# Transcriptome Assembly

## BIOINF 3005 /7160: Transcriptomics Applications

Dr Terry Bertozzi

South Australian Museum
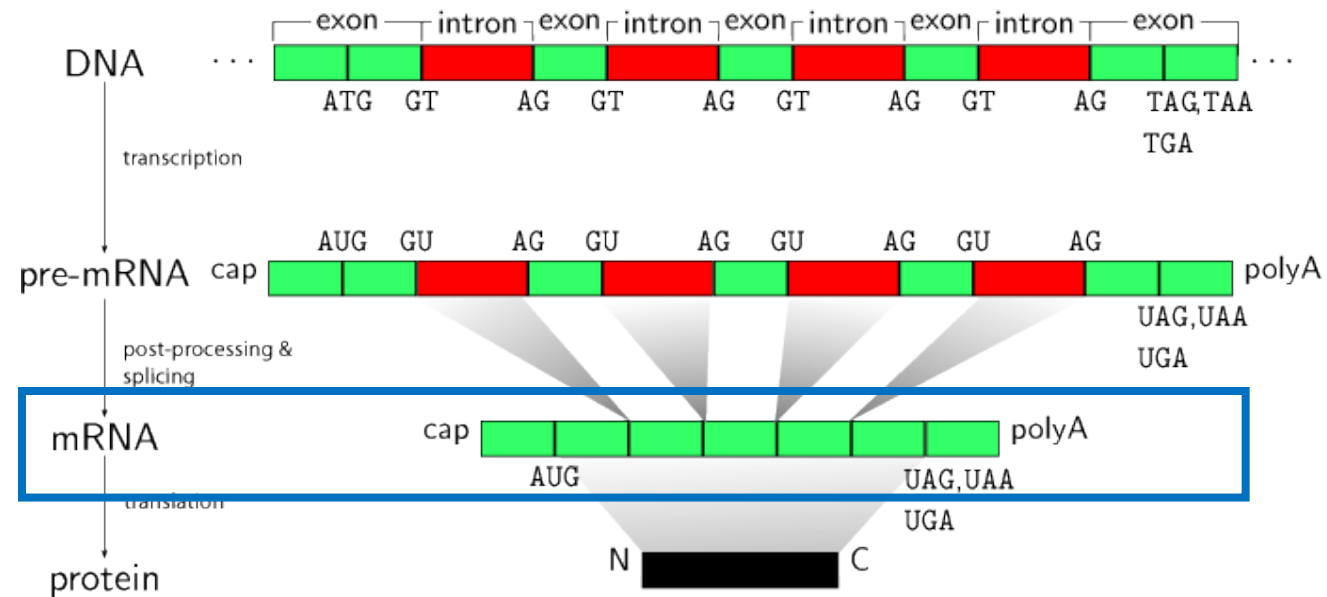
University of Adelaide

June 1, 2020

# Outline

- Transcriptome recap

- Genome vs Transcriptome assembly

- Transcriptome assembly
  - Short read assembly
  - Long read methods
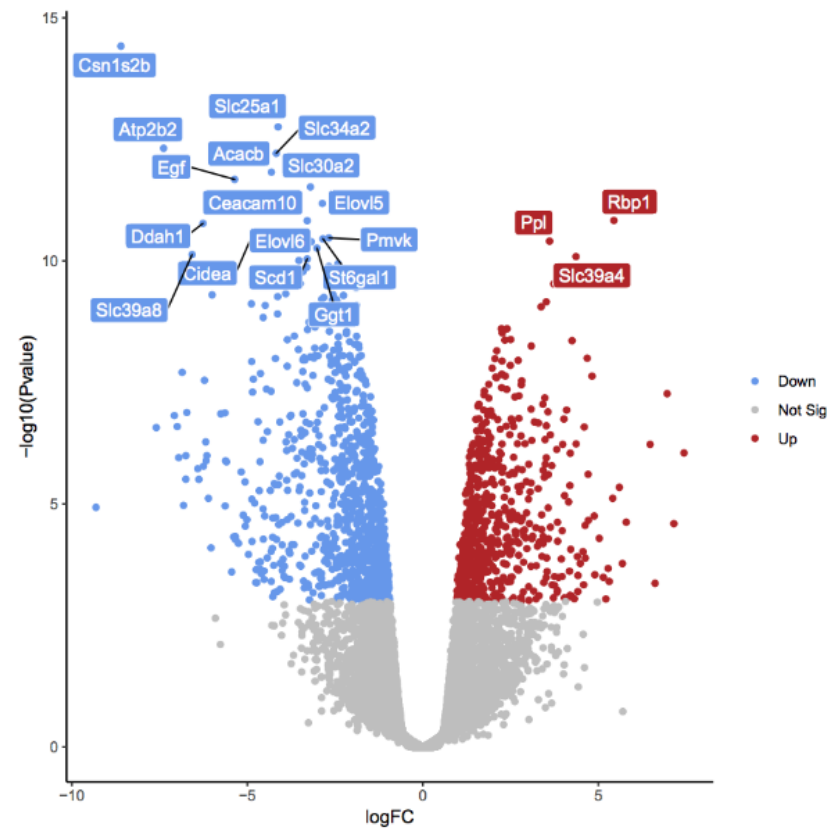  - Guided assembly

