# BULK RNA-SEQ: UPSTREAM ANALYSIS

**Presenter: Duy Dao**

-03/15/2025-

# TABLE OF CONTENTS

**RNA SEQ: UPSTREAM ANALYSIS**

**RNASEQ: UPSTREAM ANALYSIS WORKFLOW**

# RAW DATA PROCESSING

## SEQUENCE QUALITY CONTROL (FASTQC)

> **FASTQC Summary**



"FASTQC is a useful tool to check sequences quality."



**Basic Statistics**

| Measure | Value |
|---------|-------|
| Filename | NIST7035_TAAGGCGA_L001_R1_001.fastq.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 20203002 |
| Total Bases | 2 Gbp |
| Sequences flagged as poor quality | 0 |
| Sequence length | 101 |
| %GC | 49 |

## READ TRIMMING & FILTERING

usadellab/
**Trimmomatic**



👥 2 Contributors    ⊙ 25 Issues    ☆ 131 Stars    ⑂ 56 Forks

This program does adaptive quality trimming, head and tail crop, and adaptor removal.

*Check QC → Trim → Check QC again.*



Trimming:
- Quality trimming
- Adapter trimming.

5

https://home.cc.umanitoba.ca/~psgendb/doc/TrimmomaticManual.pdf

**Quality problems**

- Quality problems typically originate either in the sequencing itself or in the preceding library preparation.

- They include low-confidence bases, sequence-specific bias, 3′/5′ positional bias, polymerase chain reaction (PCR) artifacts, untrimmed adapters, and sequence contamination.

- These problems can seriously affect mapping to reference, assembly, and expression estimates, but luckily many of them can be corrected for by filtering, trimming, error correction, or bias correction.

- Some problems cannot be corrected for, but you should at least be aware of them when interpreting results.

**ALIGNMENT**

*Workflow for a RNA-seq analysis.*

https://hbctraining.github.io/Intro-to-rnaseq-hpc-O2/lessons/03_alignment.html

**Basic alignment (Contiguous Alignment / Non-spliced Alignment)**

In contiguous alignment, sequences are aligned continuously without any gaps or interruptions.

**Problem when using basic alignment to map RNA-seq data**



- Unmapped reads due to intron splicing.

https://www.rna-seqblog.com/introduction-and-application-of-transcriptome-sequencing/

**Contiguos Alignment vs Splice-Aware Alignment**



- Contiguous Aligners: BWA, Bowtie2,...
- Spliced Aligners: HiSAT2, TopHat, STAR,...

# RNA-SEQ: ALIGNMENT / MAPPING

**Compare the Alignment of DNA-seq and RNA-seq**

**Compare the Alignment of DNA-seq and RNA-seq**

**Purpose**

**DNA-seq**
- Variant calling
- Genome Coparative Analysis
- …

**RNA-seq**
- Gene expression analysis
- Identifying Differentially expressed genes
- Studying Alternative splicing events
- Characterizing transcript isoforms.
- …

Spliced-aware alignment algorithms employ various strategies to handle splice junctions, such as:

- **Split reads:** Allows for precise alignment across the splice junctions.

- **Novel splice junction detection:** Detect previously unknown splicing events, providing insights into alternative splicing patterns and transcriptome complexity.

- **Splice junction annotation:** Aligners may utilize existing splice junction annotations, such as those obtained from databases or previous studies, to guide the alignment process.

