## SAMPLE PROJECT ON RNA-SEQUENCING

Name- Manjushri Kalpande

**SRR504382 our read id  - paired end (Study: Muscleblind-Like 2 mediated alternative splicing in the developing bain by mRNA sequencing)**

.

```
Files included in this directory:

mm39.fa.gz - "Soft-masked" assembly sequence in one file.
    Repeats from RepeatMasker and Tandem Repeats Finder (with period
    of 12 or less) are shown in lower case; non-repeating sequence is
    shown in upper case.
```

1) Downloaded a genome sequence of mouse



2) Next we created a new folder in sra data to store our prefetch our id and store fastq file

3) FASTQC output of both reads

1st read (good quality)



2nd read:: (ok but few place have low coverage)

4) Then we move the fastq files to place where we want to do BWA

```
mv: target 'altsplice' is not a directory
mkalpande@ManjushriK:~/sra_data/altsplice$ mv SRR504382_1.fastq SRR504382_2.fastq /home/mkalpande/bwa_index_files/sample\ al
tsplice
```

5) Next we give permissions to file, gunzip it as follows:

```
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ ls -lrt
total 9597536
-rw-r--r-- 1 mkalpande mkalpande        172 Dec 14 17:39 mm39.fa.gz:Zone.Identifier
-rw-r--r-- 1 mkalpande mkalpande  870543764 Dec 14 17:39 mm39.fa.gz
-rw-r--r-- 1 mkalpande mkalpande 4478654890 Dec 14 18:50 SRR504382_1.fastq
-rw-r--r-- 1 mkalpande mkalpande 4478654890 Dec 14 18:50 SRR504382_2.fastq
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ chmod +x mm39.fa.gz
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ ls -lrt
total 9597536
-rw-r--r-- 1 mkalpande mkalpande        172 Dec 14 17:39 mm39.fa.gz:Zone.Identifier
-rwxr-xr-x 1 mkalpande mkalpande  870543764 Dec 14 17:39 mm39.fa.gz
-rw-r--r-- 1 mkalpande mkalpande 4478654890 Dec 14 18:50 SRR504382_1.fastq
-rw-r--r-- 1 mkalpande mkalpande 4478654890 Dec 14 18:50 SRR504382_2.fastq
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ unzip mm39.fa.gz
Archive:  mm39.fa.gz
  End-of-central-directory signature not found.  Either this file is not
  a zipfile, or it constitutes one disk of a multi-part archive.  In the
  latter case the central directory and zipfile comment will be found on
  the last disk(s) of this archive.
note:  mm39.fa.gz may be a plain executable, not an archive
unzip:  cannot find zipfile directory in one of mm39.fa.gz or
        mm39.fa.gz.zip, and cannot find mm39.fa.gz.ZIP, period.
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ ls
SRR504382_1.fastq  SRR504382_2.fastq  mm39.fa.gz  mm39.fa.gz:Zone.Identifier
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ gunzip mm39.fa.gz
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ ls
SRR504382_1.fastq  SRR504382_2.fastq  mm39.fa  mm39.fa.gz:Zone.Identifier
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$
```

6) Now will do indexing using bwa index -p SRR82 file.fa
   Here, we are using **-p as prefix** as we can use same ref for other reads also.

```
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ ls
SRR504382_1.fastq  SRR504382_2.fastq  mm39.fa  mm39.fa.gz:Zone.Identifier
mkalpande@ManjushriK:~/bwa_index_files/sample_altsplice$ bwa index -p SRR82 mm39.fa
[bwa_index] Pack FASTA... 12.00 sec
[bwa_index] Construct BWT for the packed sequence...
[BWTIncCreate] textLength=5456444902, availableWord=395935328
[BWTIncConstructFromPacked] 10 iterations done. 99999990 characters processed.
[BWTIncConstructFromPacked] 20 iterations done. 199999990 characters processed.
[BWTIncConstructFromPacked] 30 iterations done. 299999990 characters processed.
[BWTIncConstructFromPacked] 40 iterations done. 399999990 characters processed.
```

```
[BWTIncConstructFromPacked] 590 iterations done. 5402330102 chara
[BWTIncConstructFromPacked] 600 iterations done. 5427525990 chara
[BWTIncConstructFromPacked] 610 iterations done. 5449916134 chara
[bwt_gen] Finished constructing BWT in 614 iterations.
[bwa_index] 2193.71 seconds elapse.
[bwa_index] Update BWT... 14.17 sec
[bwa_index] Pack forward-only FASTA... 12.57 sec
[bwa_index] Construct SA from BWT and Occ... 1259.67 sec
[main] Version: 0.7.17-r1198-dirty
[main] CMD: bwa index -p SRR82 mm39.fa
[main] Real time: 3501.289 sec; CPU: 3492.117 sec
```

