054. We supplied SLOPPR with the 13 known SLs and designated *Tmu*-
631 SL1, *Tmu*-SL4 and *Tmu*-SL12 as SL2-type. SLOPPR detected a relatively high SL *trans*-splicing rate of
632 35 % and clustered the 13 SLs into the expected groups comprising *Tmu*-SL1, *Tmu*-SL4 and *Tmu*-SL12

633 versus all other SLs (Figure 3). Using a relaxed SL2:SL1 threshold of two, SLOPPR predicted 718 operons
634 comprising 1,492 operonic genes with a median intergenic distance of 418 bp. This is larger than the c.
635 100 bp in *C. elegans* but is consistent with the observed elevated intergenic distance among manually
636 curated *Tmu* benchmark operons [59] and also with the considerably elevated intergenic distance among
637 non-operonic genes (5,519 bp compared with c. 3,500 bp observed in all other species in this study;
638 Supplemental Table S2). These results echo those we obtained in *T. spiralis* [19] and demonstrate that
639 SLOPPR allows to identify subfunctionalisation among SLs that may correspond to SL1 and SL2-type
640 *trans*-splicing. *Tmu*-SL1, *Tmu*-SL4 and *Tmu*-SL12 are very likely used to resolve polycistronic RNAs
641 in this organism and SLOPPR has predicted plausible candidate operons that warrant future curation
642 efforts.

643 Validation of SLOPPR in other eukaryotes

644 Having established that SLOPPR is a powerful method for predicting operons in organisms that use
645 specialised SLs to resolve downstream operonic genes (SL2-type SLs), we finally aimed to illustrate that
646 SLOPPR is also able to infer operons in organisms that lack SL specialisation. Here we demonstrate this
647 ability in two tunicates, *Ciona intestinalis* and *Oikopleura dioica*, both of which have only a single SL
648 that resolves operons but is also added to monocistronic genes.

649 In such situations, the SL must be designated as SL2-type such that all genes that receive the SL
650 are classed as operonic; this set of genes will contain *bona fide* operonic genes but will also contain all
651 monocistronic genes that receive the SL. Therefore, these initial candidate operonic genes must be filtered
652 by intergenic distance to partition out true operonic genes. SLOPPR can be configured to either use
653 a user-supplied cutoff if the expected intergenic distances are known, or to bisect the distribution
654 of intergenic distances empirically into two groups using K-means clustering and retaining those
655 genes with short distances. By exploring several parameter combinations, specificity and sensitivity in
656 partitioning out operons can be traded off (Figure 4).

657 *Ciona intestinalis*

658 The tunicate *C. intestinalis* splices a single SL to downstream operonic genes and infrequently to upstream
659 operonic genes, consistent with operons in nematodes [21, 39]. Using short intergenic distances (<100 bp)
660 as the sole criterion, a total of 1,310 operons comprising 2,909 genes have been predicted [21, 39]. These
661 operons are predominantly dicistronic and have extremely small intergenic distances, often lacking an
662 intergenic region altogether [21, 39], similar to the rare SL1-dependent operons observed in *C. elegans*
663 [71]. The genome annotations take SL *trans*-splicing into account and define gene boundaries correctly
664 between poly(A) and *trans*-splicing sites [21].

665 We obtained the KH genome assembly, the KH gene models (2013) and KH operon annotations (2013;
666 containing 1,328 operons) from the Ghost database (http://ghost.zool.kyoto-u.ac.jp/download_kh.html).
667 We used the same 13 RNA-Seq libraries from bioproject PRJNA376667 that contained disproportionately
668 more evidence of SL *trans*-splicing using SLIDR than other tested libraries. SLOPPR detected an overall
669 SL *trans*-splicing rate of 51 %, close to the 58 % expectation [65], although the libraries varied considerably
670 in genome alignment rate and SL *trans*-splicing rate as observed earlier (Supplemental Table S2). We first
671 ran SLOPPR with the same default configuration as for the nematodes, ignoring intergenic distances
672 after operon inference. Indeed, this first run predicted a vastly inflated set of 3,594 operons, of which 1,196
673 (33 %) matched reference operons and 90 % of the reference operons were detected. The contamination