**STAR** (Spliced Transcripts Alignment to a Reference)

**STAR alignment strategy**



Seed searching…

Clustering, stitching, and scoring

Input → Output

Input
- *R1.fastq*
- *R2.fastq*
- *genome.fa*
- *genome.gtf*

Output
- *SAM/BAM*
- *SJ file*
- *Log file*

**Alignment statistics & Utilities for manipulating alignment files.**

```
                        Started job on |    Jun 18 14:12:08
                    Started mapping on |    Jun 18 14:12:13
                           Finished on |    Jun 18 14:12:43
  Mapping speed, Million of reads per hour |    666.81

                  Number of input reads |    5556750
              Average input read length |    124
                          UNIQUE READS:
           Uniquely mapped reads number |    4987609
                Uniquely mapped reads % |    89.76%
                  Average mapped length |    124.32
              Number of splices: Total |    244266
       Number of splices: Annotated (sjdb) |    236101
             Number of splices: GT/AG |    243487
             Number of splices: GC/AG |    63
             Number of splices: AT/AC |    11
       Number of splices: Non-canonical |    705
             Mismatch rate per base, % |    0.08%
                  Deletion rate per base |    0.01%
                 Deletion average length |    1.35
                 Insertion rate per base |    0.00%
                Insertion average length |    1.07
                     MULTI-MAPPING READS:
      Number of reads mapped to multiple loci |    270940
           % of reads mapped to multiple loci |    4.88%
      Number of reads mapped to too many loci |    30963
           % of reads mapped to too many loci |    0.56%
                          UNMAPPED READS:
  Number of reads unmapped: too many mismatches |    0
     % of reads unmapped: too many mismatches |    0.00%
            Number of reads unmapped: too short |    266450
               % of reads unmapped: too short |    4.80%
               Number of reads unmapped: other |    788
                  % of reads unmapped: other |    0.01%
                        CHIMERIC READS:
```
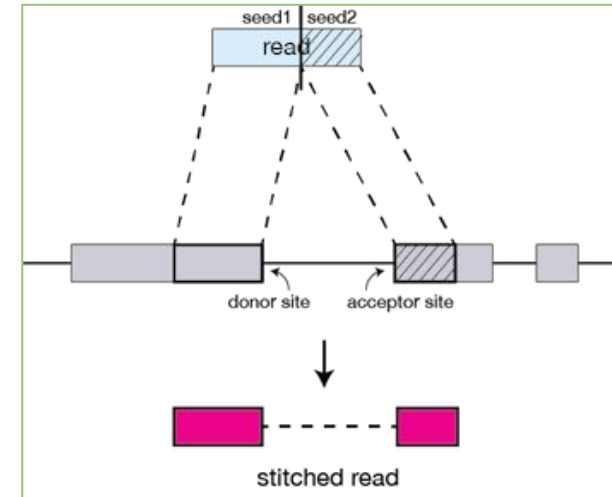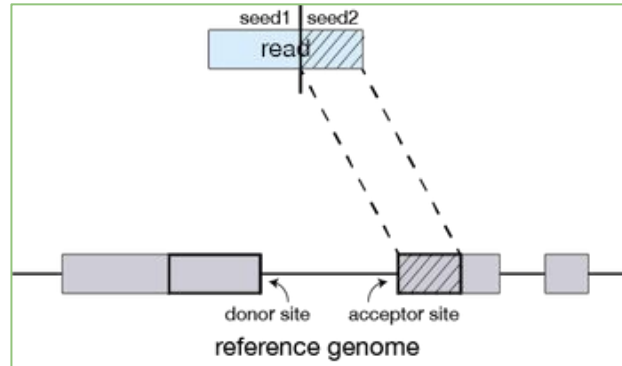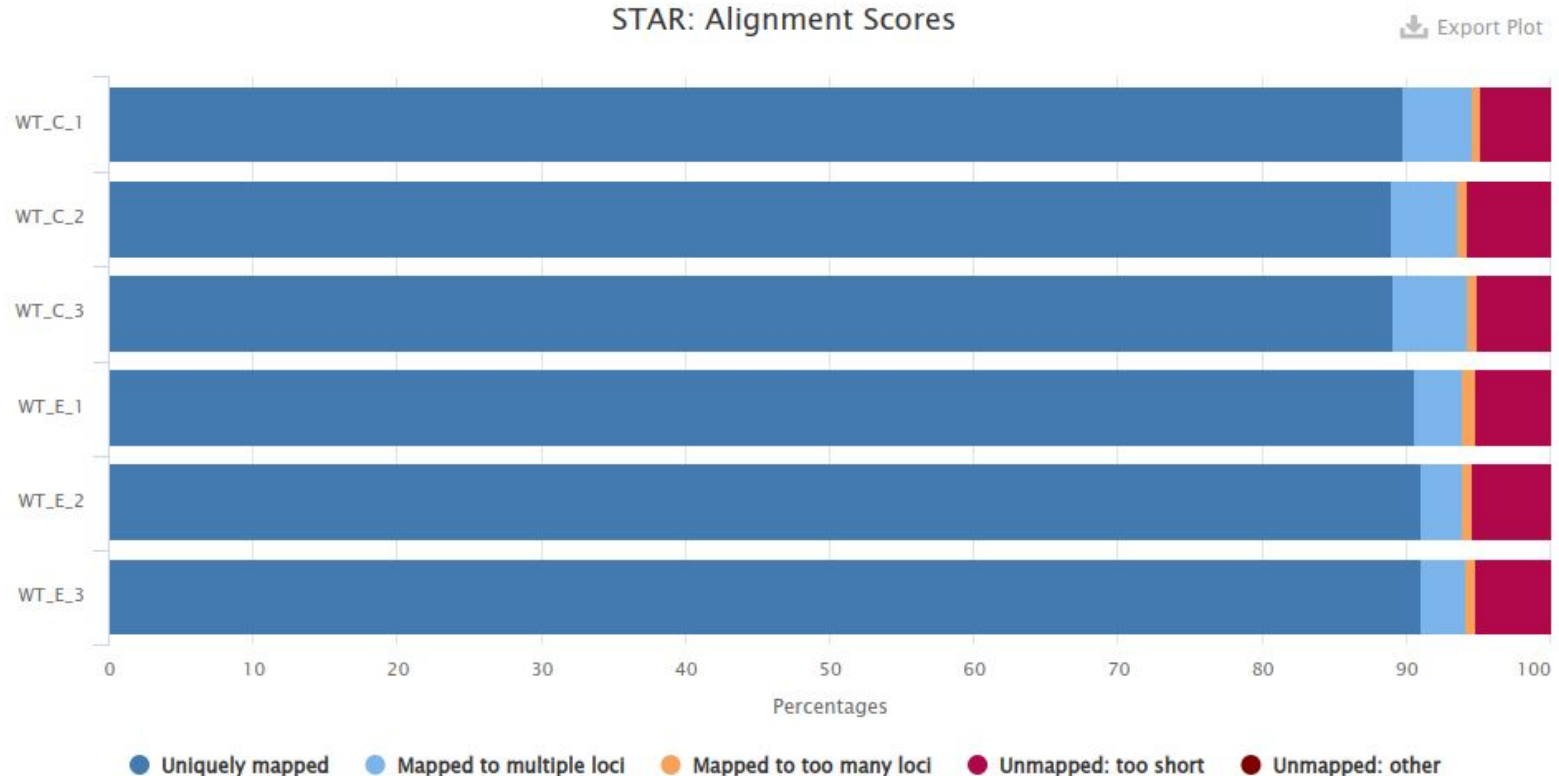
Alignment statistics (log file)

