

1 **Genome annotation of *Caenorhabditis briggsae* by TEC-RED identifies new exons,**
2 **paralogs, and conserved and novel operons**

3

4 Nikita Jhaveri^{1§}, Wouter van den Berg^{1§}, Byung Joon Hwang^{2,3}, Hans-Michael Muller², Paul W.
5 Sternberg², and Bhagwati P. Gupta^{1*}

6

7 ¹Department of Biology, McMaster University, Hamilton, Ontario L8S 4K1, Canada

8 ²Division of Biology and Biological Engineering, California Institute of Technology, Pasadena,
9 California 91125, USA.

10

11 [§]Equal first authors

12

13 ³Current address: Department of Molecular Bioscience, College of Biomedical Science,
14 Kangwon National University, Chuncheon, South Korea.

15

16

17 * Author for correspondence: guptab@mcmaster.ca; +1-905-525-9140 x26451.

18

19

20 Running head: *C. briggsae* genome annotation and operons

21

22 Keywords: Nematode, *C. briggsae*, *Trans*-splicing, Spliced leader, Operons, Paralog, Genome
23 annotation

24

25

26 ORCIDs:

27 Bhagwati P. Gupta 0000-0001-8572-7054

28 Paul W. Sternberg 0000-0002-7699-0173

29

1 **ABSTRACT**

2 The nematode *Caenorhabditis briggsae* is routinely used in comparative and evolutionary
3 studies involving its well-known cousin *C. elegans*. The *C. briggsae* genome sequence has
4 accelerated research by facilitating the generation of new resources, tools, and functional studies
5 of genes. While substantial progress has been made in predicting genes and start sites,
6 experimental evidence is still lacking in many cases. Here, we report an improved annotation of
7 the *C. briggsae* genome using the *Trans*-spliced Exon Coupled RNA End Determination (TEC-
8 RED) technique. In addition to identifying the 5' ends of expressed genes, we have discovered
9 operons and paralogs. In summary, our analysis yielded 10,243 unique 5' end sequence tags with
10 matches in the *C. briggsae* genome. Of these, 6,395 were found to represent 4,252 unique genes
11 along with 362 paralogs and 52 previously unknown exons. These genes included 14 that are
12 exclusively trans-spliced in *C. briggsae* when compared with *C. elegans* orthologs. A major
13 contribution of this study is the identification of 493 operons, of which two-thirds are fully
14 supported by tags. In addition, two SL1-type operons were discovered. Interestingly,
15 comparisons with *C. elegans* showed that only 40% of operons are conserved. Of the remaining
16 operons, 73 are novel, including 12 that entirely lack orthologs in *C. elegans*. Further analysis
17 revealed that four of the 12 novel operons are conserved in *C. nigoni*. Altogether, the work
18 described here has significantly advanced our understanding of the *C. briggsae* system and
19 serves as a rich resource to aid biological studies involving this species.

20

21

22 **INTRODUCTION**

23 Nematodes are a mainstay in fundamental biological research. While most work has been based
24 on *C. elegans* over the last half a century since its proposed role as a model organism (BRENNER
25 1974), the close relative *C. briggsae* offers many of the same advantages in carrying out studies.
26 Despite diverging roughly 20-30 million years ago (CUTTER 2008), the two species exhibit
27 similar behavioral, developmental, and morphological processes including a hermaphroditic
28 mode of reproduction (GUPTA *et al.* 2007). Moreover, many of the experimental techniques and
29 protocols developed to manipulate *C. elegans* can be adopted to *C. briggsae* with minimal to no
30 modification (BAIRD AND CHAMBERLIN 2006; GUPTA *et al.* 2007). These features make *C.*
31 *briggsae* - *C. elegans* an ideal pair for comparative and evolutionary studies.

1
2 The genome of *C. briggsae* was sequenced many years ago and revealed extensive genomic and
3 genic conservation with *C. elegans* (STEIN *et al.* 2003). Subsequent work reported the assembly
4 of genomic fragments into chromosomes and improved gene predictions (HILLIER *et al.* 2007;
5 ROSS *et al.* 2011). While a diverse array of techniques have been applied to improve the
6 annotation of the *C. elegans* genome (HWANG *et al.* 2004; SPIETH AND LAWSON 2006; HILLIER *et*
7 *al.* 2009; SALEHI-ASHTIANI *et al.* 2009; ALLEN *et al.* 2011), a similar approach is lacking for *C.*
8 *briggsae*. The current *C. briggsae* genome annotation is largely based on homology with the *C.*
9 *elegans* genome. More analysis that uses experimental data gathered directly from *C. briggsae* is
10 needed to improve gene identification and gene models. To this end, we used *trans*-spliced exon
11 coupled RNA end determination (TEC-RED) (HWANG *et al.* 2004), a technique based on
12 exploiting the phenomenon of spliced leader (SL) *trans*-splicing which has been observed in
13 nematodes and several other phyla including platyhelminths, chordates and trypanosomes
14 (LASDA AND BLUMENTHAL 2011).

15
16 The advantage of TEC-RED compared to other genome annotation techniques like EST (MARRA
17 *et al.* 1998) and SAGE (VELCULESCU *et al.* 1995) is that it is capable of identifying transcripts of
18 most expressed genes, and uniquely allows for the identification of 5' transcript start sites and
19 alternative transcripts with different 5' ends of a gene. The approach is based on two principles:
20 one, a short sequence from the 5' end of a transcript can be used to uniquely identify the
21 initiation site of the transcript, and two, the 5' ends of mRNAs are spliced to common leader
22 sequences known as spliced leader (SL) sequences. The SL *trans*-splicing process involves
23 replacing the outtron of a pre-mRNA with a 22 nucleotide SL sequence donated by a 100-
24 nucleotide small ribonucleoprotein (snRNP) (BLUMENTHAL 2005; ALLEN *et al.* 2011). *C.*
25 *elegans* and *C. briggsae* both have two types of spliced leader sequences: SL1 and SL2 (QIAN
26 AND ZHANG 2008; BLUMENTHAL *et al.* 2015).

27
28 We recovered well over 120,000 5' end tags from sequencing reactions representing 10,243
29 unique tags (7,234 for SL1; 3,009 for SL2) with matches in the *C. briggsae* genome. The tags
30 were analyzed using WormBase release WS276 and it was found that more than 60% could be
31 aligned to exons curated in WormBase (www.wormbase.org). Most of the tags were found to

1 have unique hits in the genome and identified a total of 4,252 genes. Other tags with high
2 confidence hits to unannotated regions or to multiple locations of the genome identified 52 novel
3 exons and 362 paralog genes, respectively. The novel exons could either represent previously
4 unknown genes or new exons of existing genes. The paralogs define 133 sets of two or more
5 genes. Of these sets, 21 were confirmed as exact matches with known paralogs in WormBase.
6 The rest could potentially be new paralogous pairs that need further validation. While the
7 majority of the genes discovered by tags confirmed 5' ends of genes listed in WormBase, there
8 are many for which 5' ends indicated by tags differ from current gene models, suggesting the
9 need to revise existing annotations.

10

11 A comparison of the splicing pattern of *C. briggsae* genes with *C. elegans* revealed some
12 changes. Specifically, 14 genes are spliced to leader sequences in *C. briggsae* but their *C.*
13 *elegans* orthologs lack such splicing information. We also investigated the presence of operons.
14 It was reported earlier that 96% of *C. elegans* operons are conserved in *C. briggsae* based on
15 collinearity (STEIN *et al.* 2003). Our analysis revealed a total of 1,199 operons including 493 for
16 which splicing identities of two or more genes are reported in this study. Of these operons, 334
17 are fully supported by tags. Comparison of the latter with *C. elegans* revealed that 40% are
18 conserved, the largest of which is composed of seven genes. Another 38% are termed partially
19 conserved since gene sets do not fully correspond to any of the operons in *C. elegans*. The
20 remaining are novel, i.e., consisting of divergent genes as well as genes whose *C. elegans*
21 orthologs are not reported in operons. Of the divergent operons, four were found to be conserved
22 in a closely related sister species, *C. nigoni*. Lastly, two SL1-type operons have been identified.
23 Overall, the results presented in this study have substantially improved the annotation of the *C.*
24 *briggsae* genome by identifying the 5' ends of a large number of genes as well as discovering
25 novel operons, new exons, and paralogs. The findings strengthen the utility of *C. briggsae* as a
26 model organism and serve as a platform to accelerate comparative and evolutionary studies
27 involving nematodes and other metazoans.

28

29

30 **MATERIALS AND METHODS**

31 *Generation of tags*

1 We followed the protocol described earlier for *C. elegans* (HWANG *et al.* 2004). Briefly, the steps
2 involved purification of poly(A) RNA from the wild type *AF16* mixed stage animals, RT-PCR to
3 generate cDNA, amplification of cDNAs using biotin-attached primers homologous to SL1 and
4 SL2 sequences carrying mismatches to create *BpmI* restriction enzyme site (Supplementary
5 Tables 1-3), digestion of amplified cDNAs using *BpmI* to produce short fragments (termed “5’
6 tags”), ligation of tags to adaptor DNA sequences, and sequential ligation of DNA to create
7 concatenated products. The ligated DNA pieces were cloned into a vector and sequenced.

8

9 ***5' Tag sequence analysis and exon identification***

10 We wrote several Perl scripts to analyze the tags and genes. A flowchart is provided in
11 Supplementary figure 1. Briefly, tags were collected and assigned unique tag IDs. Tag locations
12 in the genome were determined by comparing the tag sequence to WS176 and WS276 genome
13 files, where orientation and chromosome location for each tag was noted. Subsequently the
14 splice sequence for each tag was obtained by finding the seven bases directly upstream of each
15 location where a tag matched on the genome.

16

17 The criteria to identify tag matches to exonic regions were described earlier (HWANG *et al.*
18 2004). These included ‘same orientation of the tag as that of the corresponding exon’, ‘distance
19 to the first ATG’, ‘a minimum distance to the nearest in-frame stop codon’ and ‘presence of a
20 splice acceptor sequence following the tag’. The latter was scored on how well they fit the
21 consensus splice site ‘TTTTCAG’ (BLUMENTHAL AND STEWARD 1997). In cases where tags had
22 multiple matches, we applied stricter splice acceptor site criteria. Perfect consensus sequence
23 was given the highest weight. Sites having mismatches were assigned lower weights with
24 priority given to conserved bases. While this approach resulted in most tags identifying unique
25 exons, a small number still showed multiple matches and were used to search for potential
26 paralogs (see below).

27

28 Each tag was used to find the nearest ATG of an open reading frame (ORF), i.e., the proposed
29 start of a coding sequence (CDS). This ATG location was compared to known coordinates of
30 start sites of nearest exons as annotated in WS176 and WS276 genome annotation (gff3) files.

31 This was done using coordinates of annotated CDS. Two broad categories of exon matches were

1 identified based on tags that had unique matches: one, where the 5' end corresponded to the start
2 of a known exon (first exon: 1a, internal exon: 1b) and two, matches for which the 5' end
3 differed from a nearest exon (Figure 1). Depending on the distance between the 5' end and the
4 exon, the second category of matches were further divided into two sub-categories. These
5 consisted of exons that were either within 20 bp from the 5' end ('minor misprediction') or
6 further away ('major misprediction'). The major misprediction sub-category also includes
7 matches where 5' ends were more than 3 kb away and may define brand new exons of existing
8 genes as well as potentially new, previously unknown genes.

9

10 ***Manual curation of genes***

11 We found that 75 tag-matched genomic regions in the WS276 gff3 file had no known
12 genes/exons within 3kb downstream of the matched ATG. The surrounding chromosomal
13 regions of these matches were confirmed by manually searching the WormBase genome browser
14 for presence of annotated exons. Of the 75, 21 were found to be false positives due to incorrect
15 script calls. Two were excluded from analysis because the genes are not assigned to any
16 chromosomes. The remaining 52 matches may represent novel exons.

17

18 ***Analysis of intergenic regions and operons***

19 The distance between two genes, termed 'intergenic region' (IGR) was determined based on the
20 distance from the end of the 3' UTR of an upstream gene to 5' start of the coding sequence of the
21 immediate downstream gene. Graphs were generated using Graphpad Prism 7.0 and Microsoft
22 Excel. Genes having IGR >5000 bp (257) were excluded from the analysis. For pairs of genes
23 where the second gene is located within the first gene, IGR length is calculated as a negative
24 value. Intercistronic regions (ICRs) were calculated in the same manner. The ICR analysis was
25 done only for genes identified to be part of an operon.

26

27 To identify genes that could be present in operons, all genes *trans*-spliced with SL2 or SL1/SL2
28 and present downstream of an SL1-spliced gene were categorized into a single operon model
29 along with the upstream SL1 spliced gene. We based our assumption of genes being in an operon
30 together on this pattern of splice leader sequences. If the splicing of the first upstream gene was
31 unknown, the operon models were termed 'non-tag-supported' whereas those models in which

1 the identity of the first upstream gene was known were termed ‘tag-supported’. We compared the
2 ‘tag-supported’ operon models to those in *C. elegans* (WormBase) to determine how well
3 operons are conserved. Based on the conservation of genes, the ‘tag-supported’ operons were
4 classified into Exact match, Partial match, and Novel.

