

## **Why do we do RNA-seq analysis?**

- 1- Quantify expression Genome-wide in a single assay.
- 2- Finding novel genes and splice variants (isoforms of known genes).
- 3- Compare genes and transcripts expression under two or more conditions.

### **A) Materials and Software:**

- Bowtie2 (<http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>)
- Samtools (<http://samtools.sourceforge.net>)
- Tophat2 (<http://ccb.jhu.edu/software/tophat/index.shtml>)
- STAR (<https://github.com/alexdobin/STAR>)
- Cufflinks (<http://cole-trapnell-lab.github.io/cufflinks/>)
- Rstudio (<https://www.rstudio.com>)
- CummeRbund package
- DESeq2 package

### **B) Datasets and Genome reference**

- Reference Genome: Homo\_Sapiens release  
([http://ftp.ensembl.org/pub/release-75/fasta/homo\\_sapiens/dna/](http://ftp.ensembl.org/pub/release-75/fasta/homo_sapiens/dna/))
- SRA: ERP001304 (RNA-seq datasets)  
(<http://www.ncbi.nlm.nih.gov/sra/?term=ERP001304>)

## **Results**

### A) Quality Control Results for Health samples

- 1) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103422\\_C02\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103422_C02_fastqc.html)
- 2) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103423\\_C04\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103423_C04_fastqc.html)
- 3) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103424\\_C05\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103424_C05_fastqc.html)
- 4) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103428\\_C25\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103428_C25_fastqc.html)

### B) Quality Control Results for Schizophrenia Disease samples

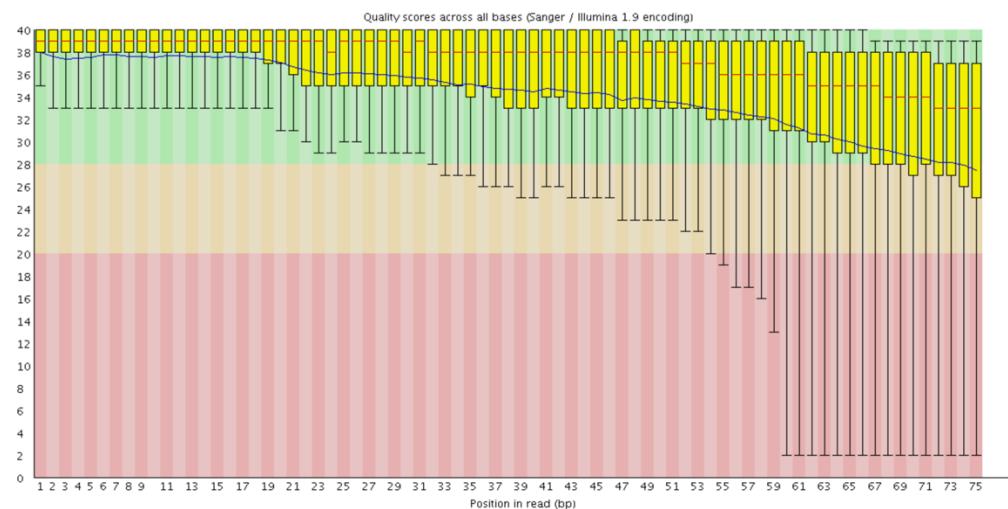
- 1) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103431\\_S02\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103431_S02_fastqc.html)
- 2) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103432\\_S04\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103432_S04_fastqc.html)
- 3) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103433\\_S05\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103433_S05_fastqc.html)
- 4) [file:///Users/sultanaharbi/Desktop/Schizophrenia\\_analysis/ERR103436\\_S23\\_fastqc.html](file:///Users/sultanaharbi/Desktop/Schizophrenia_analysis/ERR103436_S23_fastqc.html)



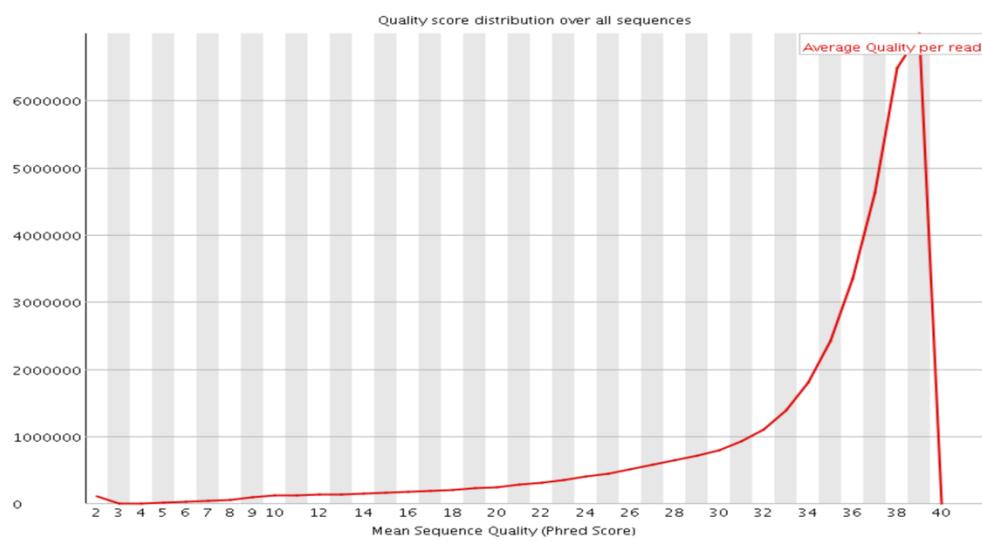
### Basic Statistics

Measure	Value
Filename	ERR103436_S23.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	36526047
Sequences flagged as poor quality	0
Sequence length	75
%GC	49