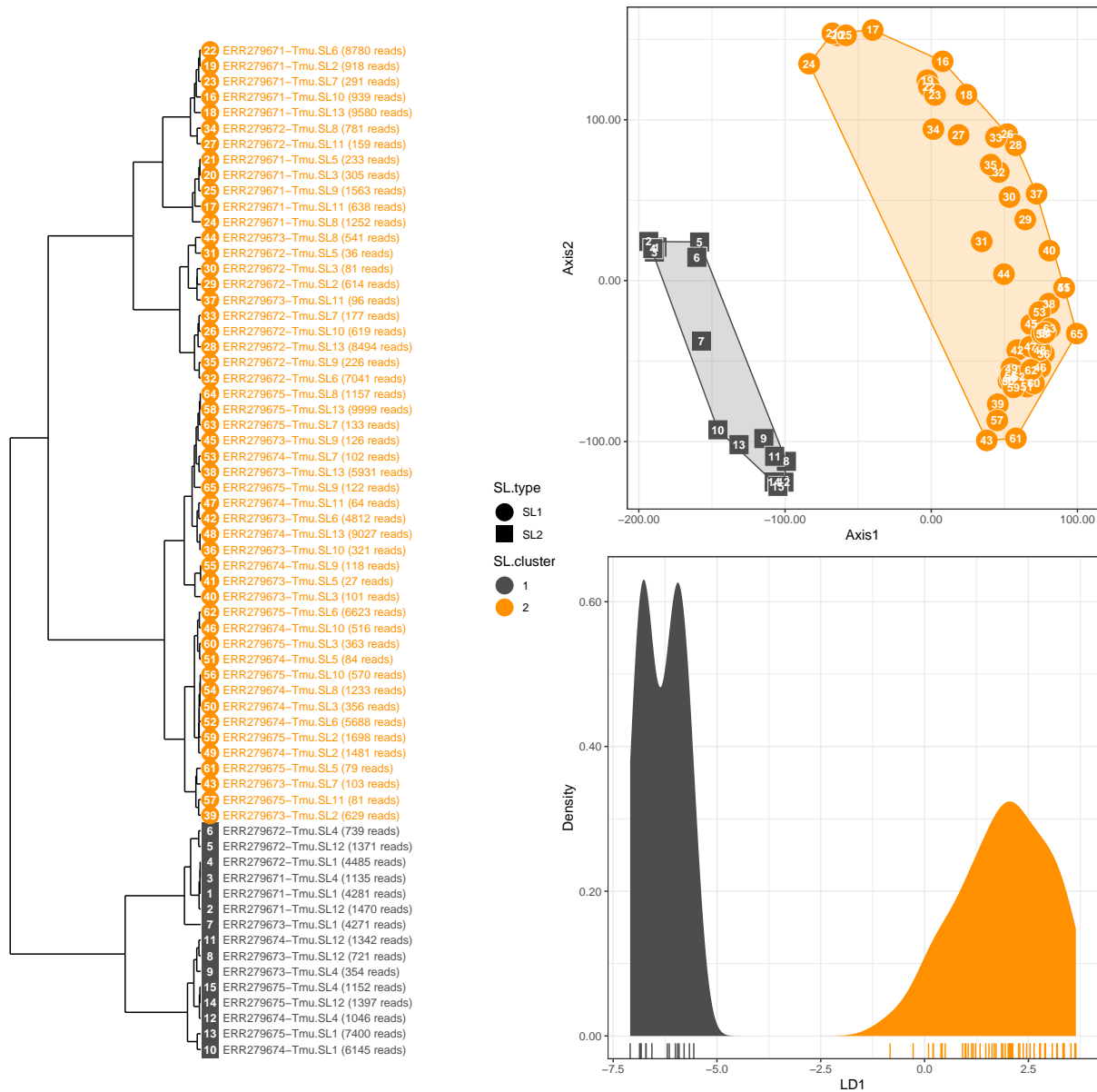


Figure 3: Genome-wide spliced-leader (SL) *trans*-splicing patterns among 13 SLs and five RNA-Seq libraries in *Trichuris muris*. Top right: generalized PCA of normalised genome-wide SL read counts. Symbol shape represents *a priori* SL type (circle: SL1; square: SL2) and colour represents cluster membership inferred via K-means clustering (dark grey: cluster1; orange: cluster2). Numbers inside symbols refer to library identifiers as detailed in the dendrogram on the left (hierarchical Ward's clustering of PCA eigenvectors). Bottom right: linear discriminant analysis between the two clusters, highlighting complete cluster differentiation. *Tmu*-SL1, *Tmu*-SL4 and *Tmu*-SL12 are correctly identified as distinct from all other SLs, confirming their functional specialisation as SL2-type SLs.

