# *C. briggsae* genome annotation and comparative analysis with *C. elegans* using RNA-Seq data

**Final oral examination for the degree of**
**Master of Science**

**SHINTA THIO**

**Monday, April 6th, 2020**

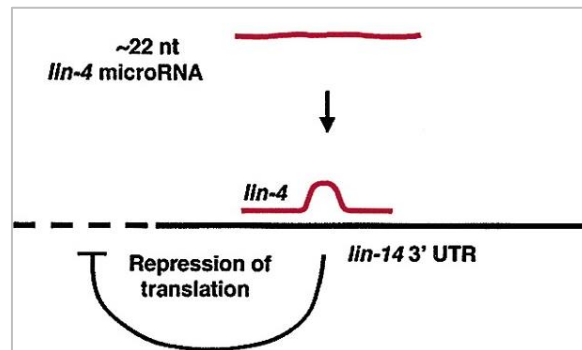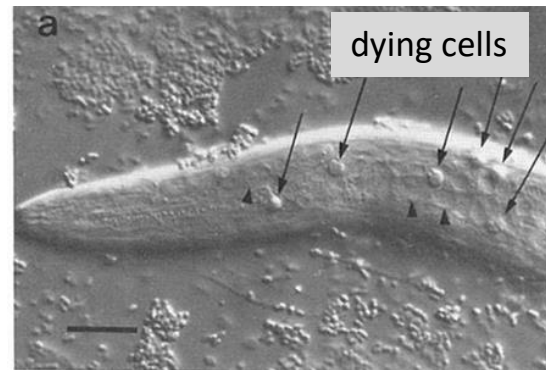**C. elegans:**
- **Small**: ~1mm long
- **Simple body plan**: 959 cells
- **Compact genome size**: ~100 Mbp, first multicellular genome sequenced[1]
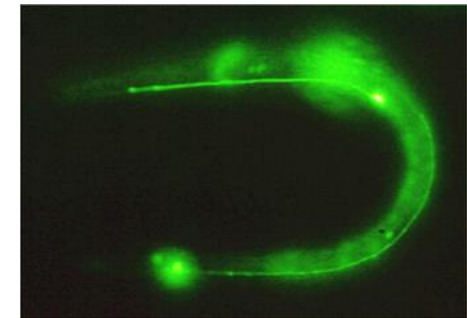
**Key discoveries using *C. elegans*:**

1. microRNA[2]



2. Apoptosis pathway[3]



dying cells

Nomarski photomicrograph, newly hatched *ced-1* larva. Bar=10u

3. Cell visualization in living organism using GFP[5]



[1]The C. elegans Genome Sequencing Consortium, 1998
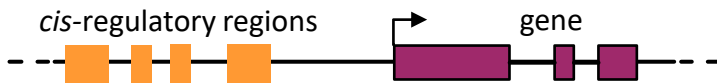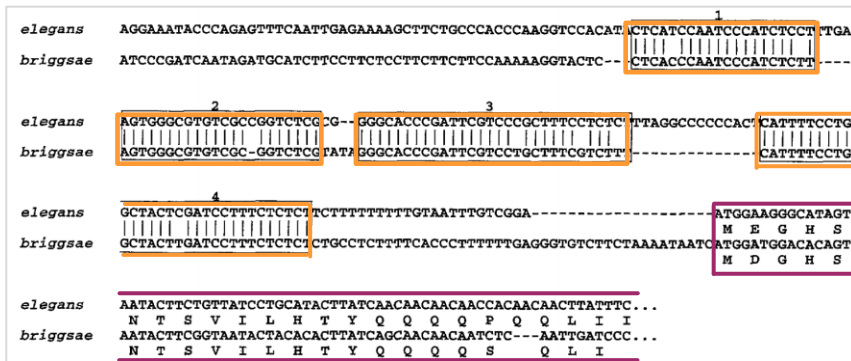[2]Lee et al. , 1993
[3]Ellis & Horvitz, 1986
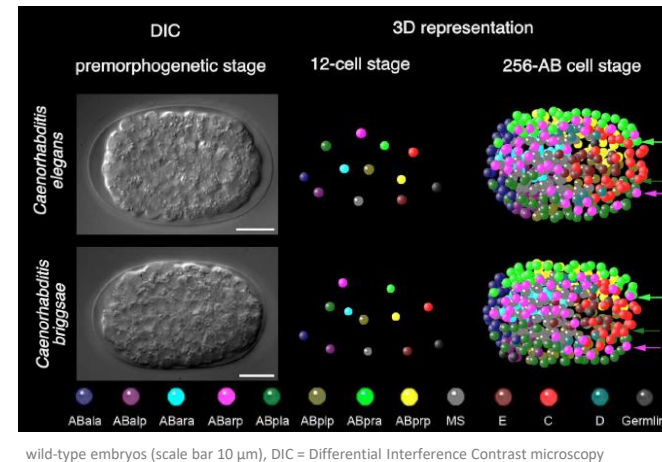[4]Chalfie et al., 1994

*C. elegans* is an ideal model organism to investigate biology
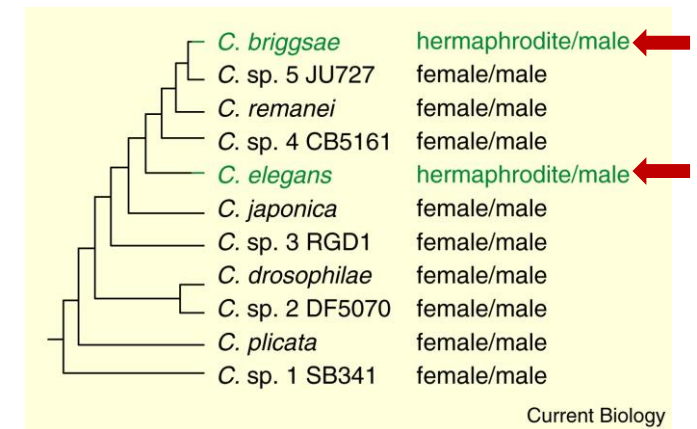
# *C. briggsae* facilitates *C. elegans* research

## 1. Improving genome annotation (ortholog-based)[1-7]



*cis*-regulatory regions          gene

## 2. Understanding embryonic development[8-10]



wild-type embryos (scale bar 10 µm), DIC = Differential Interference Contrast microscopy

## 3. Understanding evolution of hermaphroditism[11]
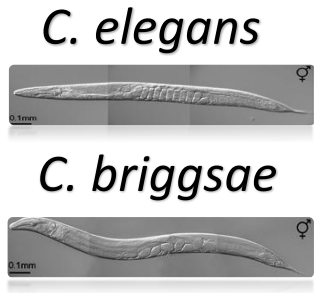


*C. briggsae* is a powerful comparative tool to improve the understanding of *C. elegans*
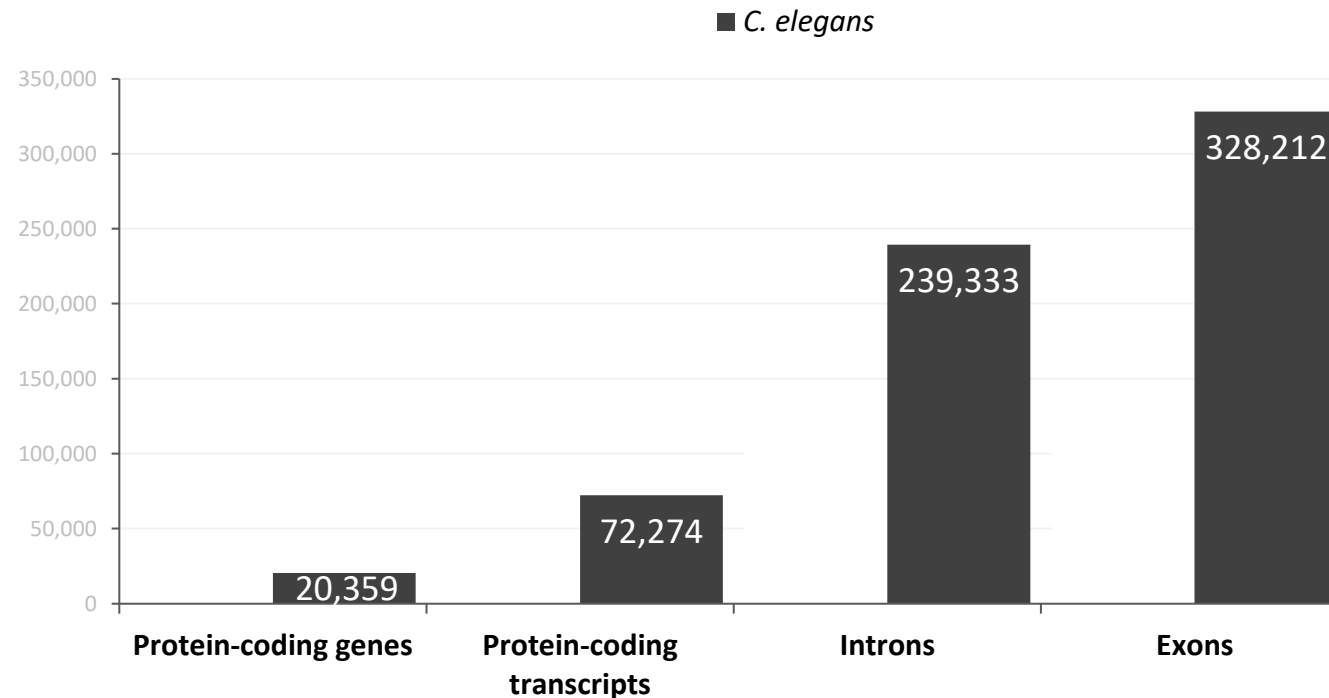
[1]Clark et al., 1995, [2]The C. elegans Genome Sequencing Consortium, 1998, [3]Chen et al., 2006, [4]Gaudet et al., 2004, [5]Kirouac and Sternberg, 2003, [6]Stein et al, 2003, [7]Uyar et al, 2012, [8]Sulston et al., 1983, [9]Zhao et al., 2008, [10]Memar et al., 2019, [11]Kiontke et al., 2004

# Current *C. briggsae* annotation is inadequate

*C. elegans*

*C. briggsae*

| Life cycle | Reproduction | Number of cells (adult ⚥) | Chromosome | Genome size |
|---|---|---|---|---|
| Adult → Eggs → L1 → L2 → L3 → L4 → Young adult → Adult | Hermaphrodite & male (<0.2%) | 959 cells | 6 | 100 Mbp |
| | | | | 108 Mbp |

■ *C. elegans*

- Protein-coding genes: 20,359
- Protein-coding transcripts: 72,274
- Introns: 239,333
- Exons: 328,212

[1]Douglas, M., 2018
[2]WormBase WS250

# Current *C. briggsae* annotation is inadequate

*C. elegans*

*C. briggsae*

| Life cycle | Reproduction | Number of cells (adult ♀) | Chromosome | Genome size |
|---|---|---|---|---|
| Young adult → Adult → Eggs → L1, L2, L3, L4 | Hermaphrodite & male (<0.2%) | 959 cells | 6 | 100 Mbp |
| | | | | 108 Mbp |

**Legend:** ■ *C. briggsae*  ■ *C. elegans*

| Category | *C. briggsae* | *C. elegans* |
|---|---|---|
| Protein-coding genes | 21,814 | 20,359 |
| Protein-coding transcripts | 22,475 | 72,274 |
| Introns | 107,848 | 239,333 |
| Exons | 121,849 | 328,212 |