Its done now.

7) Now we will do bwa mem for mapping of fastq files and ref sequence.
   bwa mem SRR82 SRR504382_1.fastq SRR504382_2.fastq > SR82output.sam

8) For sam to bam we will use following command:
   samtools view -1 -bS SR82output.sam > SR82output.bam

**where-  samtools view**: Starts the conversion process.

**-1**: Specifies the input file is in SAM format.

   **-b**: Converts the output to BAM format.

   **-S:** Sorts the alignments by reference coordinates before converting to BAM. This is crucial for efficient downstream analyses.

9) Next we will do <u>sorting</u> of bam file followed by <u>indexing</u>..

   <u>samtools sort -T temp -o sorted_SR82output.bam SR82output.bam</u>

   - samtools sort orders the alignments by **chromosomal coordinates**, meaning reads from the beginning of the first chromosome will appear first in the file, followed by reads from the beginning of the second chromosome….
   - -T This command tells Samtools to use the prefix temp for the temporary files and write the final sorted BAM file
   - -o Specifies the path and filename for the output BAM file containing the sorted alignments.



10) Next we will do indexing of sorted bam file
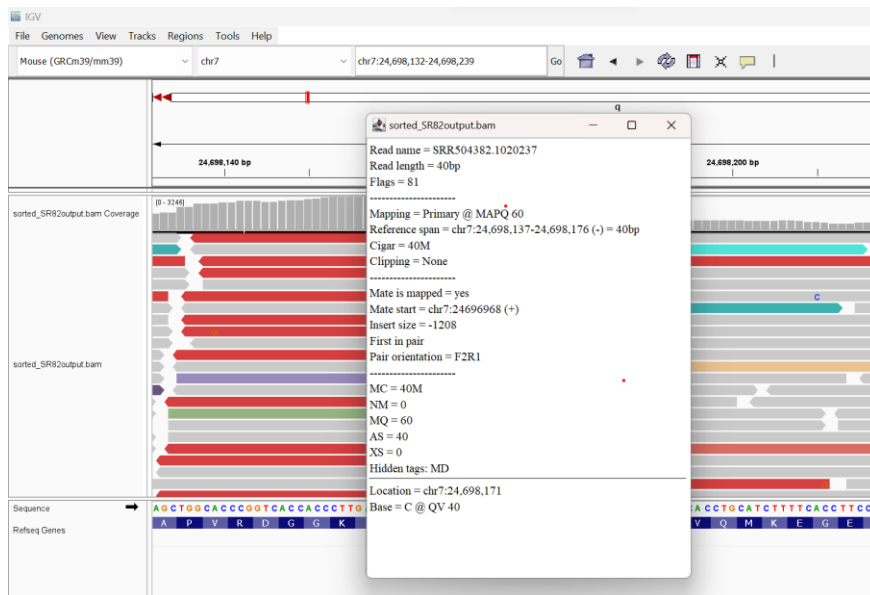
   <u>samtools index sorted_SR82output.bam</u>

   - **Samtools index**-The index file acts as a map, allowing tools to quickly find specific regions of the genome within the BAM file, improves the performance of many downstream analysis tasks, such as variant calling, read counting, and gene expression analysis.
   - **.bai is index file**
   - Next we will view our data in IGV