5
6 We examined the enrichment of germline genes in *C. briggsae* high confidence operons. For
7 this, *C. elegans* orthologs were identified and researched for association with germline function
8 (WANG *et al.* 2009). The significance of overlap was tested by the hypergeometric probability
9 test. Next, to identify processes related to genes in operons, Gene Ontology (GO) (ASHBURNER
10 *et al.* 2000) analysis was carried out for all operon genes. A similar analysis was conducted for
11 genes present in *C. elegans* operons using a published data set (ALLEN *et al.* 2011).

12 13 ***Paralog analysis***

14 A total of 203 tags had multiple hits in the genome. Since many of these consisted of overlapping
15 sequences, we retained only the longest tags. This filtering step narrowed down the count to 158.
16 The genes identified by these tags were compared with annotated paralogs in WormBase. The
17 matches allowed us to place the predicted paralogs into three different categories.

18 19 ***Uniquely spliced C. briggsae genes***

20 To identify the genes that are trans splicing in *C. briggsae* but not in *C. elegans*, we used datasets
21 published by two different groups that together constitute the most complete collection of *trans*-
22 spliced genes in *C. elegans* (ALLEN *et al.* 2011; TOURASSE *et al.* 2017). Initial comparisons with
23 the Allen *et al.* dataset revealed 198 genes that are present only in our analysis. The number was
24 further reduced to 14 genes when compared with the Tourasse *et al.* study (Supplementary data
25 file 3).

26 27 28 **RESULTS**

29 **Overview of the TEC-RED method in *C. briggsae***

30 To implement the TEC-RED approach to identify transcripts, we first isolated *C. briggsae*
31 mRNAs containing an SL1 or SL2 sequence at their 5' ends. A total of 121,189 5' tags (91,733

1 for mRNA with an SL1 and 29,456 for mRNA with an SL2 spliced leader sequence) were
2 recovered from DNA sequencing reactions. These tags represent 14,678 different sequences, of
3 which 10,400 (71%) are for SL1 and 4,278 (29%) for SL2 sites. More than two-thirds of all tags
4 found matches in the genome (10,243, 70%), of which 46% are unique, i.e., matching only once
5 and others matching multiple times (Table 1). The proportions were similar for both spliced
6 leader categories, demonstrating no bias in the experimental protocol. The remaining 4,435 tags
7 (30%) had no match, likely due to reasons such as sequencing errors, gaps in the genome
8 sequence, and incorrect sequence assembly.

9

10 **Exon validations and predictions in *C. briggsae* based on 5' tag matches**

11 After filtering the matches (see Methods), 62.5% of all tags (6,395 of 10,243) were retained for
12 further analysis (Table 2). Next, we determined the locations of these tags relative to annotated
13 exons in WormBase. Most of the tags (6,192, 96.8%) matched uniquely to one exon, with a
14 small number having multiple matches (203, 3.2%) (Supplementary data file 1). For both SL1
15 and SL2 tags, roughly 80% of the matches correspond to known 5' ends of annotated genes
16 (Category 1a), providing support to existing gene models in WormBase. Less than one percent of
17 the tags matched to internal exons (Category 1b), suggesting an alternate 5' end of the
18 corresponding genes. The remaining tags identified start sites that differed from current
19 WormBase gene models and were categorized as mispredicted genes. In most of these cases
20 (roughly three-quarters of all mispredictions) the nearest exon was more than 20 bp away. This
21 leads us to suggest that, particularly in these cases, existing gene models may need to be revised.
22 These exons may define new 5' ends of known genes as well as novel, previously unidentified
23 genes. More experiments are needed to investigate these possibilities. As expected, both types of
24 tags, i.e., with unique and multiple hits have a similar distribution of categories (Figure 2,
25 Supplementary data file 1).

26

27 **Identification of genes based on tag matches**

28 Next, we compiled a list of *C. briggsae* genes based on exons identified by unique tags. A total
29 of 4,252 genes were recovered by SL1 and SL2 tags (Supplementary data file 2). Almost two-
30 thirds of the genes (65%) are spliced with SL1 and 18% with SL2. Another 18% of exons
31 matched with both SL1 and SL2 tags (SL1/SL2), suggesting the genes are part of hybrid operons

1 (Allen et al., 2011) (Table 3, Supplementary data file 2). Based on their genomic locations, these
2 genes are roughly evenly distributed on the chromosomes except for Chr. X which had the
3 lowest gene count. However, the trend was different for gene density with Chr. III being the
4 densest chromosome and Chr. X the sparsest (Supplementary table 4). Whether the uneven
5 distribution is by chance or a characteristic of *trans*-spliced genes in *C. briggsae* remains to be
6 seen. A tiny fraction of genes (0.1%) is located in unmapped genomic fragments.

7
8 The recovery of *C. briggsae* genes prompted us to examine evolutionary changes in *trans*-
9 splicing. A comparison with *C. elegans* studies (ALLEN *et al.* 2011; TOURASSE *et al.* 2017)
10 revealed 14 genes that appear to be uniquely spliced to leader sequences in *C. briggsae* but not in
11 *C. elegans* (Supplementary data file 3).

12
13 Next, we searched for transcripts resulting from *cis*-splicing of the *C. briggsae* genes. Almost
14 95% of the curated genes (4,025 of 4,252) were found to be associated with unique tag
15 sequences, i.e., 5' ends matched to just one exon, providing support for the presence of a single
16 transcript for these genes (Table 4). In the majority of cases (82%, 3,288 of 4,025), the tag-
17 identified 5' ends matched with a known first exon (category 1a tags). Less than one percent of
18 the tags identify 5' ends that match with internal exons (category 1b). The remaining genes
19 (18%) consist of exons belonging to minor and major misprediction categories. The rest of the
20 genes (5%, 227 of 4,252) identified by tags consist of those that produce multiple transcripts
21 (Table 5). In 84% of these cases, at least one 5' end identified by tags matched with the first
22 exon (category 1a). Five of the genes were alternatively spliced using internal exons as the 5'
23 start site (category 1b). Most of the genes consisted of at least one major mispredicted exon,
24 suggesting that genes with multiple splice variants require further validation.

25
26 As mentioned above, 203 tags had multiple matches in the genome. Further analysis narrowed
27 down the set to 158 unique sequences (see Methods). We reasoned that these tags may represent
28 paralogs and performed searches in WormBase. The analysis identified 133 potential paralog
29 sets consisting of 362 genes. These sets fall into three distinct categories (Supplementary data
30 file 10). Paralogs that fully matched with WormBase annotation were termed 'Exact Match' (21
31 paralogous sets, 42 genes). The other sets matched only partially or did not match to paralog sets

1 recorded on WormBase (Partial Match: 66 sets, 174 genes; No Match: 46 sets, 146 genes). It is
2 worth mentioning that about half of the genes in the No Match category have no paralogous
3 information available, whereas the remaining half have paralogs in WormBase but these did not
4 match with our analysis. To further validate the paralogous relationships, we determined
5 chromosomal locations of the genes. Gene duplications arising from mechanisms such as
6 slipped-strand mispairing can cause the creation of paralogous genes in adjacent stretches of
7 sequence on the same chromosome (LEVINSON AND GUTMAN 1987). In *C. elegans*, paralogs
8 originating from gene duplications are more than twice as likely to be present on the same
9 chromosome and tend to be located closely together (SEMPLE AND WOLFE 1999). Additionally,
10 studies in humans and other higher eukaryotes have revealed that intergenic distances between
11 paralogous genes are smaller than random gene pairs on the same chromosome (IBN-SALEM *et*
12 *al.* 2017). Of the paralog sets identified in this study, 63% (84 sets) were present on the same
13 chromosome including 35% (16 sets) that belong to the No Match category. The IGR analysis
14 revealed that the distances in five cases are less than 10 kb (Supplementary table 5), which is
15 more than 500-fold shorter than the average distance between a random pair of genes on the
16 same chromosome (5.58 +/- 0.89 Mb in *C. elegans*) (LEE AND SONNHAMMER 2003).

17

18 **Validations of TEC-RED-identified transcripts**

19 We took three different approaches to validate subsets of TEC-RED predictions with the goal of
20 demonstrating the usefulness of the technique in improving gene identification and gene models.
21 One approach involved comparing different categories of tag-identified exons between two
22 WormBase releases. As described above, a significant number of exons are categorized as minor
23 and major mispredictions (22%, 943 of 4,252; see Tables 4 and 5). We hypothesized that
24 mispredicted exons may be confirmed with improvements in genome annotation. To test this
25 hypothesis, 1a category of transcripts were compared with those reported in an older WormBase
26 release (WS176). The analysis involved SL1 spliced transcripts belonging to category 1a (2,143
27 (Table 4). As expected, a vast majority of the genes (74%, 1,583) are in category 1a in both
28 releases, providing support for these gene models (Figure 3A, 3D, Supplementary data file 4).
29 The next two largest categories consist of genes that are mispredicted (11.7%, 250 genes) and
30 newly predicted, i.e., absent in WS176 (13.2%, 286 genes). Few genes (0.5%, 11) have start sites
31 that correctly match with internal 5' ends of internal exons. The rest (0.5%, 14 genes) could not

1 be uniquely placed into a single category since these had multiple tag matches in the older
2 annotation. Roughly similar results were obtained by analyzing 1a category of SL2 spliced and
3 SL1/SL2 spliced genes (182 genes, 115 genes, respectively) (Table 4; Figure 3B-D;
4 Supplementary data file 4). Altogether, 858 annotation improvements are supported by our
5 analysis. The demonstrated improvements in gene identification and genome annotation as
6 observed in WS276 prove the accuracy of our 5' start site determination method.

7
8 The second type of validation focused on a subset of the major misprediction category of genes
9 whose 5' ends mapped more than 3 kb away from nearest exons. Most of these (94%, 49 of 52)
10 are in intergenic regions (Supplementary data file 5). 37% (19 of 52) of the exons are supported
11 by RNA sequencing reads (WormBase), providing proof of accuracy to our method
12 (Supplementary figure 2). These novel exons are likely to either belong to nearby existing genes
13 or define brand new genes.

14
15 The last set of validations consisted of comparisons with *C. elegans* gene models. In this case
16 category 1b of single and multiple transcripts (Tables 4 and 5 respectively) were manually
17 examined. The results showed that 38% of newly discovered 5' ends (6 single transcript and two
18 multiple transcripts) are supported by *C. elegans* orthologs (Supplementary figure 3,
19 Supplementary data file 6), providing further support to our analysis. We took a similar approach
20 to analyze a subset of transcripts in the major mispredictions category. Of the 10% of such
21 predictions that were tested, 34% (17 of 50) are supported by WormBase *C. elegans* gene
22 models. With this success rate, another 115 of the remaining single transcript genes of the major
23 misprediction category are likely to be validated. Overall, the 5' tag analysis serves as a rich
24 resource to improve the *C. briggsae* genome annotation.

25 26 **Discovery of operons**

27 The identification of genes based on unique tag matches in *C. briggsae* allowed us to search for
28 operons. In *C. elegans* it has been shown that the first gene in an operon is SL1 spliced (CONRAD
29 *et al.* 1991), whereas downstream genes are spliced either with SL2, SL2 variants or both SL1
30 and SL2 (BLUMENTHAL 2005). Genes that are both SL1 and SL2 spliced cause the operon they

1 are part of to be considered as ‘hybrid operon’. Ultimately, global analysis of *trans*-splicing in *C.*
2 *briggsae* will reveal all operons and operon genes.

3
4 Our data suggests the existence of up to 1,199 *C. briggsae* operons (Table 6, Supplementary data
5 file 7). These include 334 operons that are fully supported by tags, i.e., we were able to
6 determine the splicing pattern of every gene, with operons ranging from two to seven genes The
7 remaining 865 operons (ranging between two to six genes) are categorized as ‘Predicted
8 operons’ since the splicing identity of the first gene in these cases remains to be determined. In
9 this set, the predicted operons that contain three or more genes (159) are large enough to be
10 labeled as bona fide operons. Added together with the 334 fully supported operons, this allows
11 us to report at least 493 operons in *C. briggsae* with sufficient certainty.

12
13 In *C. elegans*, operon genes tend to be very closely spaced, typically having an ICR of less than
14 1 kb (ALLEN *et al.* 2011; BLUMENTHAL *et al.* 2015). To examine whether the same is true in *C.*
15 *briggsae*, we calculated ICRs and found that a vast majority of the genes (78%) are separated by
16 less than 200 bp (Figure 4). We also determined intergenic distances (IGRs) for SL1 and
17 SL1/SL2 hybrid spliced genes discovered in our study. The results suggested that the IGR to the
18 nearest gene upstream of SL2-spliced genes is smaller compared to those spliced with SL1 and
19 SL1/SL2. While the SL2-spliced genes have a median distance of 180 bp, the medians of SL1
20 and SL1/SL2 spliced genes are 4,631 bp and 1,242 bp, respectively (Figure 5A). Furthermore, as
21 we would expect, genes with larger IGRs are more likely spliced with SL1 than SL2 or SL1/SL2
22 (Figure 5B, Supplementary data file 8) and are thus less likely to be part of the same operon.