**Extensive studies** in *C. elegans*, **limited studies** in *C. briggsae*.

*C. briggsae*

**Computational**
- *Ab initio* gene finding
- Sequence conservation

**Experimental**
- ESTs & Protein-based comparisons
- RNA Sequencing (RNA-Seq, **2 libraries**)

*C. elegans*

**Computational**
- *Ab initio* gene finding
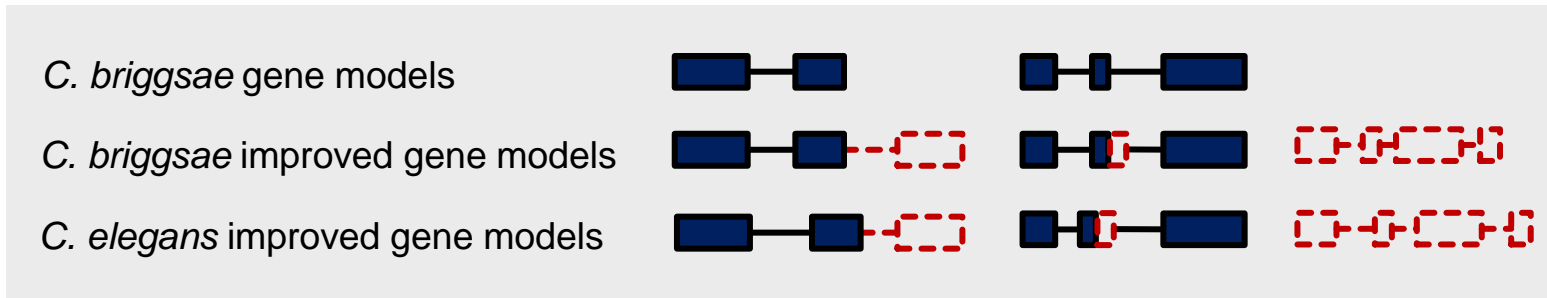- Homology-based gene prediction

**Experimental**
- Expressed sequence tags (ESTs)
- Open reading frame sequence tags (OSTs)
- Serial analysis of gene expression (SAGE)
- Rapid Amplification of cDNA ends (RACE)
- Trans-spliced exon coupled RNA end determination (TEC-RED)
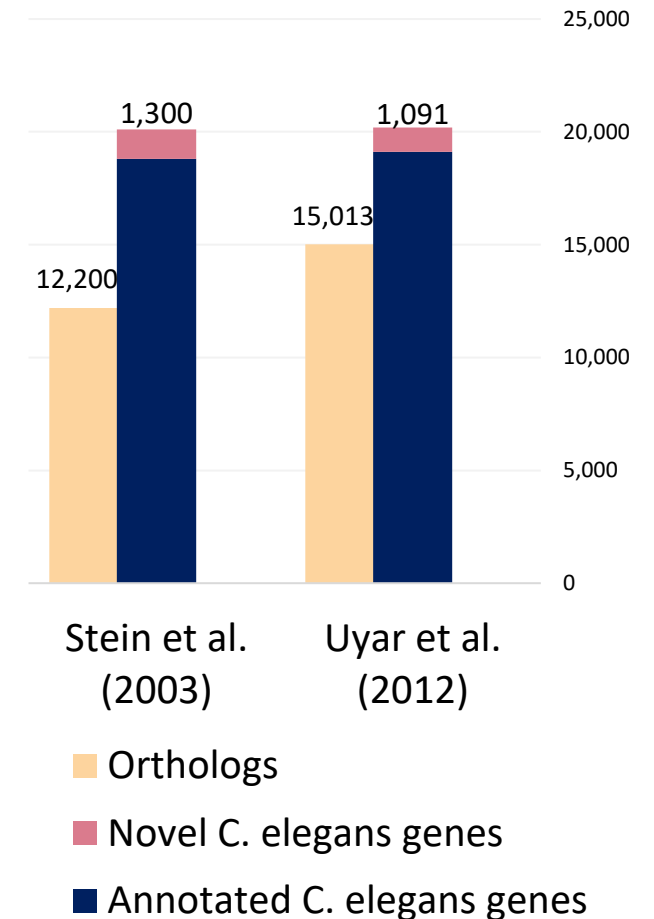- RNA Sequencing (RNA-Seq, **~800 libraries**)

Hypothesis #1: **The *C. briggsae* genome annotation is incomplete and can be improved using additional RNA-Seq data.**

# Complete genome annotation is essential for comparative genomics

- Protein-coding sequences and *cis*-regulatory regions are usually highly conserved

- A more complete or accurate genome annotation can increase the quality of the annotation of its close relative



*C. briggsae* gene models

*C. briggsae* improved gene models

*C. elegans* improved gene models

Hypothesis #2: **Using the improved *C. briggsae* annotation, we can find additional orthologous relationships with *C. elegans* that were previously missed and additional *C. elegans* gene models.**



Stein et al. (2003)

Uyar et al. (2012)

1,300    1,091
15,013
12,200

- Orthologs
- Novel C. elegans genes
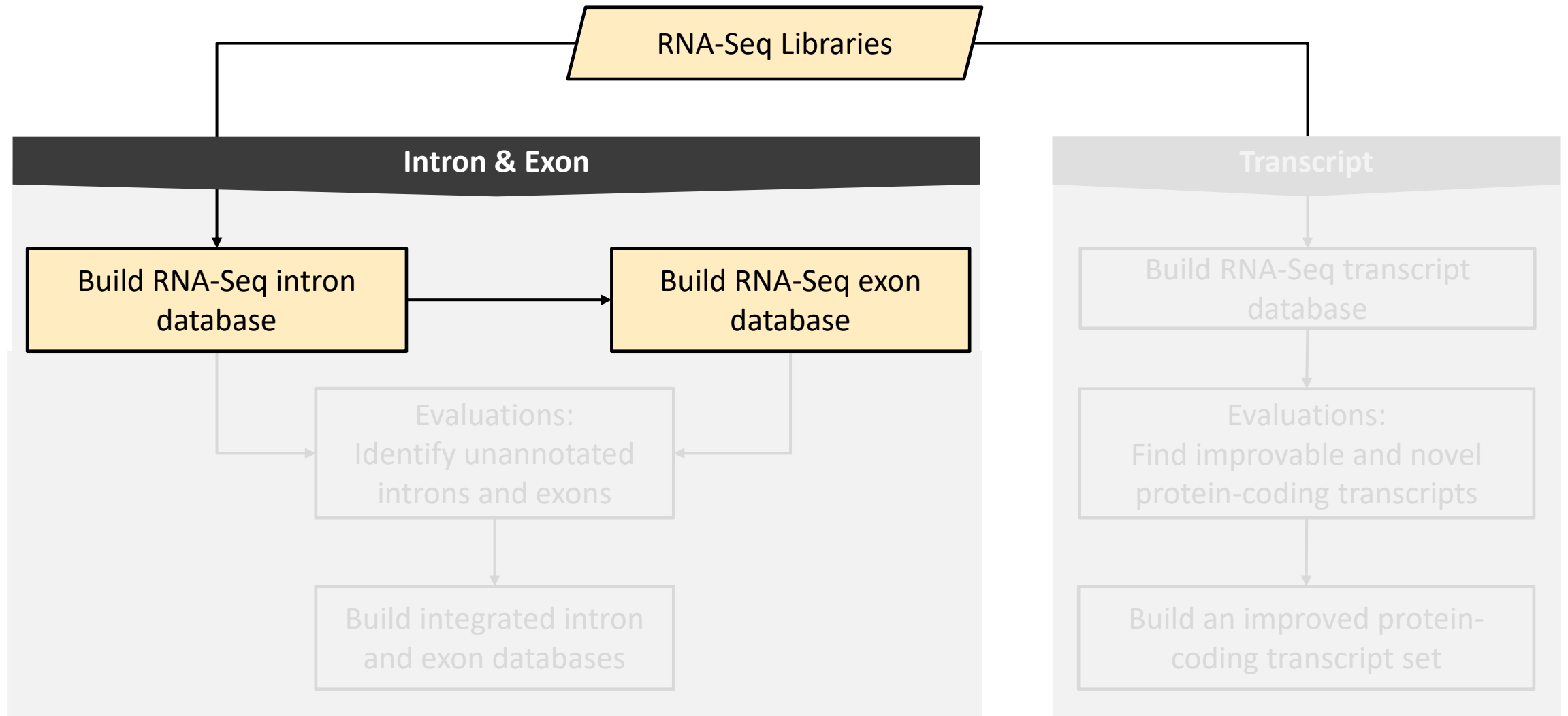- Annotated C. elegans genes

*C. briggsae* genome annotation is incomplete which limits its utility as a comparative platform for *C. elegans.*

- Hypothesis #1: The different number of molecular features observed between the two species is due to the incomplete annotation of *C. briggsae* genome. The *C. briggsae* annotation can be improved using RNA-Seq data.

- Hypothesis #2: Using the improved *C. briggsae* annotation, we can find additional orthologous relationships with *C. elegans* that were previously missed and additional *C. elegans* gene models.

# Specific Aims

- <u>Aim #1:</u> Improve the *C. briggsae* genome annotation at the intron, exon, transcript levels.

- <u>Aim #2:</u> Find additional orthologous relationships between *C. briggsae* and *C. elegans* and improve *C. elegans* genome annotation at the transcript level.

# Aim 1: Improve *C. briggsae* genome annotation

RNA-Seq Libraries

**Intron & Exon**

Build RNA-Seq intron database → Build RNA-Seq exon database

Evaluations:
Identify unannotated introns and exons

Build integrated intron and exon databases

**Transcript**

Build RNA-Seq transcript database

Evaluations:
Find improvable and novel protein-coding transcripts

Build an improved protein-coding transcript set

# Method: Building RNA-Seq intron and exon databases

**Data Selection & pre-processing**

1. Obtain public (NCBI) and in-house RNA-Seq data

11 publicly available, 2 in-house PE libraries (174M read pairs)

2. Remove excess rRNA reads

BBDuk[1] – 0.64% read pairs removed

3. Remove adapter and low-quality reads
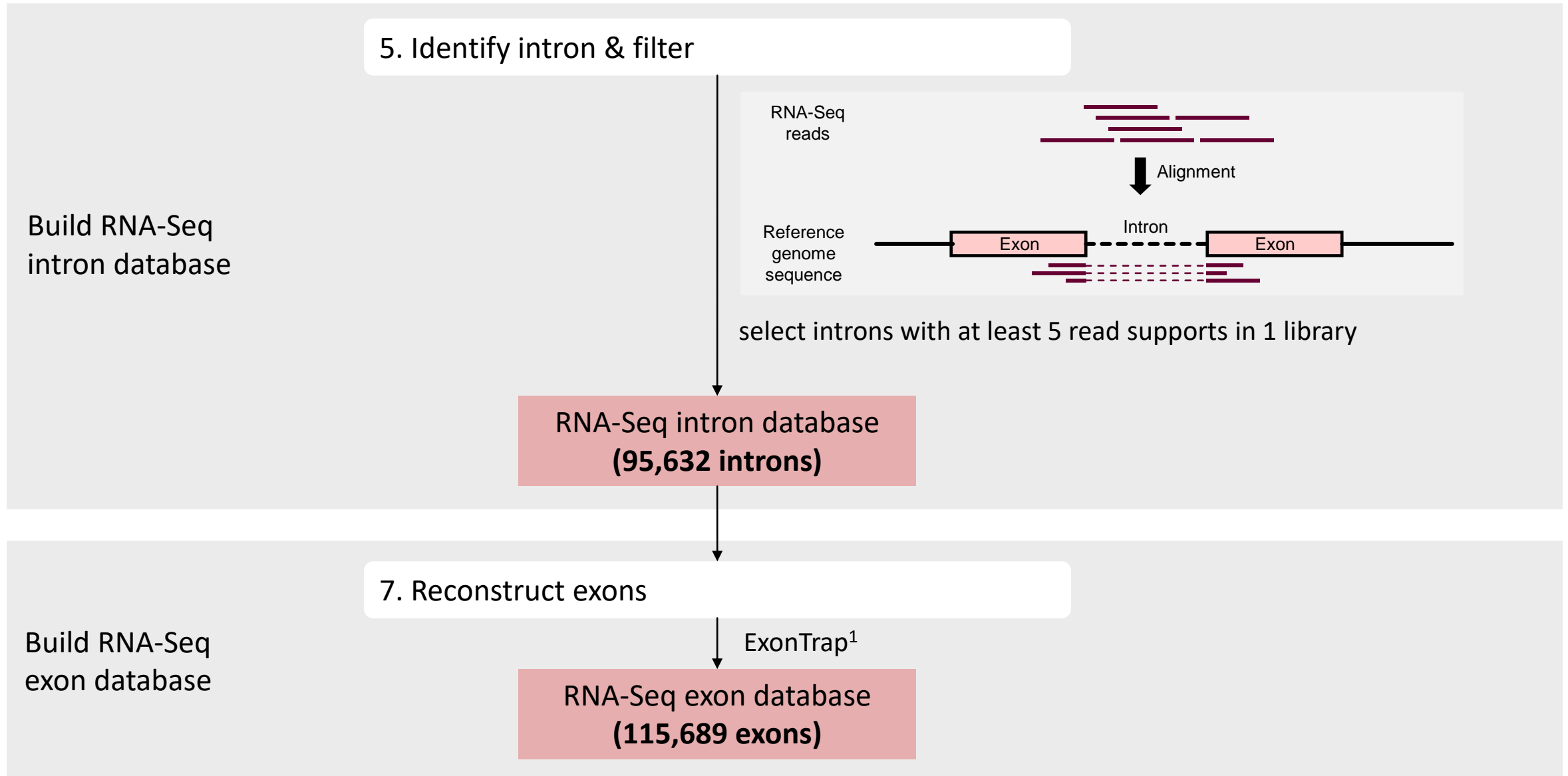
Trimmomatic[2] – 31.26% read pairs removed