**SJ.out.tab**

| chrI | 12728 | 12823 | 1 | 1 | 0 | 0 | 1 | 13 |
|------|-------|-------|---|---|---|------|----|----|
| chrI | 87388 | 87500 | 1 | 1 | 1 | 70 | 0 | 32 |
| chrI | 128525 | 129021 | 2 | 2 | 0 | 1 | 0 | 26 |
| chrI | 142254 | 142619 | 1 | 1 | 1 | 1820 | 0 | 32 |
| chrI | 142254 | 143349 | 1 | 1 | 0 | 1 | 0 | 22 |
| chrI | 151007 | 151096 | 2 | 2 | 1 | 4 | 0 | 32 |
| chrI | 206383 | 206517 | 1 | 1 | 0 | 2 | 1 | 25 |
| chrII | 5120 | 5335 | 2 | 2 | 0 | 0 | 1 | 30 |
| chrII | 45645 | 45977 | 1 | 1 | 0 | 1087 | 2 | 31 |
| chrII | 47059 | 47146 | 2 | 2 | 1 | 11 | 0 | 23 |
| chrII | 60194 | 60697 | 2 | 2 | 1 | 2960 | 0 | 32 |
| chrII | 89133 | 89440 | 2 | 2 | 0 | 81 | 0 | 31 |
| chrII | 110421 | 110505 | 2 | 2 | 1 | 88 | 0 | 32 |
| chrII | 110880 | 110948 | 1 | 1 | 1 | 23 | 0 | 32 |
| chrII | 125155 | 125270 | 1 | 1 | 1 | 67 | 0 | 32 |
| chrII | 142750 | 142846 | 2 | 2 | 1 | 181 | 0 | 32 |
| chrII | 142754 | 142846 | 2 | 2 | 0 | 1 | 0 | 26 |
| chrII | 167650 | 230011 | 2 | 2 | 0 | 0 | 17 | 12 |
| chrII | 168425 | 168808 | 1 | 1 | 1 | 308 | 0 | 31 |
| chrII | 170621 | 170804 | 1 | 1 | 0 | 6 | 0 | 27 |
| chrII | 170677 | 170804 | 1 | 1 | 1 | 82 | 0 | 32 |