23

24 Tag-supported operons

25 We examined the conservation of tag-supported operons in *C. elegans*. The analysis of orthologs
26 helped define three distinct categories (Supplementary data file 7). The two largest categories are
27 named ‘exact match’ and ‘partial match’ operons (40% and 38%, respectively). Exact match
28 operons consist entirely of *C. elegans* orthologs, whereas in partial match operons only some of
29 the genes are conserved. The remaining one-fifth of operons define a third category, named
30 ‘novel’ (73). While a majority of these (61, 18%) consist of conserved genes whose orthologs are

1 not present in *C. elegans* operons, others (12, 4%) consist of divergent, *C. briggsae*-specific
2 genes.

3
4 The largest *C. briggsae* operons (CBROPX0001) consists of seven genes, six of which
5 (*CBG25571*, *CBG03062*, *CBG25572*, *CBG03061*, *CBG03060*, *CBG03059*) are conserved in *C.*
6 *elegans* and are part of the orthologous operon CEOP2496. The 5th gene in CBROPX0001
7 (*CBG25573*) does not appear to have a *C. elegans* ortholog. Syntenic alignments revealed that
8 *CBG25573* is conserved in *C. brenneri*, suggesting that the gene may have been lost in the *C.*
9 *elegans* lineage (Supplementary figure 4). While we did not recover a tag for *Cbr-rpb-6*
10 (*CBG03063*), whose ortholog is the first gene in CEOP2496, we hypothesize that it is part of *C.*
11 *briggsae* operon CBROPX0001 based on the distance from its neighbor *CBG25571* (195 bp)
12 (Figure 6). More experiments are needed to confirm if *Cbr-rpb-6* is the eighth gene in
13 CBROPX0001.

14
15 A few operons were manually updated. For example, CBROP0002 and CBROPX0002 were split
16 based on homology information in *C. elegans*, resulting in four different operons: CBROP0002A
17 (*CBG02635*, *CBG02634*), CBROP0002B (*CBG02633*, *CBG02632*), CBROP0132 (*CBG01778*,
18 *CBG31146*, *CBG01779*), and CBROP0133 (*CBG01783*, *CBG01784*). In a different case,
19 CBROPX0007 is predicted to consist of four genes (*CBG03212*, *CBG03213*, *CBG03214*, and
20 *CBG03215*) (Supplementary Figure 5). The *C. elegans* orthologs of these genes constitute two
21 distinct operons (CEOP2396 and CEOP2749) (Figure 7). While the ICR between *CBG03213* and
22 *CBG03214* is larger than 2 kb, all downstream genes in CBROPX0007 are either SL2 or
23 SL1/SL2 spliced. Further experiments are needed to validate the structure of CBROPX0007.
24 Table 7 lists the updated numbers of operons in each category.

25
26 We also analyzed partially conserved operons in some detail. While all of these contain *C.*
27 *elegans* orthologs, their structures are not conserved. Specifically, the number of genes or some
28 of the orthologs in corresponding operons differ between the two species (Supplementary data
29 file 7). Of the 127 such operons, 82 contain two or more conserved genes including 58 (70% of
30 83) with less than 1 kb ICR between every gene. One such operon (CBROPX0003) consists of

1 five genes (Figure 8). While the *C. elegans* operon CEOP1484 contains orthologs of all of these,
2 CEOP1484 encompasses three additional genes.

3
4 Our tag searches identified 73 novel operons (Supplementary data file 7). A majority of these
5 (61, 84%) consist of a mix of conserved genes and those that lack orthologs in *C. elegans*. It is
6 important to point out that none of the conserved genes are part of *C. elegans* operons. The other
7 12 (17%) operons consist entirely of genes that lack orthology in *C. elegans*. In seven of these
8 cases, ICRs are less than 1 kb, providing further support to the operon structures (Table 8).

9
10 To investigate whether the 12 novel operons are unique to *C. briggsae* or might be conserved,
11 the analysis was extended to *C. nigoni*, a sister species to *C. briggsae* (WOODRUFF *et al.* 2010).
12 For this, we manually searched 5' upstream regions of orthologs with $\geq 90\%$ sequence similarity
13 to 5' tags and splice acceptor sites of *C. briggsae* genes. The sequence searches revealed that
14 four of the operons have orthologs in the same genomic order with highly similar splice site
15 sequences and small ICRs (≥ 1200 bp), suggesting that they are conserved (Supplementary data
16 file 11).

17 18 Predicted (Non-tag-supported) operons

19 We report 865 predicted operons (Supplementary data file 7). While the downstream genes in
20 these cases are spliced either with SL2 or SL1/SL2, the splicing status of the upstream gene is
21 unknown. Most, if not all, of these are predicted to be genuine operons, especially those that are
22 larger, i.e., consist of more than two genes. A comparison with *C. elegans* of 159 operons
23 containing three or more genes revealed that 26 (16%) are fully conserved. A couple of examples
24 include CBROPX0206 (5 genes) (Figure 9A, C) and CBROPX0207 (5 genes) (Figure 9D). The
25 corresponding *C. elegans* operons are CEOP4500 (six genes) (Figure 9B, C) and CEOP5248
26 (seven genes) (Figure 9E). Comparison of genes in CBROPX0206 and CEOP4500 revealed that
27 these share four orthologs. We observed two additional differences between CBROPX0206 and
28 CEOP4500: the order of genes has changed and CBROPX0206 includes *CBG26297* which
29 appears to lack a *C. elegans* ortholog (Figure 9C). Given that *CBG06240* and *CBG36241* are
30 immediately upstream of CBROPX0206 and their orthologs are part of CEOP4500, the *C.*
31 *briggsae* operon may be extended to include both these genes. However, we have excluded these

1 from our operon model in the absence of corresponding TEC-RED tags. The second example,
2 CBROPX0207, contains five genes, all of which have orthologs in CEOP5248. However, the *C.*
3 *elegans* operon contains two additional genes (*ZK856.16* and *ZK856.19*) which are not conserved
4 in *C. briggsae*.

6 SL1-type operons

7 We also found two operons that contain only SL1-spliced genes. These genes are positioned
8 directly adjacent to one another, with no space between them. The SL1-type operons have been
9 described previously in *C. elegans* and shown to lack ICR (WILLIAMS *et al.* 1999). One of the *C.*
10 *briggsae* SL1-type operons consists of two genes: CBROP0134 (*CBG16825*, *Cbr-vha-*
11 *11/CBG16826*). Its *C. elegans* ortholog, CEOP4638, also consists of two genes. Another SL1-
12 type operon identified by our study is CBROPX0001. Its *C. elegans* ortholog is CEOP2496.
13 Interestingly, CBROPX0001 and CEOP2496 consist of more than two genes (Figure 6). In the
14 case of CEOP2496, the first two genes (*rpb-6* and *dohh-1*) are known to be spliced exclusively
15 with SL1 (defined as SL1 operon) whereas the remaining downstream genes with SL2 or both
16 SL1 and SL2.

17
18 There is also a potential SL1-type operon consisting of *CBG03984* and *CBG03983*. These two
19 genes have a single base pair IGR (Figure 10). Interestingly, the *C. elegans* orthologs, *F23C8.6*
20 and *F23C8.5* (SL1 and SL1/SL2 spliced, respectively) are part of one operon, CEOP1044, with
21 an ICR of more than 400 bp (Allen *et al.*, 2011). More work is needed to determine whether the
22 *C. briggsae* genes are indeed part of an SL1-type operon.

24 *C. briggsae* operons show enrichment of germline genes and highly expressed growth genes

25 Studies in *C. elegans* and *P. pacificus* have reported that germline genes are overrepresented in
26 operons (REINKE AND CUTTER 2009; SINHA *et al.* 2014). We did a gene-association study in *C.*
27 *briggsae* to examine a similar possibility. The results revealed a significant enrichment of
28 germline genes in high confidence operons ($p < 7.40E-98$) (Supplementary data file 9).

29
30 In addition to investigating germline genes, we performed GO term analysis of operon genes and
31 found enrichment of terms associated with metabolic and biosynthesis processes. The pattern of

1 enrichment was similar to what was observed with a *C. elegans* operon dataset (Supplementary
2 data file 9). We also found enrichment of growth-related genes, as in *C. elegans*, specifically,
3 female gamete generation (GO:0007292), embryo development ending in birth or egg hatching
4 (GO:0009792), reproduction (GO:0000003) and embryo development (GO:0009790)
5 (ZASLAVER *et al.* 2011). It is important to point out that while GO terms are similar in both
6 species, *C. briggsae* operon genes associated with specific processes are not always the orthologs
7 of *C. elegans* gene sets. We therefore conclude that functions of operon genes are conserved
8 even if specific genes are not.

9

10 **DISCUSSION**

11 This paper reports major improvements in the genome annotation of *C. briggsae*. We recovered
12 10,243 unique 5' end tags with matches in the genome, of which 6,395 correspond to SL1 and
13 SL2 spliced exons and provide support to the existence of 4,252 unique *trans*-spliced genes.
14 Another 362 genes have been identified as paralogs, including 42 for which the paralogous
15 relationship is supported by WormBase annotation. We also report 52 novel exons that may
16 define new genes or exons of existing genes. Figure 11 provides a global overview of sequences
17 identified in our analysis.

18

19 In *C. elegans*, 84% of all genes are spliced to leader sequences (TOURASSE *et al.* 2017). If the
20 percentage is comparable in *C. briggsae*, then our work has resulted in the identification of
21 roughly one-quarter of all *trans*-spliced genes in this species. Further analysis has revealed that
22 two-thirds of all *C. briggsae* genes are spliced with SL1, while the rest are split evenly between
23 SL2 and SL1/SL2 hybrid sequences (65% SL1, 18% SL2 and 18% SL1/SL2). Assuming that the
24 TEC-RED method is unbiased in regard to the recovery of SL1 and SL2 spliced transcript tags,
25 the proportion of spliced genes in *C. briggsae* differs from those in *C. elegans* as reported by
26 Allen *et al.* (ALLEN *et al.* 2011), (82% SL1, 12% SL2 and 8% SL1/SL2). Additionally, 14 genes
27 were found to be spliced to leader sequences only in *C. briggsae* and not in *C. elegans*. More
28 work is needed to determine if *trans*-splicing of these genes has indeed diverged between the two
29 species.

30

1 Our analysis revealed that most of the genes identified by unique tag matches are represented by
2 a single transcript (94.8%) and very few (5.2%) by multiple transcripts. Studies in *C. elegans*
3 have reported roughly 18% of genes giving rise to multiple isoforms (WANG *et al.* 2010; SPIETH
4 *et al.* 2014), although this number is predicted to be as high as 25% (RAMANI *et al.* 2011;
5 ZAHLER 2012). Considering this, along with the fact that our experiments captured only a partial
6 set of all spliced genes, the actual proportion of genes with multiple transcripts in *C. briggsae* is
7 likely to be much higher. Among other things, it was found that 77.8% of the genes in our study
8 have 5' start sites that match with those annotated by WormBase. The remaining ones were
9 considered mispredictions, most of which were major mispredictions (15.7%) as 5' start sites in
10 these cases map anywhere between 20 bp to 3 kb away from known locations. We also found 52
11 new, previously unreported exons that map more than 3 kb upstream to the nearest exon of
12 existing genes, and potentially include some that define 5' start sites of new genes.

13

14 Several approaches were taken to validate tag-based gene models. One involved comparing
15 results with those obtained using an older genome annotation (gff) file. The findings revealed
16 that a total of 858 genes for which 5' ends were correctly annotated in WS276 were mispredicted
17 or absent in the older version, which demonstrates that our data can help improve start sites of
18 many *C. briggsae* genes. Another approach involved comparing 5' ends of some of the genes
19 with those of *C. elegans* orthologs. Of the 21 alternate start sites and 50 major mispredicted start
20 sites analyzed, 38% and 34%, respectively, are supported by *C. elegans* transcripts. Finally, we
21 examined the 52 newly discovered exons and found that 37% of these are supported by RNA-seq
22 data in WormBase. The above three validations provide significant support to the accuracy of
23 our analysis of expressed transcripts in *C. briggsae*.

24

25 The identification of genes spliced with leader sequences in *C. briggsae* allowed us to curate
26 operons and study their conservation. Even though the operon-based organization of genes in *C.*
27 *elegans* and *C. briggsae* is similar to those found in bacteria and archaea, work in *C. elegans* has
28 shown that worm operons have no ancestral relationship with prokaryotes and appear to have
29 evolved independently within the nematode phylum (BLUMENTHAL 2004; QIAN AND ZHANG
30 2008). We identified a total of 1,199 operons, of which 28% consist entirely of tag-supported
31 genes. Of the remaining operons with partial tag support, 159 contain three or more genes.