**Alignment**

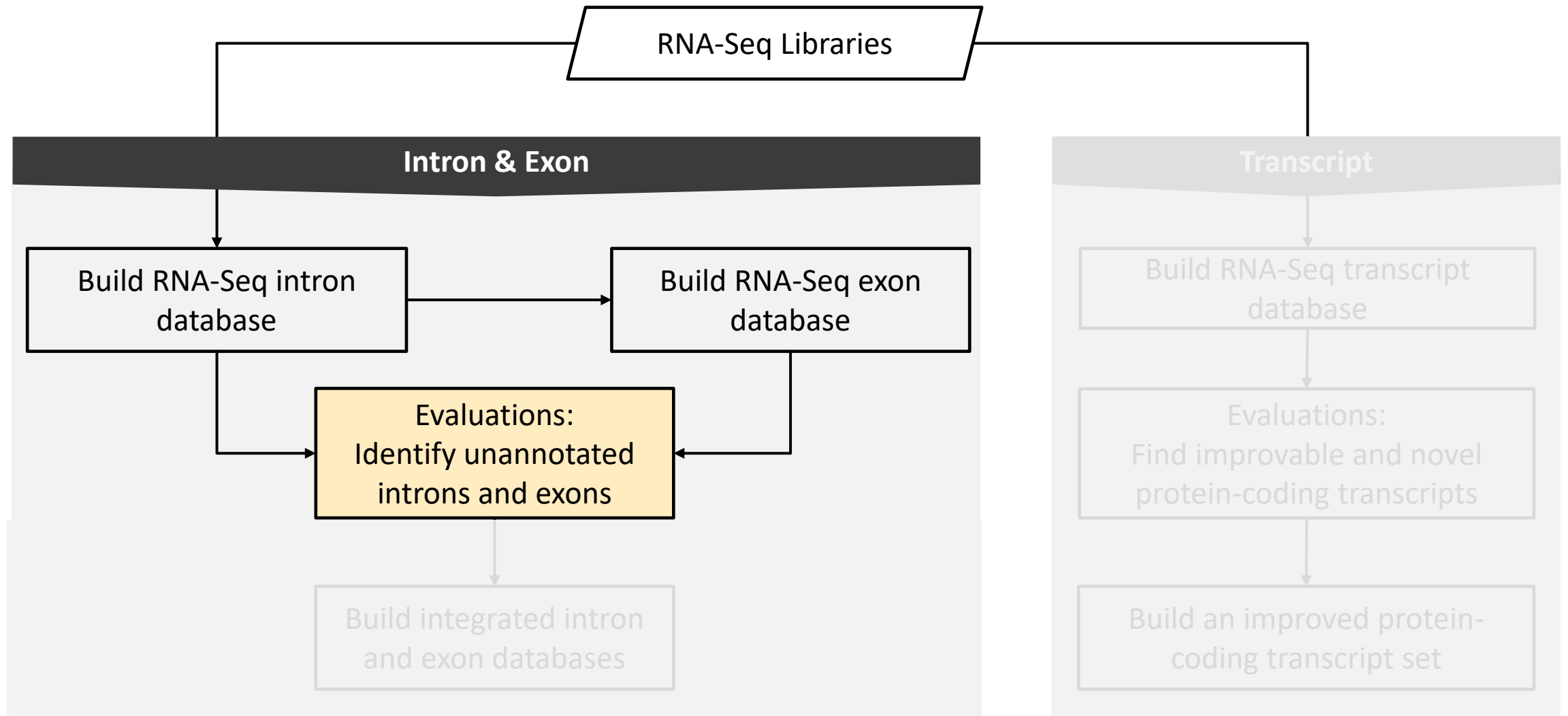4. Alignment to reference genome & filtering

STAR[3] – WS254 ref genome, filtered 10.36% multimapped alignments

[1]https://jgi.doe.gov/data-and-tools/bbtools/bb-tools-user-guide/bbduk-guide; [2]Bolger et al., 2014; [3]Dobin et al., 2013

# Method: Building RNA-Seq intron and exon databases

**5. Identify intron & filter**

**Build RNA-Seq intron database**

RNA-Seq reads

Alignment

Reference genome sequence

Exon   Intron   Exon

select introns with at least 5 read supports in 1 library

**RNA-Seq intron database (95,632 introns)**

**7. Reconstruct exons**

**Build RNA-Seq exon database**

ExonTrap[1]

**RNA-Seq exon database (115,689 exons)**

# Aim 1: Improve *C. briggsae* genome annotation

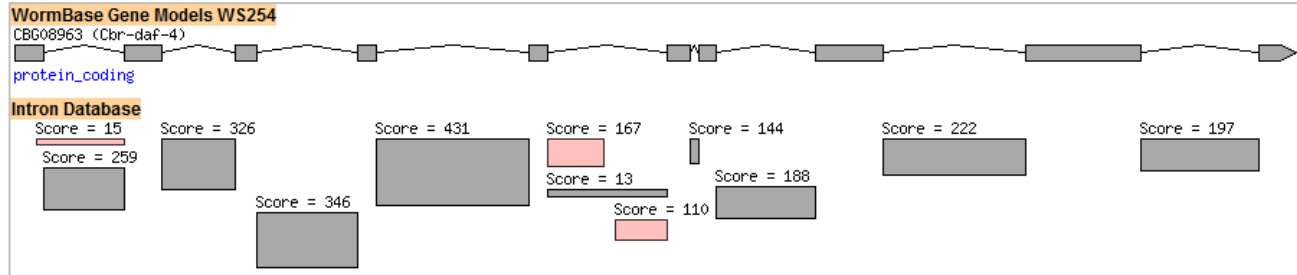**Intron**

RNA-Seq
95,632

74,972
(73%)

WormBase
103,314

All WormBase introns are
supported by RNA-Seq introns



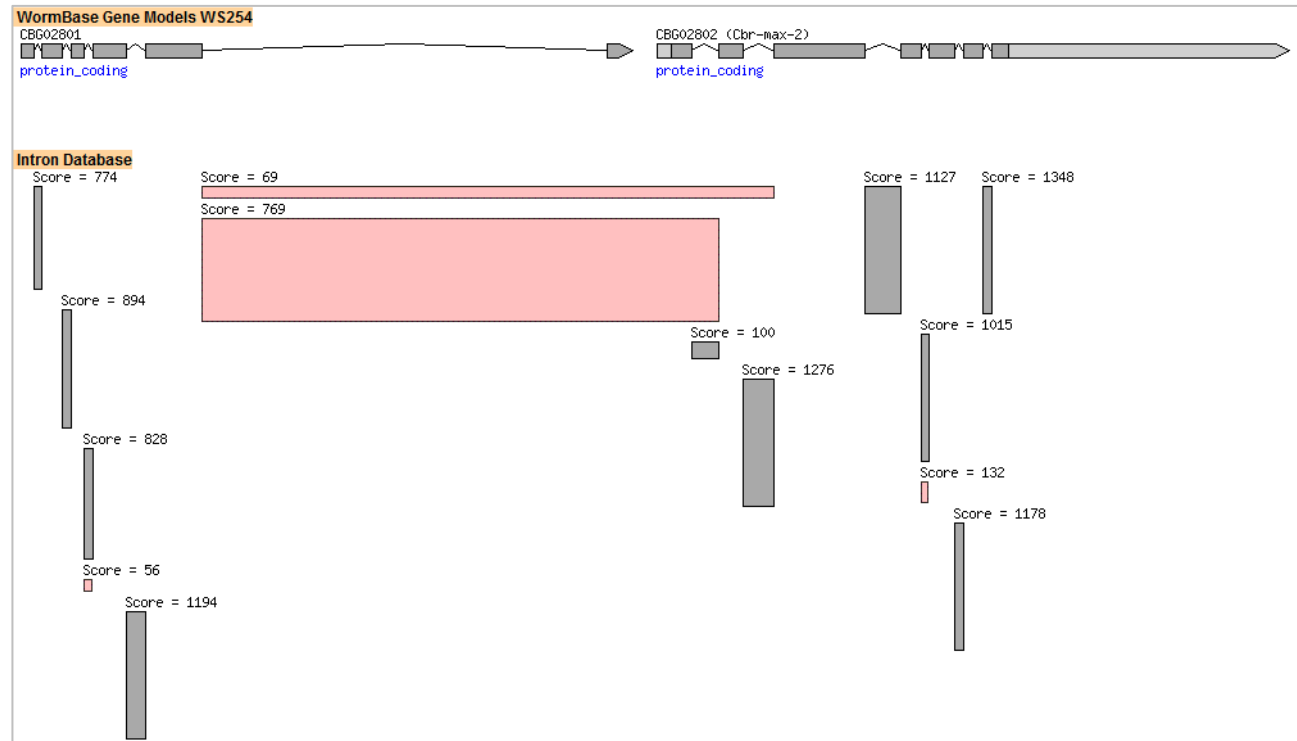**Pictured:** *Cbr-unc-18* (uncoordinated), an ortholog of *C. elegans unc-18*

Validating ~¾ of WormBase introns demonstrates the utility of our intron database

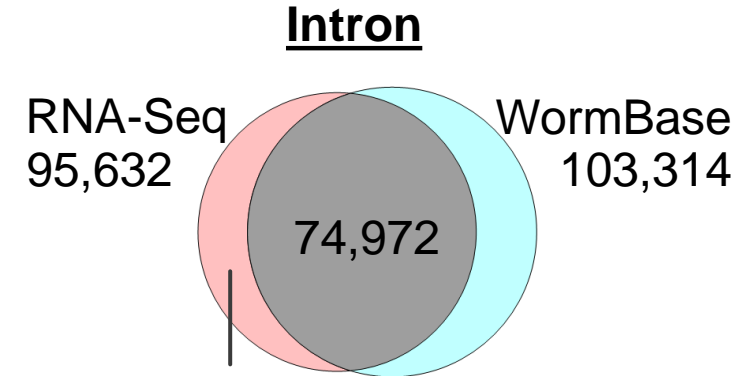WormBase = nematode information resource/database

14

**Pictured:** *Cbr-daf-4* (abnormal dauer formation), an ortholog of *C. elegans daf-4*



**Pictured:** *Cbr-max-2* (motor axon guidance), an ortholog of *C. elegans max-2*
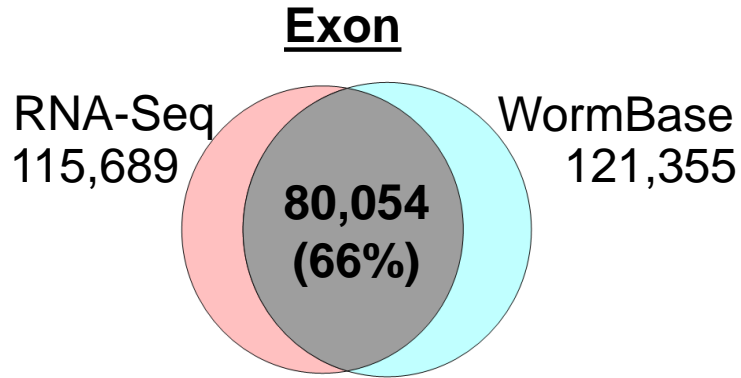
**Intron**



RNA-Seq 95,632    WormBase 103,314

74,972

**RNA-Seq specific**
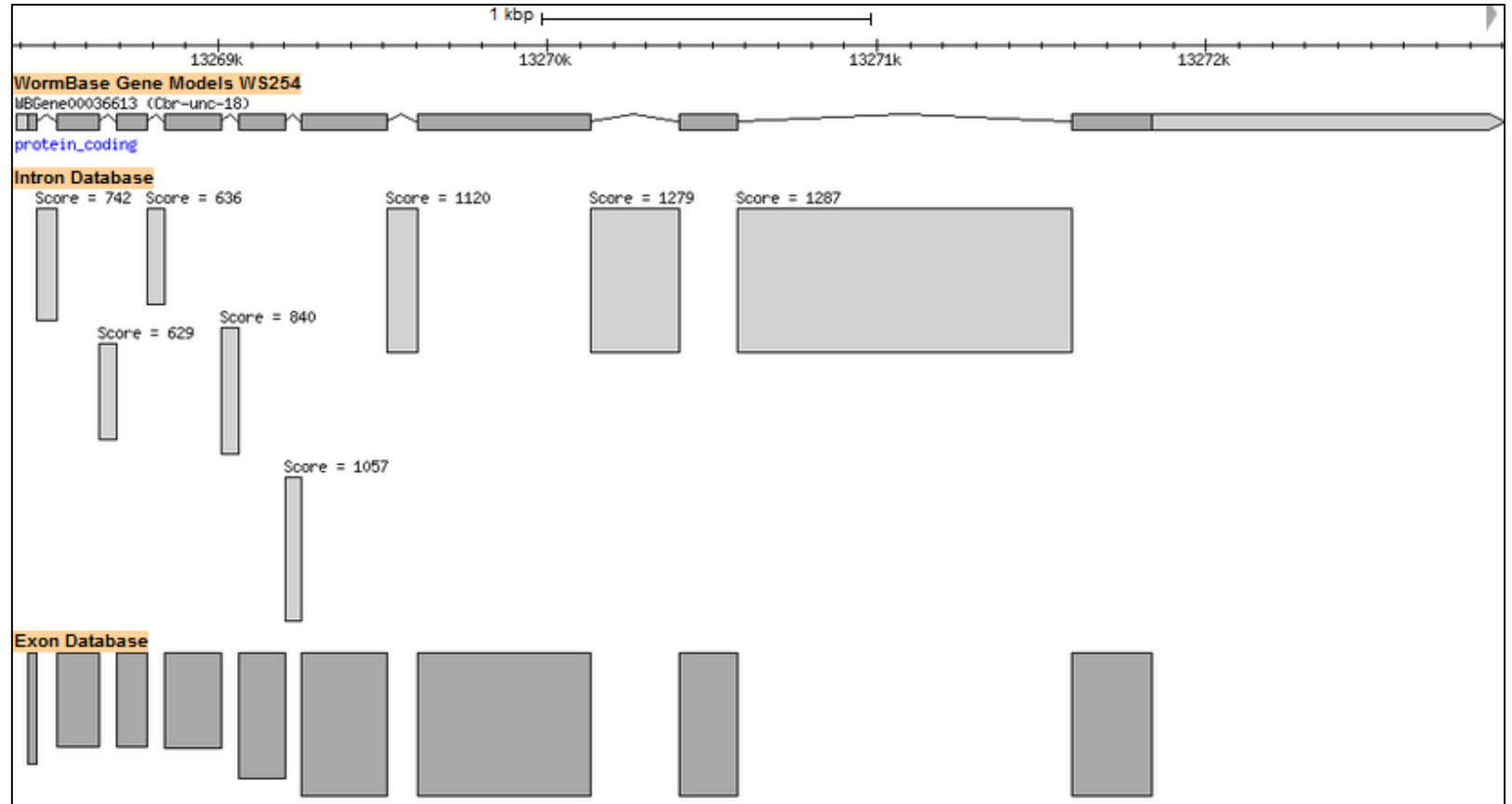**20,660 (22%)**
9,516 protein-coding genes affected

**Novel introns suggest gene model modifications and novel genes**
- 73% located internal of existing genes (top left)
- 12% extending existing genes
- 1% merging existing genes (bottom left)
- 14% of introns did not map to existing genes