- Transcriptome evaluation

# Transcriptome

- **The set of all RNA transcripts, including coding and non-coding, in an individual or a population of cells.**

- mRNA

- lncRNA

- tRNA

- rRNA

- Small RNAs (e.g. miRNA, siRNA)

- Typically avoid rRNA
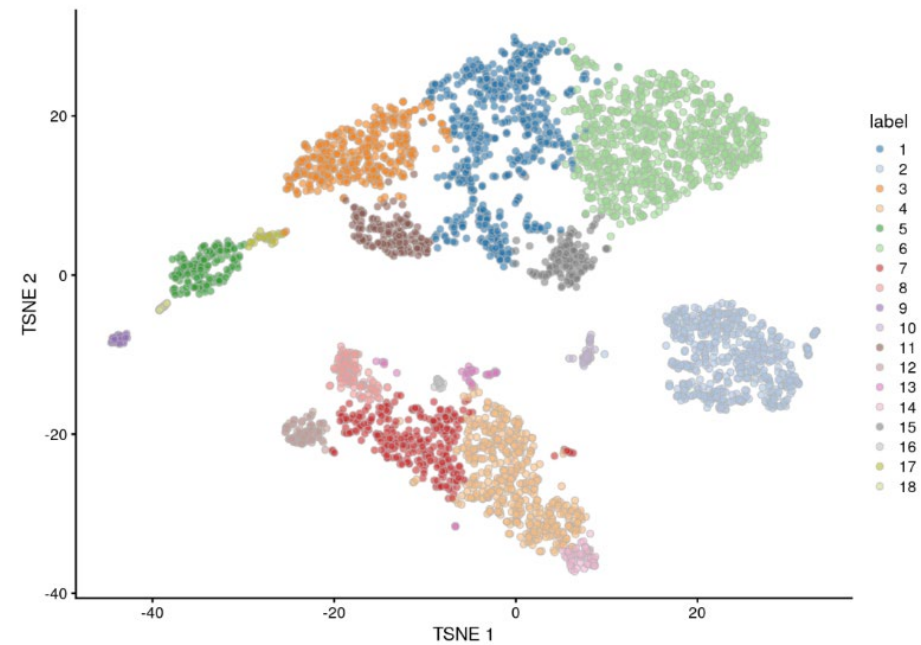  - >80% total RNA
  - ~5% mRNA

# Transcriptome

- Differential expression

- Single cell clustering



Image from https://galaxyproject.github.io
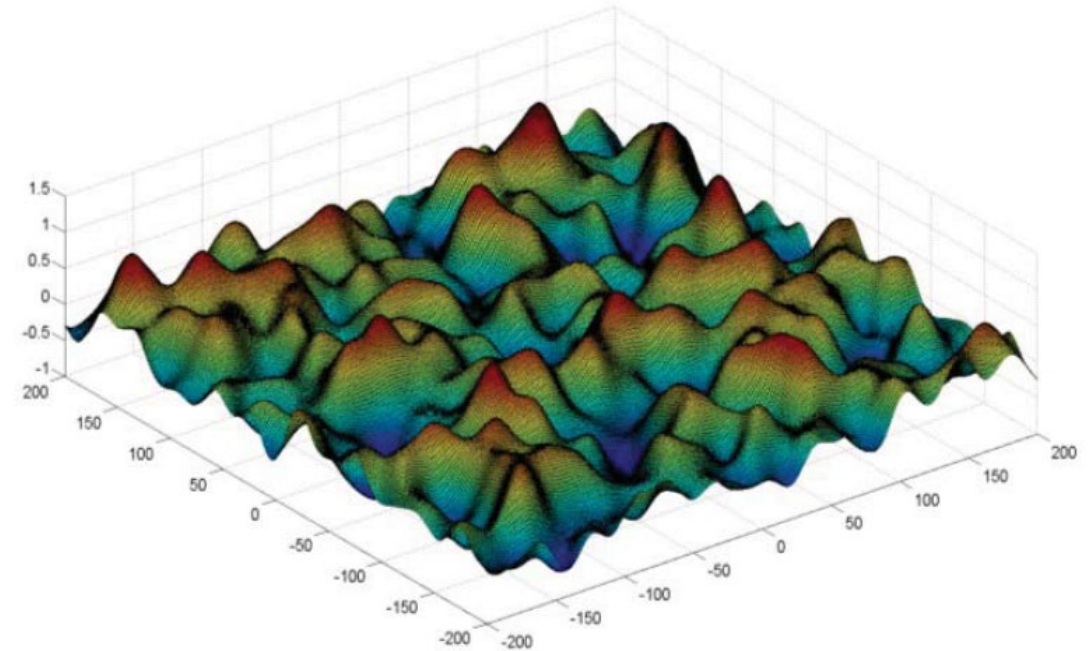
# Genome vs Transcriptome assembly
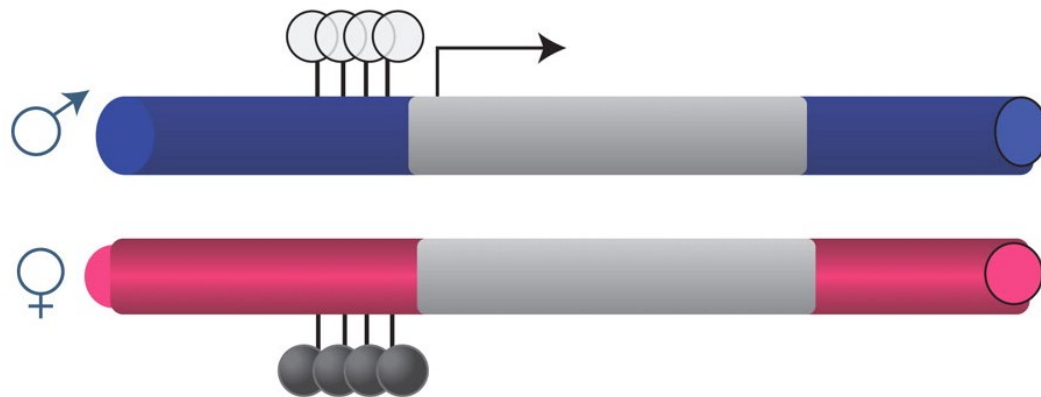
**Genome**
- Uniform coverage