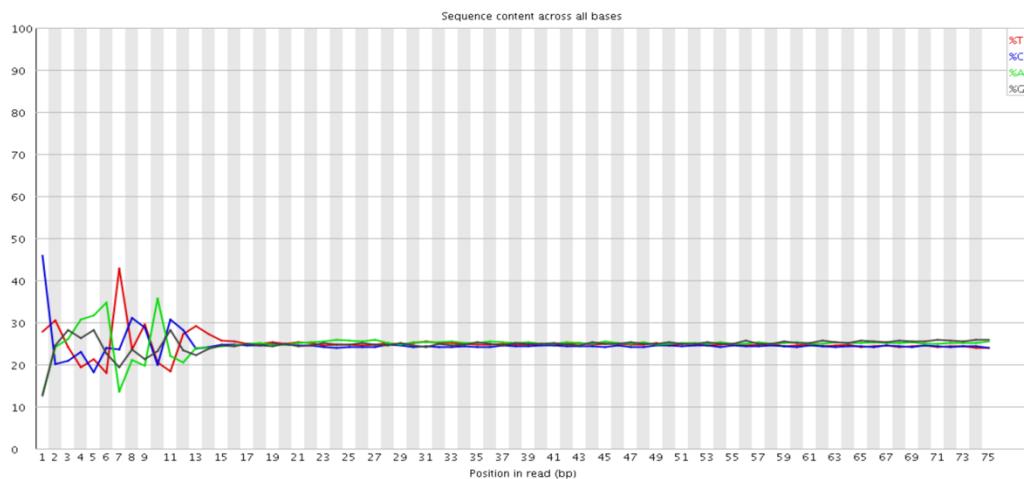
### Per base sequence quality



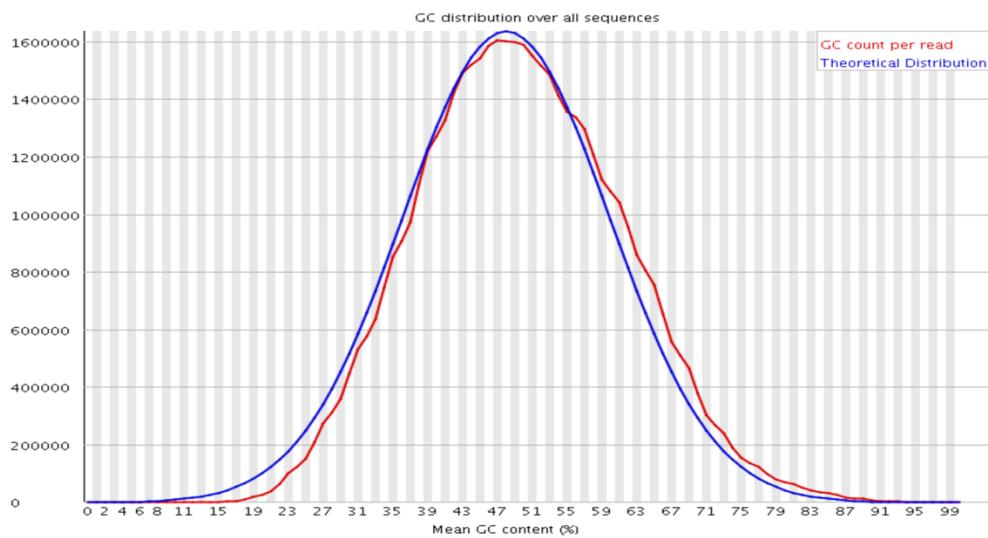
### Per sequence quality scores



### ✖ Per base sequence content



### ✓ Per sequence GC content



### C) Tophat2 Results (aligning results)

```
sultanaharbi@sULTANs-iMac:~$ head /Users/sultanaharbi/Desktop/RNA_seq_Analysis/SULTAN_BIN/ERR103422_C02/align_summary.txt
Reads:
Input      : 14880770
Mapped     : 10505850 (70.6% of input)
of these:   985733 ( 9.4%) have multiple alignments (578 have >20)
70.6% overall read mapping rate.
```

```
sultanaharbi@sULTANs-iMac:~$ head /Users/sultanaharbi/Desktop/RNA_seq_Analysis/SULTAN_BIN/ERR103423_C04/align_summary.txt
Reads:
Input      : 31819015
Mapped     : 28779389 (90.4% of input)
of these:   2594075 ( 9.0%) have multiple alignments (1640 have >20)
90.4% overall read mapping rate.
```

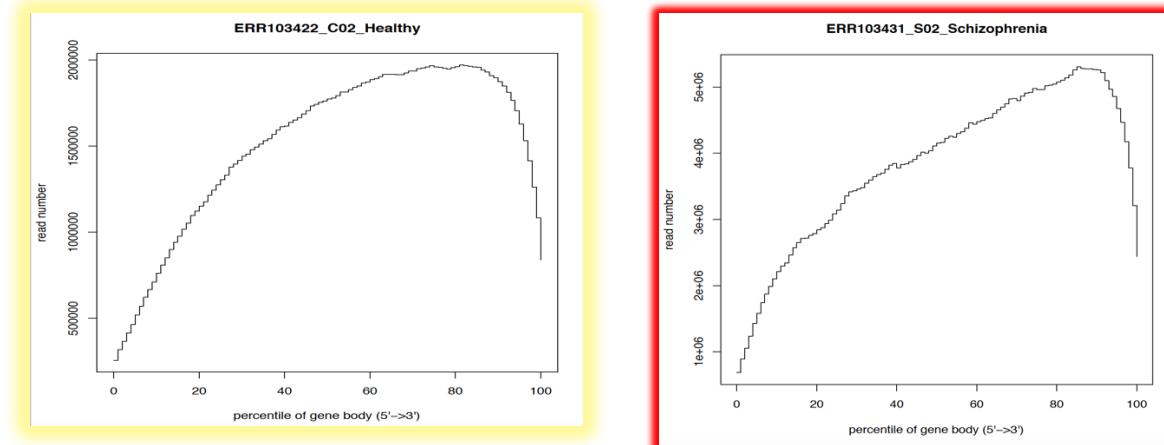
### D) Calculating the distribution of reads to different genomic feature types

Group	Total_bases	Tag_count	Tags/Kb
CDS_Exons	36860499	6149726	16.59
5'UTR_Exons	3073977	8035	22.85
3'UTR_Exons	5786939	2718461	47.03
Introns	1453678930	542622	0.37
TSS_up_5kb	2046080	64	0.25
TSS_up_5kb	141452483	18410	0.13
TSS_up_10kb	252939815	25671	0.10
TES_down_1kb	33557126	20268	0.68
TES_down_1kb	145246303	1526	0.43
TES_down_10kb	255858394	77862	0.30