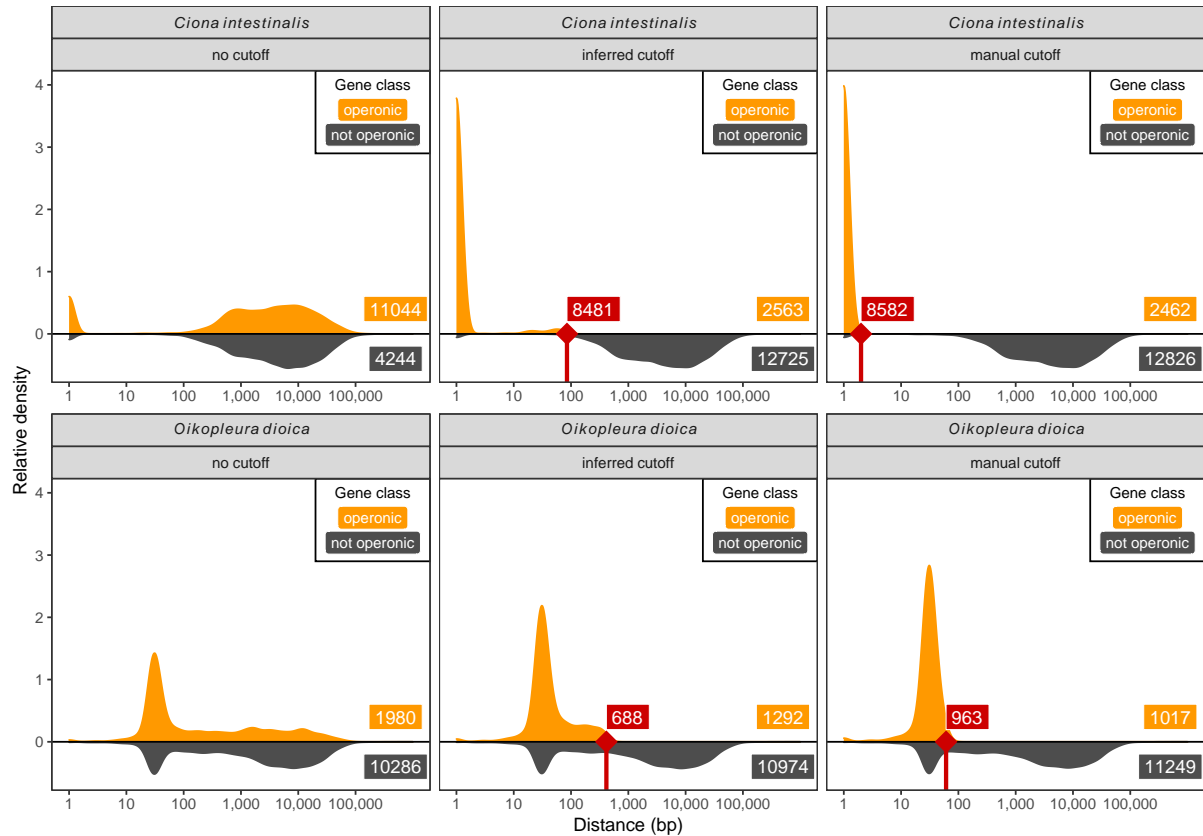


Figure 4: Separation of operonic genes from monocistronic genes in the absence of specialised spliced leaders (SLs), illustrated with SLOPPR data from the tunicates *Ciona intestinalis* (top panels) and *Oikopleura dioica* (bottom panels). All panels display distributions of distances between operonic and non-operonic genes, and labels provide gene numbers. Left panels: in both organisms, a single SL is added to monocistronic and operonic genes, causing SLOPPR to incorrectly designate monocistronic SL-receiving genes with large intergenic distances as operonic. Middle panels: an optimal distance cutoff for operonic genes is inferred via K-means clustering, and genes at or above the cutoff (red notches at 85 and 414 bp respectively) are re-classified as monocistronic non-operonic (red labels). Right panels: a lower manual cutoff (1 bp and 60 bp respectively; red notches at 2 and 61 bp) further reduces the set of genes retained as operonic. Note the peak of tightly-spaced non-operonic genes in the *O. dioica* panels; these genes are likely operonic genes but no SL evidence was obtained from the RNA-Seq data.

674 with monocistronic genes inflated the intergenic distances (median 2,287 bp) but a distinct set of
675 genes had very low intergenic distances, likely representing true operons (Figure 4).

676 We then partitioned out true operonic genes by re-running SLOPPR with automatic inference of the
677 optimal intergenic distance cutoff and also switching off addition of upstream operonic genes. These
678 strict settings require all genes in an operon to be SL *trans*-spliced and all operons to be at least dicistronic.
679 SLOPPR predicted only 856 operons with a median intergenic distance of 1 bp (inferred cutoff: 68 bp).
680 Of these operons, 823 (96%) matched reference operons, indicating high specificity, but only 62 % of the
681 reference operons were detected (Supplemental Table S2). Therefore, we re-ran the same analysis but re-
682 instating upstream genes (i. e., allowing them not to be SL *trans*-spliced). This time, SLOPPR predicted
683 1,172 operons, of which 1,100 (94%) matched reference operons and represented a considerably improved
684 83 % of reference operons. The median intergenic distance was again 1 bp (inferred cutoff: 84 bp),
685 consistent with the notion that many operons in this organism have no intergenic regions [21, 39]. To
686 quantify this proportion of operons, we re-ran the analysis enforcing a maximum intergenic distance of
687 1 bp. This yielded 1,128 operons, indicating that only 44 operons had intergenic regions (Supplemental
688 Table S2).

689 Overall, SLOPPR predicted most of the previously proposed operons and also novel operons comprising
690 up to seven genes instead of six as previously reported [21]. This analysis demonstrates the flexibility that
691 SLOPPR brings to operon prediction, by first identifying SL *trans*-splicing patterns and then filtering
692 candidate operons by intergenic distance. This approach enables extraction of operon sets with varying
693 stringency and biological characteristics.

694 *Oikopleura dioica*

695 Like *C. intestinalis*, the tunicate *O. dioica* possesses only a single SL that is *trans*-spliced to both mono-
696 cistronic genes and genes in operons, where upstream genes are not required to be SL *trans*-spliced [20, 72].
697 At least 39 % of genes are SL *trans*-spliced and 58 % of SL *trans*-spliced transcripts originate from oper-
698 ons [72]. A total of 1,765 operons comprising 5,005 genes have been predicted via short intergenic
699 distances of at most 60 bp [49]. We downloaded genome assembly GCA_000209535.1 (V3) and genome
700 annotations from OikoBase (<http://oikoarrays.biology.uiowa.edu/Oiko/>), and operon annotations from
701 the Genoscope Oikopleura Genome Browser (<https://wwwdev.genoscope.cns.fr/oikopleura/>). Curiously,
702 both the genome and operon annotations also contain entries that reference an inferior, much smaller,
703 assembly containing much shorter contigs (GCA_000209555.1). After removing these redundant entries
704 from both annotation sets we were left with the expected 1,765 instead of 2,971 operons.

705 We downloaded four unstranded 2x90 bp libraries from bioproject PRJNA269316 and 16 stranded
706 2x100 bp libraries from bioproject PRJDB5668, representing various life stages. Similarly to *C. in-*
707 *testinalis*, we observed poor and highly variable background alignment rates (17-69 %) but large numbers
708 of SL reads (5.8 million in total). However, these reads covered only 9 % of genes, which is much lower
709 than the expected 39 % (Supplemental Table S2). In default mode, SLOPPR predicted 885 operons
710 with median intergenic distance of 57 bp. Of these, 644 (73 %) matched reference operons. As in *C.*
711 *intestinalis*, the operons were contaminated with monocistronic genes having much larger intergenic
712 distances (median of 2,178 bp) (Figure 4). We therefore re-ran SLOPPR with inference of the optimal
713 intergenic distance cutoff and obtained 577 operons, of which 521 (90 %) matched reference operons.
714 The median intergenic distance was reduced to 33 bp, but the inferred cutoff was still fairly high at
715 413 bp (Figure 4). We thus re-ran the analysis with the same hard cutoff of 60 bp that was used to
716 predict the 1,765 reference operons [49] and were left with 464 operons (median intergenic distance
717 of 31 bp), of which 454 (98%) matched reference operons (Supplemental Table S2). We also tested the

718 effect of enforcing SL *trans*-splicing at upstream genes (-u option) across the same three analysis runs
719 and obtained a much more stringent set of 111-165 operons of which 106-143 (87-95 %) matched reference
720 operons (Supplemental Table S2).