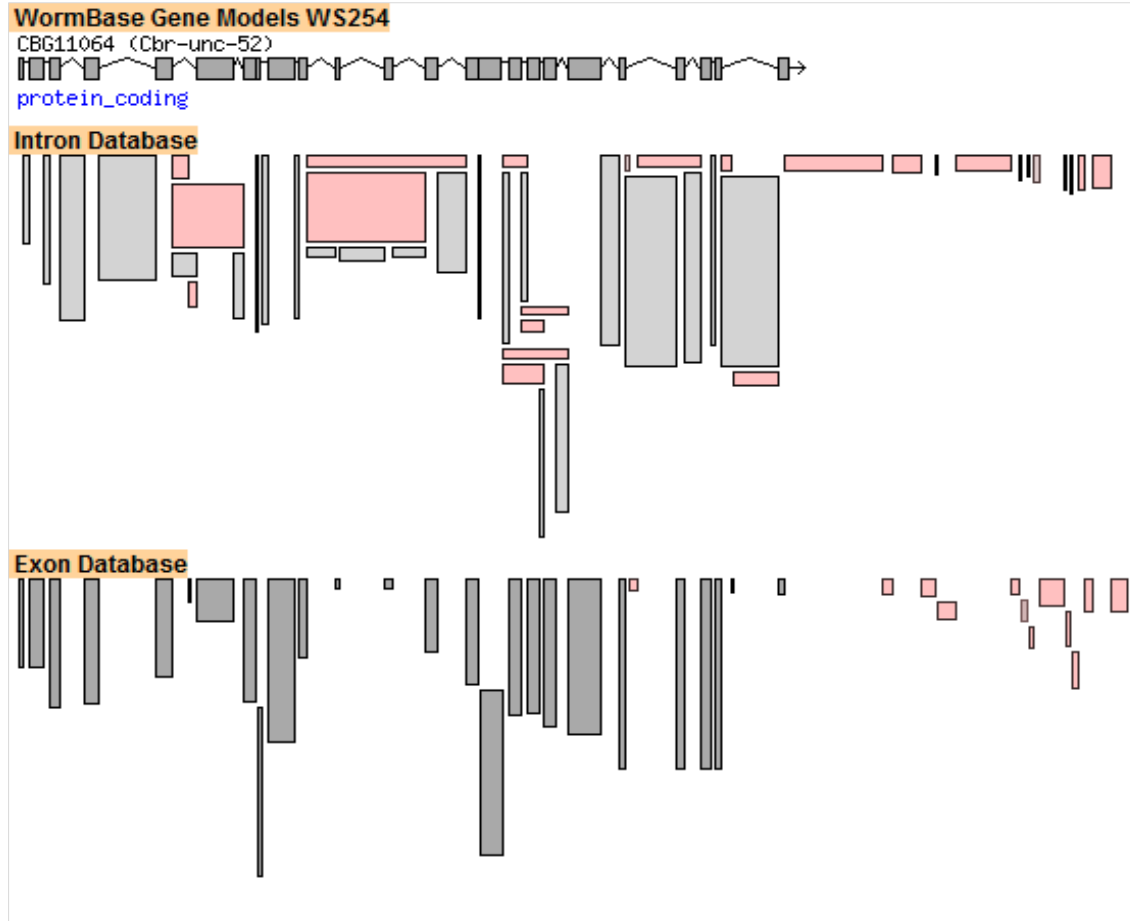# RNA-Seq exons validated 66% of WormBase exons



**Exon**

RNA-Seq
115,689

WormBase
121,355

**80,054
(66%)**

All WormBase exons are
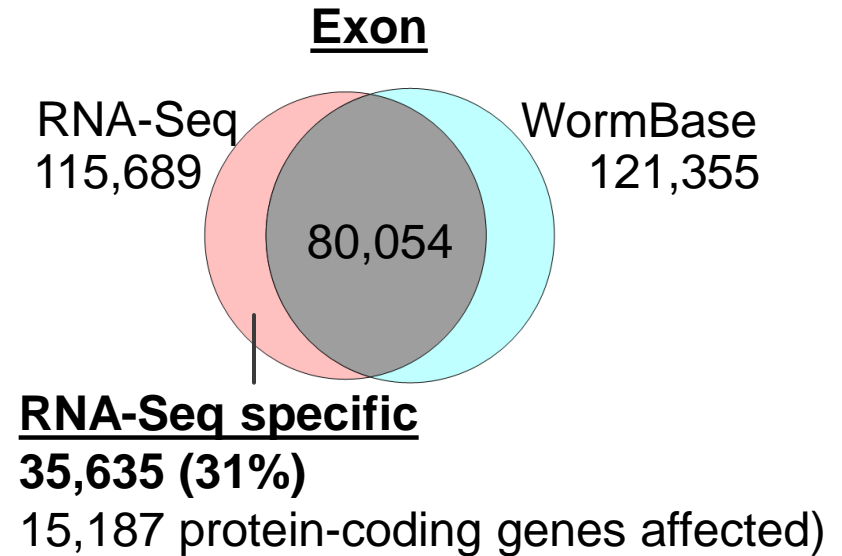supported by RNA-Seq exons

**Pictured:** *Cbr-unc-18* (<u>unc</u>oordinated), an ortholog of *C. elegans unc-18;* ~46% of WormBase transcripts are completely validated by our RNA-Seq introns and exons

Validating ⅔ of WormBase exons demonstrates the utility of our exon database
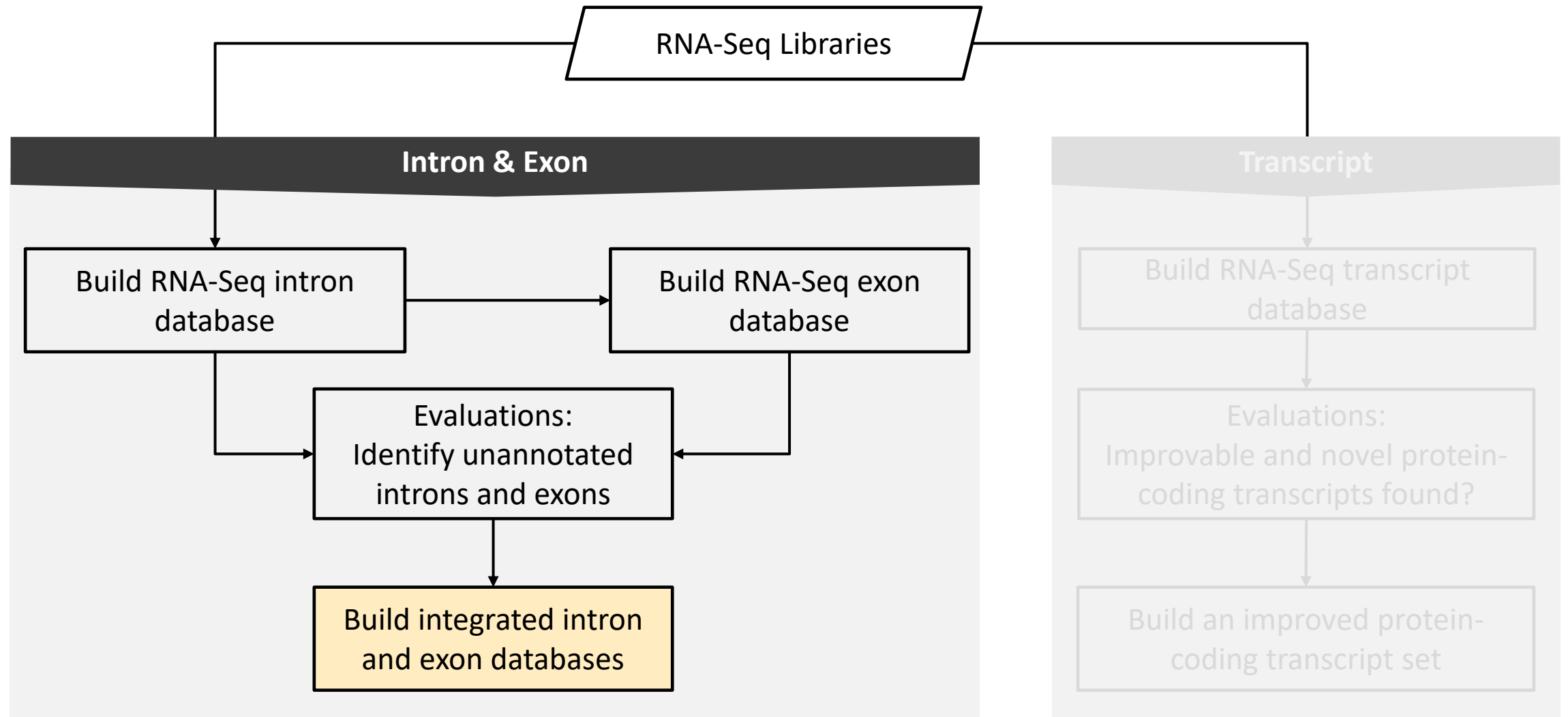
**Pictured:** *Cbr-unc-52* (underline{uncoordinated}), an ortholog of *C. elegans unc-52*; 31% of WormBase transcripts are partially validated by our RNA-Seq introns and exons