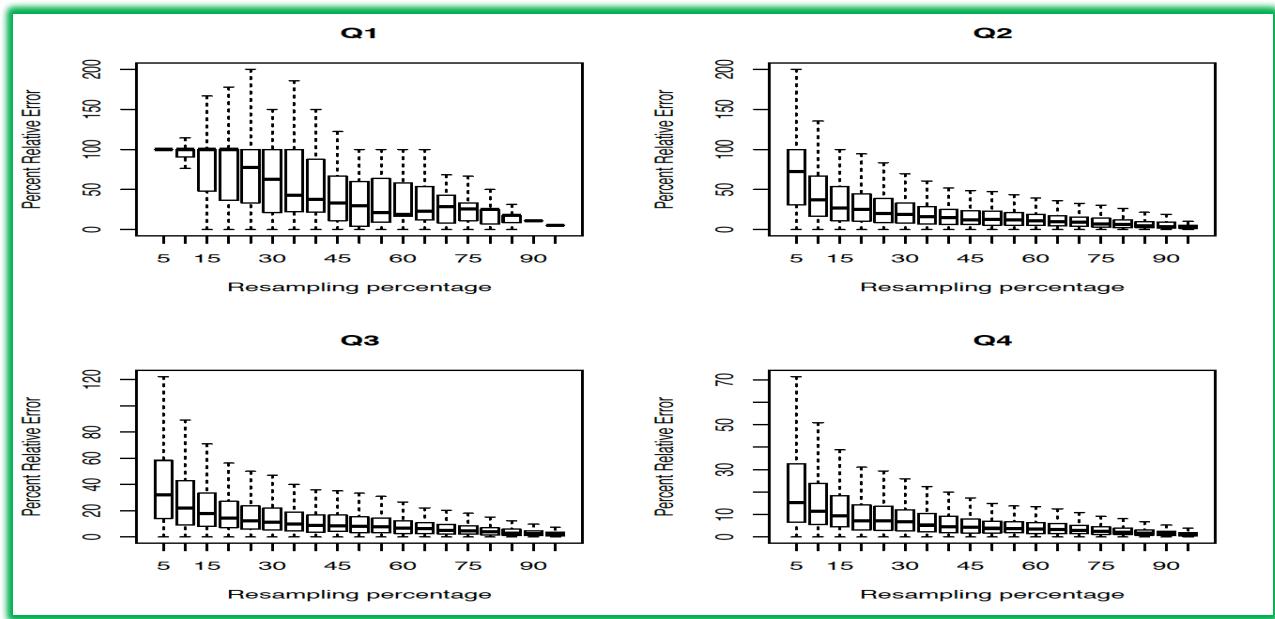
  

Group	Total_bases	Tag_count	Tags/Kb
CDS_Exons	36860499	12374885	335.72
5'UTR_Exons	35730397	1456686	40.77
3'UTR_Exons	5786934	9346534	156.23
Introns	1453670930	1303662	0.25
TSS_up_1kb	31788498	24125	0.76
TSS_up_5kb	141452483	57281	0.48
TSS_up_10kb	252939815	7807	0.31
TES_down_1kb	33557126	66840	1.97
TES_down_1kb	145246303	186425	1.28
TES_down_10kb	255858394	232561	0.91

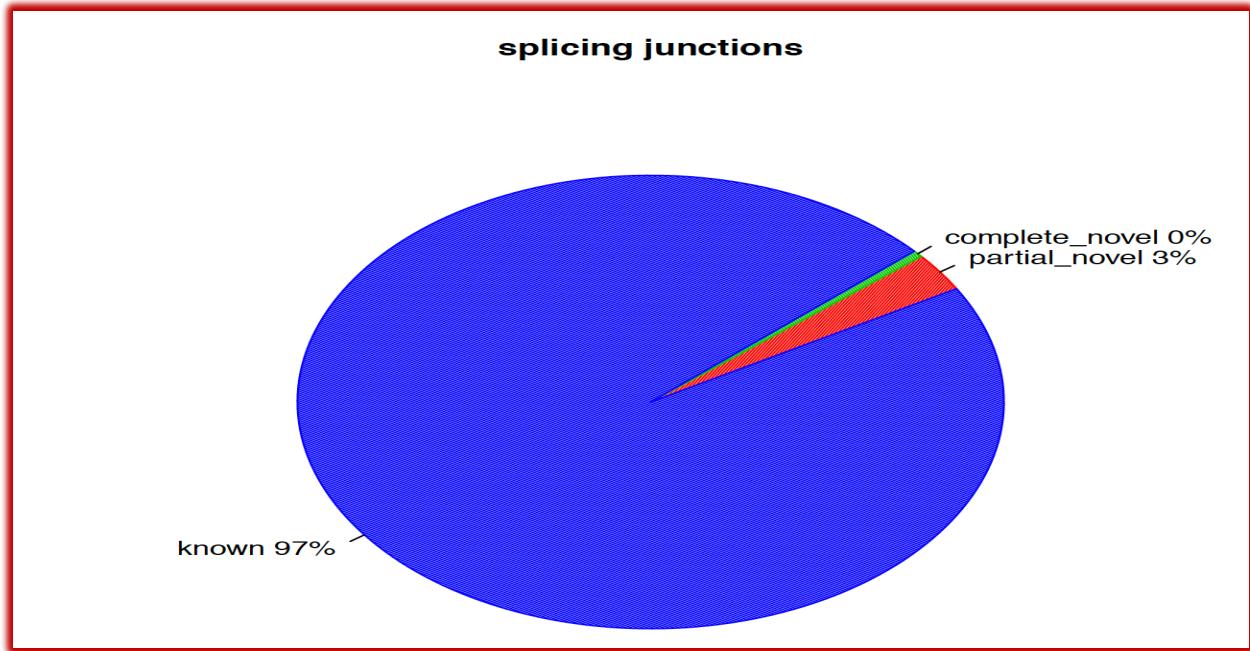
### E) Checking the uniformity of transcripts along 3' and 5' ends