721 These results indicate that SLOPPR can discriminate operonic genes from monocistronic genes receiving
722 the same SL, identify the vast majority of previously described operons and also predict a small number
723 of novel operons. One limitation with this dataset was that only few genes received the SL, which was a
724 consistent observation across all libraries tested from several bioprojects and suggests that more RNA-Seq
725 data would need to be generated to fully characterise the SL *trans*-splicing landscape in this organism.

726 Conclusions

727 We have created two computational pipelines that fill a long-standing gap in our ability to identify and
728 quantify SL *trans*-splicing and eukaryotic operons in any species where RNA-Seq data and reference
729 genomes/transcriptomes are available. SLIDR is a more sensitive, specific and efficient SL discovery
730 pipeline than SLFinder [33], able to uncover a wealth of untapped SL diversity. SLOPPR is the first
731 universal pipeline to predict operons from SL *trans*-splicing events, closing this important gap left by
732 existing SL quantification pipelines [34, 35]. We have demonstrated here and elsewhere [19] that SLOPPR
733 identifies both *bona fide* and novel operons, blazing the trail for routine operon prediction in any organism
734 with SL *trans*-splicing. Importantly, SLOPPR exploits biological replicates to infer subfunctionalisation
735 among SLs and to moderate noise in SL quantification, which lays a foundational framework for developing
736 a new field of eco-evolutionary “SL-omics” investigating differential SL usage and *trans*-splicing levels
737 among biological replicates, experimental groups or wild populations.

738 A fundamental limitation of both SLIDR and SLOPPR is that they were designed for traditional RNA-
739 Seq data where sequencing error is low but only a small fraction of reads originate from the 5' end of
740 the transcript containing the SL. Most RNA-Seq library preparation methods also show considerable loss
741 of coverage at the 5' end, which often limits SL detection to a short c. 10 bp portion at typically <1
742 % of reads [19, 35, 73]. This means that SLOPPR in particular is likely to underestimate the extent of
743 SL *trans*-splicing and operonic gene organisation unless huge amounts of sequencing data are available
744 [34] or specialised SL-enrichment library preparation methods are used [18, 22, 29]. However, our SLIDR
745 analysis on *Hydra vulgaris* vividly demonstrates that SLs at nearly 100 % of all genes can be detected
746 from RNA-Seq data if coverage is sufficient.

747 We decided to build these pipelines on RNA-Seq data because a wealth of datasets already exists for many
748 species, which continues to grow rapidly. We are thus, for the first time, in the position to investigate
749 SL *trans*-splicing systematically throughout the tree of life without needing to generate novel sequence
750 data. Nevertheless, a powerful future avenue for capturing the full 5' end of transcripts is direct RNA or
751 cDNA sequencing on the Oxford NanoPore or PacBio long-read platforms [74, 75]. This would require
752 much less sequencing effort because the full molecule is sequenced instead of a short random fraction.
753 SLIDR and SLOPPR could easily be expanded to accept long-read data but would require tailored error-
754 tolerant screening methods to accommodate the higher error rate of NanoPore reads. As these long-read
755 transcriptomics datasets become more commonplace, we expect SL-omics to become a routine molecular
756 tool for uncovering the causes and consequences of this enigmatic source of molecular diversity.

757 Availability and requirements

758 Project name: SLIDR and SLOPPR

759 Project home page: <https://github.com/wenzelm/slidr-sloppr>

760 Operating system(s): Linux

761 Programming language: BASH, R

762 Other requirements: CUTADAPT (tested v2.3), GFFREAD (tested v0.11.4), HISAT2 (tested v2.1.0), BOWTIE2
763 (tested v2.3.5), SAMTOOLS (tested v1.9), BEDTOOLS (tested v2.28.0), SEQTK (tested v1.3), CDHIT (tested
764 v4.8.1), VSEARCH (tested v2.4.3), BLASTN (tested v2.9.0), SUBREAD (tested v1.6.2), VIENNARNA (tested
765 v2.4.14), R (tested v3.6.0), R packages *glmpca*, *data.table*, *parallel*, *ggdendro*, *MASS*, *ggplot2*

766 License: Creative Commons Attribution License (CC BY 4.0)

767 Any restrictions to use by non-academics: None

768 Abbreviations

769 **SL** spliced leader

770 **UTR** untranslated region

771 **CPM** counts per million reads

772 **PCA** principal component analysis

773 Declarations

774 Ethics approval and consent to participate

775 Not applicable

776 Consent for publication

777 Not applicable

778 Availability of data and materials

779 All datasets analysed in this study are publicly available from NCBI (<https://www.ncbi.nlm.nih.gov/>),
780 ENA (<https://www.ebi.ac.uk/ena/browser/home>), WormBase (<https://wormbase.org/>), Ghost database
781 (http://ghost.zool.kyoto-u.ac.jp/download_kh.html), OikoBase (<http://oikoarrays.biology.uiowa.edu/Oiko/>)
782 and Genoscope (<https://wwwdev.genoscope.cns.fr/oikopleura/>). Accession numbers are detailed in the
783 main text and in supplemental materials. All data generated in this study are included in this published
784 article and its supplemental material.

785 Competing interests

786 The authors declare that they have no competing interests.

787 Funding

788 This work was supported by the Biotechnology and Biological Sciences Research Council [BB/J007137/1
789 to JP and BM, and BB/T002859/1 to BM and JP]. The funding body had no role in the design of the
790 study and collection, analysis, and interpretation of data and in writing the manuscript.