1 Combined with the fully tag-supported operons, this totals to 493 operons in *C. briggsae* with a
2 high degree of confidence. Comparison of tag-supported operons with *C. elegans* revealed that
3 134 (40%) are conserved, with the remainder being partially conserved (127, 39%) and novel
4 (73, 21%). A subset of novel operons (12, 17%) consists entirely of genes that lack *C. elegans*
5 orthologs. Further comparisons with *C. nigoni* revealed that four of the twelve are likely to be
6 conserved, suggesting that these might have arisen in the common ancestor shared between *C.*
7 *briggsae* and *C. nigoni*. Whether the remaining eight are unique to *C. briggsae* requires more
8 analysis. Along with the above-mentioned operons, we also uncovered two conserved SL1-type
9 operons. Together, these data demonstrate that while many of the operons are conserved, there
10 are substantial differences between the two species. The findings represent the first
11 comprehensive analysis of operons in *C. briggsae*.

12

13 In conclusion, the data presented in this study has significantly improved the annotation of the *C.*
14 *briggsae* genome by validating existing gene models, refining start sites of many genes,
15 identifying novel gene exons, alternate transcripts, and by providing a comprehensive analysis of
16 operons and paralogous gene sets. While the majority of the genes and operons are conserved in
17 *C. elegans*, our work has also revealed substantial differences between the two species. The
18 improvements to the genome annotation reported here are expected to strengthen *C. briggsae* as
19 a model for comparative and evolutionary studies.

20

21

22 **DATA AVAILABILITY**

23 The data underlying this article are available in the article and in its online supplementary
24 material.

25

26

27 **ACKNOWLEDGEMENTS**

28 We thank WormBase for assistance with some aspects of data analysis, Mary Ann Allen and
29 Tom Blumenthal for discussions on *C. elegans* operons, and members of the Gupta lab for
30 feedback on the manuscript. This work was supported by grants to BPG (NSERC Discovery) and

1 PWS (U24-HG002223). PWS was an Investigator with the HHMI, which partially supported this
2 work.
3

1 **REFERENCES**

- 2
- 3 Allen, M. A., L. W. Hillier, R. H. Waterston and T. Blumenthal, 2011 A global analysis of *C.*
4 *elegans* trans-splicing. *Genome Res* 21: 255-264.
- 5 Ashburner, M., C. A. Ball, J. A. Blake, D. Botstein, H. Butler *et al.*, 2000 Gene ontology: tool
6 for the unification of biology. The Gene Ontology Consortium. *Nat Genet* 25: 25-29.
- 7 Baird, S. E., and H. M. Chamberlin, 2006 *Caenorhabditis briggsae* methods. *WormBook*: 1-9.
- 8 Blumenthal, T., 2004 Operons in eukaryotes. *Brief Funct Genomic Proteomic* 3: 199-211.
- 9 Blumenthal, T., 2005 Trans-splicing and operons. *WormBook*: 1-9.
- 10 Blumenthal, T., P. Davis and A. Garrido-Lecca, 2015 Operon and non-operon gene clusters in
11 the *C. elegans* genome. *WormBook*: 1-20.
- 12 Blumenthal, T., and K. Steward, 1997 RNA Processing and Gene Structure in *C. elegans II*,
13 edited by nd, D. L. Riddle, T. Blumenthal, B. J. Meyer and J. R. Priess, Cold Spring
14 Harbor (NY).
- 15 Brenner, S., 1974 The genetics of *Caenorhabditis elegans*. *Genetics* 77: 71-94.
- 16 Conrad, R., J. Thomas, J. Spieth and T. Blumenthal, 1991 Insertion of part of an intron into the 5'
17 untranslated region of a *Caenorhabditis elegans* gene converts it into a trans-spliced gene.
18 *Mol Cell Biol* 11: 1921-1926.
- 19 Cutter, A. D., 2008 Divergence times in *Caenorhabditis* and *Drosophila* inferred from direct
20 estimates of the neutral mutation rate. *Mol Biol Evol* 25: 778-786.
- 21 Gupta, B. P., R. Johnsen and N. Chen, 2007 Genomics and biology of the nematode
22 *Caenorhabditis briggsae*. *WormBook*, ed. The *C. elegans* Research Community
23 doi/10.1895/wormbook.1.136.1, <http://www.wormbook.org>.
- 24 Hillier, L. W., R. D. Miller, S. E. Baird, A. Chinwalla, L. A. Fulton *et al.*, 2007 Comparison of
25 *C. elegans* and *C. briggsae* Genome Sequences Reveals Extensive Conservation of
26 Chromosome Organization and Synteny. *PLoS Biol* 5: e167.
- 27 Hillier, L. W., V. Reinke, P. Green, M. Hirst, M. A. Marra *et al.*, 2009 Massively parallel
28 sequencing of the polyadenylated transcriptome of *C. elegans*. *Genome Res* 19: 657-666.
- 29 Hwang, B. J., H. M. Muller and P. W. Sternberg, 2004 Genome annotation by high-throughput 5'
30 RNA end determination. *Proc Natl Acad Sci U S A* 101: 1650-1655.
- 31 Ibn-Salem, J., E. M. Muro and M. A. Andrade-Navarro, 2017 Co-regulation of paralog genes in
32 the three-dimensional chromatin architecture. *Nucleic Acids Res* 45: 81-91.
- 33 Lasda, E. L., and T. Blumenthal, 2011 Trans-splicing. *Wiley Interdiscip Rev RNA* 2: 417-434.

- 1 Lee, J. M., and E. L. Sonnhammer, 2003 Genomic gene clustering analysis of pathways in
2 eukaryotes. *Genome Res* 13: 875-882.
- 3 Levinson, G., and G. A. Gutman, 1987 Slipped-strand mispairing: a major mechanism for DNA
4 sequence evolution. *Mol Biol Evol* 4: 203-221.
- 5 Marra, M. A., L. Hillier and R. H. Waterston, 1998 Expressed sequence tags--ESTablishing
6 bridges between genomes. *Trends Genet* 14: 4-7.
- 7 Qian, W., and J. Zhang, 2008 Evolutionary dynamics of nematode operons: easy come, slow go.
8 *Genome Res* 18: 412-421.
- 9 Ramani, A. K., J. A. Calarco, Q. Pan, S. Mavandadi, Y. Wang *et al.*, 2011 Genome-wide
10 analysis of alternative splicing in *Caenorhabditis elegans*. *Genome Res* 21: 342-348.
- 11 Reinke, V., and A. D. Cutter, 2009 Germline expression influences operon organization in the
12 *Caenorhabditis elegans* genome. *Genetics* 181: 1219-1228.
- 13 Ross, J. A., D. C. Koboldt, J. E. Staisch, H. M. Chamberlin, B. P. Gupta *et al.*, 2011
14 *Caenorhabditis briggsae* recombinant inbred line genotypes reveal inter-strain
15 incompatibility and the evolution of recombination. *PLoS Genet* 7: e1002174.
- 16 Salehi-Ashtiani, K., C. Lin, T. Hao, Y. Shen, D. Szeto *et al.*, 2009 Large-scale RACE approach
17 for proactive experimental definition of *C. elegans* ORFeome. *Genome Res* 19: 2334-
18 2342.
- 19 Semple, C., and K. H. Wolfe, 1999 Gene duplication and gene conversion in the *Caenorhabditis*
20 *elegans* genome. *J Mol Evol* 48: 555-564.
- 21 Sinha, A., C. Langnick, R. J. Sommer and C. Dieterich, 2014 Genome-wide analysis of trans-
22 splicing in the nematode *Pristionchus pacificus* unravels conserved gene functions for
23 germline and dauer development in divergent operons. *RNA* 20: 1386-1397.
- 24 Spieth, J., and D. Lawson, 2006 Overview of gene structure. *WormBook*: 1-10.
- 25 Spieth, J., D. Lawson, P. Davis, G. Williams and K. Howe, 2014 Overview of gene structure in
26 *C. elegans*. *WormBook*: 1-18.
- 27 Stein, L. D., Z. Bao, D. Blasiar, T. Blumenthal, M. R. Brent *et al.*, 2003 The genome sequence of
28 *Caenorhabditis briggsae*: a platform for comparative genomics. *PLoS Biol* 1: E45.
- 29 Tourasse, N. J., J. R. M. Millet and D. Dupuy, 2017 Quantitative RNA-seq meta-analysis of
30 alternative exon usage in *C. elegans*. *Genome Res* 27: 2120-2128.
- 31 Velculescu, V. E., L. Zhang, B. Vogelstein and K. W. Kinzler, 1995 Serial analysis of gene
32 expression. *Science* 270: 484-487.

- 1 Wang, F., S. Huang and L. Ma, 2010 *Caenorhabditis elegans* operons contain a higher proportion
2 of genes with multiple transcripts and use 3' splice sites differentially. *PLoS One* 5:
3 e12456.
- 4 Wang, X., Y. Zhao, K. Wong, P. Ehlers, Y. Kohara *et al.*, 2009 Identification of genes expressed
5 in the hermaphrodite germ line of *C. elegans* using SAGE. *BMC Genomics* 10: 213.
- 6 Williams, C., L. Xu and T. Blumenthal, 1999 SL1 trans splicing and 3'-end formation in a novel
7 class of *Caenorhabditis elegans* operon. *Mol Cell Biol* 19: 376-383.
- 8 Woodruff, G. C., O. Eke, S. E. Baird, M. A. Felix and E. S. Haag, 2010 Insights into species
9 divergence and the evolution of hermaphroditism from fertile interspecies hybrids of
10 *Caenorhabditis* nematodes. *Genetics* 186: 997-1012.
- 11 Zahler, A. M., 2012 Pre-mRNA splicing and its regulation in *Caenorhabditis elegans*.
12 *WormBook*: 1-21.
- 13 Zaslaver, A., L. R. Baugh and P. W. Sternberg, 2011 Metazoan operons accelerate recovery from
14 growth-arrested states. *Cell* 145: 981-992.
15
16

1 **LIST OF TABLES**

2

3 **Table 1: Overview of SL1 and SL2 5' tags identified in the study.**

4

	Total unique tags	Matches in genome	Unique hits	Multiple hits
All	14,678	10,243	4,753	5,490
SL1	10,400	7,234	3,281	3,953
SL2	□ 4,278	3,009	1,472	1,537

5

6

1 **Table 2: Breakdown of tag matches into different categories.**

2 The numbers include both unique and multiple hits. Tag matches termed as ‘Others’ are those
3 that cannot be placed uniquely into any of the main categories.

4

Category of tag matches	SL1	SL2	Total
1a	3,537	1,542	5,079 (79.4%)
1b	20	3	23 (0.3%)
Minor misprediction	245	91	336 (5.2%)
Major misprediction	639	291	930 (14.5%)
Others	22	5	27 (0.4%)
Total	4,463	1,932	6,395

5

6

1 **Table 3: Breakdown of genes by spliced leader sequences.**

2

	Number of genes
SL1 type	2,750 (65%)
SL2 type	743 (18%)
SL1/SL2 type	759 (18%)
Total	4,252

3

4

1 **Table 4. Genes supported by the presence of a single 5' end (single transcript).**

2 Numbers refer to genes identified by SL1, SL2 and SL1/SL2 tags. The genes have been divided
3 further into various categories based on distance from the nearest exon (see figure 1). Novel
4 exons and potential paralogs are excluded.

5

	All	SL1	SL2	SL1/SL2
Matching first exon (1a)	3,288	2,143 (65.2%)	558 (17.0%)	587 (17.9%)
Matching internal exon (1b)	16	14 (87.5%)	1 (6.2%)	1 (6.2%)
Minor misprediction of first or internal exon	204	146 (71.6%)	34 (16.7%)	24 (11.7%)
Major misprediction of first or internal exon	517	357 (69%)	127 (24.6%)	33 (6.4%)
Total	4,025	2,660	720	645

6

7

1 **Table 5: Genes supported by the presence of multiple 5' ends.**

2 Numbers refer to genes identified by SL1, SL2 and SL1 and SL2 tags. These genes have been
3 divided further into various categories based on distance from the nearest exon (see figure 1).
4 Genes for which exons belong to multiple categories are grouped as 'Others'. Novel exons and
5 potential paralogs are excluded.

6

	All	SL1	SL2	SL1/SL2
Matching first exon (1a) and matching internal exon (1b)	5	2 (25%)	0	3 (75%)
Matching internal exons (1b)	0	0	0	0
Matching first exon (1a) and minor misprediction of one or more internal exons	39	14 (36%)	7 (18%)	18 (46%)
Matching first exon (1a) and major misprediction of one or more internal exons	145	57 (39%)	14 (10%)	74 (51%)
Others	6	2 (33%)	0	4 (67%)
All mispredicted exons (minor and major)	32	16 (50%)	2 (6%)	14 (44%)
Total	227	91	23	113

7

8

1 **Table 6: Breakdown of *C. briggsae* operons based on the number of genes present.**

2 Operons are placed into two broad categories, those consisting entirely of genes with known
3 spliced leader sequences (Tag-supported) and others where the splice leader identity of the first
4 gene is unknown (Predicted).