**Transcriptome**
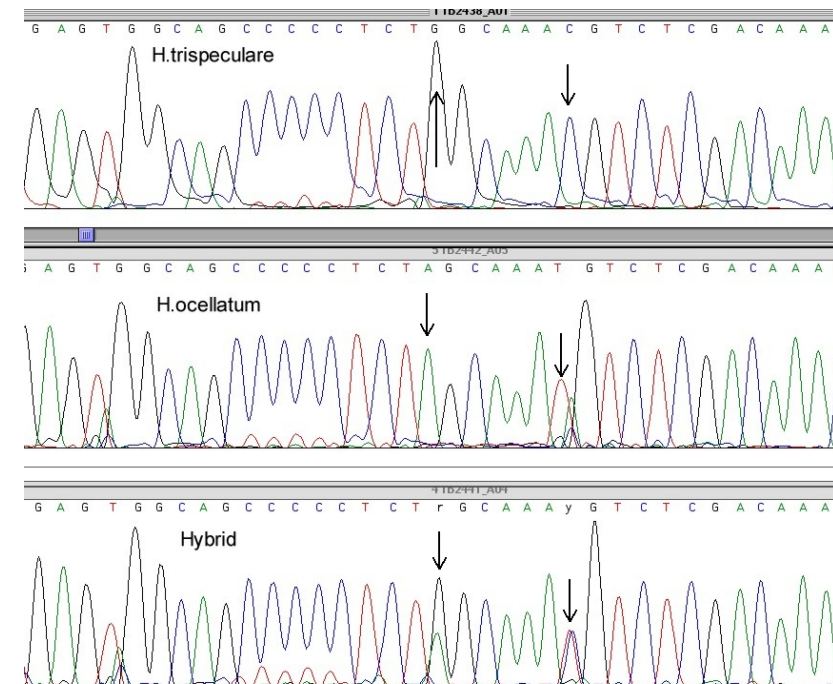- Exponentially distributed coverage

# Uneven coverage



- Cell/tissue compartmentalisation/function
  - Rare transcripts

- Allele specific expression
  - X-inactivation (Xist: X-inactive specific transcript)
  - Full / partial imprinting





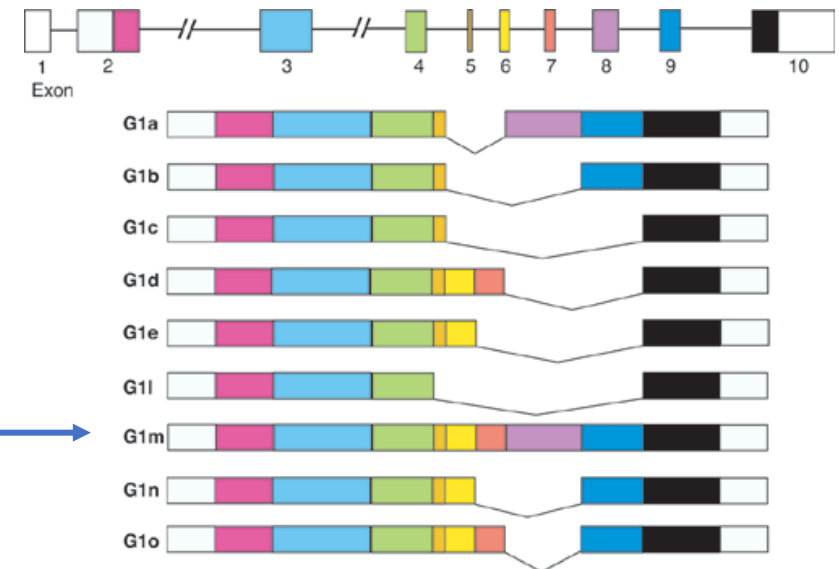Bartolomei 2011 Cold Spring Harbor Perspectives in Biology