The SJ.out.tab contains filtered splice junctions detected in the mapping

**Alignment statistics & Utilities for manipulating alignment files.**



STAR: Alignment Scores

**Alignment statistics & Utilities for manipulating alignment files.**

**SAM format**

| QNAME | FLAG | RNAME | POS | MAPQ | CIGAR | RNEXT | PNEXT | TLEN | SEQ |
|---|---|---|---|---|---|---|---|---|---|

```
HWI-D00119:50:H7AP8ADXX:1:1101:1322:2086   113   chr17   74243430    60   86M   =   74243430    0     GCCTGGGACTGGCCAAGCGCTCCAAGTGTCTCCTGGCCTCGTGAGTATCCCTACCCGTGGACCTGGGACAAAGGG
HWI-D00119:50:H7AP8ADXX:1:1101:1322:2086   177   chr17   74243430    60   86M   =   74243430    0     GCCTGGGACTGGCCAAGCGCTCCAAGTGTCTCCTGGCCTCGTGAGTATCCCTACCCGTGGACCTGGGACAAAGGG
HWI-D00119:50:H7AP8ADXX:1:1101:1294:2087   65    chr1    49798791    60   73M   =   49798791    0     CCAAACTCTGGAAGTATTTTATTTTATTTAAAATATTTTCAGCATCTTTGTAAATTAAGTGTTAATCACCATG
HWI-D00119:50:H7AP8ADXX:1:1101:1294:2087   129   chr1    49798791    60   73M   =   49798791    0     CCAAACTCTGGAAGTATTTTATTTTATTTAAAATATTTTCAGCATCTTTGTAAATTAAGTGTTAATCACCATG
HWI-D00119:50:H7AP8ADXX:1:1101:1406:2159   65    chr18   9733113 60   86M   =   9733113 0       CCATGCTATAAATCTTTATTTTTTCCACTTTGGTTTGAGTGGGATCACTTCCATTGGAGCTAATCCAGAAATTGCTAATTCAGTGG   HGI
HWI-D00119:50:H7AP8ADXX:1:1101:1406:2159   129   chr18   9733113 60   86M   =   9733113 0       CCATGCTATAAATCTTTATTTTTTCCACTTTGGTTTGAGTGGGATCACTTCCATTGGAGCTAATCCAGAAATTGCTAATTCAGTGG   HGI
HWI-D00119:50:H7AP8ADXX:1:1101:1303:2183   65    chr6    132570953   60   86M   =   132570953   0     TTATTATGATAATTCTGTATGGTAACATATTTCTTGTGGCTAGACGACAGGCGAAAAAGATAGAAAATACTGGTA
HWI-D00119:50:H7AP8ADXX:1:1101:1303:2183   129   chr6    132570953   60   86M   =   132570953   0     TTATTATGATAATTCTGTATGGTAACATATTTCTTGTGGCTAGACGACAGGCGAAAAAGATAGAAAATACTGGTA
HWI-D00119:50:H7AP8ADXX:1:1101:1413:2244   81    chr19   56204243    0    86M   =   56195531    -8798 TTGCCGGCTCATGCAGTTGTTGGACTTTCAATTTTCTGAGGAGCTTCTATCTGTCTCCCATATCGGCTGCTCTA
HWI-D00119:50:H7AP8ADXX:1:1101:1413:2244   161   chr19   56195531    0    86M   =   56204243    8798  GGGATGTAAATTAGAGCAGCCGATATGGGAGACAGTATAGAAGGTCCTCAGAAAATTGAAAGTCCAACAACTGCA
HWI-D00119:50:H7AP8ADXX:1:1101:1372:2246   113   chr19   50402338    60   86M   =   50402338    0     CGCCGCCCTACGAGGCCAACGTCGACTTTGAGATCCGGTACGGCCTCTGCCTCACTTCTCCGGCCTCTATCCCCAC
HWI-D00119:50:H7AP8ADXX:1:1101:1372:2246   177   chr19   50402338    60   86M   =   50402338    0     CGCCGCCCTACGAGGCCAACGTCGACTTTGAGATCCGGTACGGCCTCTGCCTCACTTCTCCGGCCTCTATCCCCAC
HWI-D00119:50:H7AP8ADXX:1:1101:1613:2099   113   chr1    172749268   27   86M   =   172749268   0     TGGTATTCTCGCCTCACTGGTAACTCAACCTGCGGATGTTATCAAAACTCATATGCAGCTTTACCCACTGAAGTT
HWI-D00119:50:H7AP8ADXX:1:1101:1613:2099   177   chr1    172749268   27   86M   =   172749268   0     TGGTATTCTCGCCTCACTGGTAACTCAACCTGCGGATGTTATCAAAACTCATATGCAGCTTTACCCACTGAAGTT
HWI-D00119:50:H7AP8ADXX:1:1101:1644:2101   113   chr1    76700193    46   86M   =   76700193    0     GGAAGCTGGCATCACTGAGAAGGTTGTTTTCGAGCAGACAAAGGTCATCGCAGATAACGTGAAGGACTGGGAGCAA
HWI-D00119:50:H7AP8ADXX:1:1101:1644:2101   177   chr1    76700193    46   86M   =   76700193    0     GGAAGC
HWI-D00119:50:H7AP8ADXX:1:1101:1657:2113   65    chr8    85661599    0    86M   =   85654762    -6838 GCTGT
HWI-D00119:50:H7AP8ADXX:1:1101:1657:2113   129   chr8    85654762    0    86M   =   85661599    6838  GCTGT
HWI-D00119:50:H7AP8ADXX:1:1101:1777:2104   113   chr19   18784989    60   86M   =   18784989    0     GGAGC
HWI-D00119:50:H7AP8ADXX:1:1101:1777:2104   177   chr19   18784989    60   86M   =   18784989    0     GGAGC
HWI-D00119:50:H7AP8ADXX:1:1101:1830:2136   65    chr1    227731993   60   86M   =   227731993   0     CAGGG
HWI-D00119:50:H7AP8ADXX:1:1101:1830:2136   129   chr1    227731993   60   86M   =   227731993   0     CAGGG
HWI-D00119:50:H7AP8ADXX:1:1101:2021:2071   65    chr17   18541691    0    86M   =   18541691    0     ACCAA
HWI-D00119:50:H7AP8ADXX:1:1101:2021:2071   129   chr17   18541691    0    86M   =   18541691    0     ACCAA
HWI-D00119:50:H7AP8ADXX:1:1101:2027:2162   65    chr2    39673045    60   86M   =   39673045    0     GTGCT
HWI-D00119:50:H7AP8ADXX:1:1101:2027:2162   129   chr2    39673045    60   86M   =   39673045    0     GTGCT
HWI-D00119:50:H7AP8ADXX:1:1101:2340:2087   65    chr2    232386840   60   86M   =   232386840   0     GCGCC
HWI-D00119:50:H7AP8ADXX:1:1101:2340:2087   129   chr2    232386840   60   86M   =   232386840   0     GCGCC
                                                  chr4    9326689 0    86M   =   9331435 4747  CGAGGGAGGGCCAGGAGATCC
                                                  chr4    9331435 0    86M   =   9326689 -4747 CGAGGGAGGGCCAGGAGATCC
```