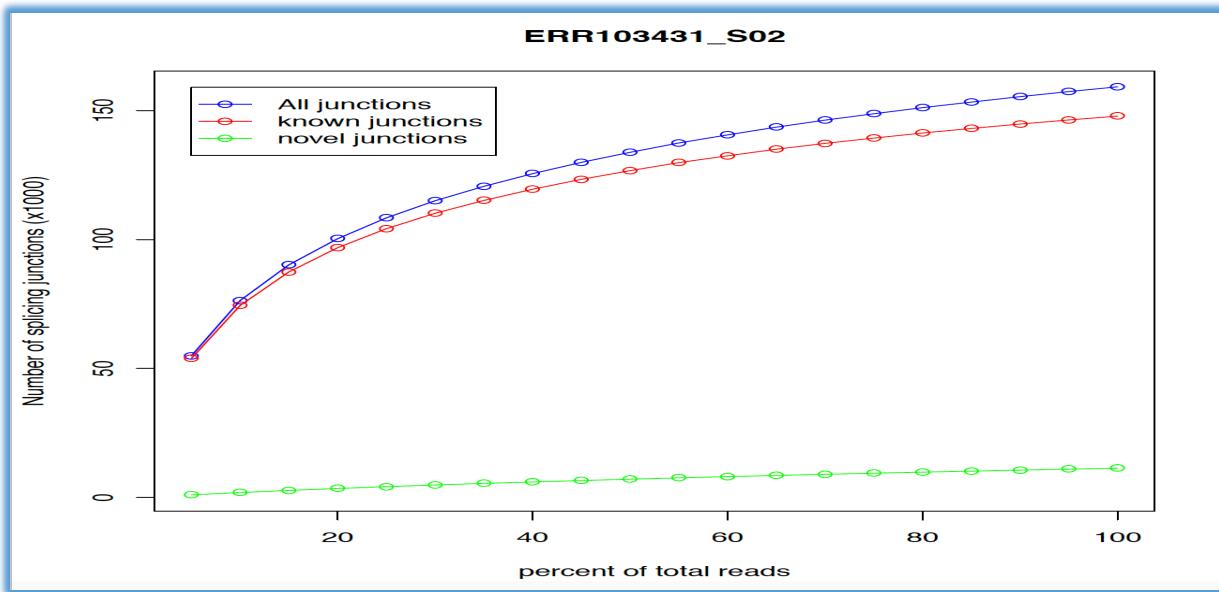
F) Calculating abundance in RPKM units



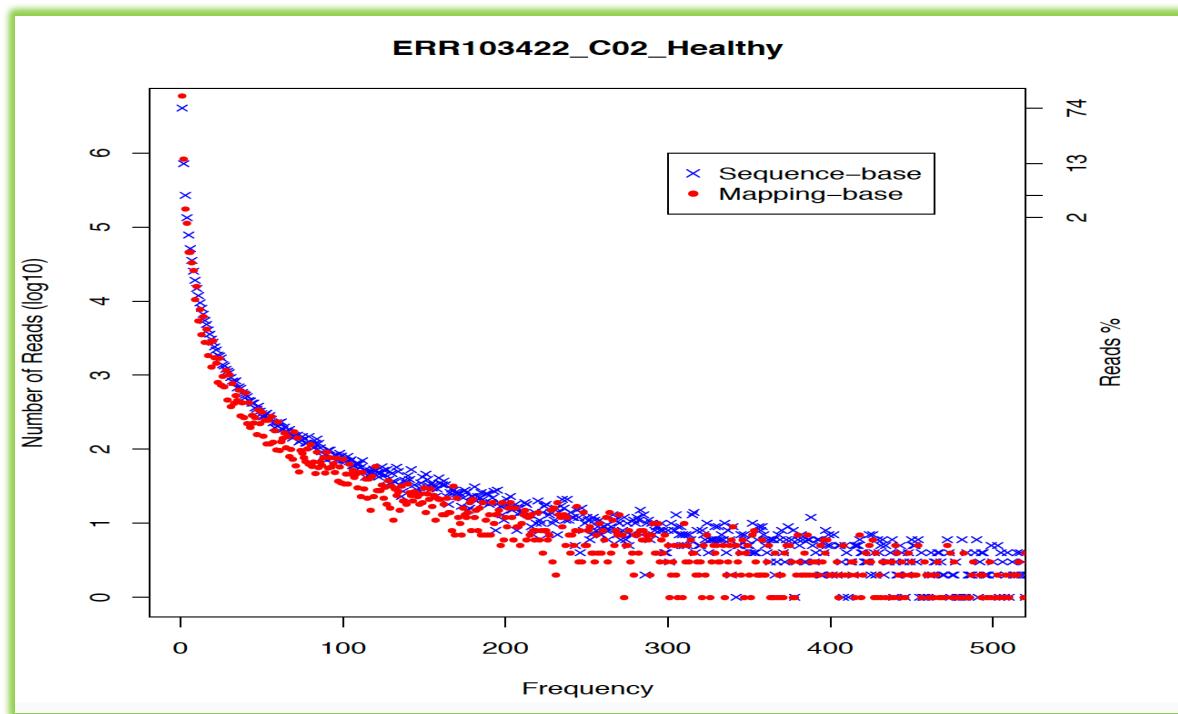
G) Detecting splicing junctions as novel, partially and known