# Genome vs Transcriptome assembly

**Genome**

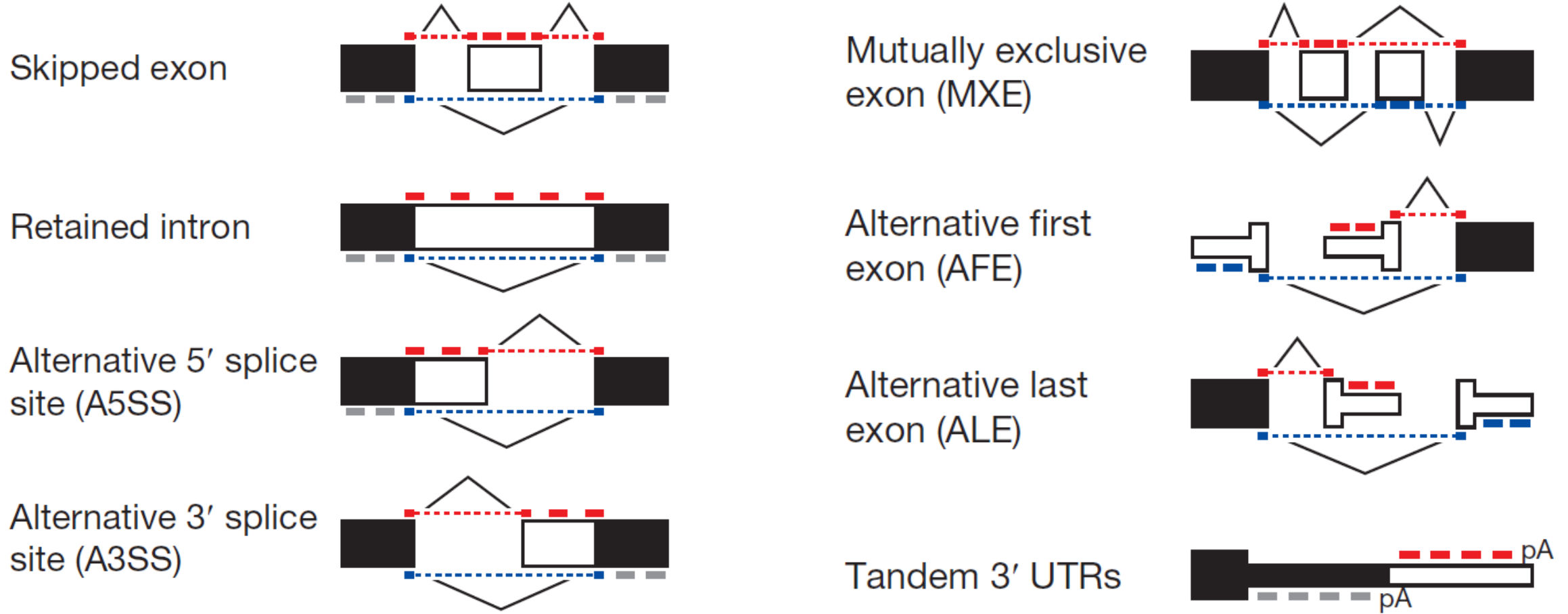- Uniform coverage

- Two contigs / locus

**Transcriptome**

- Exponentially distributed coverage

- Multiple contigs / locus
  - isoforms



Canonical transcript →