791 Authors' contributions

792 JP and BM acquired funding for the project, conceived the research and managed the research activity.
793 MW designed and implemented the computational pipelines with major contributions by JP. MW carried
794 out all data analyses and prepared all display items. MW wrote the manuscript with major contributions
795 by JM and BM. All authors read and approved the final manuscript.

796 Acknowledgements

797 The authors thank Bernadette Connolly for helpful discussions and Andreea Marin for testing the
798 pipelines. The authors acknowledge the support of the Maxwell computer cluster funded by the University
799 of Aberdeen.

800 References

- 801 [1] Nicholas A Stover, Michelle S Kaye, and Andre RO Cavalcanti. Spliced leader trans-splicing. *Current*
802 *Biology*, 16(1):R8–R9, 2006.
- 803 [2] Erika L Lasda and Thomas Blumenthal. Trans-splicing. *Wiley Interdisciplinary Reviews: RNA*, 2
804 (3):417–434, 2011.
- 805 [3] Richard E Sutton and John C Boothroyd. Evidence for trans splicing in trypanosomes. *Cell*, 47(4):
806 527–535, 1986.
- 807 [4] Bart Cuypers, Malgorzata A Domagalska, Pieter Meysman, Géraldine de Muylder, Manu Vanaer-
808 schot, Hideo Imamura, Franck Dumetz, Thomas Wolf Verdonckt, Peter J Myler, Gowthaman Ra-
809 masamy, et al. Multiplexed spliced-leader sequencing: A high-throughput, selective method for
810 rna-seq in trypanosomatids. *Scientific reports*, 7(1):1–11, 2017.
- 811 [5] George Cherian Pandarakalam, Michael Speake, Stuart McElroy, Ammar Alturkistani, Lucas
812 Philippe, Jonathan Pettitt, Berndt Müller, and Bernadette Connolly. A high-throughput screen for
813 the identification of compounds that inhibit nematode gene expression by targeting spliced leader
814 trans-splicing. *International Journal for Parasitology: Drugs and Drug Resistance*, 10:28–37, 2019.
- 815 [6] Kenneth EM Hastings. Evolutionary origin of sl-addition trans-splicing: still an enigma. *Trends in*
816 *genetics*, 21(4):240–247, 2005.
- 817 [7] Vassilis Douris, Maximilian J Telford, and Michalis Averof. Evidence for multiple independent origins
818 of trans-splicing in metazoa. *Molecular biology and evolution*, 27(3):684–693, 2010.
- 819 [8] Romain Derelle, Tsuyoshi Momose, Michael Manuel, Corinne Da Silva, Patrick Wincker, and Evelyn
820 Houliston. Convergent origins and rapid evolution of spliced leader trans-splicing in metazoa: insights
821 from the ctenophora and hydrozoa. *RNA*, 16(4):696–707, 2010.

- 822 [9] Mainá Bitar, Mariana Boroni, Andréa Mara Macedo, Carlos Renato Machado, and Glória Regina
823 Franco. The spliced leader trans-splicing mechanism in different organisms: molecular details and
824 possible biological roles. *Frontiers in genetics*, 4:199, 2013.
- 825 [10] Thomas Blumenthal. Operons in eukaryotes. *Briefings in Functional Genomics*, 3(3):199–211, 2004.
- 826 [11] John Spieth, Glenn Brooke, Scott Kuersten, Kristi Lea, and Thomas Blumenthal. Operons in *c.*
827 *elegans*: polycistronic mrna precursors are processed by trans-splicing of sl2 to downstream coding
828 regions. *Cell*, 73(3):521–532, 1993.
- 829 [12] Thomas Blumenthal. Trans-splicing and polycistronic transcription in *caenorhabditis elegans*. *Trends*
830 *in Genetics*, 11(4):132–136, 1995.
- 831 [13] Thomas Blumenthal, Donald Evans, Christopher D Link, Alessandro Guffanti, Daniel Lawson, Jean
832 Thierry-Mieg, Danielle Thierry-Mieg, Wei Lu Chiu, Kyle Duke, Moni Kiraly, et al. A global analysis
833 of *caenorhabditis elegans* operons. *Nature*, 417(6891):851, 2002.
- 834 [14] Kwang-Zin Lee and Ralf J Sommer. Operon structure and trans-splicing in the nematode *pristionchus*
835 *pacificus*. *Molecular biology and evolution*, 20(12):2097–2103, 2003.
- 836 [15] David B Guiliano and Mark L Blaxter. Operon conservation and the evolution of trans-splicing in
837 the phylum nematoda. *PLoS Genet*, 2(11):e198, 2006.
- 838 [16] Neale Harrison, Andreas Kalbfleisch, Bernadette Connolly, Jonathan Pettitt, and Berndt Müller. Sl2-
839 like spliced leader rnas in the basal nematode *prionchulus punctatus*: New insight into the evolution
840 of nematode sl2 rnas. *RNA*, 16(8):1500–1507, 2010.
- 841 [17] Bora Uyar, Jeffrey SC Chu, Ismael A Vergara, Shu Yi Chua, Martin R Jones, Tammy Wong, David L
842 Baillie, and Nansheng Chen. Rna-seq analysis of the *c. briggsae* transcriptome. *Genome research*,
843 22(8):1567–1580, 2012.
- 844 [18] Amit Sinha, Claudia Langnick, Ralf J Sommer, and Christoph Dieterich. Genome-wide analysis of
845 trans-splicing in the nematode *pristionchus pacificus* unravels conserved gene functions for germline
846 and dauer development in divergent operons. *RNA*, 20(9):1386–1397, 2014.
- 847 [19] Marius Wenzel, Christopher Johnston, Berndt Müller, Jonathan Pettitt, and Bernadette Connolly.
848 Resolution of polycistronic rna by sl2 trans-splicing is a widely-conserved nematode trait. *RNA*,
849 pages rna-076414, 2020.
- 850 [20] Philippe Ganot, Torben Kallesøe, Richard Reinhardt, Daniel Chourrout, and Eric M Thompson.
851 Spliced-leader rna trans splicing in a chordate, *oikopleura dioica*, with a compact genome. *Molecular*
852 *and cellular biology*, 24(17):7795–7805, 2004.
- 853 [21] Yutaka Satou, Katsuhiko Mineta, Michio Ogasawara, Yasunori Sasakura, Eiichi Shoguchi, Keisuke
854 Ueno, Lixy Yamada, Jun Matsumoto, Jessica Wasserscheid, Ken Dewar, et al. Improved genome
855 assembly and evidence-based global gene model set for the chordate *ciona intestinalis*: new insight
856 into intron and operon populations. *Genome biology*, 9(10):R152, 2008.
- 857 [22] Mariana Boroni, Michael Sammeth, Sandra Grossi Gava, Natasha Andressa Nogueira Jorge, An-
858 dréa Mara Macedo, Carlos Renato Machado, Marina Moraes Mourão, and Glória Regina Franco.
859 Landscape of the spliced leader trans-splicing mechanism in *schistosoma mansoni*. *Scientific Reports*,
860 8(1):1–14, 2018.