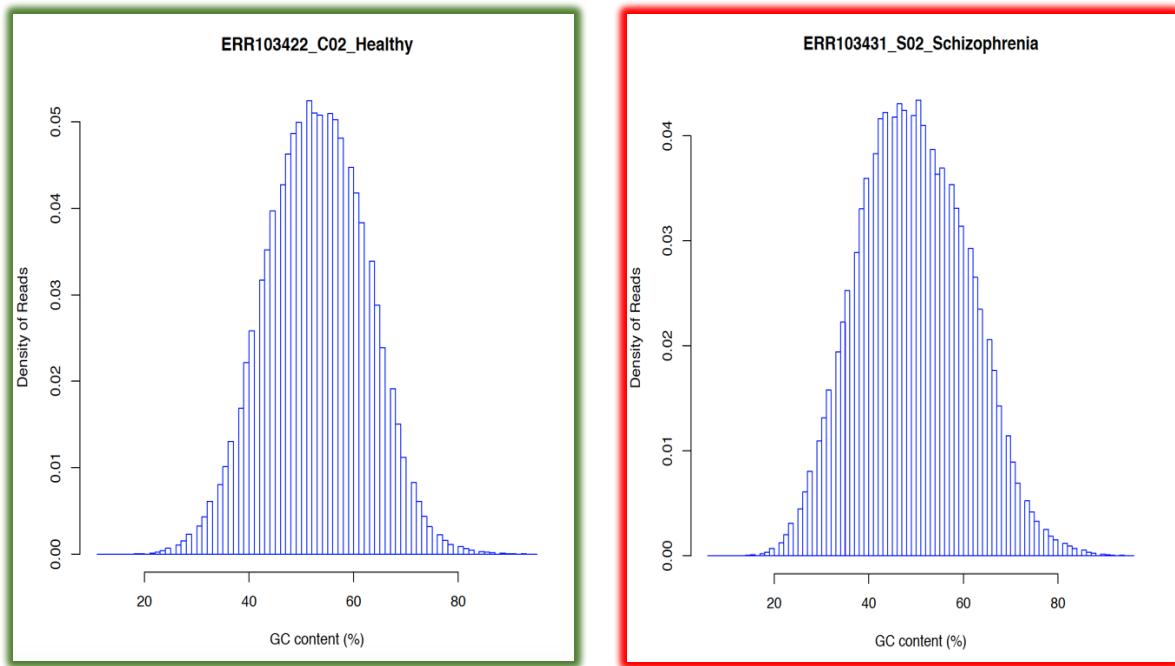
K) Detecting junctions in each subset



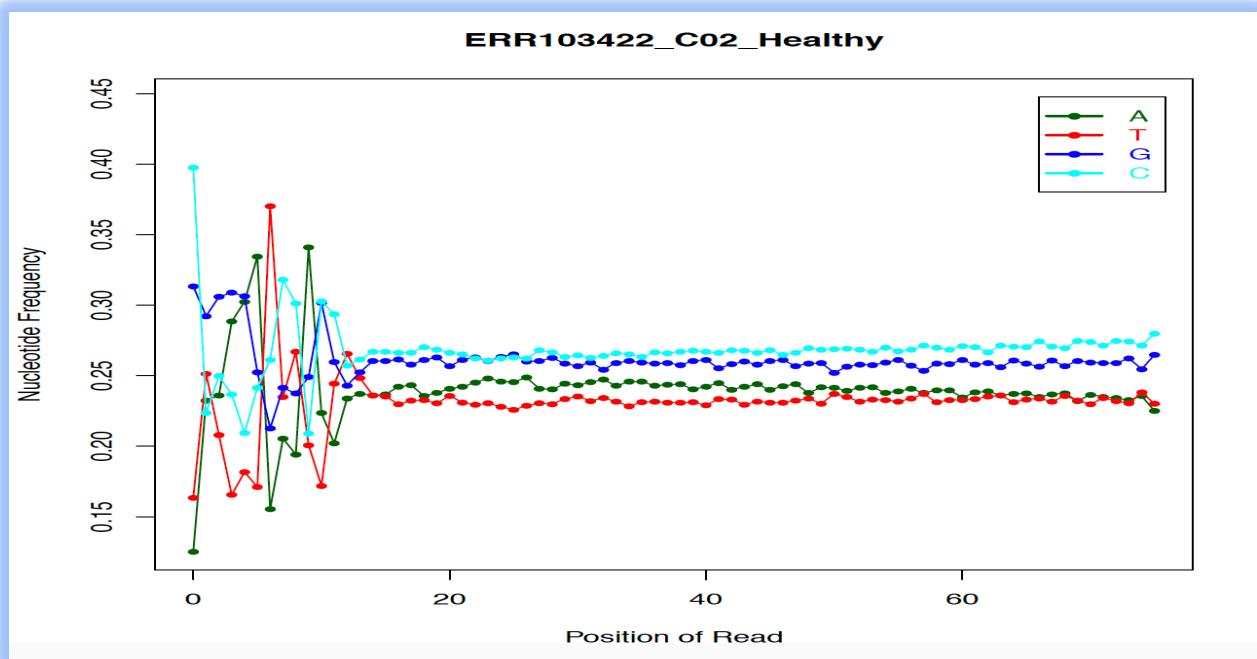
L) Checking Read Duplication



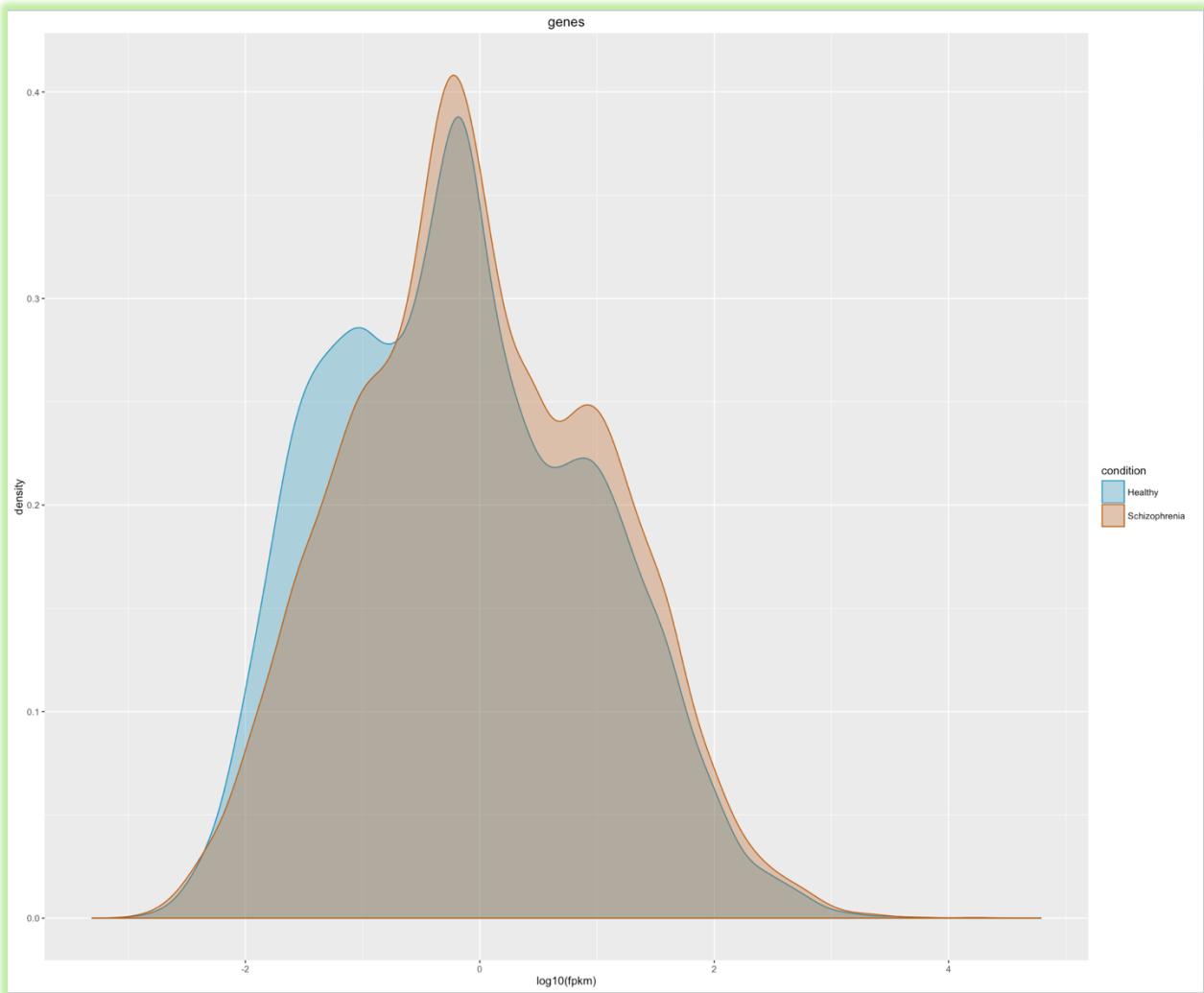
M) Distribution of read GC content (%):



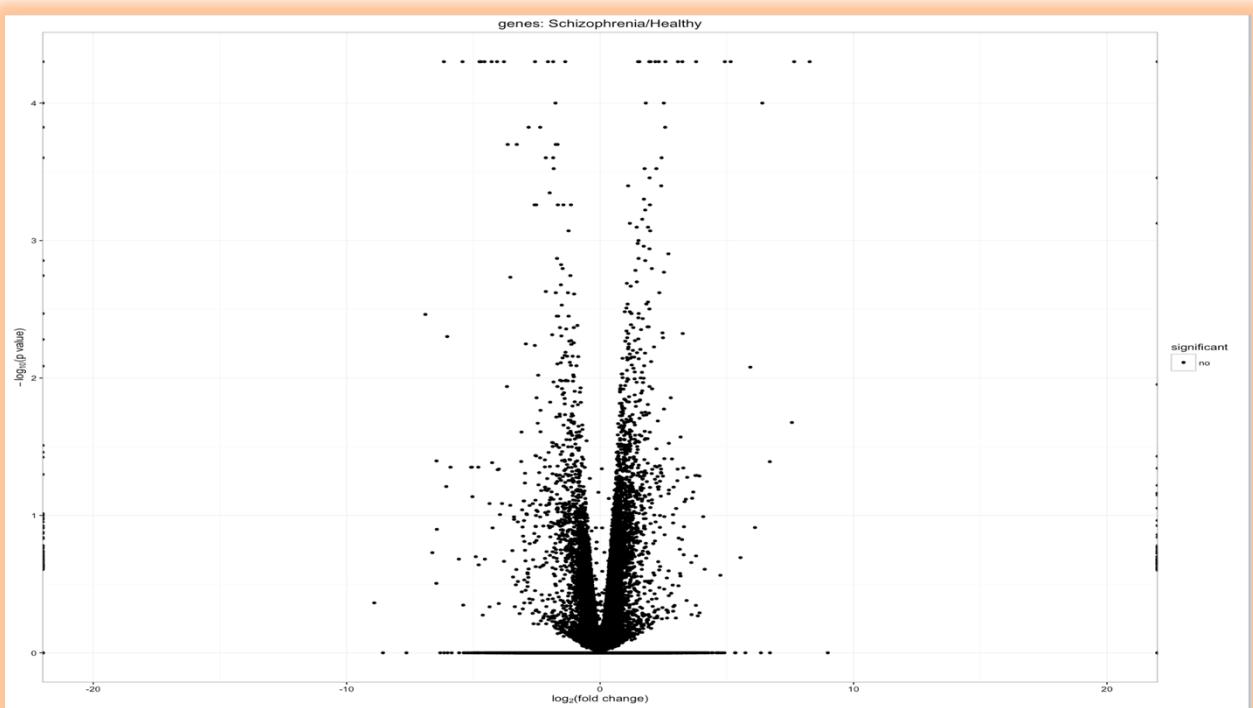
N) Read - Nucleotide vs Cycle (Phred base score vs. position in read):