For input in IGV we selected file that is (sorted_SR82output.bam)



We selected this and enter the nucleotide number along with chr number in igv as cigar string is showing different

---- we perform alignment using HISAT2 as well----

Build reference genome- **hisat_mouse** -is the index file name

## hisat2-build mm39.fa hisat_mouse

```
Headers:
    len: 2654621783
    gbwtLen: 2654621784
    nodes: 2654621784
    sz: 663655446
    gbwtSz: 663655447
    lineRate: 6
    offRate: 4
    offMask: 0xfffffff0
    ftabChars: 10
    eftabLen: 0
    eftabSz: 0
    ftabLen: 1048577
    ftabSz: 4194308
    offsLen: 165913862
    offsSz: 663655448
    lineSz: 64
    sideSz: 64
    sideGbwtSz: 48
    sideGbwtLen: 192
    numSides: 13826156
    numLines: 13826156
    gbwtTotLen: 884873984
    gbwtTotSz: 884873984
    reverse: 0
    linearFM: Yes
Total time for call to driver() for forward index: 00:43:20
```

Files generated during indexing is-

```
(venv) mkalpande@ManjushriK:~/hisat2-2.2.1/hisat2_index_files$ ls
chr11.fa              hisat_mouse.2.ht2  hisat_mouse.4.ht2  hisat_mouse.6.ht2  hisat_mouse.8.ht2  mm39.fa
hisat_mouse.1.ht2  hisat_mouse.3.ht2  hisat_mouse.5.ht2  hisat_mouse.7.ht2  mm39.build
```

2) Running hisat2 so need to use the command-

```
mkalpande@ManjushriK:~/hisat2-2.2.1/hisat2_index_files$ hisat2 -p 4 -x /home/mkalpande/hisat2-2.2.1/hisat2_index_files/hisat_mouse -1
/home/mkalpande/hisat2-2.2.1/hisat2_index_files/SRR504382_1.fastq -2 /home/mkalpande/hisat2-2.2.1/hisat2_index_files/SRR504382_2.fas
tq -S hisat_alignment.sam
21196948 reads; of these:
  21196948 (100.00%) were paired; of these:
    5683168 (26.81%) aligned concordantly 0 times
    14548882 (68.64%) aligned concordantly exactly 1 time
    964898 (4.55%) aligned concordantly >1 times
    ----
    5683168 pairs aligned concordantly 0 times; of these:
      1380294 (24.29%) aligned discordantly 1 time
    ----
    4302874 pairs aligned 0 times concordantly or discordantly; of these:
      8605748 mates make up the pairs; of these:
        4166215 (48.41%) aligned 0 times
        3695763 (42.95%) aligned exactly 1 time
        743770 (8.64%) aligned >1 times
90.17% overall alignment rate
mkalpande@ManjushriK:~/hisat2-2.2.1/hisat2_index_files$ ls
SRR504382_1.fastq  hisat_alignment.sam  hisat_mouse.6.ht2  hisat_mouse.6.ht2  hisat_mouse_indexes  venv
SRR504382_2.fastq  hisat_mouse.1.ht2    hisat_mouse.4.ht2  hisat_mouse.7.ht2  mm39.build
chr11.fa           hisat_mouse.2.ht2    hisat_mouse.5.ht2  hisat_mouse.8.ht2  mm39.fa
```