**Exon**

RNA-Seq
115,689

WormBase
121,355

80,054

**RNA-Seq specific**
**35,635 (31%)**
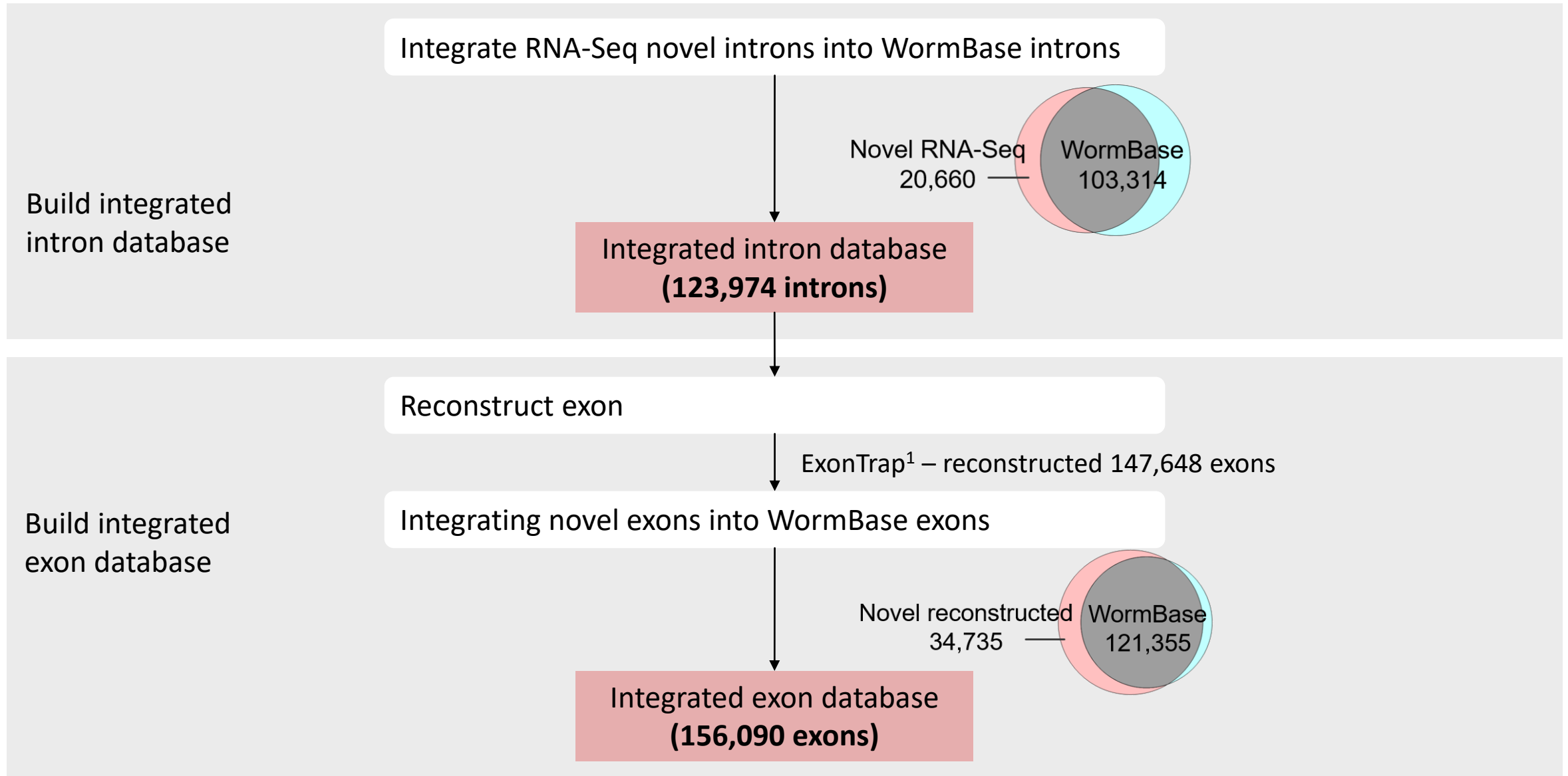15,187 protein-coding genes affected)

**Novel exons suggest gene model modifications and novel genes**
- 76% located internal of existing genes
- 14% extending existing genes (left)
- <1% merging existing genes
- 9% of exons did not map to existing genes
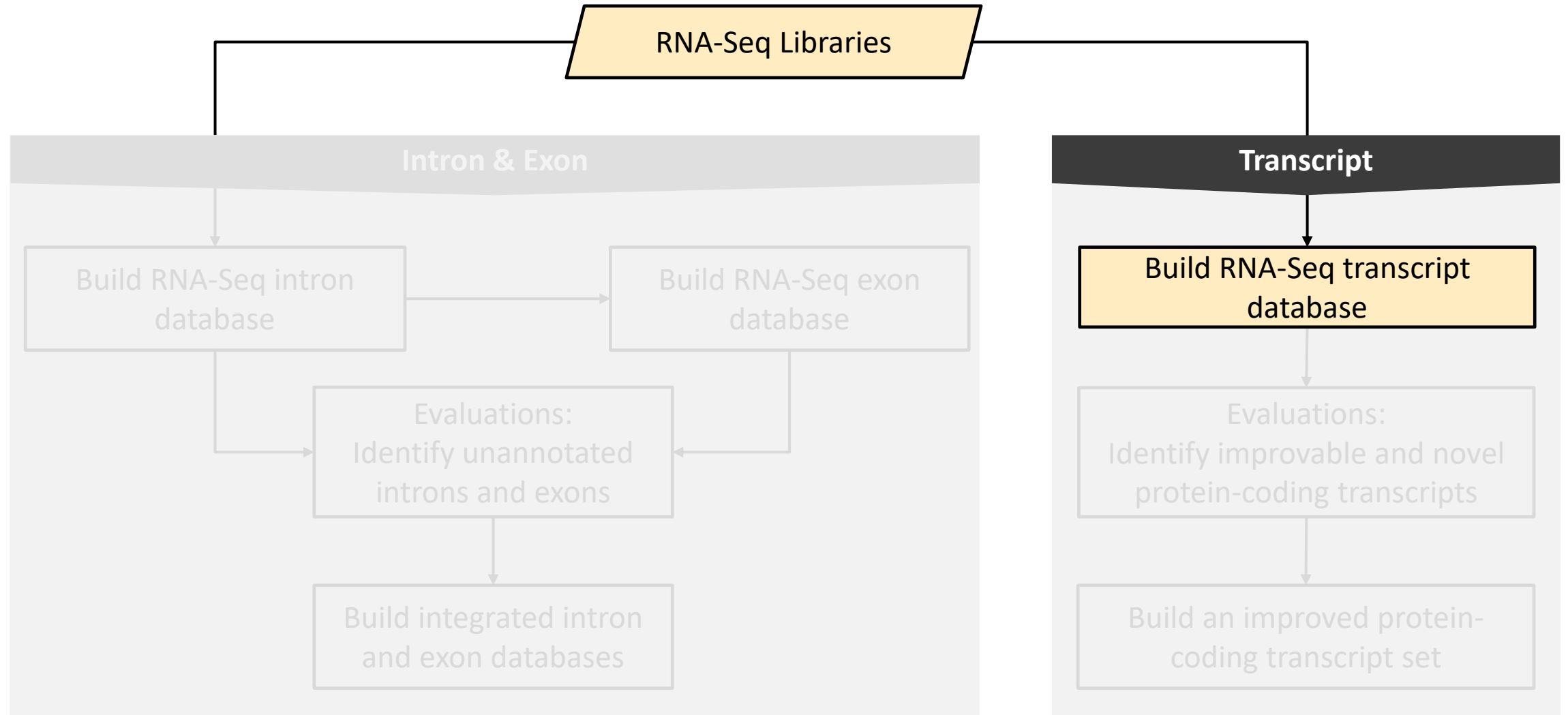
# Aim 1: Improve *C. briggsae* genome annotation

RNA-Seq Libraries

**Intron & Exon**

Build RNA-Seq intron database

Build RNA-Seq exon database

Evaluations: Identify unannotated introns and exons

Build integrated intron and exon databases

**Transcript**

Build RNA-Seq transcript database

Evaluations: Improvable and novel protein-coding transcripts found?

Build an improved protein-coding transcript set

# Method: Building improved intron and exon databases

**Build integrated intron database**

Integrate RNA-Seq novel introns into WormBase intrors

Novel RNA-Seq
20,660 —— WormBase
103,314

Integrated intron database
**(123,974 introns)**

**Build integrated exon database**

Reconstruct exon

ExonTrap[1] – reconstructed 147,648 exons

Integrating novel exons into WormBase exons

Novel reconstructed
34,735 —— WormBase
121,355

Integrated exon database
**(156,090 exons)**

- **Evidence that the *C. briggsae* annotation is incomplete and can be improved using 13 RNA-Seq libraries**

  - Identified 20,660 novel introns and 35,635 novel exons

  - Validated 73% and 66% WormBase introns and exons

- **Built improved intron and exon databases that serve as a more complete annotation at the intron and exon level**

  - Intron database consisting of 123,974 introns

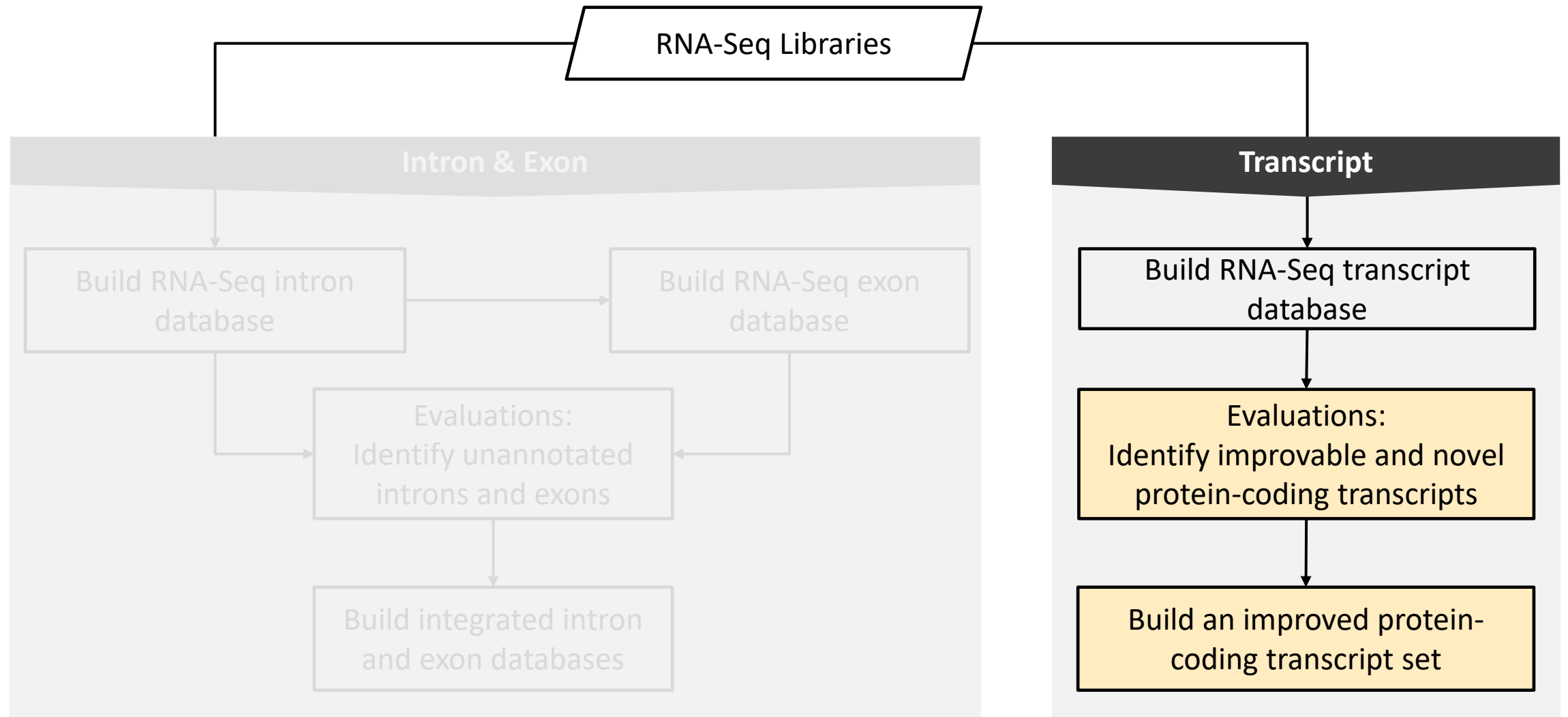  - Exon database consisting of 150,690 exons

# Aim 1: Improve *C. briggsae* genome annotation



RNA-Seq Libraries

**Intron & Exon**

Build RNA-Seq intron database

Build RNA-Seq exon database

Evaluations:
Identify unannotated introns and exons

Build integrated intron and exon databases

**Transcript**

Build RNA-Seq transcript database

Evaluations:
Identify improvable and novel protein-coding transcripts

Build an improved protein-coding transcript set

Build transcript database supported by introns and exons

1. *De-novo* and genome-guided transcript assembly

Cufflinks[1], StringTie[2], TransABySS[3]

2. Select transcripts supported by integrated introns and exons and merge transcripts

Gffcompare[4] – 29,352 fully supported transcripts

Predict coding regions

3. Identify candidate coding transcripts

TransDecoder[5]

Protein-coding transcripts **(24,705 transcripts)**

RNA-Seq Libraries

**Intron & Exon**

Build RNA-Seq intron database

Build RNA-Seq exon database

Evaluations:
Identify unannotated introns and exons

Build integrated intron and exon databases

**Transcript**

Build RNA-Seq transcript database

Evaluations:
Identify improvable and novel protein-coding transcripts

Build an improved protein-coding transcript set

WormBase transcript model

complete match

1. Complete match (WB confirmed)
2. 3' extension
3. 5' extension
4. 5' & 3' extension
5. Intron overlapping internal exon
6. Introns overlapping intron
7. Alternative donor (5'ss)
8. Alternative acceptor (3'ss)
9. Alternative donor & acceptor
10. Merged
11. Complex changes
12. Novel
13. Other (single-exon, partial)

partial match

no match

partial match

Criteria:

- Both assembled transcript and WB transcript have identical intron chains



Pictured: transcript CBG00984.2 of *Cbr-usp-14* (ubiquitin specific protease), ortholog of *C. elegans' usp-14.* All introns in WormBase CBG00984 transcript are observed in the assembled transcript

Criteria:

- All introns in the annotated transcript are observed in the assembled transcript

- One or more additional introns are found extending 3' of the transcript



The transcript we identified shares similar structure with its *C. elegans'* ortholog

Criteria:

- All introns in the annotated transcripts are observed in the assembled transcripts

- One more additional intron is found overlapping an annotated internal exon of existing transcripts

- Both transcripts have the same leftmost and rightmost intron boundaries



Intron with high read support

Share similar structure with its *C. elegans'* ortholog

Criteria:

- Has not been annotated previously

- No overlapping genes in the genomic regions



Pictured: Transcript MERGE_00026402 (nCBG00109). Previously not annotated by WormBase

WormBase transcript model

1. Complete match (WB confirmed)
2. 3' extension
3. 5' extension
4. 5' & 3' extension
5. Intron overlapping internal exon
6. Introns overlapping intron
7. Alternative donor (5'ss)
8. Alternative acceptor (3'ss)
9. Alternative donor & acceptor
10. Merged
11. Complex changes
12. Novel
13. Other (single-exon, partial)

partial match

no match

- Using intron chain comparison method, a transcript is discarded when either WormBase or assembled transcript is single-exon (i.e., no intron)

WormBase transcript

Assembled transcript (Case 1)

Assembled transcript (Case 2)



Case 2

Single-exon

Multi-exon

Pictured: transcript CBG23430.2 of *Cbr-bnc-1* (basonuclin-1 zinc finger protein homolog), ortholog of *C. elegans' bnc-1*

# Identified 6,285 candidate protein-coding transcripts

| # | Category | Diagram | Protein-coding transcripts | Protein-coding genes |
|---|----------|---------|---------------------------|---------------------|
| 1. | Complete match (WB confirmed) | | 8,080 (32.71%) | 8,055 |
| 2. | 3' extension | | 316 (1.28%) | 287 |
| 3. | 5' extension | | 753 (3.05%) | 687 |
| 4. | 5' & 3' extension | | 26 (0.11%) | 25 |
| 5. | Intron overlapping internal exon | | 358 (1.45%) | 332 |
| 6. | Introns overlapping intron | | 217 (0.88%) | 205 |
| 7. | Alternative donor (5'ss) | | 777 (3.15%) | 746 |
| 8. | Alternative acceptor (3'ss) | | 882 (3.57%) | 810 |
| 9. | Alternative donor & acceptor | | 346 (1.40%) | 327 |
| 10. | Merging 2 or more genes | | 206 (0.83%) | 116 |
| 11. | Complex changes | | 2,245 (9.09%) | 1,517 |
| 12. | Novel | | 159 (0.64%) | 159 |
| 13. | Other – Single-exon (no intron) | | 120 (0.49%) | 95 |
| | Other – Partial | | 10,220 (41.37%) | 5,304 |

# Method: Building an improved transcript database

| Category | | Diagram | Protein-coding transcripts | With proper start & stop codon |
|---|---|---|---|---|
| 2. | 3' extension | | 316 | 284 |
| 3. | 5' extension | | 753 | 692 |
| 4. | 5' & 3' extension | | 26 | 23 |
| 5. | Intron overlapping internal exon | | 358 | 341 |
| 6. | Introns overlapping intron | | 217 | 197 |
| 7. | Alternative donor (5'ss) | | 777 | 717 |
| 8. | Alternative acceptor (3'ss) | | 882 | 818 |
| 9. | Alternative donor & acceptor | | 346 | 284 |
| 10. | Merged | | 206 | 179 |
| 11. | Complex changes | | 2,245 | 2,015 |
| 12. | Novel | | 159 | 104 |
| **Total** | | | 6,285 (100.00%) | **5,654 (89.96%)** |

Integrated protein-coding transcripts
**(28,129 transcripts)**

# Transcript summary

- **Evidence that the *C. briggsae* annotation is incomplete and can be improved using 13 RNA-Seq libraries**

  - Identified 24,705 protein-coding transcripts, including 5,654 candidate protein-coding transcripts to improve *C. briggsae* annotation

- **Generated an improved *C. briggsae* transcript database consisting of 28,129 transcripts (25% higher than current annotation)**

# Discovery power is proportional to RNA-Seq data quantity

- Using 13 *C. briggsae* RNA-Seq libraries, we identified thousands of introns, exons, and transcripts. We hypothesized that data availability does limit the discovery of features.