**Alignment section**

SAM format: 11 mandatory fields

**Table 9.1.** Mandatory fields of the SAM Format.

| Col | Field | Description | Example |
|---|---|---|---|
| 1 | QNAME | Query template NAME | read_1 |
| 2 | FLAG | Bitwise FLAG | 0 |
| 3 | RNAME | Reference sequence NAME | chrE |
| 4 | POS | Left-most mapping POSition (1-based) | 11 |
| 5 | MAPQ | MAPping Quality | 37 |
| 6 | CIGAR | CIGAR string | 10M |
| 7 | RNEXT | Ref. name of the mate or NEXT read | * |
| 8 | PNEXT | Position of the mate or NEXT read | 0 |
| 9 | TLEN | Observed Template LENgth | 0 |
| 10 | SEQ | Segment SEQuence | ACGCATACTG |
| 11 | QUAL | Base QUALity string | DIGAFHHBCA |

*Note:* Each line in the alignment section of a SAM file comprises 11 mandatory fields.
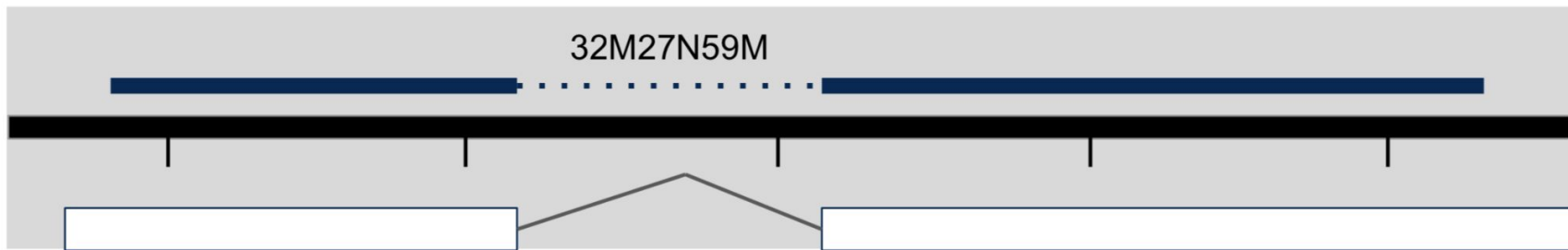
**Alignment statistics & Utilities for manipulating alignment files.**