5

	No. of operons	Operons consisting of					
		2 genes	3 genes	4 genes	5 genes	6 genes	7 genes
Tag-supported operons	334	263	54	14	2	0	1
Predicted operons	865	706	125	26	7	1	0

6

7

1 **Table 7: Tag-supported operons in *C. briggsae*.**

2 Exact match operons are conserved between *C. briggsae* and *C. elegans*. Partially conserved
3 operons may contain some but not all orthologs that are part of corresponding *C. elegans*
4 operons. Novel operons may contain *C. elegans* orthologs and divergent, *C. briggsae*-specific,
5 genes.

6

Operon type	Number (% of total)
Fully conserved operons (Exact match)	134 (40.1%)
Partially conserved operons (Partial match)	127 (38%)
Novel operons	73 (21.9%)
- consisting of both divergent genes as well as orthologs that are not part of <i>C. elegans</i> operons	61 (18.3%)
- consisting entirely of divergent genes	12 (3.6%)
Total	334

7

8

- 1 **Table 8: Novel *C. briggsae* operons identified in this study with ICRs of less than 1 kb.**
2 None of the genes in these operons have orthologs in *C. elegans*. The numbers in brackets refer
3 to ICR.
4

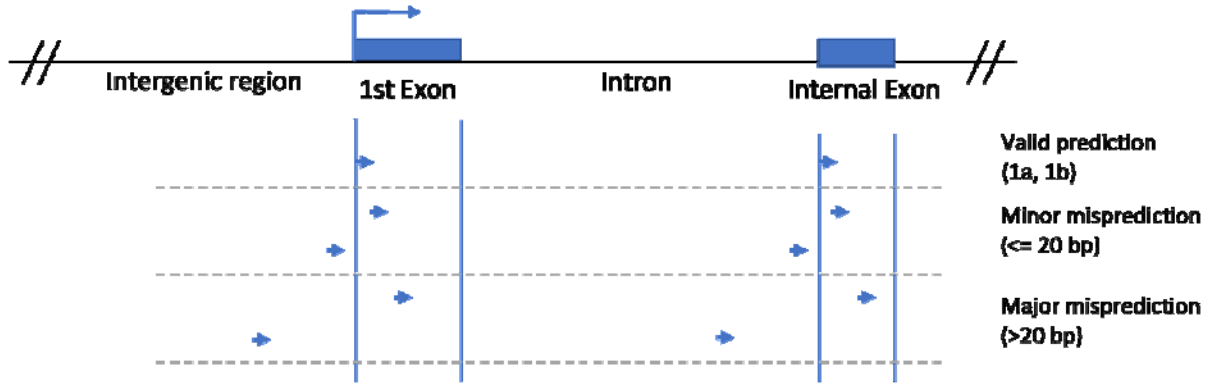
<i>C. briggsae</i> operon	No. of genes	Gene names (ICR)
CBROPX0130	3	<i>CBG30062</i> (172) <i>CBG25686</i> (105) <i>CBG25687</i>
CBROPX0131	3	<i>CBG27303</i> (533) <i>CBG27302</i> (116) <i>CBG27301</i>
CBROPX0140	2	<i>CBG11551</i> (162) <i>CBG31489</i>
CBROPX0129	3	<i>CBG21606</i> (235) <i>CBG30457</i> (493) <i>CBG21605</i>
CBROPX0139	2	<i>CBG30329</i> (76) <i>CBG30328</i>

- 5
6

1 **LIST OF FIGURES**

2

3 **Figure 1: Representative model of locations of tag sequences within the genome.** Three
4 broad categories of matches are: valid prediction (termed 1a and 1b), minor misprediction, and
5 major misprediction.



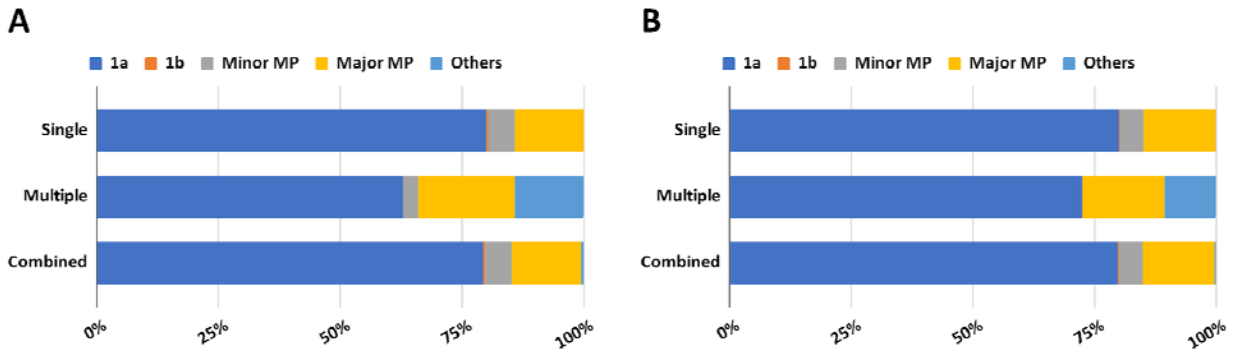
6

7

8

1 **Figure 2. Proportion of tags belonging to different categories.** The majority of SL1 (A) and
2 SL2 tags (B) have single (unique) hits in the genome and belong to category 1a, i.e., predicted 5'
3 ends match with WormBase gene models. Minor MP: Minor misprediction, Major MP: Major
4 misprediction, Others: mix category of matches.

5

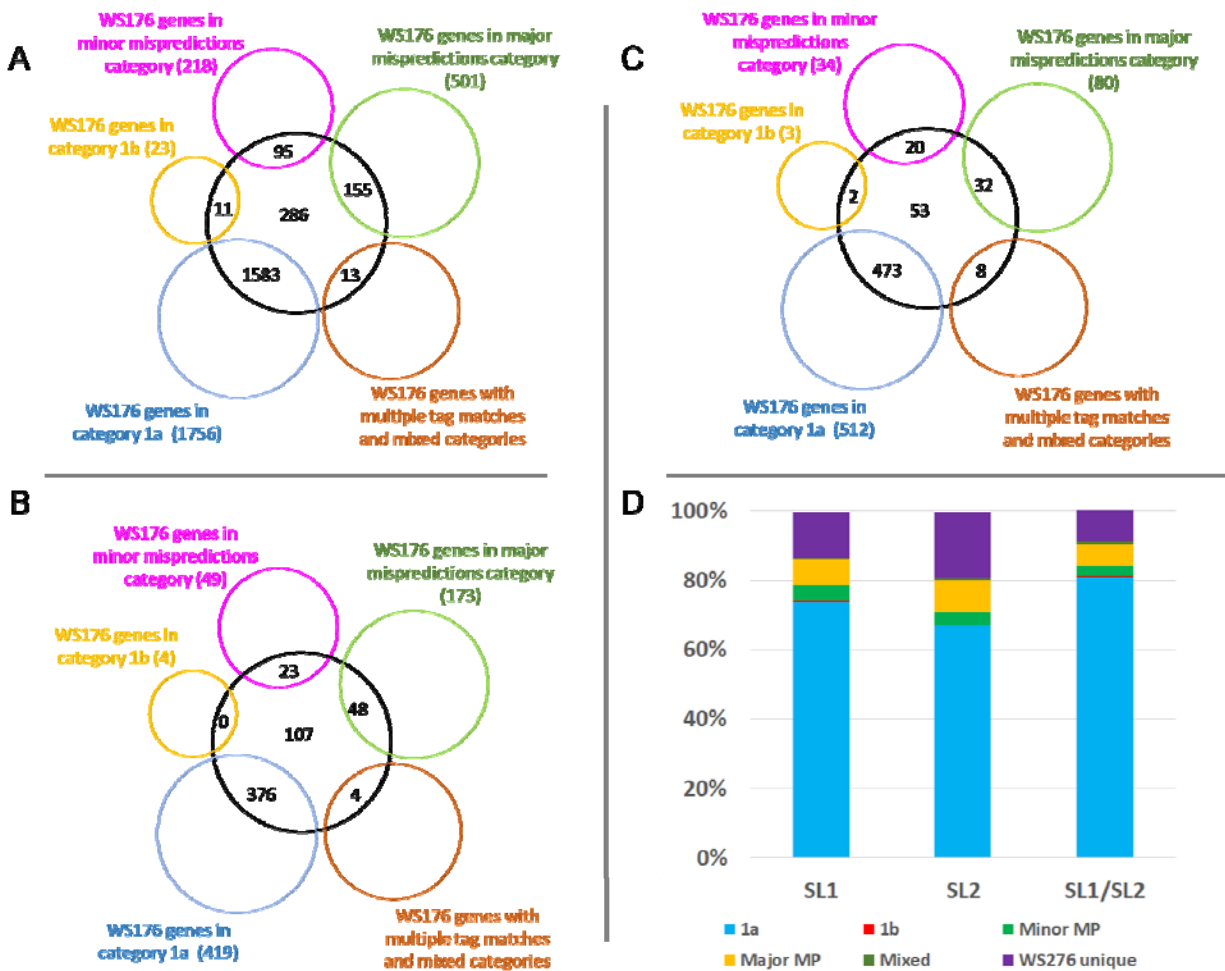


6

7

1 **Figure 3: Reclassification of genes from WS176 categories to Category 1a in WS276.** Only
 2 single transcript genes were compared. (A-C) Venn diagrams, with WS276 genes of category 1a
 3 in black circles and WS176 genes of various categories in colored circles. Numbers in
 4 overlapping circles represent genes of a given category in WS176 that are annotated as 1a type in
 5 WS276. Numbers in the middle of black circles (non-overlapping) represent genes that are
 6 unique to WS276 analysis (A, 286 or 13.2% of SL1-spliced; B, 107 or 19.0% of SL2-spliced; C,
 7 53 or 9% of SL1/SL2 hybrid-spliced) whereas those in brackets next to colored circles are total
 8 genes identified by tag searches in WS176. (D) Histogram showing the proportion of genes with
 9 matching 5' ends in WS276 (category 1a) that overlap with various categories in the WS176
 10 analysis.

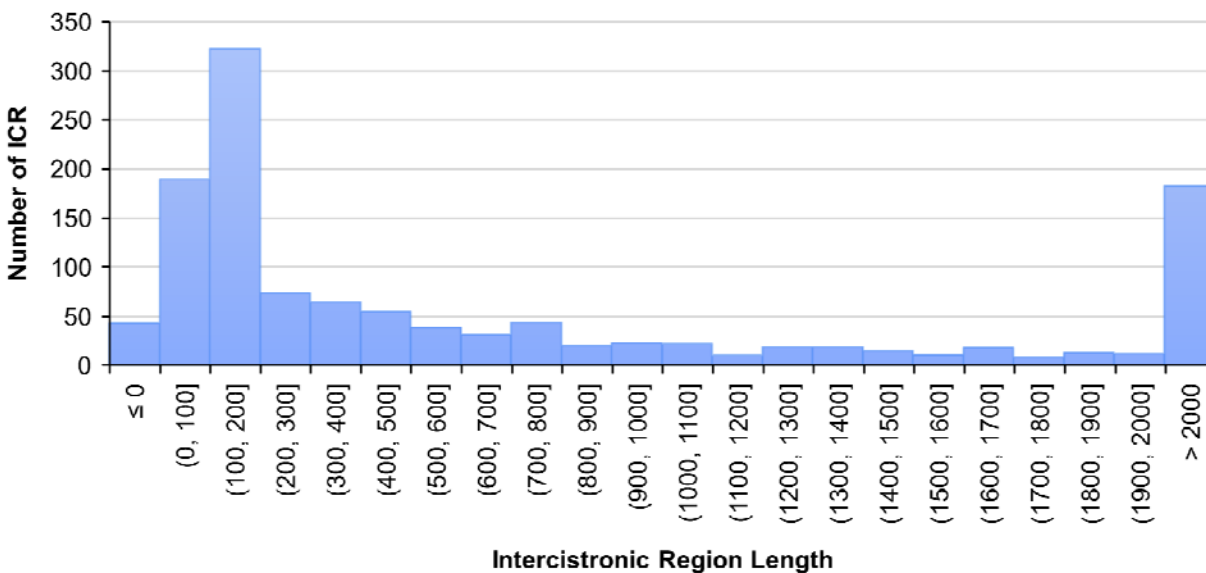
11



12

13

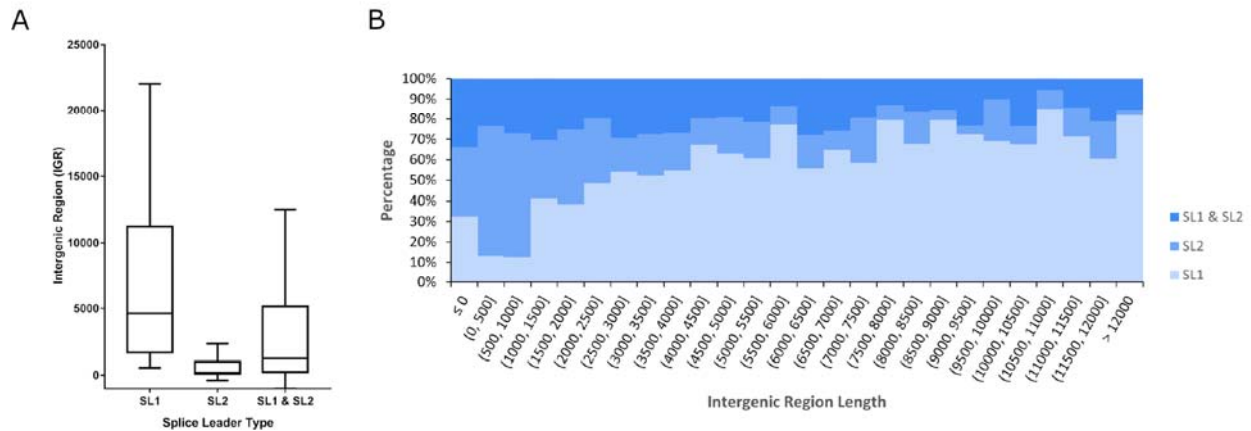
- 1 **Figure 4: Frequencies of ICR lengths between SL2 and hybrid-spliced genes in operons.**
- 2 ICRs are sorted in bins of 100 nucleotides. For pairs of genes where the second gene is within
- 3 the first gene, ICR is calculated as a negative value. For bin sizes, round brackets indicate
- 4 exclusive bound, square brackets indicate inclusive bounds. Genes with larger than 2kb ICRs are
- 5 shown as a single peak.