## Results of Differential Expression analysis

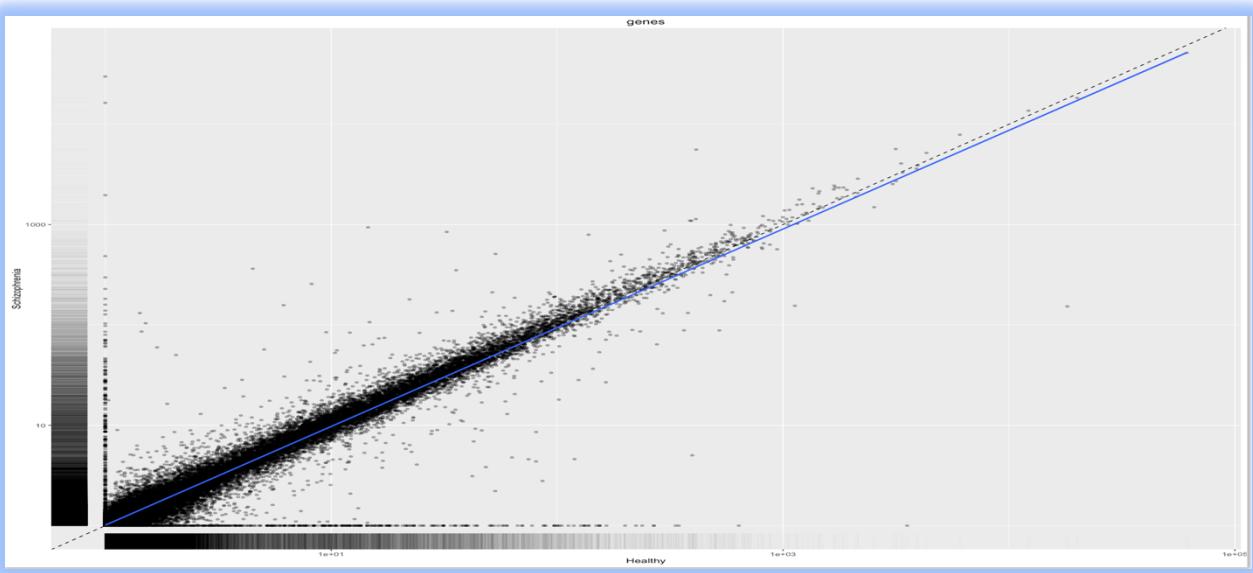


- 1) A plot showing the expression levels for each sample, using the `csDensity` Function (CummeRbund package). Command line: `> csDensity(genes(cuffdata))`



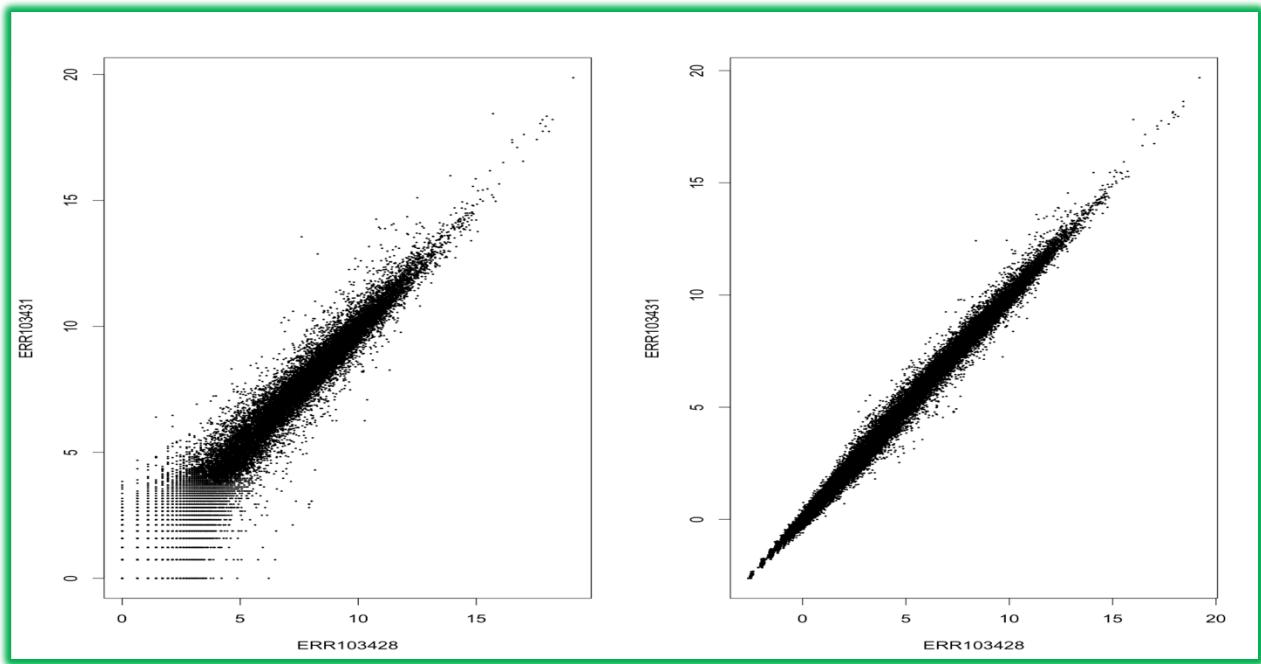
2) A Volcano plot showing differentially expressed genes across the two samples