**CIGAR string with "N"**

The "N" in the CIGAR string represents a stretch of skipped reference bases (also known as introns or gaps) in a sequence alignment.
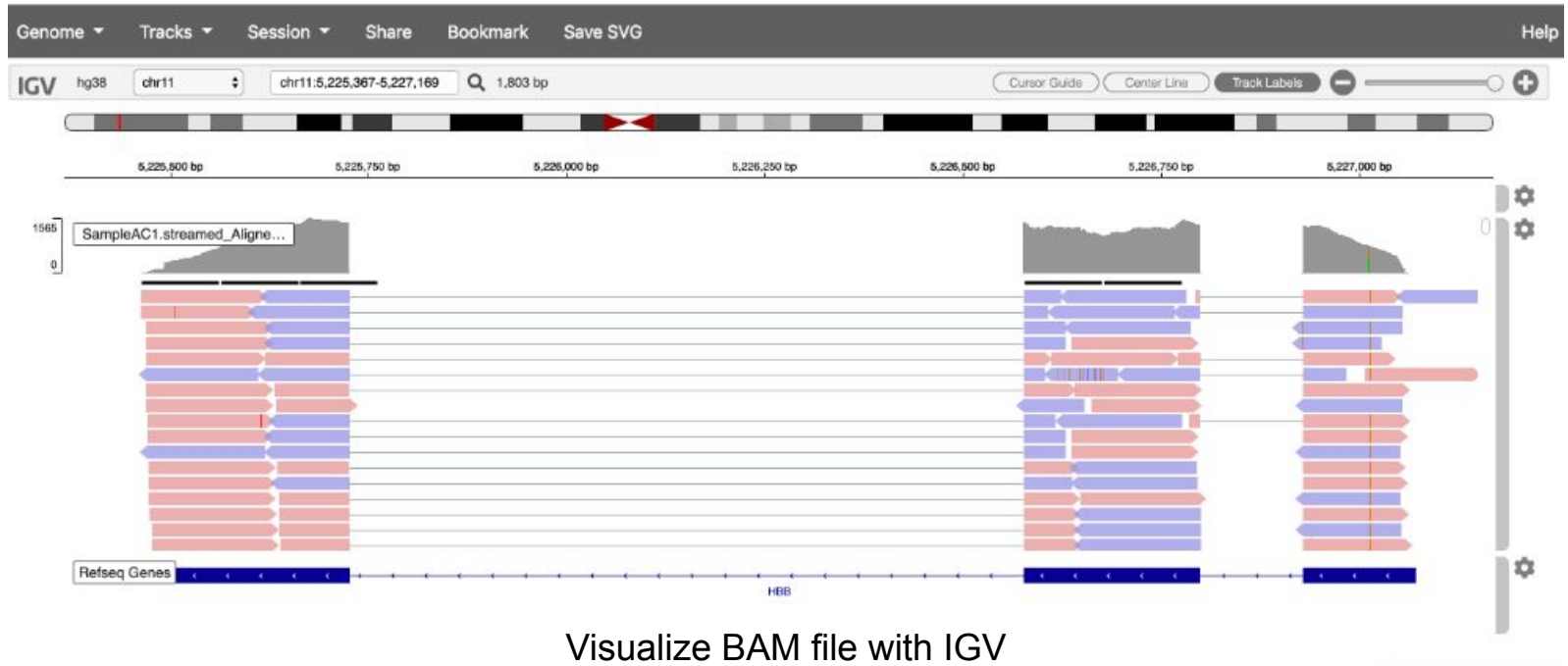
It indicates that the read aligns to the reference genome, but there is a region of the reference sequence that is not covered by the read.

**Splicing:**



32M27N59M

https://ucdavis-bioinformatics-training.github.io/2019_March_UCSF_mRNAseq_Workshop/data_reduction/alignment.html

**Alignment statistics & Utilities for manipulating alignment files.**



Visualize BAM file with IGV

**ALIGNMENT QUALITY CONTROL**

# RNA-SEQ: ALIGNMENT QUALITY CONTROL

Alignment quality metrics:

- Coverage uniformity along transcripts

- Saturation of sequencing depth

- Ribosomal RNA content (rRNA)

- Read distribution between exons, introns & intergenic regions.

- …

## RSeQC: Genebody Coverage



→ Used to assess the sequencing depth and coverage across the entire length of genes

- Gene expression quantification
- Transcript isoform analysis
- Detection of gene expression biases
- Assessing RNA integrity and sample quality



*Which samples may have been degraded?*

25

**RSeQC: Junction Annotation**



Constitutive splicing

Exon skipping

Intron retention

Mutually exclusive exons

Alternative 5' splice site

Alternative 3' splice site

*Splice variation*

The junction-annotation command:

1. Search an RNA-Seq bam file for splice junctions.

2. Compare them to a gene model.

3. Output whether the found junctions are novel, partially novel, or already annotated in a gene model.

## RSeQC: Junction Annotation



Transcripts

Junctions

Annotations from regtools

'DA': Known donor-acceptor

'NDA': Novel donor-acceptor combination

'D': Known donor, novel acceptor

'A': Known acceptor, novel donor

'N': Novel acceptor and novel donor.

**splicing junctions**

complete_novel 18%

partial_novel 4%

known 78%

**splicing events**

complete_novel 7%

partial_novel 2%

known 90%

https://regtools.readthedocs.io/en/latest/commands/junctions-annotate/

## RSeQC: Junction Saturation



*A sample that reaches a plateau before getting to 100% data indicates that all junctions in the library have been detected, and that further sequencing will not yield more observations.*

Junction Saturation Analysis

- Evaluates the depth of sequencing coverage at splice junctions.

- It helps determine if sufficient sequencing depth has been achieved to capture the full diversity of splice junctions.