6

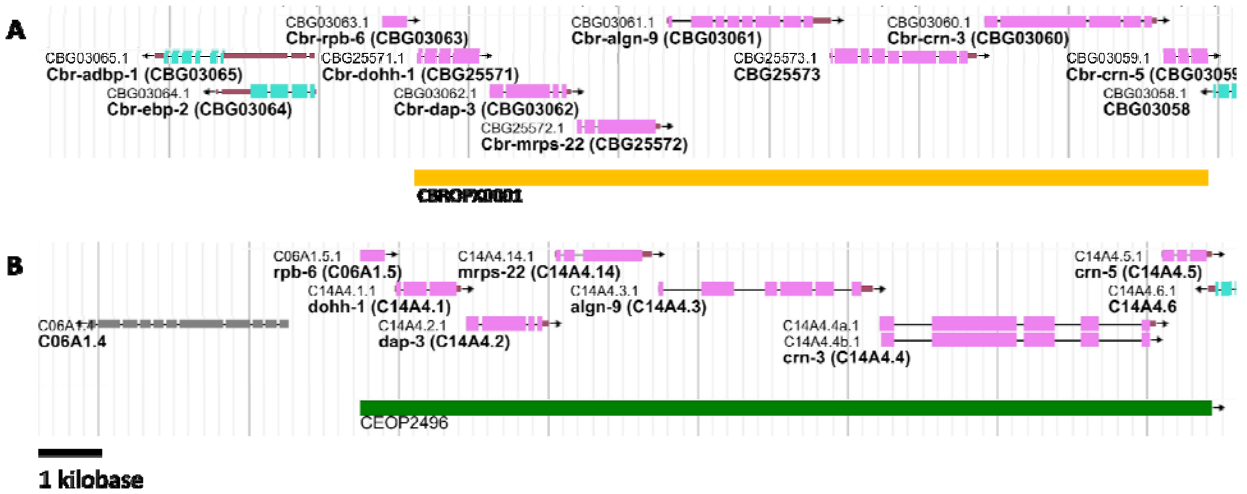
7

1 **Figure 5: A: IGRs of genes identified by tag matches.** (A) Box plots show IGRs for SL1-
2 spliced, SL2-spliced, and SL1/SL2-spliced genes. The inside line marks the median, lower and
3 upper lines represent the borders of the 25th and 75th quartile of the data sample, respectively.
4 Whiskers enclose the 10-90% range of the data. (B) 100% stacked columns of IGR length.
5 Lengths are sorted in bins of 500 nucleotides. For pairs of genes where the second gene is
6 overlapping or inside the first gene, length was calculated as a negative value. For bin sizes,
7 round brackets indicate exclusive bound, square brackets indicate inclusive bounds.



8
9

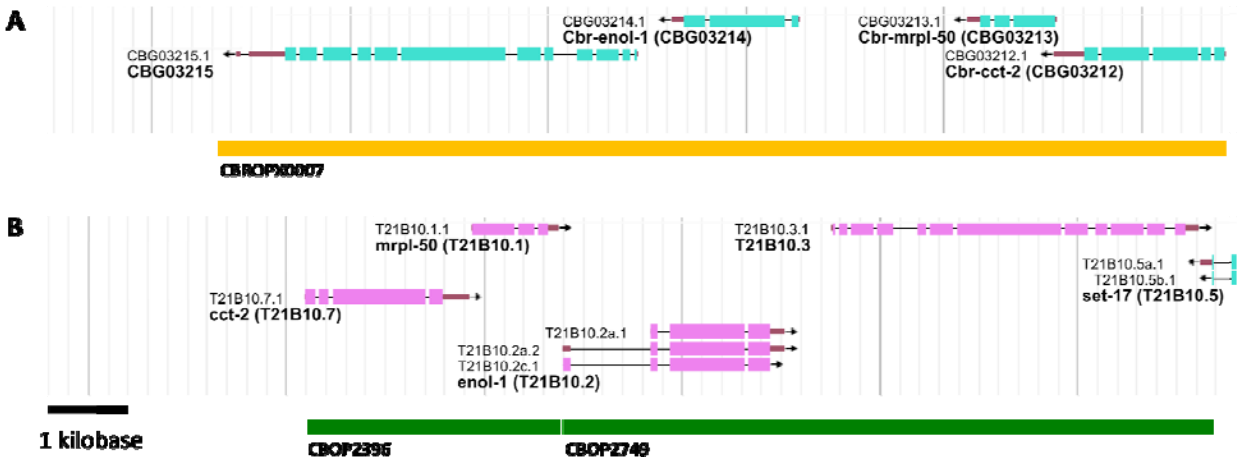
1 **Figure 6. Genomic regions of *C. briggsae* CBROPX0001 and *C. elegans* CEOP2496. (A)**
2 CBROPX0001 is proposed to contain at least seven, and possibly eight, genes depending on the
3 inclusion of *CBG03063*. **(B)** Homologous *C. elegans* operon CEOP2496 contains seven genes.
4 Genomic feature visualizations in this and subsequent figures are modified versions of images
5 obtained from WormBase JBrowse.
6



7
8

1 **Figure 7: *C. briggsae* operon CBROPX0007.** (A) A cluster of four genes that define
2 CBROPX0007. (B) The orthologs of the four genes are split between two *C. elegans* operons
3 CEOP2396 and CEOP2749.

4

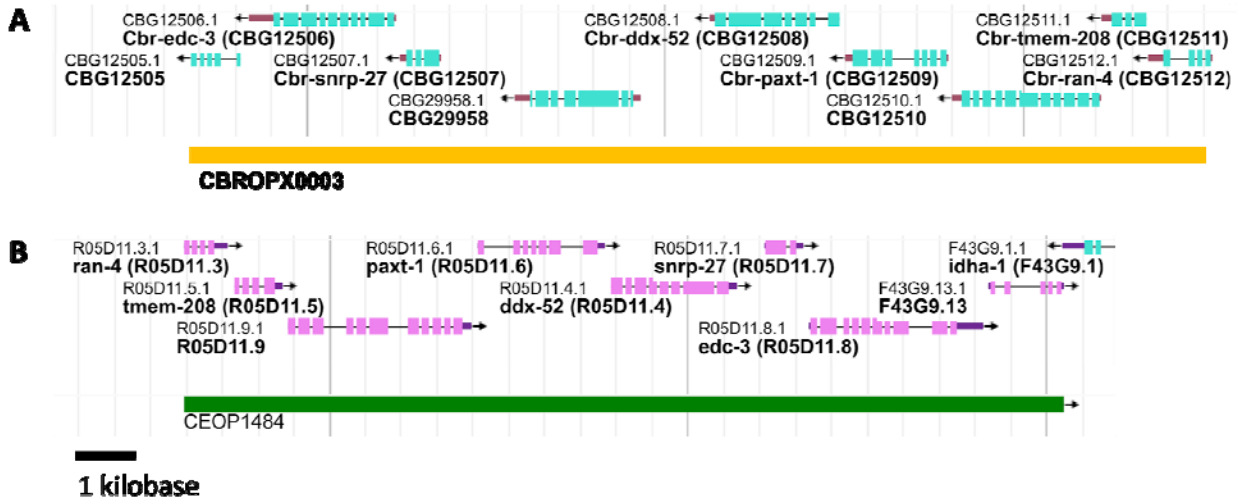


5

6

1 **Figure 8. Partially conserved operon and its *C. elegans* ortholog.** (A) CBROPX0003 is an
2 example of a partially conserved operon identified in this study. (B) CEOP1484, *C. elegans*
3 operon orthologous to *C. briggsae* operon CBROPX0003.