- Method: Applied the same pipeline on limited *C. elegans* RNA-Seq data



| 13 *C. elegans* libraries of equivalent sizes | Following the same pipeline *C. briggsae* | Intron database ? introns | Exon database ? exons | Transcript database ? transcripts |

**Introns**

| | | |
|---|---|---|
| C. briggsae (WS254) | | 103,314 |
| C. briggsae (13 libraries, 174M PE reads) | | 95,632 |
| C. elegans (13 libraries, 177M PE reads) | | **93,107** |
| C. elegans (802 libraries, 54B PE reads) | | 239,333 |

**Exons**

| | |
|---|---|
| 121,355 |
| 115,689 |
| **111,365** |
| 328,212 |

**Protein-coding transcripts**

| | |
|---|---|
| 21,814 |
| 22,110 |
| **20,402** |
| 72,274 |

## Orthology analysis

*C. briggsae* WormBase WS254 protein-coding transcripts

*C. briggsae* improved protein-coding transcripts

Ortholog assignment to *C. elegans* WormBase WS254 protein-coding transcripts

## *C. elegans* gene models improvement

Predicted *C. elegans* protein-coding transcripts

Evaluations:
- Find improvable and novel protein-coding transcripts
- Filter those supported by RNA-Seq introns

Orthology analysis before and after *C. briggsae* transcript improvement

1. Obtain peptide sequences and chose longest peptide as representatives

*C. briggsae* WS254: 21,814 longest peptides
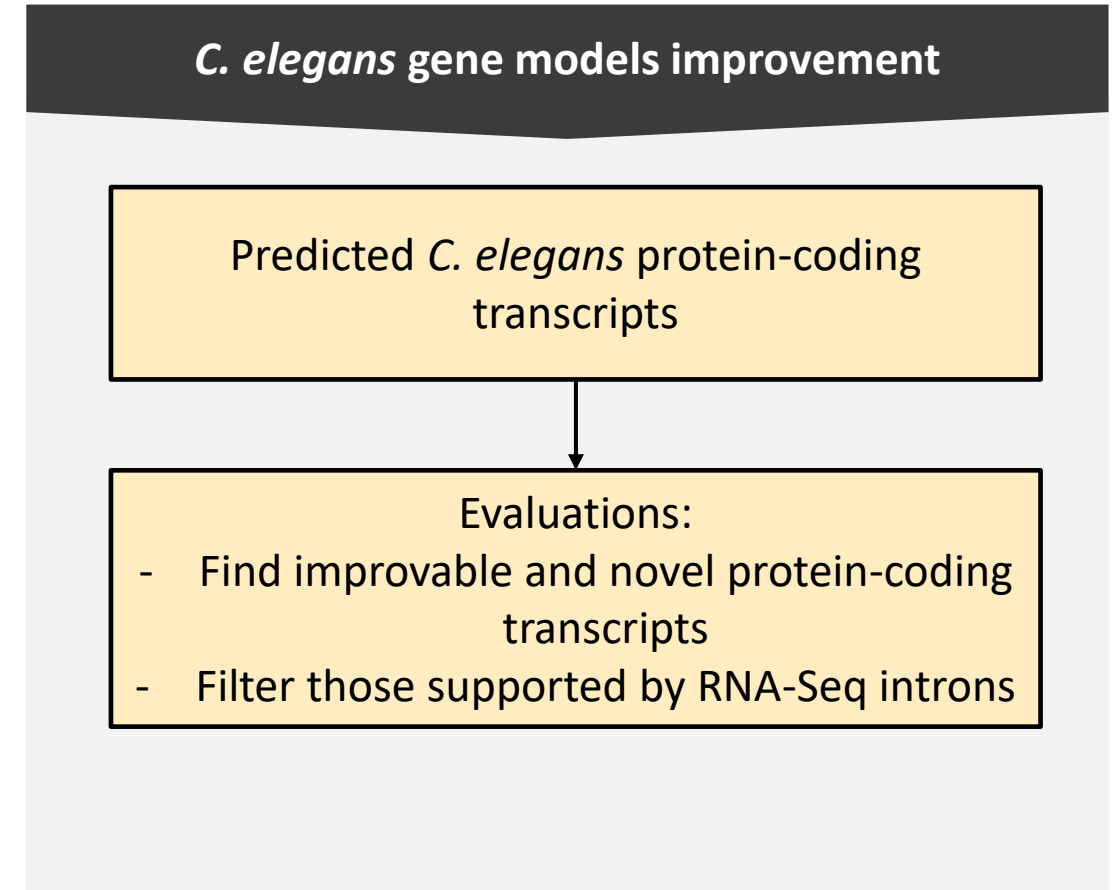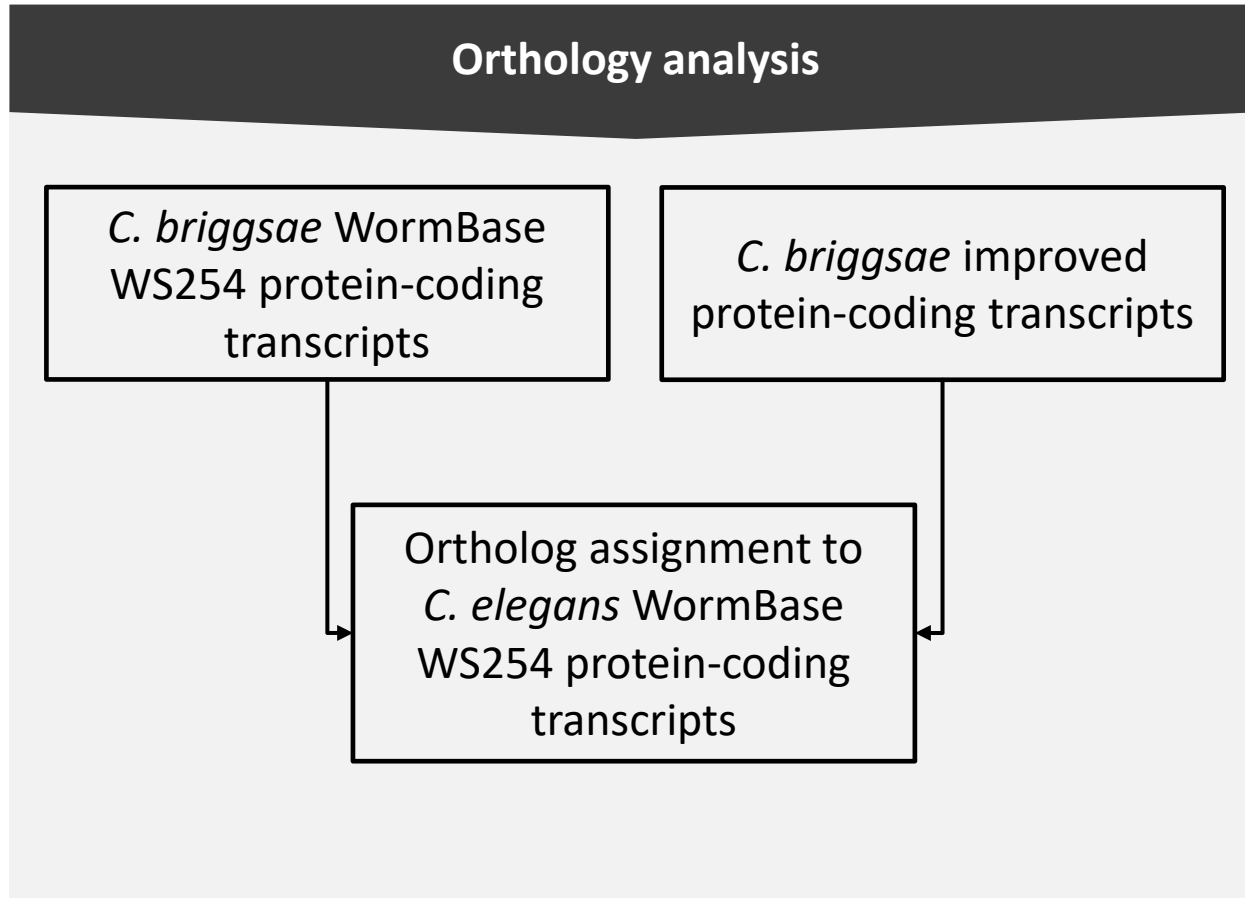*C. briggsae* improved: 21,913 longest peptides
*C. elegans* WS254: 20,254 longest peptides

2. All-vs-all BLASTP

OrthoMCL[1], E-value cutoff 1e-5

| | *C. briggsae* WS254 transcripts | *C. briggsae* improved transcripts |
|---|---|---|
| *C. elegans* WS254 transcripts | Ortholog pairs = 16,748 | Ortholog pairs = 16,880 (**32 pairs belong to novel genes**) |

[1]Li et al., 2003

# Aim 2: Find additional orthologs and improve *C. elegans* annotation

## Orthology analysis

*C. briggsae* WormBase WS254 protein-coding transcripts

*C. briggsae* improved protein-coding transcripts

Ortholog assignment to *C. elegans* WormBase WS254 protein-coding transcripts

## *C. elegans* gene models improvement

Predicted *C. elegans* protein-coding transcripts

Evaluations:
- Find improvable and novel protein-coding transcripts
- Filter those supported by RNA-Seq introns

# Method: *C. elegans* gene model improvement

**Prediction**

1. Predict *C. elegans* protein-coding transcripts using *C. briggsae* candidate protein-coding transcripts

GeMoMa[1] – predicted 4,313 protein-coding transcripts

2. Evaluation

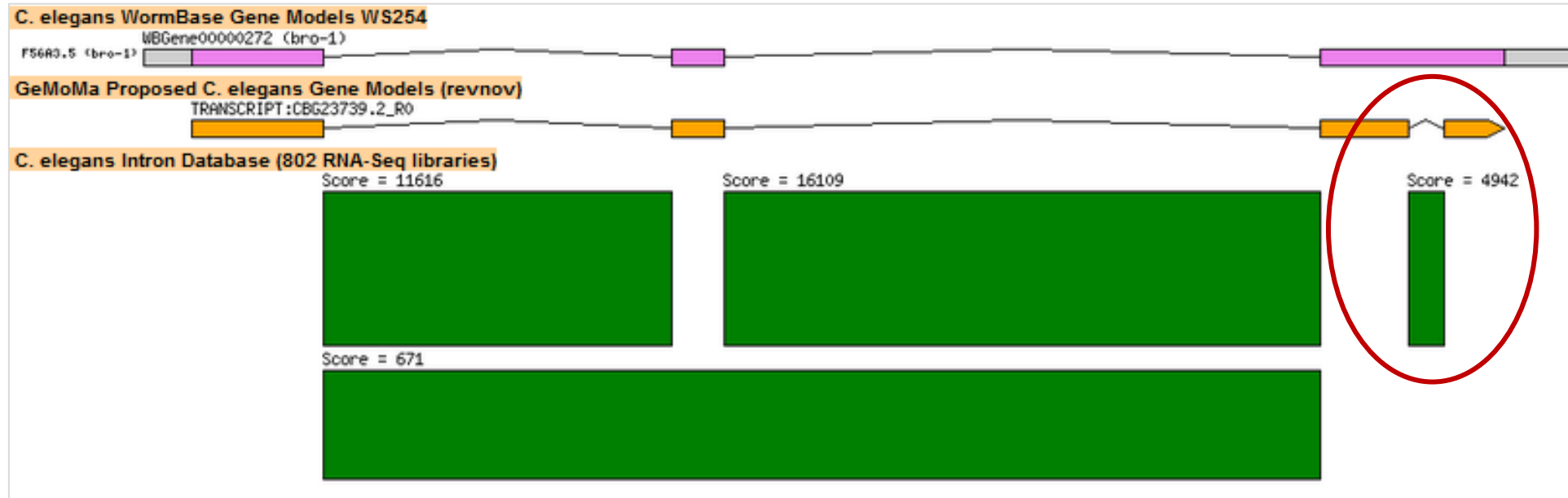categorized transcripts into 13 categories