Command line: `> csVolcano(genes(cuffdata), 'Healthy', 'Schizopherina')`

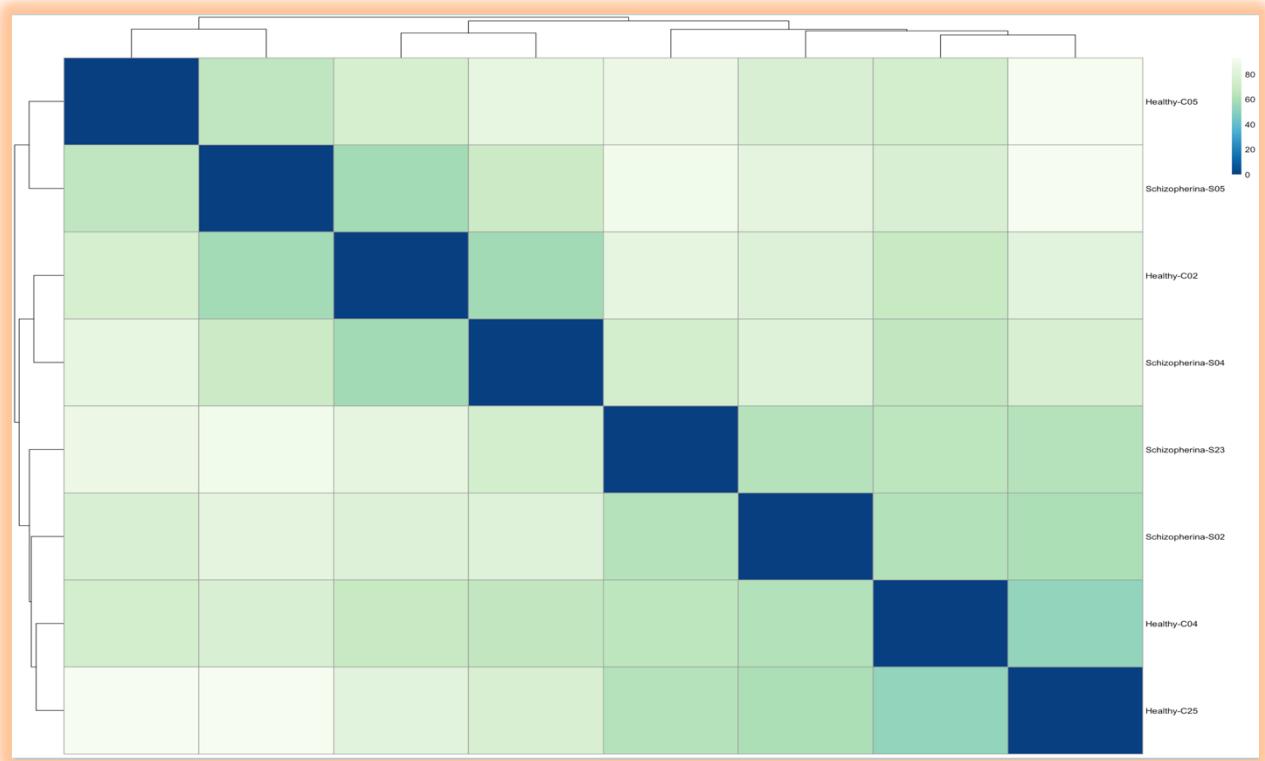


3) A scatter plot comparing the gene across the two samples Command Line:

```
> csScatter(genes(cuffdata), 'Healthy', 'Schizopherina', smooth = T, logMode= TRUE,
colorByStatus= T, hexbin=T)
```



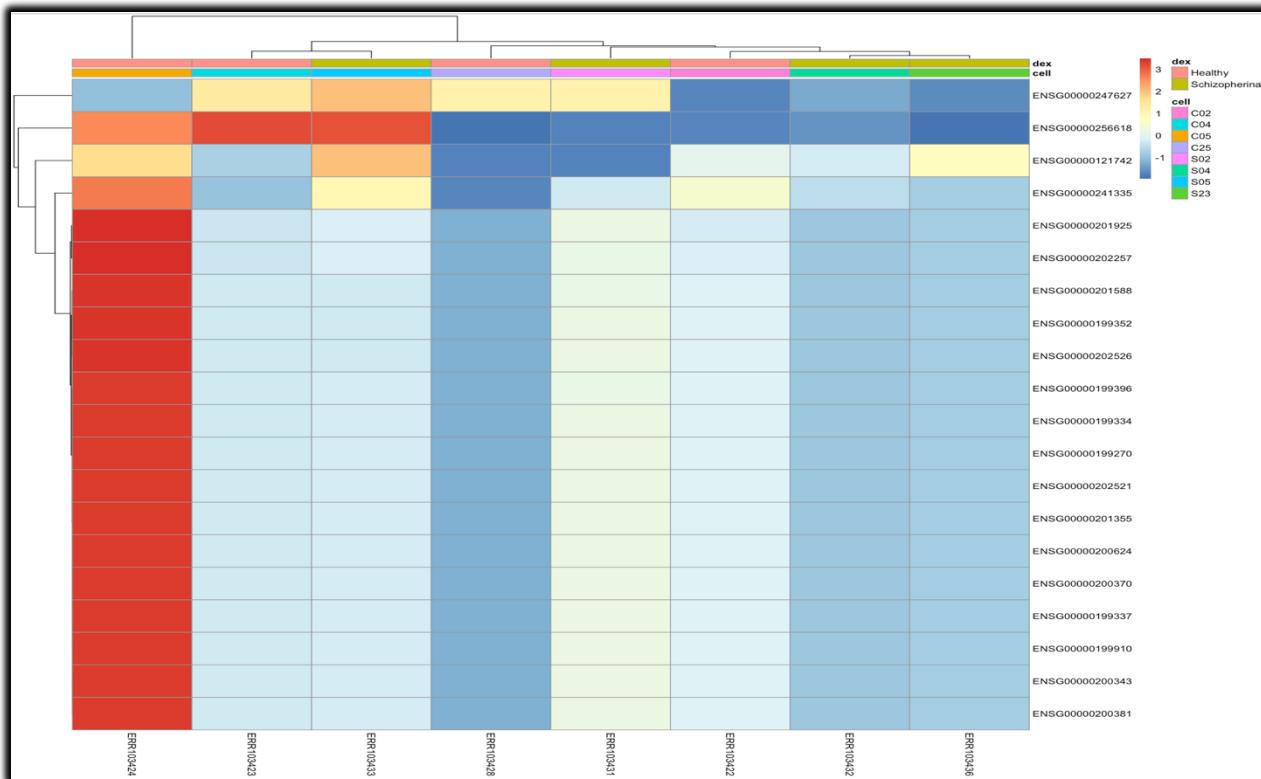
- 4) Scatterplot of transformed counts from two samples. Shown are scatterplots using the log2 transform of normalized counts (left side) and using the rlog (right side).



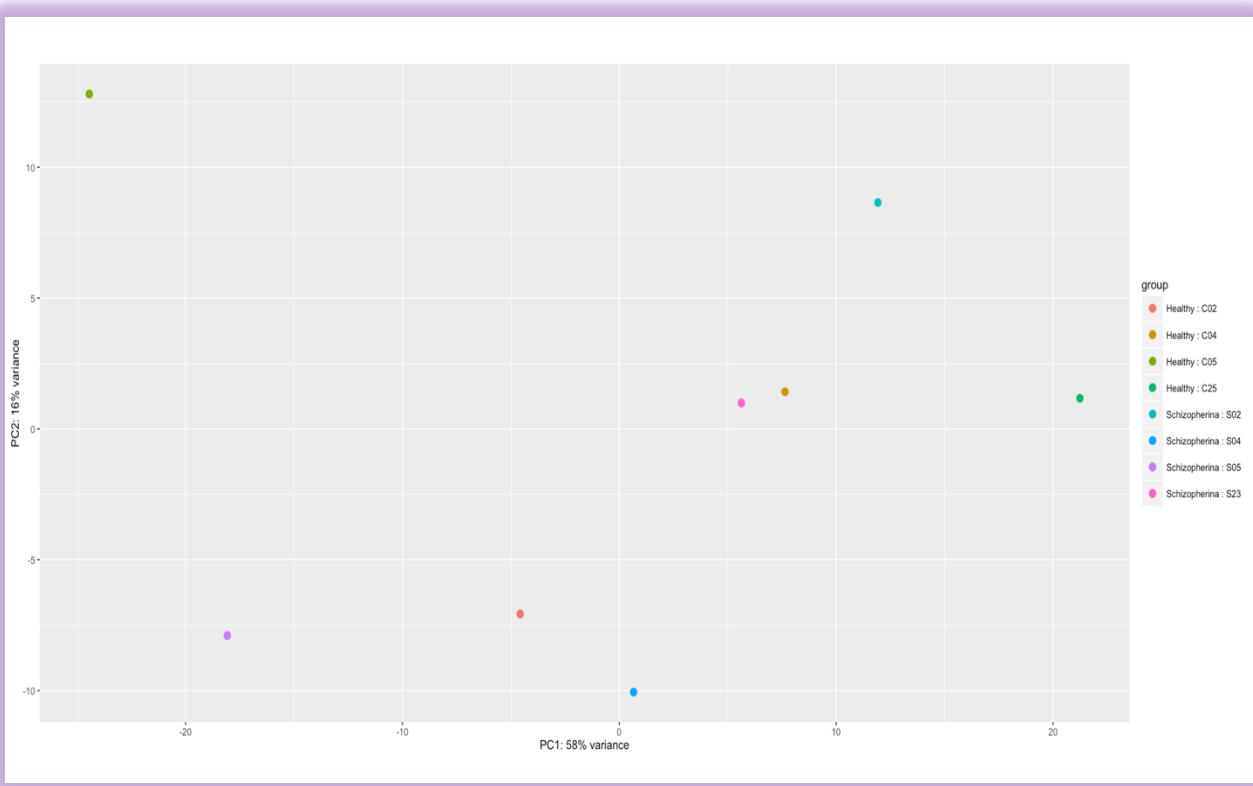
- 5) Heatmap of sample-to-sample distances using the rlog-transformed values