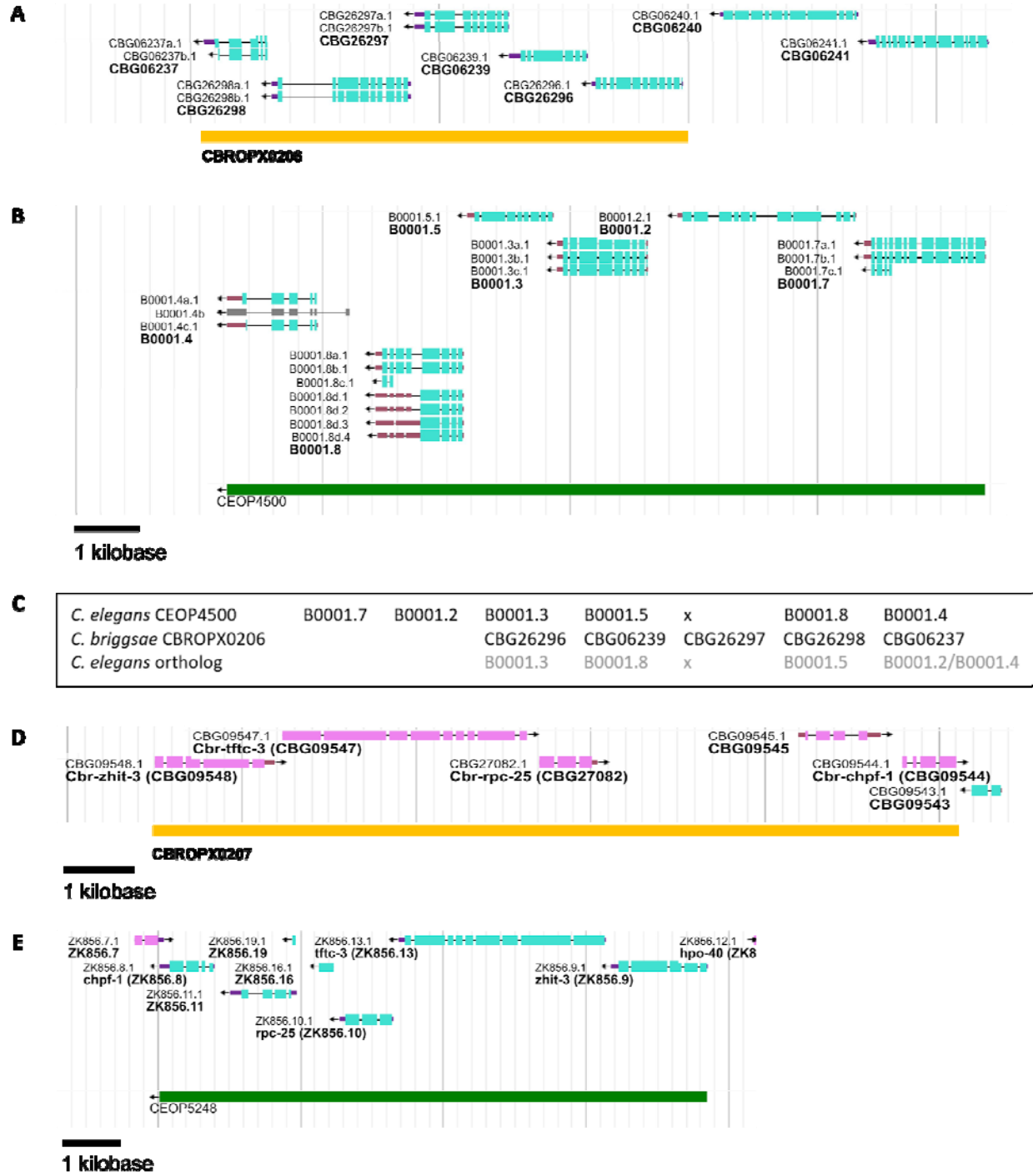
4



5

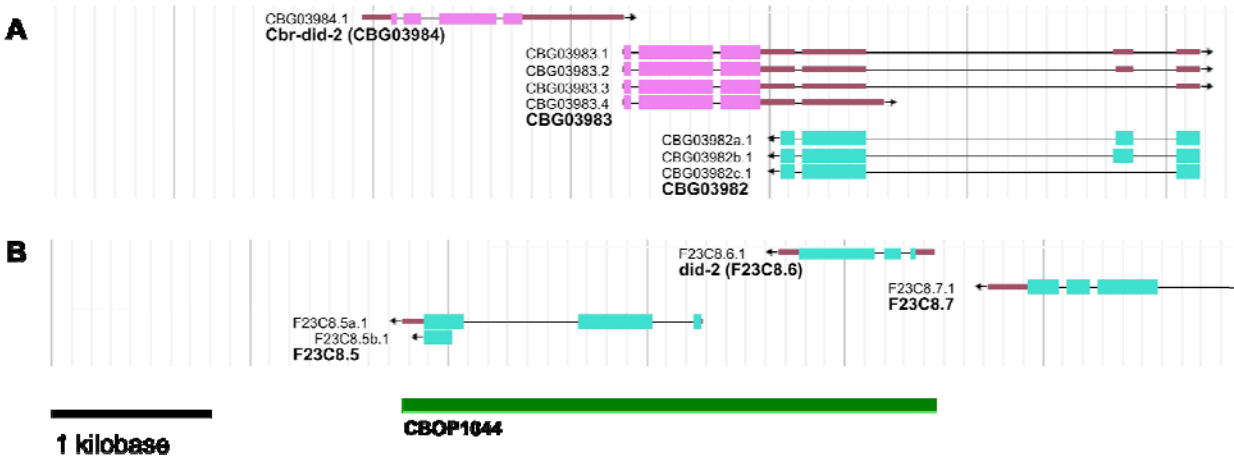
6

1 **Figure 9: Two predicted operons in *C. briggsae* along with their *C. elegans* counterparts.**
2 **(A, B)** CBROPX0206 with five genes and its orthologous operon CEOP4500 in *C. elegans*.
3 Three genes are conserved between these two operons. **(C)** The arrangement of genes in *C.*
4 *elegans* operon CEOP4500 (row 1) and *C. briggsae* operon CBROPX0206 (row 2). The *C.*
5 *elegans* orthologs of CBROPX0206 genes are shown in row 3. ‘x’ denotes a gene that is missing.
6 *CBG06237* is orthologous to both *B0001.2* and *B0001.4*. **(D, E)** CBROPX0207 with five genes
7 and its orthologous operon CEOP5428 with seven genes. All five genes of the *C. briggsae*
8 operon are conserved in CEOP5428. Two additional genes are present in CEOP5428.
9



1
2
3

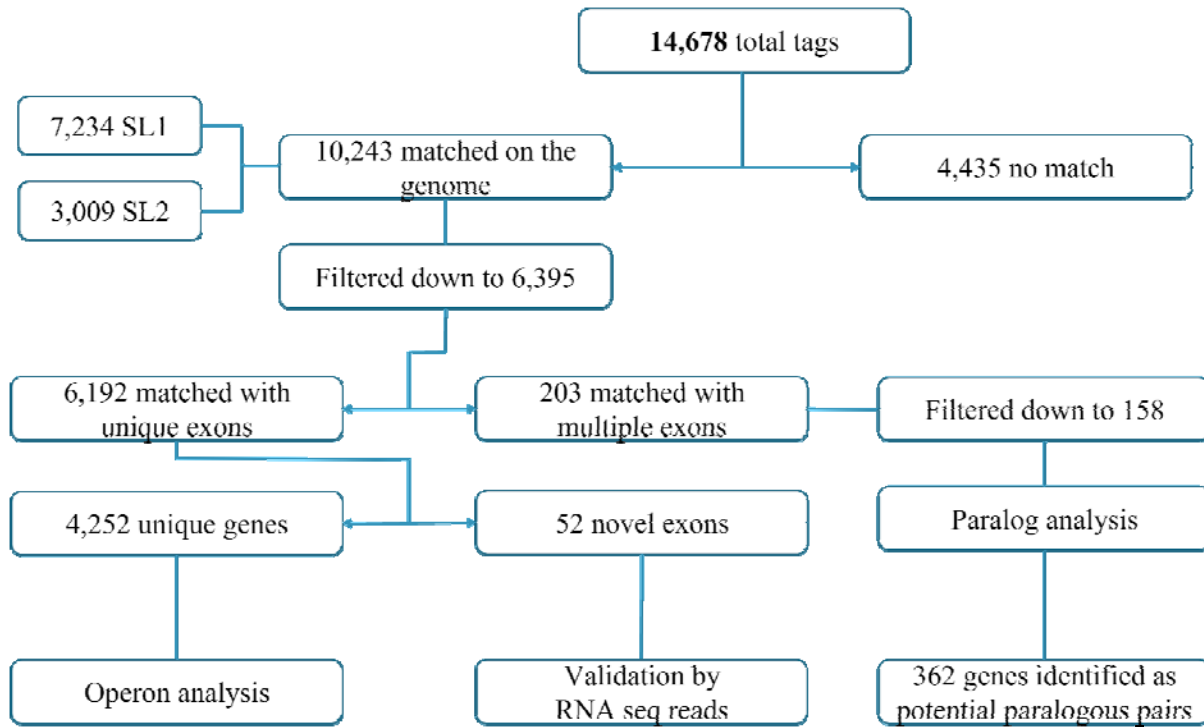
1 **Figure 10. *C. briggsae* genes with a single base pair ICR. (A)** Both *CBG03984* and *CBG03983*
2 are spliced with SL1 leader sequences and located 1 bp apart. **(B)** *C. elegans* orthologs *did-2* and
3 *F23C8.5*, respectively, are part of the operon CEOP1044.
4



5
6

1 **Figure 11. An overview of TEC-RED analysis in *C. briggsae*.** The 5' sequence tags were used
2 to identify exons and genes. Further analysis resulted in the discovery of operons, paralogs, and
3 novel exons.

4



5

6

1 **SUPPLEMENTARY MATERIALS**

2

3 **Jhaveri and van den Berg et al. Genome annotation of *Caenorhabditis briggsae* by TEC-**
4 **RED identifies new exons, paralogs, and conserved and novel operons**

5

6

7

8 **SUPPLEMENTARY DATA FILES (Microsoft Excel spreadsheets)**

9

File name	Description
Supplementary data file 1	Exons identified in our analysis
Supplementary data file 2	Unique genes identified
Supplementary data file 3	Genes uniquely spliced in <i>C. briggsae</i>
Supplementary data file 4	Validation based on overlap with WS176
Supplementary data file 5	New exons identified by our study
Supplementary data file 6	Manual curation of 1b and major mispredictions based on <i>C. elegans</i> orthologs
Supplementary data file 7	List of operons
Supplementary data file 8	Intergenic region values
Supplementary data file 9	Germline genes present in operons and GO analysis
Supplementary data file 10	Proposed paralog sets
Supplementary data file 11	Matches to <i>C. briggsae</i> novel operons in <i>C. nigoni</i>

10

11

1 **SUPPLEMENTARY TABLES**

2

3 **Supplementary Table 1: Primers used to generate Biotin-RT-PCR products**

4

Primers	Sequence (5' to 3')
RT primer	GTGATGTCTCGAGTAGTTCGAAATGGCC (T)22
5' SL1- <i>Bpm</i> I RT-PCR primer	Biotin/ AGACGCAAGGTTTAATTACCCAAGCTGGAG
5' SL2- <i>Bpm</i> I RT-PCR primer	Biotin/ AGACGCAAGGTTTAACCCAGTTACTGGAG
3' RT-PCR primer	GAGGTGATGTCTCGAGTAGTTCGAAATGGC

5

6

1 **Supplementary Table 2:** PCR primers used to generate mono-tags from the 5' biotin-adaptor
2 DNA fragments

3

Primers	Sequence (5' to 3')
5' SL1- <i>Xho</i> I primer	AGACGCAAGGTTTAATTACCCAAGCTCGAG
5' SL2- <i>Xho</i> I primer	AGACGCAAGGTTTTAACCCAGTTACTCGAG
3' for adaptor 1 (<i>Kpn</i> I)	CTATAGGGCTCAAAGATGACGAGAGGA
3' for adaptor 2 (<i>Hind</i> III)	CAAGATTCTCACGACGATGTTTCGGAGT
3' for adaptor 3 (<i>Eag</i> I)	TGAAGATTGCACAGAGGAGAGACCGCT
3' for adaptor 4 (<i>Sac</i> I)	CAGTTGGAATGAATGAAGCTATAACCAT
3' for adaptor 5 (<i>Mlu</i> I)	CTAGTATACGTTCTAGTATCAGAGGAA
3' for adaptor 6 (<i>Nhe</i> I)	TCTTGCAGTGATTAGCGTCAGTGCCTG

4

5

1 **Supplementary Table 3:** Adaptors used for ligation onto *BpmI*-digested, 5' biotin-DNA
2 fragments
3

Adapter	Sequence (5' to 3')	Sequence (3' to 5')
Adapter 1 (<i>KpnI</i>)	CTATAGGGCTCAAAGATGACGAGA GGAGGTACC	TGCTCTCCTCCATGG
Adapter 2 (<i>HindIII</i>)	CAAGATTCTCACGACGATGTTCCG AGTAAGCTT	CAAGCCTCATTCGAA
Adapter 3 (<i>EagI</i>)	TGAAGATTGCACAGAGGAGAGACC GCTCGGCCG	CTCTGGCGAGCCGGC
Adapter 4 (<i>SacI</i>)	CAGTTGGAATGAATGAAGCTATAC CATGAGCTC	GATATGGTACTCGAG
Adapter 5 (<i>MluI</i>)	CTAGTATACGTTCTAGTATCAGAG GAAACGCGT	AGTCTCCTTTGCGCA
Adapter 6 (<i>NheI</i>)	TCTTGCAGTGATTAGCGTCAGTGC CTGGCTAGC	GTCACGGACCGATCG

1 **Supplementary Table 4:** Chromosomal locations of 4,252 unique genes identified by TEC-RED. Chr: Chromosome, Un: unmapped
 2 genomic region.
 3

Chr	Total gene count	SL1 genes	Fraction	Density	SL2 genes	Fraction	Density	SL1/SL2 genes	Fraction	Density	Chr length (Mb)*
I	763	447	16.26	28.93	158	21.27	10.23	158	20.79	10.23	15.45
II	752	473	17.21	28.46	147	19.78	8.84	132	17.37	7.94	16.62
III	751	444	16.15	30.47	145	19.52	9.95	162	21.32	11.12	14.57
IV	717	447	16.26	25.57	133	17.90	7.61	137	18.03	7.84	17.48
V	749	534	19.43	27.40	105	14.13	5.39	110	14.47	5.64	19.49
X	515	400	14.55	18.57	54	7.27	2.51	61	8.03	2.83	21.54
Un	5	4			1			0			
	4252	2749		26.14	743		7.07	760		7.23	105.15

4
 5
 6 *Ross, J. A., D. C. Koboldt, J. E. Staisch, H. M. Chamberlin, B. P. Gupta et al., 2011 *Caenorhabditis briggsae* recombinant inbred line
 7 genotypes reveal inter-strain incompatibility and the evolution of recombination. PLoS Genet 7: e1002174.
 8

1 **Supplementary Table 5:** Intergenic distances of selected Category 3 genes that are less than 10
2 kb apart.

3

Adjacent genes identified by tags	IGR	BLAST alignment	<i>C. elegans</i> orthologs	<i>C. briggsae</i> gene orientation
<i>CBG25816</i> , <i>CBG00473</i>	2,903 bp	Yes	None	Opposite
<i>CBG08766</i> , <i>CBG08768a</i>	578 bp	No*	<i>F25E5.8</i> and <i>nhr-117</i>	Opposite
<i>CBG25203</i> , <i>CBG29819</i>	2,251 bp	No	<i>F59A3.2</i> and <i>ubl-5</i>	Opposite
<i>CBG26374</i> , <i>CBG05421</i>	3,564 bp	No*	None, <i>fan-1</i>	Same
<i>CBG26845</i> , <i>CBG26846</i>	8,559 bp	Yes	None	Same

4

5 *BLAST match showed some similarity in a very small 5' region.