**Evaluation**

Candidate *C. elegans* protein-coding transcripts
(1,698 transcripts)

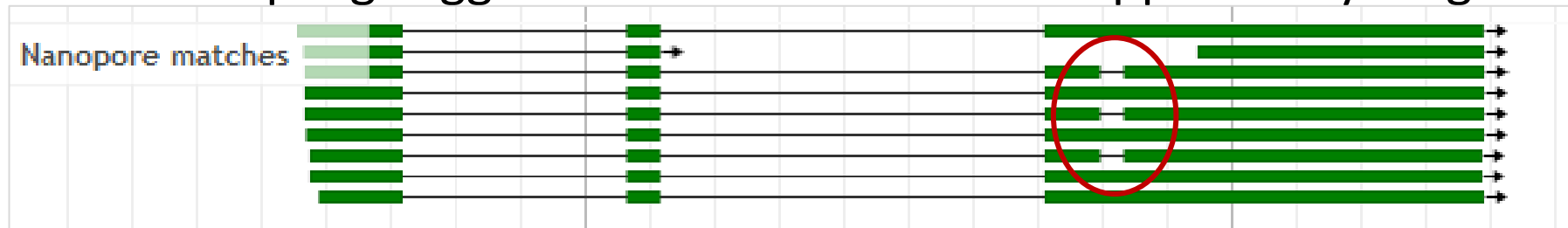used *C. elegans* RNA-Seq introns from 802 libraries[2]

Candidate *C. elegans* protein-coding transcripts
with RNA-Seq support
**(279 improvable transcripts and 2 novel genes)**

[1]Keilwagen et al., 2018, [2]Douglas, M., 2018

*C. elegans* intron database credit: Matt Douglas

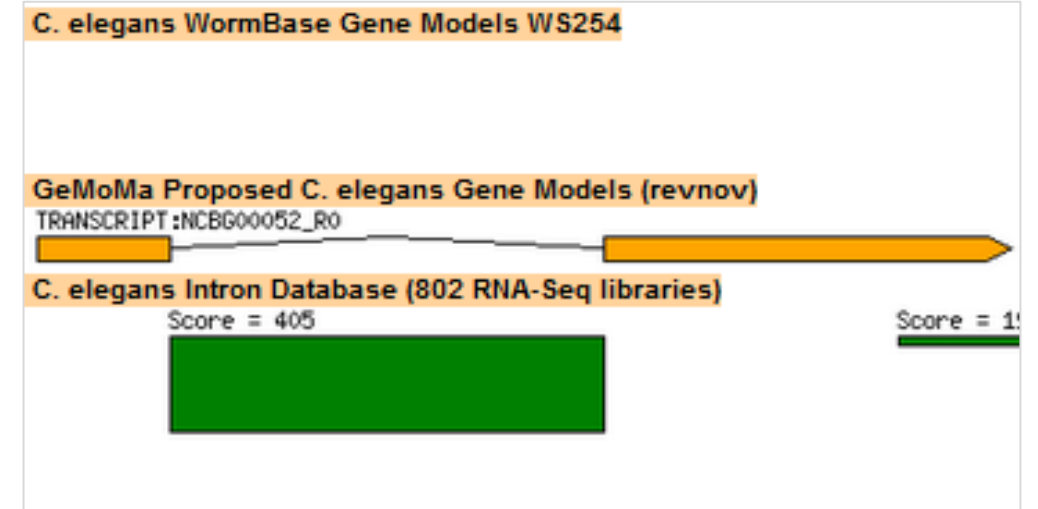- Random sampling suggested that some are also supported by long-read alignments



Long-read alignments credit: WormBase JBrowse (as of Dec 20th, 2019)

Pictured: *bro-1* (BROther (Drosophila tx factor partner) homolog), predicted to contribute to sequence-specific DNA binding activity, ortholog of human CBFB (core-binding factor subunit beta). Human ortholog of this gene are implicated in acute myeloid leukemia.

- Previously have not been annotated

- No overlapping genes in the genomic regions

- Have more than 400 reads supporting the introns

# Comparative analyses summary
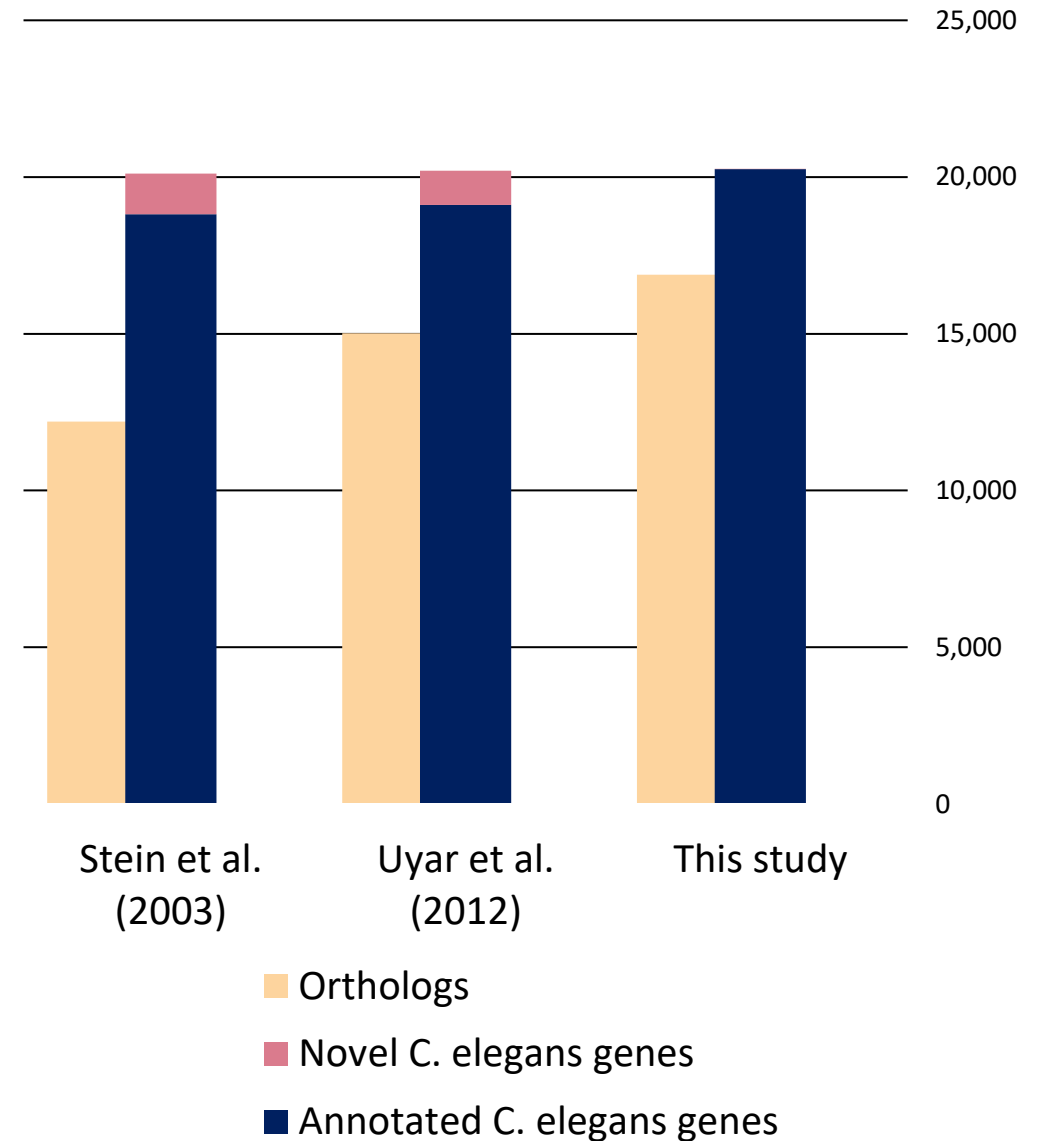
- **Homology and RNA-Seq based comparative analysis using the improved *C. briggsae* annotation revealed more ortholog relationships between *C. briggsae* and *C. elegans,* and improved *C. elegans* annotation**

  - Revealed 132 new ortholog pairs, 32 belong to *C. briggsae* novel transcripts

  - Revealed 279 transcripts and 2 novel *C. elegans* genes

# Conclusions

- RNA-Seq provides evidence to improve *C. briggsae* annotation

  - Reveals thousands of novel introns and exons, as well as hundreds of novel protein-coding transcripts

- The improved *C. briggsae* annotation together with comparative analyses reveals novel *C. briggsae–C. elegans* ortholog relationships and novel *C. elegans* protein-coding transcripts

- Despite limited data available for *C. briggsae*, the improved annotation has enhanced the utility of *C. briggsae* as a comparative platform for *C. elegans*.

# Significance

- As more RNA-Seq data becomes available, this method can be used to further refine not only *C. briggsae* annotation but also *C. elegans* annotation.