- 861 [23] Ferdinand Marlétaz, André Gilles, Xavier Caubit, Yvan Perez, Carole Dossat, Sylvie Samain, Gabor
862 Gyapay, Patrick Wincker, and Yannick Le Parco. Chaetognath transcriptome reveals ancestral and
863 unique features among bilaterians. *Genome biology*, 9(6):R94, 2008.
- 864 [24] Jarrod A Chapman, Ewen F Kirkness, Oleg Simakov, Steven E Hampson, Therese Mitros, Thomas
865 Weinmaier, Thomas Rattei, Prakash G Balasubramanian, Jon Borman, Dana Busam, et al. The
866 dynamic genome of hydra. *Nature*, 464(7288):592–596, 2010.
- 867 [25] Christian Preußner, Nicolas Jaé, and Albrecht Bindereif. mrna splicing in trypanosomes. *International*
868 *Journal of Medical Microbiology*, 302(4-5):221–224, 2012.
- 869 [26] Huan Zhang, Yubo Hou, Lilibeth Miranda, David A Campbell, Nancy R Sturm, Terry Gaasterland,
870 and Senjie Lin. Spliced leader rna trans-splicing in dinoflagellates. *Proceedings of the National*
871 *Academy of Sciences*, 104(11):4618–4623, 2007.
- 872 [27] Huan Zhang and Senjie Lin. Retrieval of missing spliced leader in dinoflagellates. *PLoS One*, 4(1):
873 e4129, 2009.
- 874 [28] Scott William Roy. Genomic and transcriptomic analysis reveals spliced leader trans-splicing in
875 cryptomonads. *Genome biology and evolution*, 9(3):468–473, 2017.
- 876 [29] Mitsuhiro Matsuo, Atsushi Katahata, Soichirou Satoh, Motomichi Matsuzaki, Mami Nomura, Ken-
877 ichiro Ishida, Yuji Inagaki, and Junichi Obokata. Characterization of spliced leader trans-splicing
878 in a photosynthetic rhizarian amoeba, paulinella micropora, and its possible role in functional gene
879 transfer. *PloS one*, 13(7):e0200961, 2018.
- 880 [30] Jonathan Pettitt, Berndt Müller, Ian Stansfield, and Bernadette Connolly. Spliced leader trans-
881 splicing in the nematode trichinella spiralis uses highly polymorphic, noncanonical spliced leaders.
882 *RNA*, 14(4):760–770, 2008.
- 883 [31] Yuelong Guo, David McK Bird, and Dahlia M Nielsen. Improved structural annotation of protein-
884 coding genes in the meloidogyne hapla genome using rna-seq. In *Worm*, volume 3, page e29158.
885 Taylor & Francis, 2014.
- 886 [32] Peter D Olson, Alan Tracey, Andrew Baillie, Katherine James, Stephen R Doyle, Sarah K Budden-
887 borg, Faye H Rodgers, Nancy Holroyd, and Matt Berriman. Complete representation of a tapeworm
888 genome reveals chromosomes capped by centromeres, necessitating a dual role in segregation and
889 protection. *bioRxiv*, 2020.
- 890 [33] Javier Calvelo, Hernán Juan, Héctor Musto, Uriel Koziol, and Andrés Iriarte. Slfinder, a pipeline for
891 the novel identification of splice-leader sequences: a good enough solution for a complex problem.
892 *BMC bioinformatics*, 21(1):1–18, 2020.
- 893 [34] Nicolas J Tourasse, Jonathan RM Millet, and Denis Dupuy. Quantitative rna-seq meta analysis of
894 alternative exon usage in c. elegans. *Genome research*, pages gr-224626, 2017.
- 895 [35] Carlo Yague-Sanz and Damien Hermand. Sl-quant: a fast and flexible pipeline to quantify spliced
896 leader trans-splicing events from rna-seq data. *GigaScience*, 7(7):giy084, 2018.
- 897 [36] Timothy W Nilsen. Evolutionary origin of sl-addition trans-splicing: still an enigma. *Trends in Ge-*
898 *netics*, 17(12):678 – 680, 2001. ISSN 0168-9525. doi: [https://doi.org/10.1016/S0168-9525\(01\)02499-](https://doi.org/10.1016/S0168-9525(01)02499-4)
899 4. URL <http://www.sciencedirect.com/science/article/pii/S0168952501024994>.