6

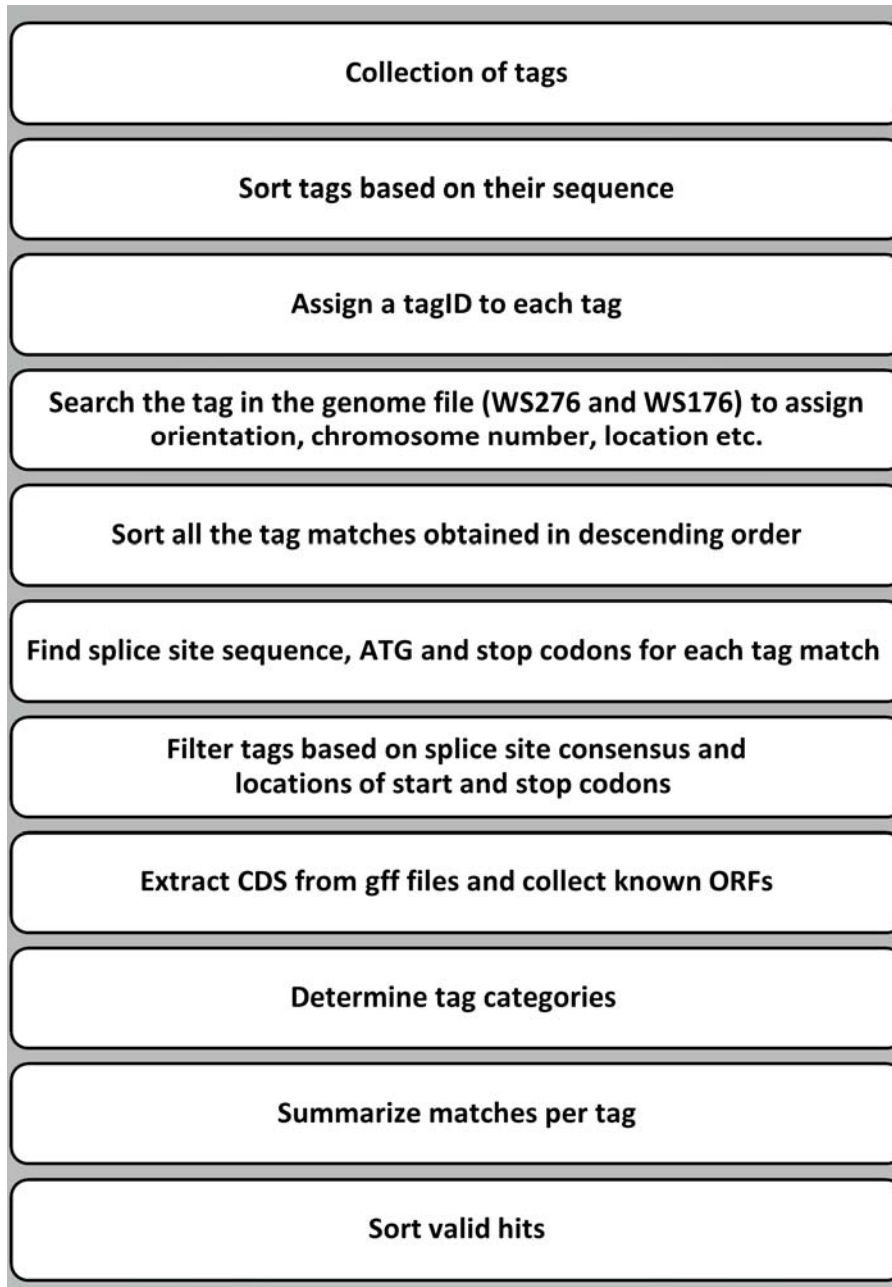
7

1 **SUPPLEMENTARY FIGURES**

2

3 **Supplementary Figure 1:** Flowchart of steps used to analyze 5' tag sequences and genes.

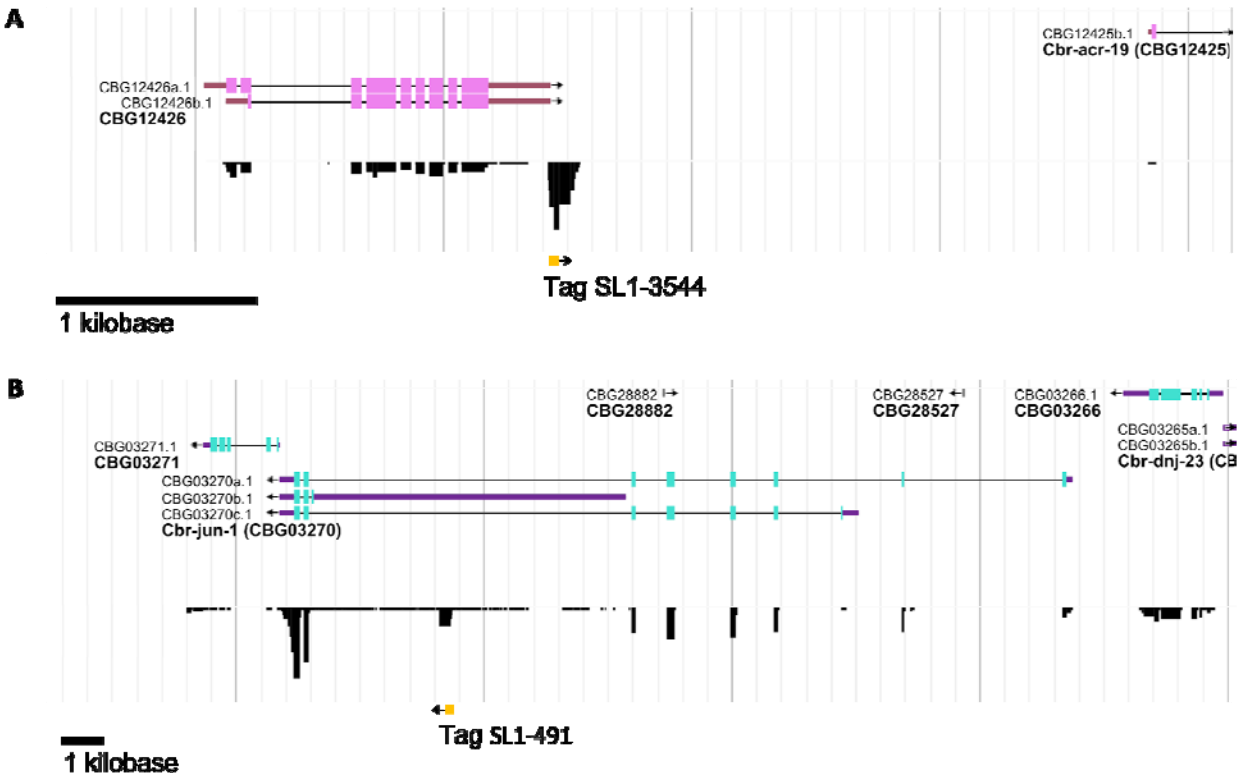
4



5

6

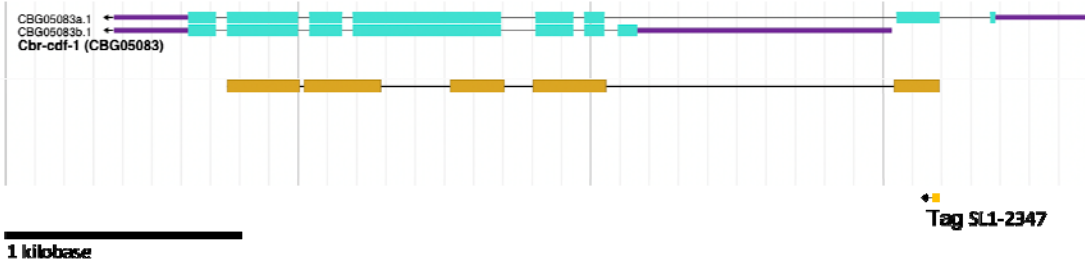
1 **Supplementary Figure 2:** Selected examples of novel exons identified by tags that are
2 supported by WormBase RNASeq data. The top track of the genome browser shows currently
3 curated genes. Second track shows alignments of short read sequences from all available
4 RNASeq projects on WormBase. The number of reads has been normalized by averaging over
5 the number of libraries. The height of reads boxes indicates the relative score of the feature. The
6 bottom track shows a TEC-RED tag binding at a genome location predicted to contain the 5'
7 start site of a new exon. (A) New exon between *CBG12426b.1* and *CBG12425b.1*. (B) A new
8 exon inside *CBG03270a.1*.
9
10



11
12

1 **Supplementary Figure 3:** An example of the 1b category in *Cbr-cdf-1* in WormBase genome
2 browser. The top track shows curated gene *Cbr-cdf-1*. The middle track shows the *C. elegans*
3 *C15B12.7a.1 (cdf-1)* gene model, which is indicated in orange. The bottom track shows a
4 category 1b TEC-RED tag binding at the 5' start site of exon 2 of *Cbr-cdf-1*. The *C. elegans*
5 gene model supports the 5' start site of a new transcript variant for *Cbr-cdf-1*.

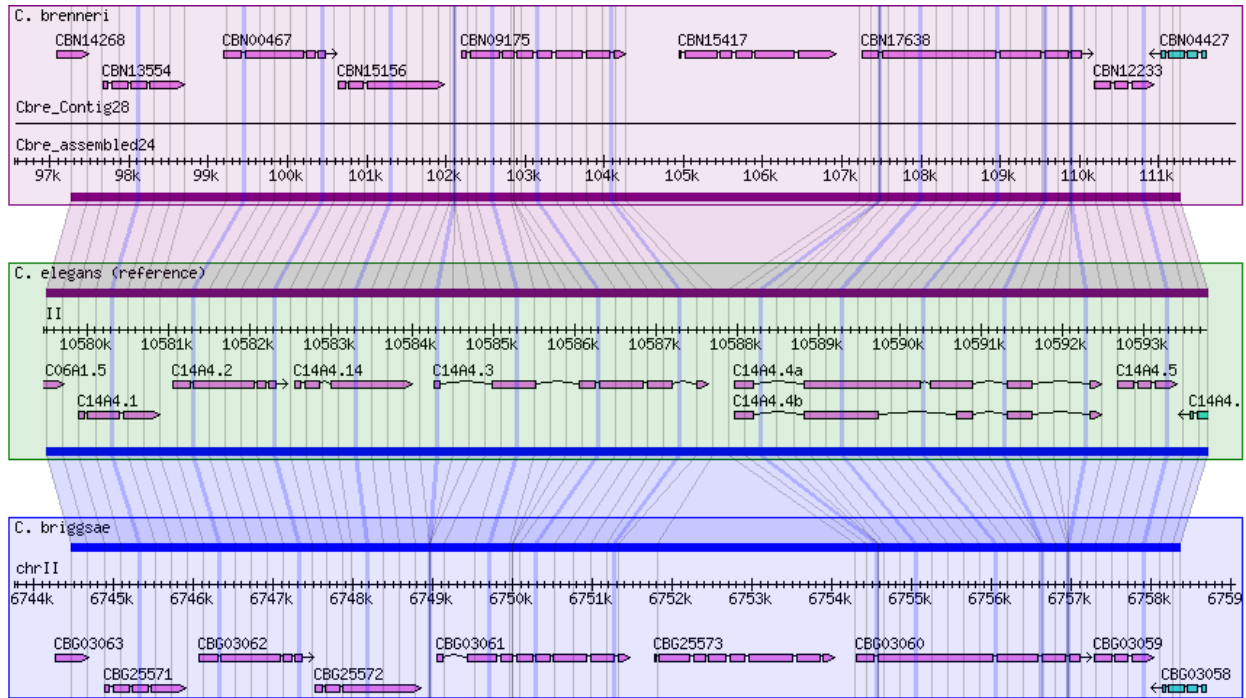
6
7



8
9

1 **Supplementary Figure 4.** A screenshot of the WormBase synteny browser showing the
2 genomic region of *C. briggsae* operon CBROPX0001, and the corresponding regions in *C.*
3 *brenneri*, *C. elegans* and *C. briggsae*.

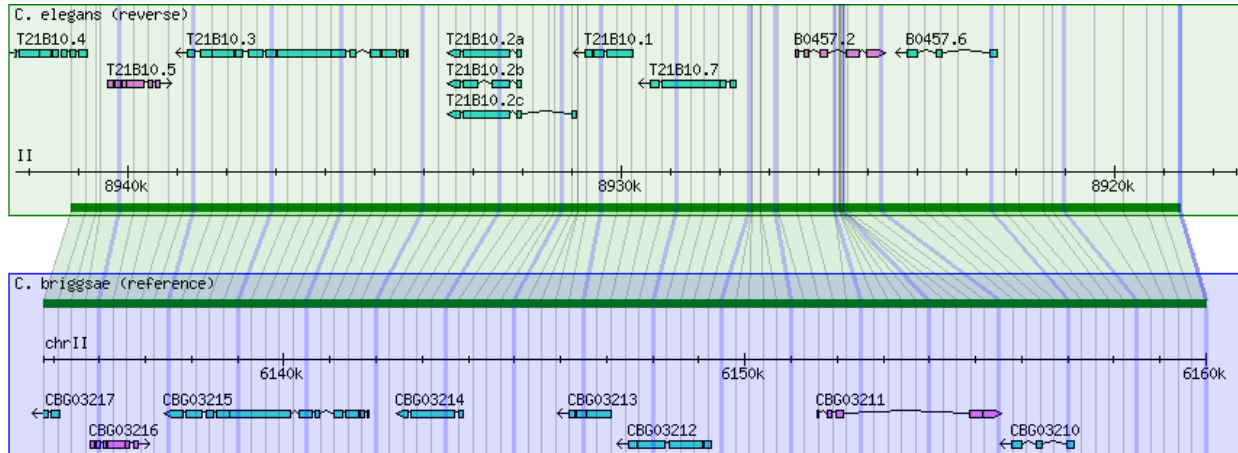
4
5



6
7

1 **Supplementary figure 5:** A screenshot of the WormBase synteny browser showing the genomic
2 region consisting of a cluster of four *C. briggsae* genes that define the CBROPX0007 operon and
3 corresponding region in *C. elegans*.

4
5



6
7