# Acknowledgements

**Senior Supervisor**

Dr. Jack Chen

**Committee Members & Examining Committee**

Dr. Fiona Brinkman

Dr. Ryan Morin

Dr. Christopher Beh

Dr. Mark Paetzel

**Chen Lab Members (past&present)**

Dr. Jiarui Li

Dr. Michelle Hu

Matthew Douglas

Marija Jovanovic

Kate Gibson

**Family and friends**

Thio, Adair, Uplands family

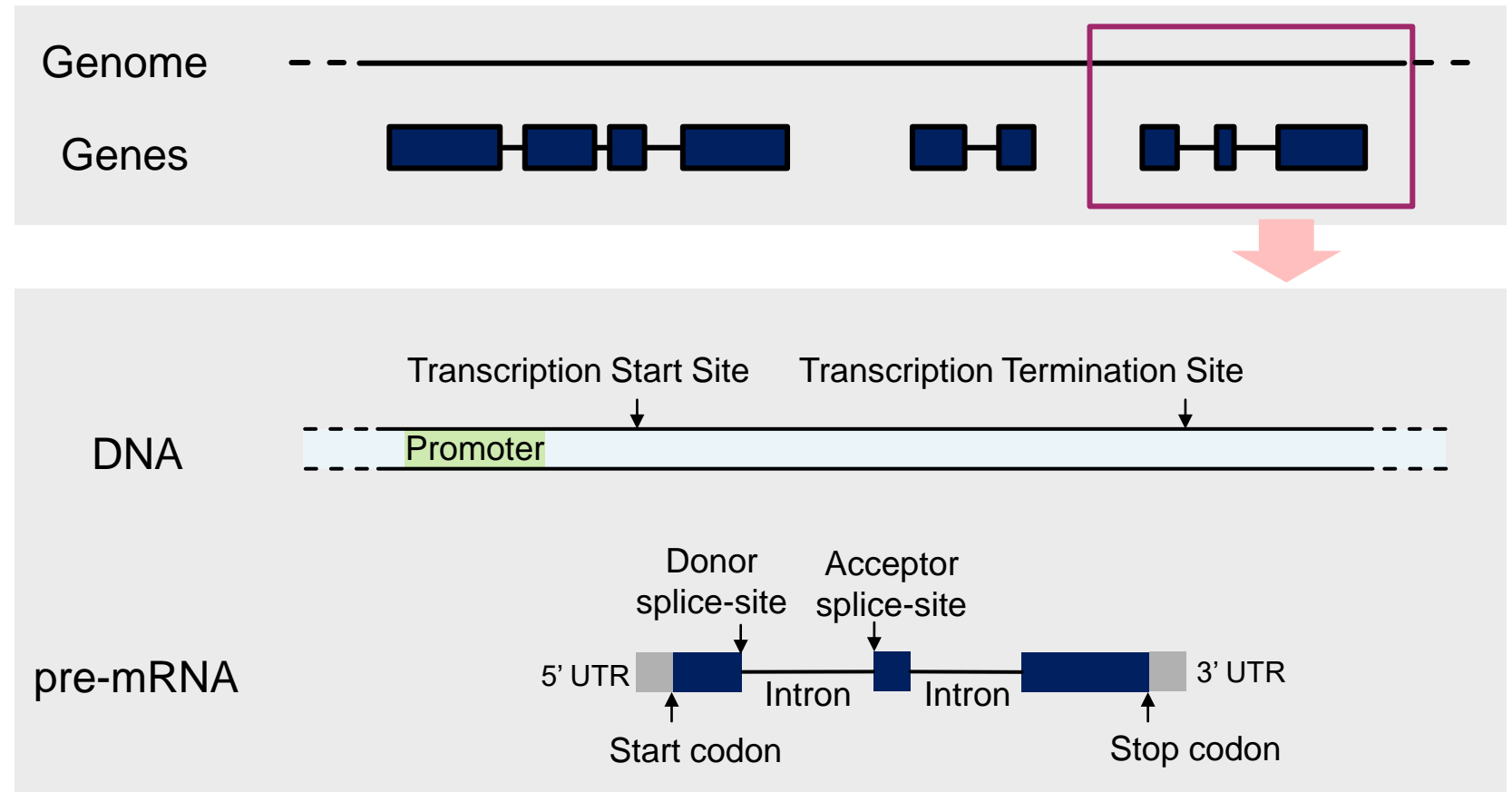MBB, SFU Omics, Indo friends

# Extras

# What is genome annotation?
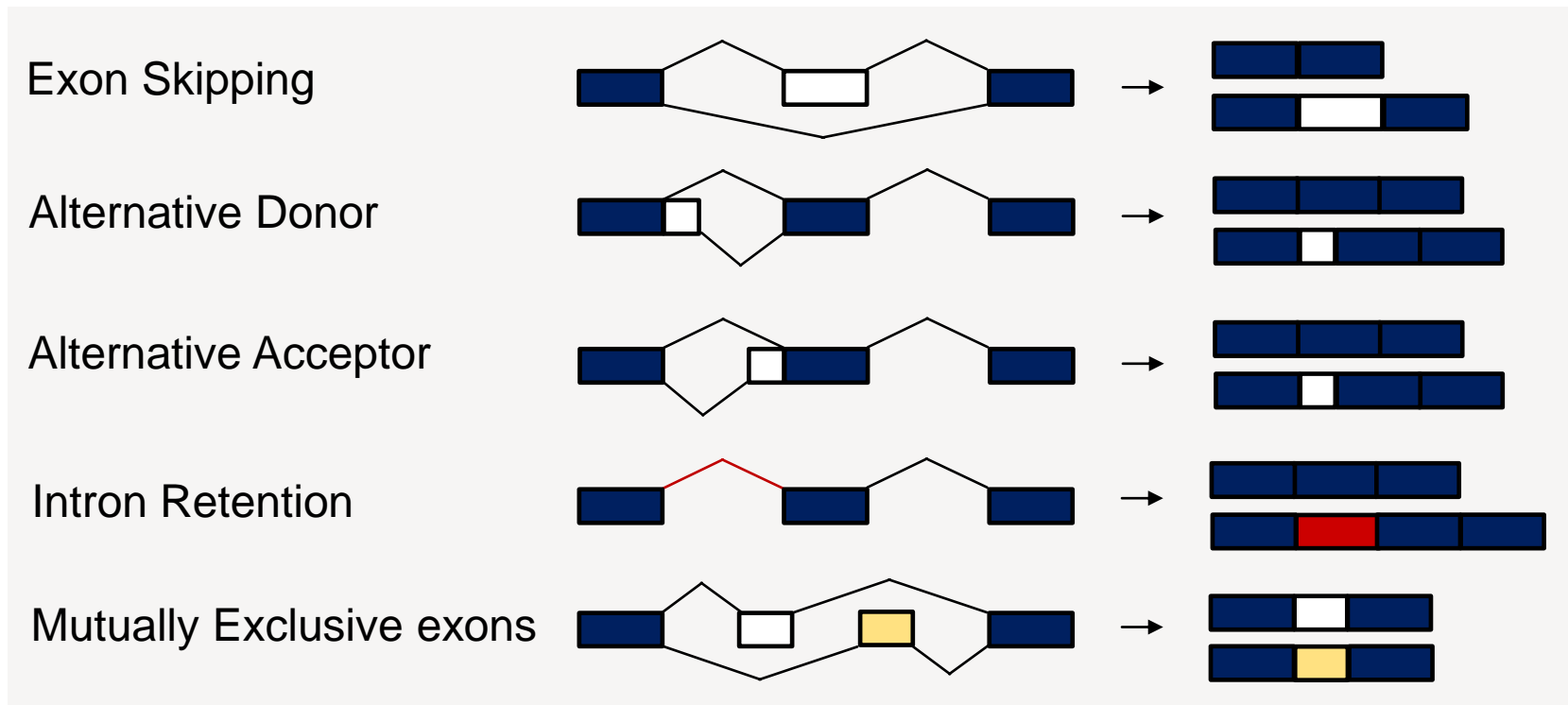
## The process of finding location and structure of all genes in the genome

- Introns
- Exons
- Splice sites
- Coding regions
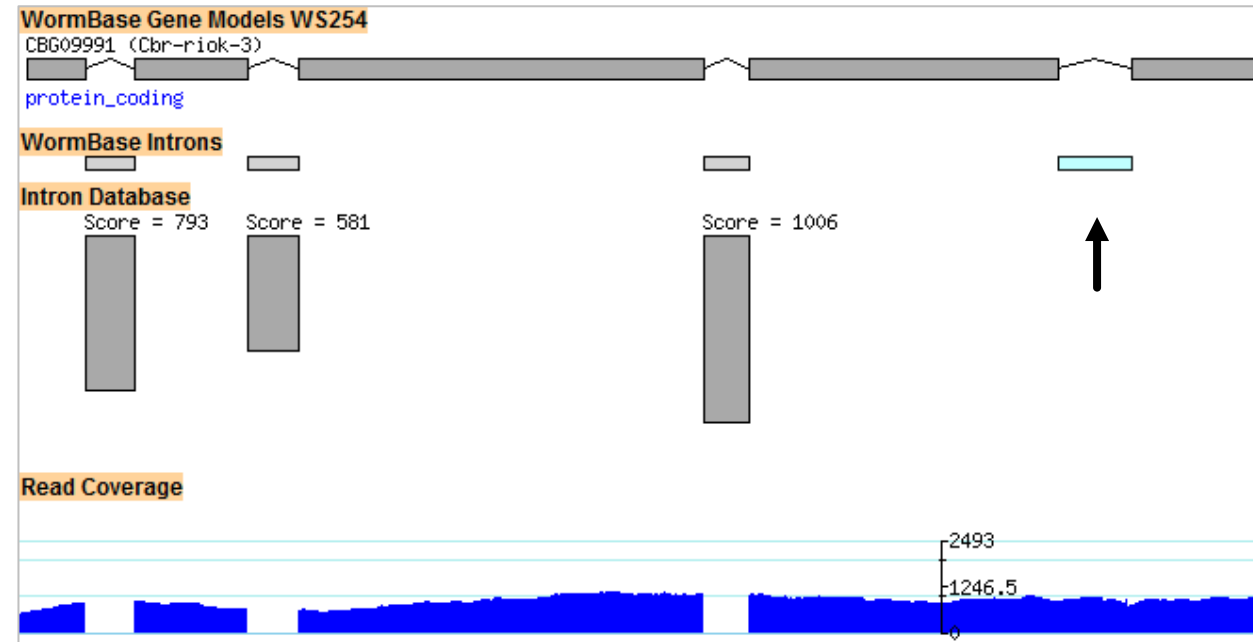- Promoter sequences
- TF binding sites
- and many more..

- Genes can have alternative splicing pathways to process pre-mRNAs into two or more mature transcripts that encode different proteins.

- This can contribute to the completeness of the annotation.

RNA-Seq
95,632

WormBase
103,314

74,972

**WormBase specific**
**28,341 (27%)**

Lack of support
(No read support) — 21%

Lack of support (despite read alignments) — 42%

Outside of intron size identification — 5%

Inadequate support (below score) — 33%

WormBase Gene Models WS254
CBG09991 (Cbr-riok-3)
protein_coding

WormBase Introns

Intron Database
Score = 793     Score = 581     Score = 1006

Read Coverage
2493
1246.5
0

RNA-Seq
115,689

WormBase
121,355

80,054

WormBase specific
41,301 (34%)

- Less than 50% of *C. briggsae* transcript models are validated by our intron and exon databases

- 31% are partially validated, which maybe due to mispredicted genes and lowly expressed transcripts

- 23% of transcripts are not validated at all.



WormBase Transcripts Evaluation

None 5,235 (23%)

Complete 10,260 (46%)

Partial 6,980 (31%)

```
Caenorhabditis briggsae Gene model confirmation status (based on the EST/mRNA evidence)
------------------------------------------------------------------
Confirmed            10331 (47.3%) Every base of every exon has transcription evidence (mRNA/EST)
Partially_confirmed   7763 (35.5%) Some, but not all exon bases are covered by transcript evidence
Predicted             3769 (17.2%) No coverage by mRNA/EST evidence
```

50

# Method: transcript-to-transcript comparison (intron chain)



Assembled protein-coding transcripts

WormBase protein-coding transcripts

**Intron chain comparison**

A) Complete intron chain match?

B) Partial intron chain match?

C) No intron chain match?

Category:
1) Match

Category:
2) 3' extension
3) 5' extension
4) 5' & 3' extension
5) Internal within exon
6) Internal within intron
7) Alternative Donor
8) Alternative Acceptor
9) Alternative Donor & Acceptor
10) Merged
11) Complex changes
13) Other (Single-exon, partial)

Category:
12) Novel

Complete match (WB confirmed) 8,080
3' extension 316
5' extension 753
5' & 3' extension 26
Intron overlapping internal exon 358
Introns overlapping intron 217
Alt. donor (5'ss) 777
Alt. acceptor (3'ss) 882
Alt. donor & acceptor 346
Merging 2 or more genes 206
Complex changes 2,245
Novel 159
Other – Single-exon (no intron) 120
Other – Partial 10,220

■ Protein-coding transcripts

52

# RNA-Seq suggests 159 novel transcript/gene



- Transcript MERGE_00026402 (nCBG00109) is previously not annotated by WormBase

- Using RNA-Seq data, we found introns and assembled a transcript suggesting putative novel gene in this genomic region

- First and second hits from BLAST result shown that the protein sequence is the most similar to *C. elegans* proteins (right)

# 34% of candidate transcripts do not start with ATG

- A functional protein-coding transcript should contain an Open Reading Frame (ORF) that begins with start codon and ends with stop codon[1].

- Limitation of TransDecoder: does not have a start-codon finding function and will include transcript from the beginning if there is no upstream in-frame stop codon at the beginning of the transcript[2,3].

| Category | | Protein-coding transcripts | Starts with ATG | Does not start with ATG |
|---|---|---|---|---|
| 2. | 3' extension | 316 | 216 | 100 |
| 3. | 5' extension | 753 | 528 | 225 |
| 4. | 5' & 3' extension | 26 | 15 | 11 |
| 5. | Intron overlapping internal exon | 358 | 246 | 112 |
| 6. | Introns overlapping intron | 217 | 144 | 73 |
| 7. | Alternative donor (5'ss) | 777 | 515 | 262 |
| 8. | Alternative acceptor (3'ss) | 882 | 559 | 323 |
| 9. | Alternative donor & acceptor | 346 | 218 | 128 |
| 10. | Merged | 206 | 119 | 87 |
| 11. | Complex changes | 2,245 | 1496 | 749 |
| 12. | Novel | 159 | 78 | 81 |
| | Total | 6,285 | 4,134 (66%) | **2,151 (34%)** |

[1]Majoros et al., 2014, [2,3]Haas, 2014, 2018

- 16,880:
  - 1894: modified, same cel, revised cbriggsae transcript (extension etc)
  - 100: new, cel exist in original, contains 4 nCBG new cbriggsae transcripts
  - 32: new from nCBG (4 redundant with one point above)
  - 80: need further analysis, maybe new new (so total new would be plus nCBG)
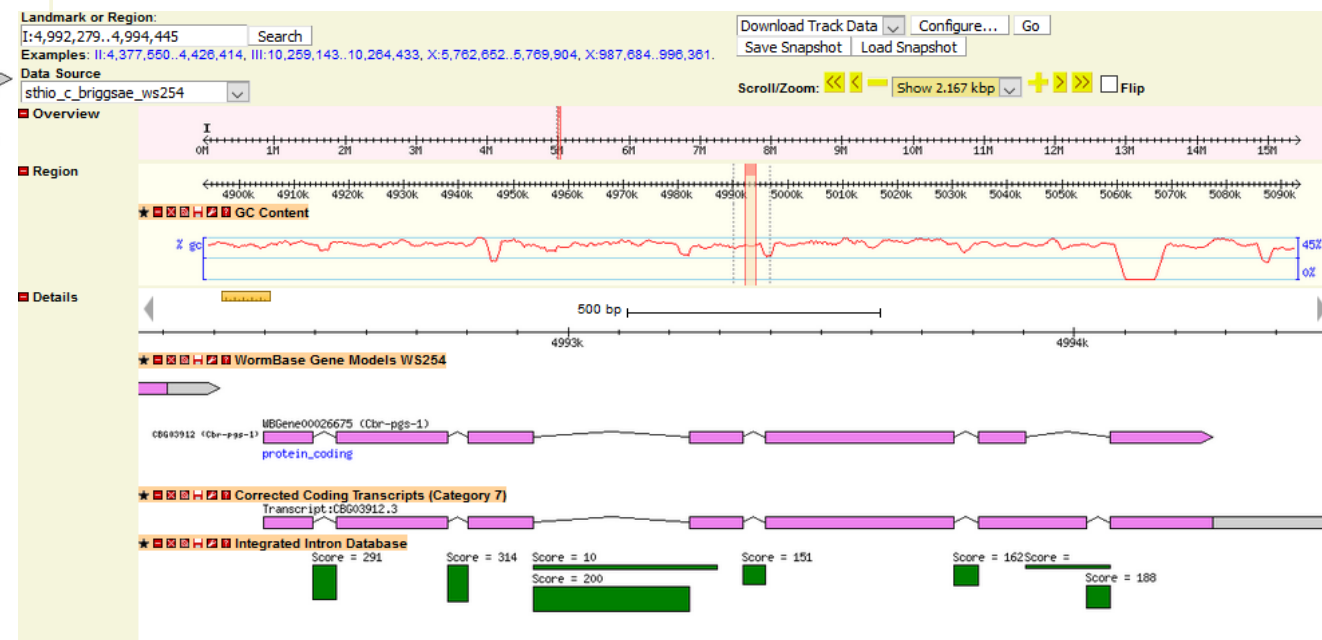
# Significance

- As more RNA-Seq data becomes available, this method can be used to further refine not only *C. briggsae* annotation but also *C. elegans* annotation.