- 900 [37] Daehwan Kim, Ben Langmead, and Steven L Salzberg. Hisat: a fast spliced aligner with low memory
901 requirements. *Nature methods*, 12(4):357, 2015.
- 902 [38] Ben Langmead and Steven L Salzberg. Fast gapped-read alignment with bowtie 2. *Nature methods*,
903 9(4):357, 2012.
- 904 [39] Yutaka Satou, Makoto Hamaguchi, Keisuke Takeuchi, Kenneth EM Hastings, and Nori Satoh. Ge-
905 nomic overview of mrna 5'-leader trans-splicing in the ascidian *ciona intestinalis*. *Nucleic acids*
906 *research*, 34(11):3378–3388, 2006.
- 907 [40] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo
908 Abecasis, and Richard Durbin. The sequence alignment/map format and samtools. *Bioinformatics*,
909 25(16):2078–2079, 2009.
- 910 [41] Torbjørn Rognes, Tomáš Flouri, Ben Nichols, Christopher Quince, and Frédéric Mahé. Vsearch: a
911 versatile open source tool for metagenomics. *PeerJ*, 4:e2584, 2016.
- 912 [42] Christiam Camacho, George Coulouris, Vahram Avagyan, Ning Ma, Jason Papadopoulos, Kevin
913 Bealer, and Thomas L Madden. Blast+: architecture and applications. *BMC bioinformatics*, 10(1):
914 421, 2009.
- 915 [43] Aaron R Quinlan and Ira M Hall. Bedtools: a flexible suite of utilities for comparing genomic
916 features. *Bioinformatics*, 26(6):841–842, 2010.
- 917 [44] Jeffrey D Thomas, Richard C Conrad, and Thomas Blumenthal. The *c. elegans* trans-spliced leader
918 rna is bound to sm and has a trimethylguanosine cap. *Cell*, 54(4):533–539, 1988.
- 919 [45] Stacey N Barnes, Rick E Masonbrink, Thomas R Maier, Arun Seetharam, Anoop S Sindhu, An-
920 drew J Severin, and Thomas J Baum. *Heterodera glycines* utilizes promiscuous spliced leaders and
921 demonstrates a unique preference for a species-specific spliced leader over *c. elegans* sl1. *Scientific*
922 *reports*, 9(1356), 2019.
- 923 [46] Ronny Lorenz, Stephan H Bernhart, Christian Höner Zu Siederdisen, Hakim Tafer, Christoph
924 Flamm, Peter F Stadler, and Ivo L Hofacker. Viennarna package 2.0. *Algorithms for molecular*
925 *biology*, 6(1):26, 2011.
- 926 [47] Donald Evans and Thomas Blumenthal. trans splicing of polycistronic *caenorhabditis elegans* pre-
927 mrnas: Analysis of the sl2 rna. *Molecular and cellular biology*, 20(18):6659–6667, 2000.
- 928 [48] R Core team. R: A language and environment for statistical computing. 2013.
- 929 [49] France Denoeud, Simon Henriët, Sutada Mungpakdee, Jean-Marc Aury, Corinne Da Silva, Henner
930 Brinkmann, Jana Mikhaleva, Lisbeth Charlotte Olsen, Claire Jubin, Cristian Cañestro, et al. Plas-
931 ticity of animal genome architecture unmasked by rapid evolution of a pelagic tunicate. *Science*, 330
932 (6009):1381–1385, 2010.
- 933 [50] Marcel Martin. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EM-*
934 *Bnet. journal*, 17(1):pp–10, 2011.
- 935 [51] Yang Liao, Gordon K Smyth, and Wei Shi. featurecounts: an efficient general purpose program for
936 assigning sequence reads to genomic features. *Bioinformatics*, 30(7):923–930, 2013.