# Alternate Splicing
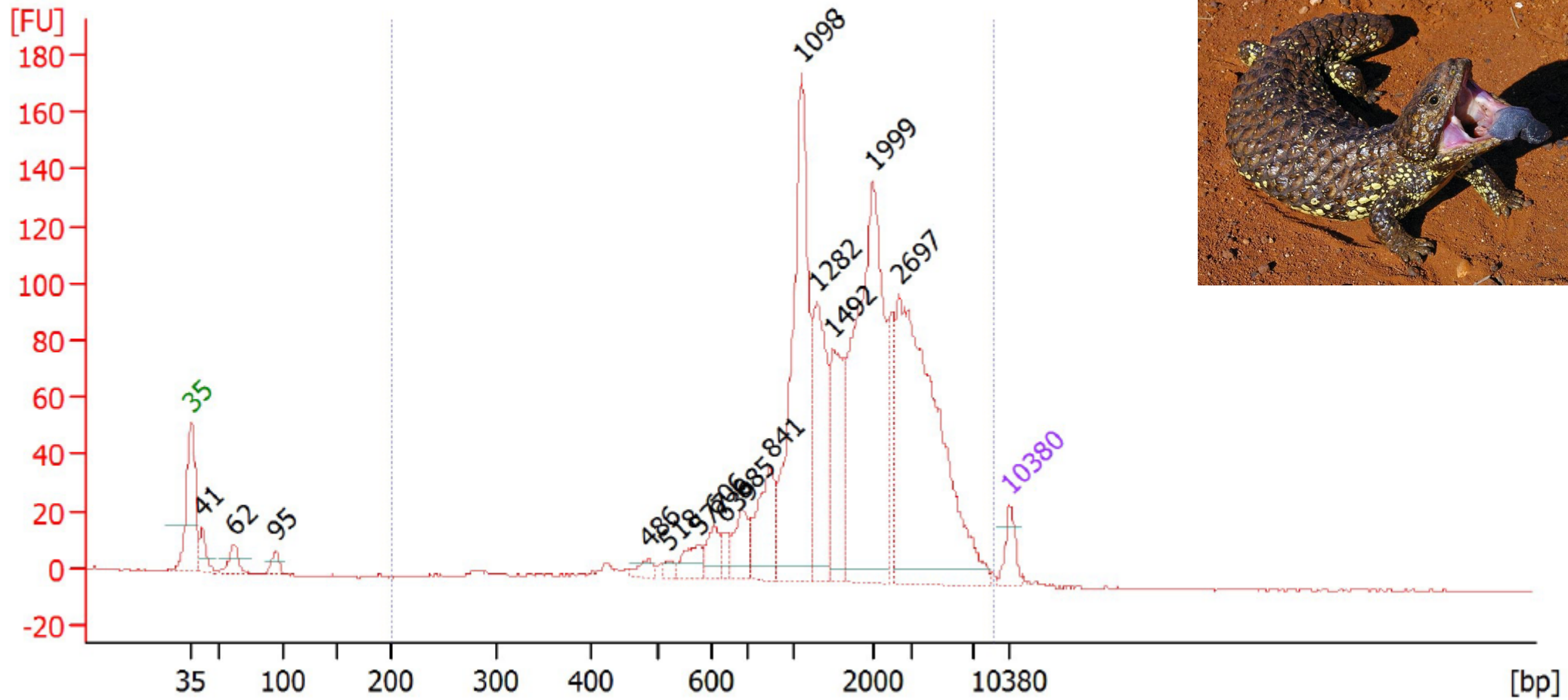
# Genome vs Transcriptome assembly

**Genome**

- Uniform coverage

- Single contig / locus

- Assemble small number of large Mb-length contigs

**Transcriptome**

- Exponentially distributed coverage