3) We converted to **BAM** file and then further sort the file .

```
        Set level of verbosity
mkalpande@ManjushriK:~/hisat2-2.2.1/hisat2_index_files$ samtools sort /home/mkalpande/hisat2-2.2.1/hisat2_index_files/hisat_alignment
.sam -o/home/mkalpande/hisat2-2.2.1/hisat2_index_files/hisat_alignment.bam
[bam_sort_core] merging from 12 files and 1 in-memory blocks...
mkalpande@ManjushriK:~/hisat2-2.2.1/hisat2_index_files$ ls
SRR504382_1.fastq  hisat_alignment.bam  hisat_mouse.2.ht2  hisat_mouse.5.ht2  hisat_mouse.8.ht2    mm39.fa
SRR504382_2.fastq  hisat_alignment.sam  hisat_mouse.3.ht2  hisat_mouse.6.ht2  hisat_mouse_indexes  venv
chr11.fa           hisat_mouse.1.ht2    hisat_mouse.4.ht2  hisat_mouse.7.ht2  mm39.build
mkalpande@ManjushriK:~/hisat2-2.2.1/hisat2_index_files$
```
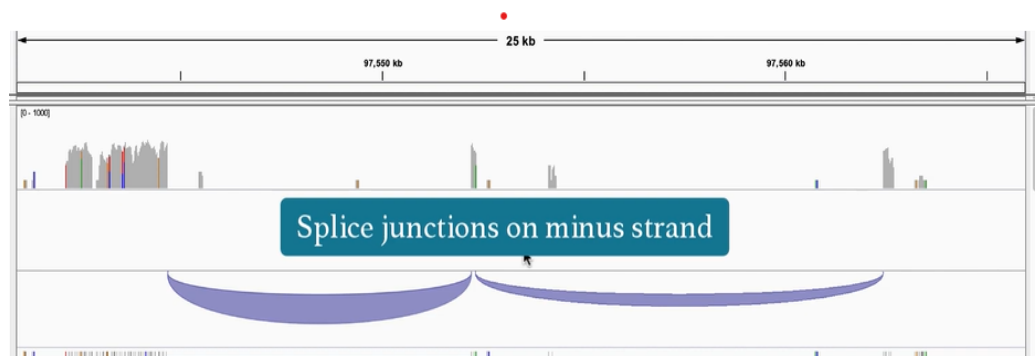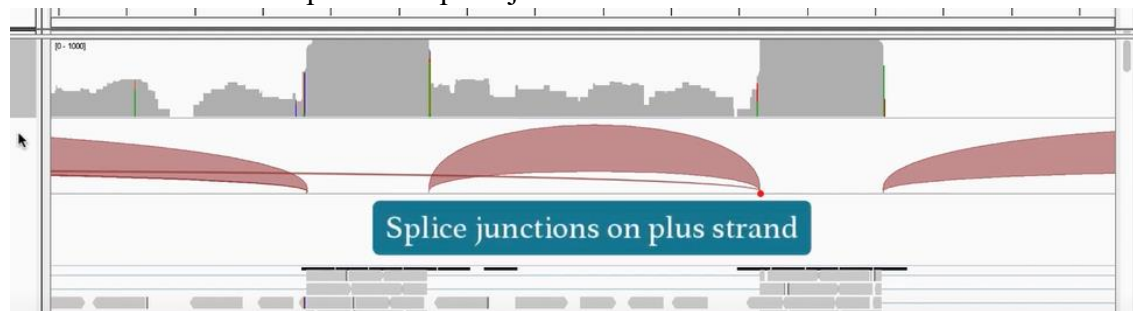
4) Next we will be using IGV to analyse the alignment
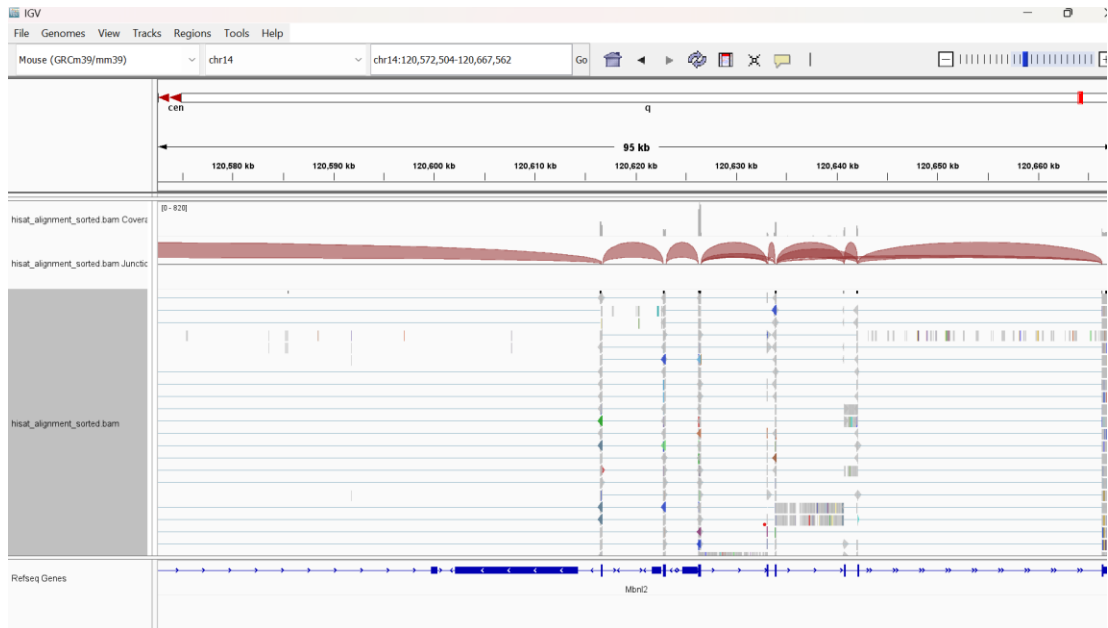- Colour alignments by read strand