- 937 [52] Cole Trapnell, Brian A Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marijke J Van Baren,
938 Steven L Salzberg, Barbara J Wold, and Lior Pachter. Transcript assembly and quantification by
939 rna-seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature*
940 *biotechnology*, 28(5):511, 2010.
- 941 [53] Mary Ann Allen, LaDeana W Hillier, Robert H Waterston, and Thomas Blumenthal. A global
942 analysis of *c. elegans* trans-splicing. *Genome research*, 21(2):255–264, 2011.
- 943 [54] F William Townes, Stephanie C Hicks, Martin J Aryee, and Rafael A Irizarry. Feature selection and
944 dimension reduction for single-cell rna-seq based on a multinomial model. *Genome biology*, 20(1):
945 1–16, 2019.
- 946 [55] F Krueger. Trim galore. *A wrapper tool around Cutadapt and FastQC to consistently apply quality*
947 *and adapter trimming to FastQ files*, 2015.
- 948 [56] Donald W Nelson and Barry M Honda. Two highly conserved transcribed regions in the 5s dna
949 repeats of the nematodes *caenorhabditis elegans* and *caenorhabditis briggsae*. *Nucleic acids research*,
950 17(21):8657–8667, 1989.
- 951 [57] Lincoln D Stein, Zhirong Bao, Darin Blasiar, Thomas Blumenthal, Michael R Brent, Nansheng
952 Chen, Asif Chinwalla, Laura Clarke, Chris Clee, Avril Coghlan, et al. The genome sequence of
953 *caenorhabditis briggsae*: a platform for comparative genomics. *PLoS Biol*, 1(2):e45, 2003.
- 954 [58] Nadine Borchert, Christoph Dieterich, Karsten Krug, Wolfgang Schütz, Stephan Jung, Alfred Nord-
955 heim, Ralf J Sommer, and Boris Macek. Proteogenomics of *pristionchus pacificus* reveals distinct
956 proteome structure of nematode models. *Genome research*, 20(6):837–846, 2010.
- 957 [59] Jonathan Pettitt, Lucas Philippe, Debjani Sarkar, Christopher Johnston, Henrike Johanna Gothe,
958 Diane Massie, Bernadette Connolly, and Berndt Müller. Operons are a conserved feature of nematode
959 genomes. *Genetics*, 197(4):1201–1211, 2014.
- 960 [60] Debjani Sarkar. *Spliced leader trans-splicing and operons in Dorylaimida (Nematoda)*. PhD thesis,
961 University of Aberdeen, 2014.
- 962 [61] Manfred G Grabherr, Brian J Haas, Moran Yassour, Joshua Z Levin, Dawn A Thompson, Ido Amit,
963 Xian Adiconis, Lin Fan, Raktima Raychowdhury, Qiandong Zeng, et al. Full-length transcriptome
964 assembly from rna-seq data without a reference genome. *Nature biotechnology*, 29(7):644, 2011.
- 965 [62] Limin Fu, Beifang Niu, Zhengwei Zhu, Sitao Wu, and Weizhong Li. Cd-hit: accelerated for clustering
966 the next-generation sequencing data. *Bioinformatics*, 28(23):3150–3152, 2012.
- 967 [63] Peng Cui, Qiang Lin, Feng Ding, Chengqi Xin, Wei Gong, Lingfang Zhang, Jianing Geng, Bing
968 Zhang, Xiaomin Yu, Jin Yang, et al. A comparison between ribo-minus rna-sequencing and poly-
969 a-selected rna-sequencing. *Genomics*, 96(5):259–265, 2010.
- 970 [64] Anna R Dahlgren, Erica Y Scott, Tamer Mansour, Erin N Hales, Pablo J Ross, Theodore S
971 Kalbfleisch, James N MacLeod, Jessica L Petersen, Rebecca R Bellone, and Carrie J Finno. Com-
972 parison of poly-a+ selection and rna depletion in detection of lncrna in two equine tissues using
973 rna-seq. *Non-coding RNA*, 6(3):32, 2020.
- 974 [65] Jun Matsumoto, Ken Dewar, Jessica Wasserscheid, Graham B Wiley, Simone L Macmil, Bruce A
975 Roe, Robert W Zeller, Yutaka Satou, and Kenneth EM Hastings. High-throughput sequence analysis
976 of *ciona intestinalis* sl trans-spliced mrnas: alternative expression modes and gene function correlates.
977 *Genome research*, 20(5):636–645, 2010.

- 978 [66] Amanda E Vandenberghe, Thomas H Meedel, and Kenneth EM Hastings. mrna 5'-leader trans-
979 splicing in the chordates. *Genes & development*, 15(3):294–303, 2001.
- 980 [67] Brendan Yeats, Jun Matsumoto, Sandra I Mortimer, Eiichi Shoguchi, Nori Satoh, and Kenneth EM
981 Hastings. Sl rna genes of the ascidian tunicates *ciona intestinalis* and *ciona savignyi*. *Zoological
982 science*, 27(2):171–180, 2010.
- 983 [68] Nicholas A Stover and Robert E Steele. Trans-spliced leader addition to mrnas in a cnidarian.
984 *Proceedings of the National Academy of Sciences*, 98(10):5693–5698, 2001.
- 985 [69] Aleksandar Rajkovic, Richard E Davis, J Neil Simonsen, and FRITZ M RoTTMAN. A spliced leader
986 is present on a subset of mrnas from the human parasite *schistosoma mansoni*. *Proceedings of the
987 National Academy of Sciences*, 87(22):8879–8883, 1990.
- 988 [70] Ashleigh B Smythe, Oleksandr Holovachov, and Kevin M Kocot. Improved phylogenomic sampling
989 of free-living nematodes enhances resolution of higher-level nematode phylogeny. *BMC evolutionary
990 biology*, 19(1):121, 2019.
- 991 [71] Carol Williams, Lei Xu, and Thomas Blumenthal. Sl1 trans splicing and 3'-end formation in a novel
992 class of *caenorhabditis elegans* operon. *Molecular and Cellular Biology*, 19(1):376–383, 1999. ISSN
993 0270-7306. doi: 10.1128/MCB.19.1.376. URL <https://mcb.asm.org/content/19/1/376>.
- 994 [72] Gemma B Danks, Martina Raasholm, Coen Campsteijn, Abby M Long, J Robert Manak, Boris
995 Lenhard, and Eric M Thompson. Trans-splicing and operons in metazoans: translational control in
996 maternally regulated development and recovery from growth arrest. *Molecular biology and evolution*,
997 32(3):585–599, 2015.
- 998 [73] Kasper D Hansen, Steven E Brenner, and Sandrine Dudoit. Biases in illumina transcriptome se-
999 quencing caused by random hexamer priming. *Nucleic acids research*, 38(12):e131–e131, 2010.
- 1000 [74] Daniel R Garalde, Elizabeth A Snell, Daniel Jachimowicz, Botond Sipos, Joseph H Lloyd, Mark
1001 Bruce, Nadia Pantic, Tigist Admassu, Phillip James, Anthony Warland, et al. Highly parallel direct
1002 rna sequencing on an array of nanopores. *Nature methods*, 15(3):201, 2018.
- 1003 [75] Daniel P Depledge, Kalanghad Puthankalam Srinivas, Tomohiko Sadaoka, Devin Bready, Yasuko
1004 Mori, Dimitris G Placantonakis, Ian Mohr, and Angus C Wilson. Direct rna sequencing on nanopore
1005 arrays redefines the transcriptional complexity of a viral pathogen. *Nature communications*, 10(1):
1006 754, 2019.