- Multiple contigs / locus

- Assemble thousands of Kb-length transcripts

# Transcript distribution
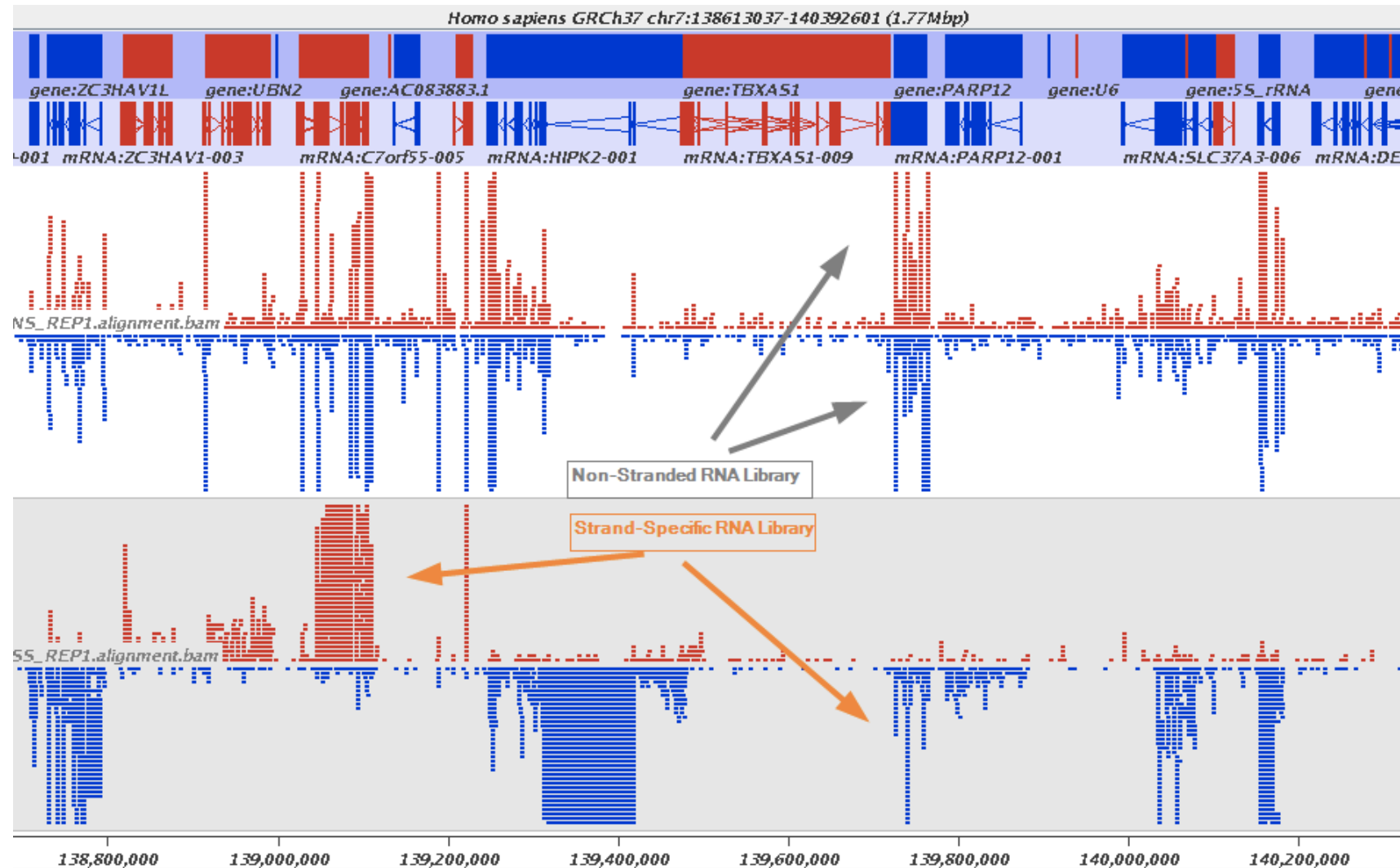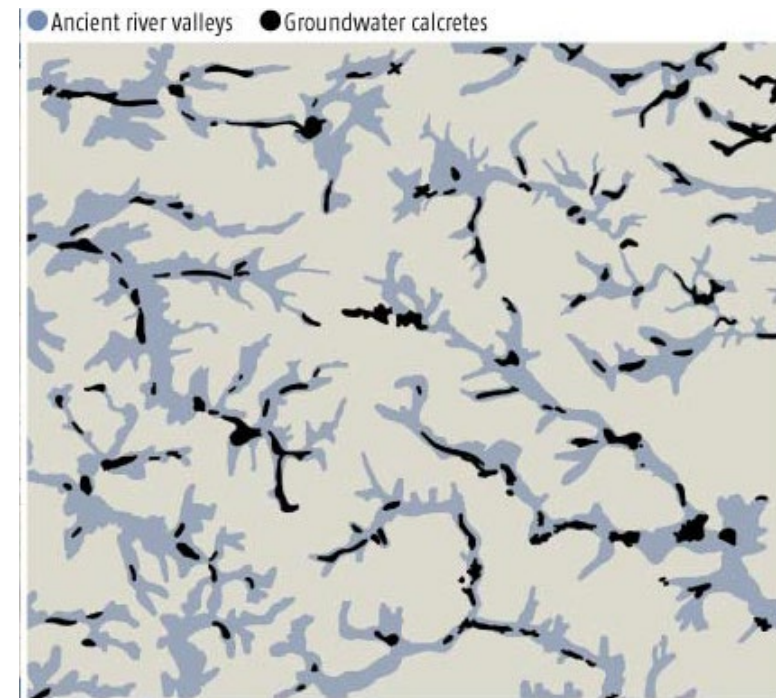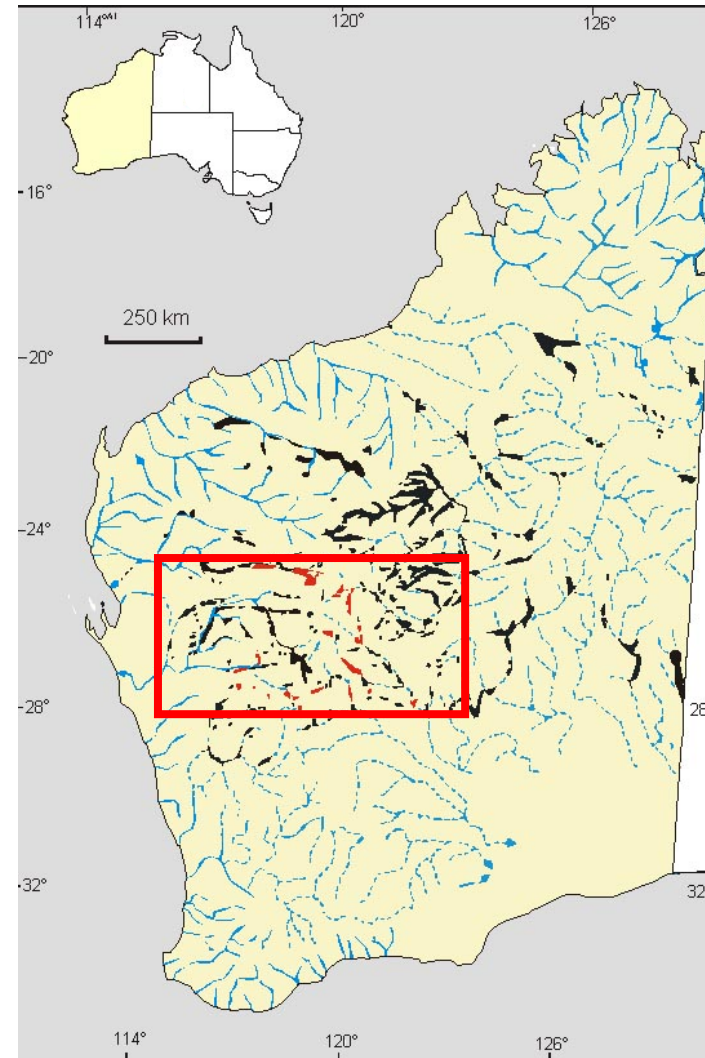
# Genome vs Transcriptome assembly