→ Guides decisions on whether additional sequencing is needed to achieve more comprehensive coverage.

→ Ensures confidence in downstream analyses (alternative splicing analysis, isoform discovery).

**RSeQC: Read Distribution**



RSeQC: Read Distribution

Calculate how mapped reads were distributed over genome feature (like CDS exon, 5'UTR exon, 3' UTR exon, Intron, Intergenic regions).

**Biotypes Count**



featureCounts Biotypes

- A good RNAseq sample should have a large portion of the reads coming from protein coding genes.

- This plot can help you spot problems with your library such as incomplete rRNA depletion.

# QUANTIFICATION

**Quantification - Read Count**



Count how many reads have mapped to each gene.

→Using the **featureCounts** tool to get the gene counts

**Input**: BAM + GTF

**Output**: Number of reads (counts) associated with each feature of interest (genes, exons, transcript, etc.).

**Counting reads with featureCounts**

- Accurate, fast and is relatively easy to use
- Counts reads that map to a single location (uniquely mapping) and follows the scheme in the figure below for assigning reads to a gene/exon.

aligned read:
start: 113217600   end: 113217650

MOV10

GTF

| chr1 unknown | exon | 113217048 | 113217252 | . | + | . | gene_id "MOV10";p_id "P5535";transcript_id "NM_001130079" |
| chr1 unknown | exon | 113217048 | 113217351 | . | + | . | gene_id "MOV10";p_id "P5535";transcript_id "NM_020963" |
| chr1 unknown | exon | 113217470 | 113217671 | . | + | . | gene_id "MOV10";p_id "P5535";transcript_id "NM_001130079" |
| chr1 unknown | CDS | 113217535 | 113217671 | . | + | 0 | gene_id "MOV10";p_id "P5535";transcript_id "NM_001130079" |
| chr1 unknown | start_codon | 113217535 | 113217537 | . | + | | gene_id "MOV10";p_id "P5535";transcript_id "NM_001130079" |

feature type                                    feature

33

**Counting reads using featureCounts**



- A read is said to overlap a feature if at least one read base is found to overlap the feature.

- For paired-end data, a fragment (or template) is said to overlap a feature if any of the two reads from that fragment is found to overlap the feature.

- If strandedness is specified, then in addition to considering the genomic coordinates it will also take the strand into account for counting.

34

*Example of a featureCounts Assignment.*

## Counting reads using featureCounts



**Output: Raw counts**

These are the "raw" counts will be used in statistical programs downstream for differential gene expression.

*featureCounts output*

# RNA-SEQ: QUANTIFICATION

**Counting reads using featureCounts**



**A table of counts**

Don't need information about the genomic coordinates, length

→ Cleaning up the featureCounts matrix

**Final output:**
A count matrix, with genes as rows and samples are columns

"Effect of the lysphosphatidylcholine analogue edelfosine on gene expression in *Saccharomyces cerevisiae*"
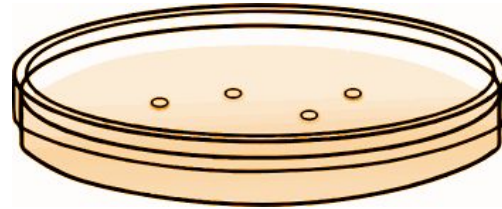
https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE227381

38

Comparative gene expression profiling analysis of RNA-seq data for S. cerevisiae cells treated with edelfosine for 60 minutes.



Control (X3)



Edelfosine treatment (X3)

```
qc_check
├── WT_C_1_R1.fastq.gz
├── WT_C_1_R2.fastq.gz
├── WT_C_2_R1.fastq.gz
├── WT_C_2_R2.fastq.gz
├── WT_C_3_R1.fastq.gz
├── WT_C_3_R2.fastq.gz
├── WT_E_1_R1.fastq.gz
├── WT_E_1_R2.fastq.gz
├── WT_E_2_R1.fastq.gz
├── WT_E_2_R2.fastq.gz
├── WT_E_3_R1.fastq.gz
├── WT_E_3_R2.fastq.gz
```

**Bulk RNA-seq Analysis**

**THANK YOU**