Red reads are forward strands & blue reads are reverse strand

- Brown colour uplifts are splice junctions

** SHASHIMI PLOT IS USED TO STUDY SPLICE JUNCTIONS



- Our shashimi plot result
- **Junction depth** refers to the number of reads that map across a specific splice junction. It is a measure of the **abundance** of that junction in the sample. A <u>higher junction depth indicates that the junction is more frequently used in the RNA transcripts.</u>



- The bottom blue lines from above image shows genes isoforms

Using Stringtie for transcript quantification:

```
(base) mkalpande@ManjushriK:~$ cd TOOLS/stringtie-2.2.1/final_stringtie_output/
(base) mkalpande@ManjushriK:~/TOOLS/stringtie-2.2.1/final_stringtie_output$ ls
alt_splice_out1.gtf  ballgown  merge.txt  merged_output3.gtf  output_transcripts2.gtf  stringe_compare
(base) mkalpande@ManjushriK:~/TOOLS/stringtie-2.2.1/final_stringtie_output$
```

- Using gffcompare generated a statistical summary file:

```
# gffcompare v0.12.9 | Command line was:
#gffcompare -R -r /home/mkalpande/mm39_genefile.gtf -o stringe_compare/str_compare /home/mkalpande/TOOLS/stringtie-2.2.1/fi
al_stringtie_output/merged_output3.gtf
#

#= Summary for dataset: /home/mkalpande/TOOLS/stringtie-2.2.1/final_stringtie_output/merged_output3.gtf
#     Query mRNAs :   153002 in    54981 loci  (124443 multi-exon transcripts)
#            (20936 multi-transcript loci, ~2.8 transcripts per locus)
# Reference mRNAs :  148859 in    54703 loci  (120749 multi-exon)
# Super-loci w/ reference transcripts:     54642
#-----------------| Sensitivity | Precision  |
        Base level:    100.0     |    98.8    |
        Exon level:     93.7     |    97.5    |
      Intron level:    100.0     |    99.6    |
Intron chain level:     99.9     |    97.0    |
   Transcript level:     99.4     |    96.7    |
        Locus level:     99.8     |    99.2    |

     Matching intron chains:   120670
     Matching transcripts:   148019
            Matching loci:    54575

        Missed exons:        0/457057  (  0.0%)
         Novel exons:      641/427887  (  0.1%)
       Missed introns:        5/289926  (  0.0%)
        Novel introns:      160/291096  (  0.1%)
         Missed loci:        0/54703  (  0.0%)
          Novel loci:      339/54981  (  0.6%)
```