**Genome**

- Uniform coverage

- Single contig / locus

- Assemble small number of large Mb-length contigs

- Double stranded

**Transcriptome**

- Exponentially distributed coverage

- Multiple contigs / locus

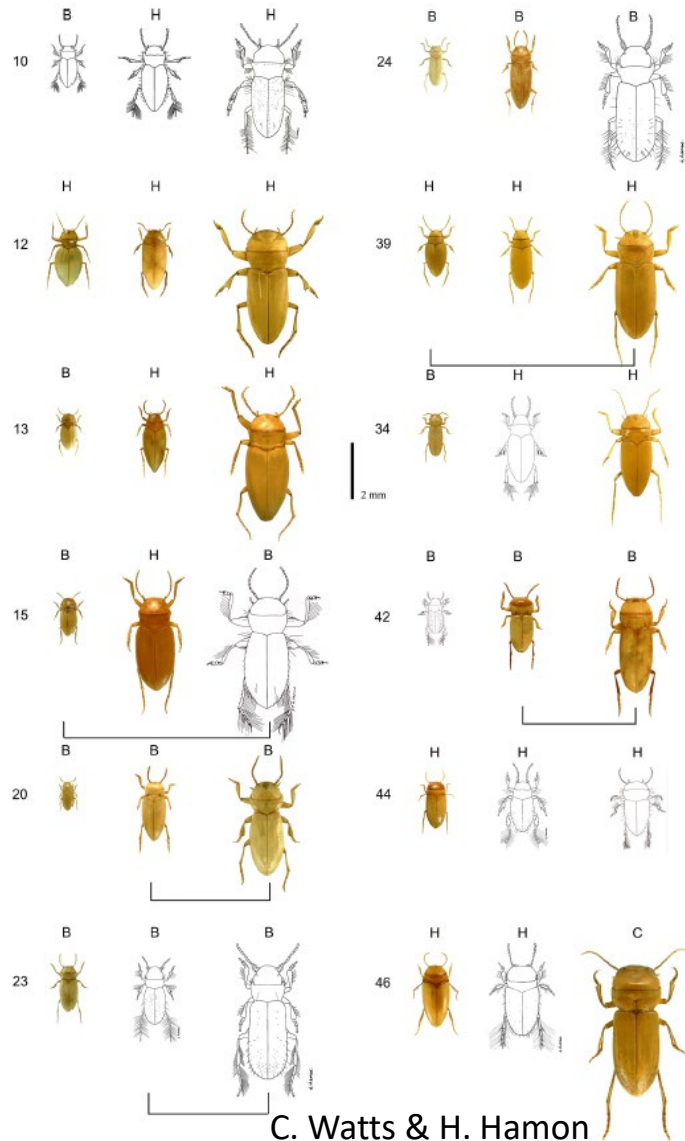- Assemble thousands of Kb-length transcripts

- Strand specific

# Strand specific expression

# Assembly tools

# Phototransduction in blind beetles



C. Watts & H. Hamon

Ancient river valleys ● Groundwater calcretes

# Phototransduction in blind beetles



Leijs *et al.* 2012 PLoS One

Surface

Subterranean

# Phototransduction in blind beetles

# Short Read Assembly

- Trinity Assembler & post assembly tools

- External required software (trinity):
  - Jellyfish (k-mer counting), samtools, Bowtie2, kallisto, salmon

- External required software (post-assembly):
  - R with Bioconductor
  - BLAST, Picard, GATK4, Hisat2, STAR
  - RSEM, express
  - Transdecoder

# Trinity

sequence

**ATGGAAGTCGCGGAATC**

7mers

ATGGAAG
TGGAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

# Trinity

- Output is a bunch of files and folders

- `trinity.fasta` contains the assembled transcripts

Transcript ID

```
>TRINITY_DN1000|c115_g5_i1 len=247 path=[31015:0-148 23018:149-246]
AATCTTTTTTGGTATTGGCAGTACTGTGCTCTGGGTAGTGATTAGGGCAAAAGAAGACAC
ACAATAAAGAACCAGGTGTTAGACGTCAGCAAGTCAAGGCCTTGGTTCTCAGCAGACAGA
AGACAGCCCTTCTCAATCCTCATCCCTTCCCTGAACAGACATGTCTTCTGCAAGCTTCTC
CAAGTCAGTTGTTCACAGGAACATCATCAGAATAAATTTGAAATTATGATTAGTATCTGA
TAAAGC
```

# UV Opsin

# Incorrectly assembled transcripts

- Multicopy gene families



Doxiadis et al 2006 PNAS

- Coverage
  - Incomplete assembly
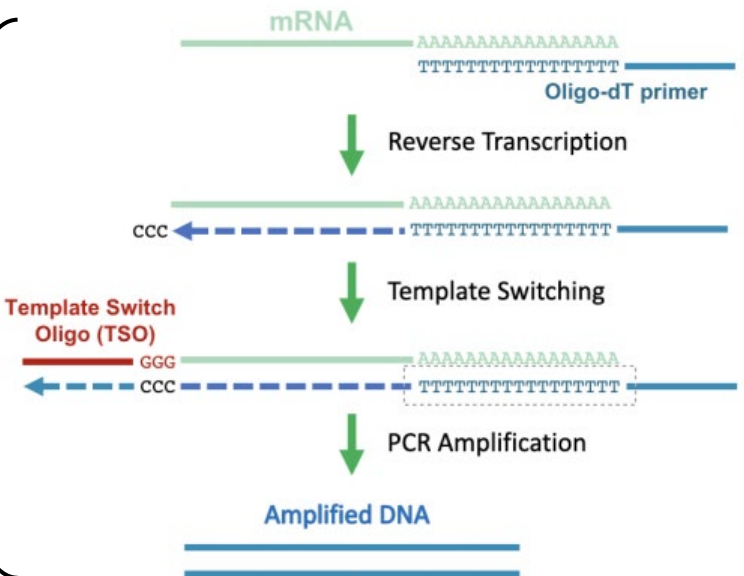  - Extreme expression

- Alternative splicing

# Long read methods

- Alternative splicing
  - 90% genes[1]; 30% ncRNA[2]
- Long read sequencing
  - No assembly required!
- Nanopore RNA sequencing
  - ✓ Native RNA sequencing
  - X Error prone reads
- Iso-Seq (Pacific Biosciences)
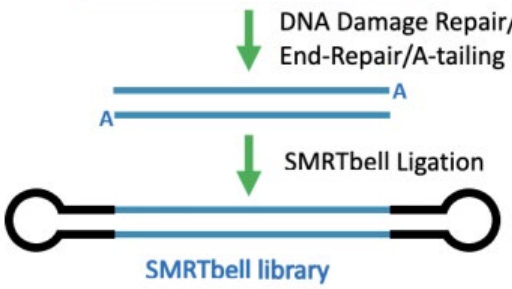  - X Native RNA sequencing
  - ✓ Highly accurate reads

[1]Wang et al 2008 Nature; [2]Cabili et al 2011 Genes & Development

# Iso-Seq

# Heat Shock Proteins (HSP70)

# Genome Guided Assembly

- Fusion transcripts



Neckles et al 2109 WIREs RNA

- RNA editing
  - A → I; C → U

DNA:    CACTGGACG
mRNA:  GIGACCUGC
Protein:   G     T     C

# Genome Guided Assembly

- Gene models

# StringTie

# StringTie



*Splice graph with heaviest path highlighted*

**Step 3**: build alternative splice graph

**Step 4**: construct flow network for path in splice graph with heaviest coverage

**Step 5**: assemble transcripts and update coverage

isoform 1

isoform 2

isoform 3

Pertea et al 2015 Nature Biotechnology

# StringTie

- General Transfer Format (GTF)

```
seqname source      feature     start   end    score  strand  frame  attributes
chrX    StringTie   transcript  281394  303355 1000    +       .       gene_id "ERR188044.1"; transcript_id "ERR188044.1.1"; reference_id "NM
chrX    StringTie   exon        281394  281684 1000    +       .       gene_id "ERR188044.1"; transcript_id "ERR188044.1.1"; exon_number "1";
...
```

https://asia.ensembl.org/info/website/upload/gff.html

- Can use current annotation (GTF)

- Also calculates coverage for expression analysis
  - Ballgown, DESeq2, edgeR

# Transcriptome Evaluation

- Benchmarking of Universal Single Copy Orthologs (BUSCO)



C: 85% [S:34%,D:51%],
F: 14%,
M: 1%,
n: 1658

C: Complete
- S: single copy;
- D: duplicated;
F: Fragmented;
M: Missing;
n: Total groups;

# Transcriptome Evaluation

- BLAST to proteome

- Quantify read support

```
76201190 reads; of these:
  76201190 (100.00%) were paired; of these:
    18166307 (23.84%) aligned concordantly 0 times
    17026716 (22.34%) aligned concordantly exactly 1 time
    41008167 (53.82%) aligned concordantly >1 times
    ----
    18166307 pairs aligned concordantly 0 times; of these:
      1769907 (9.74%) aligned discordantly 1 time
    ----
    16396400 pairs aligned 0 times concordantly or discordantly; of these:
      32792800 mates make up the pairs; of these:
        15287552 (46.62%) aligned 0 times
        3874965 (11.82%) aligned exactly 1 time
        13630283 (41.56%) aligned >1 times
89.97% overall alignment rate
```
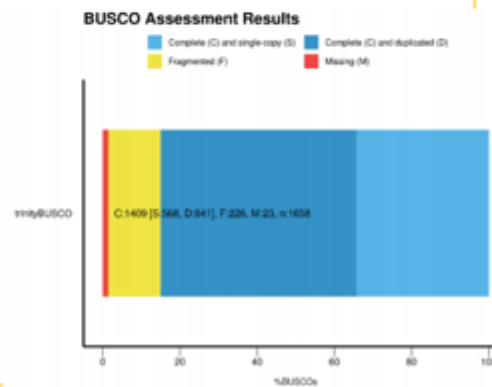
Properly paired, aligning as expected

Properly paired, wrong orientation or distance

Unmapped or SE mapping reads

# Summary

- Transcriptome assembly is hard!
  - Diverse population of RNA
  - Non-uniform coverage
  - Tissue specificity
  - Multigene families
  - Alternative splicing