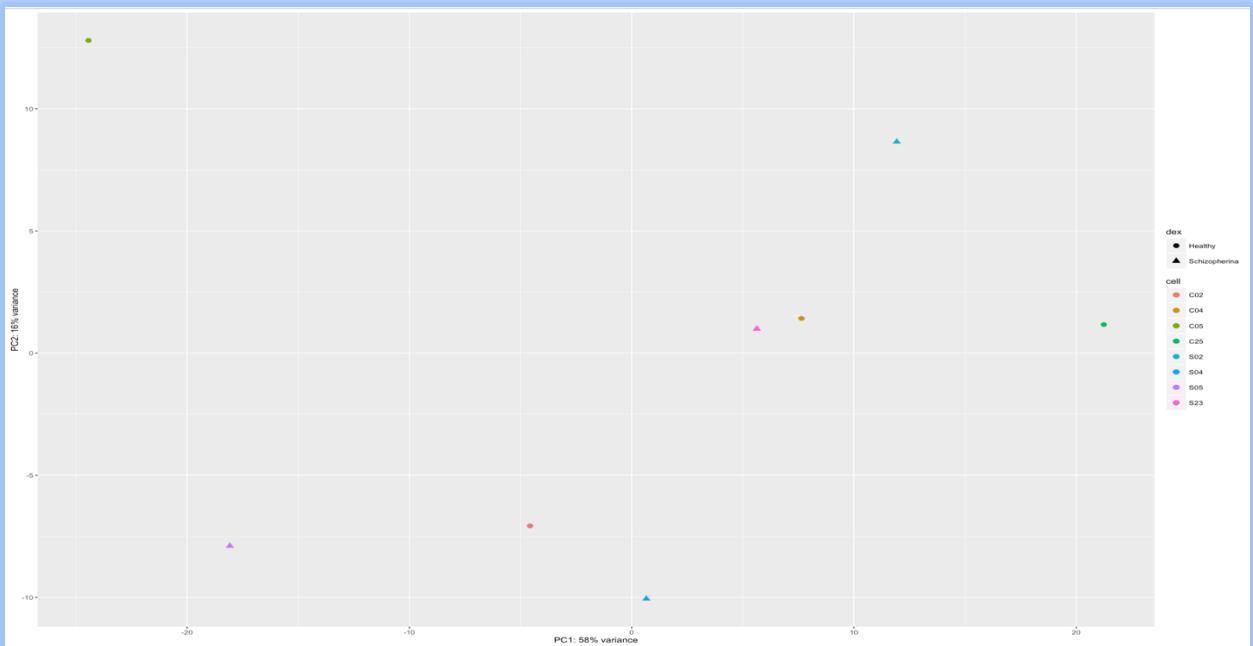
6) Heatmap of sample-to-sample distances using the Poisson Distance.



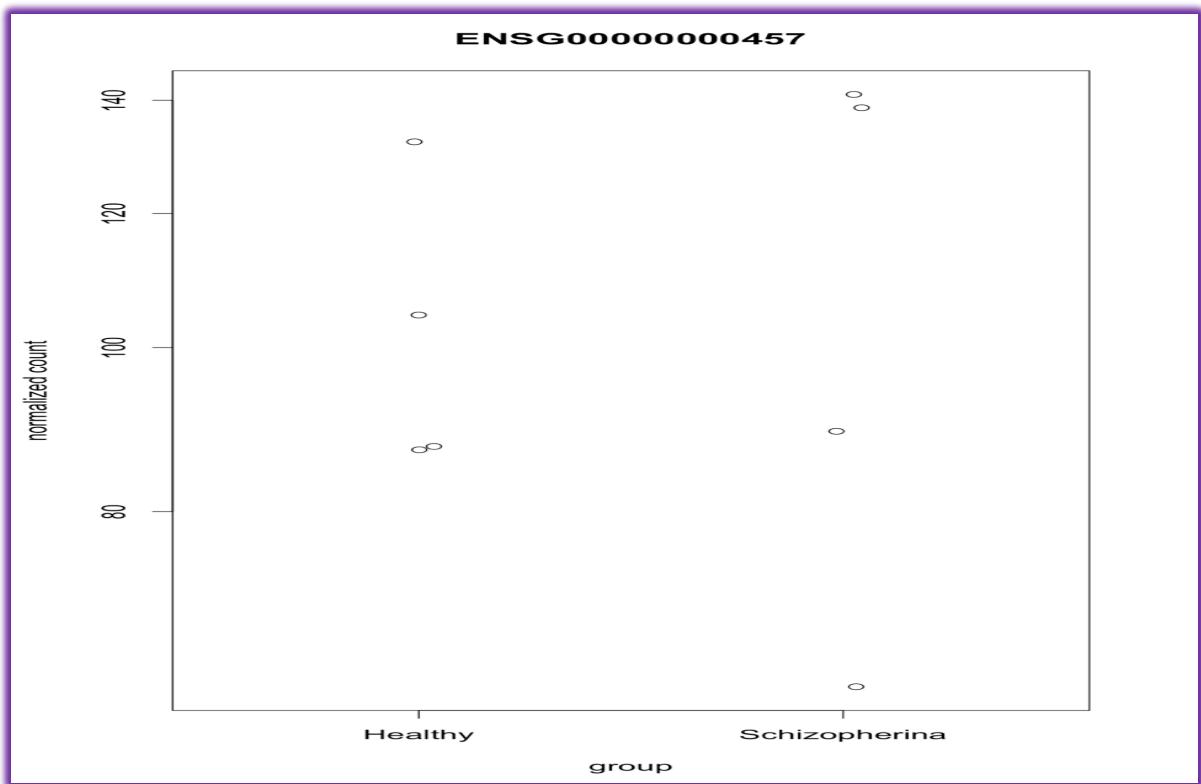
7) Heatmap of relative rlog-transformed values across samples.



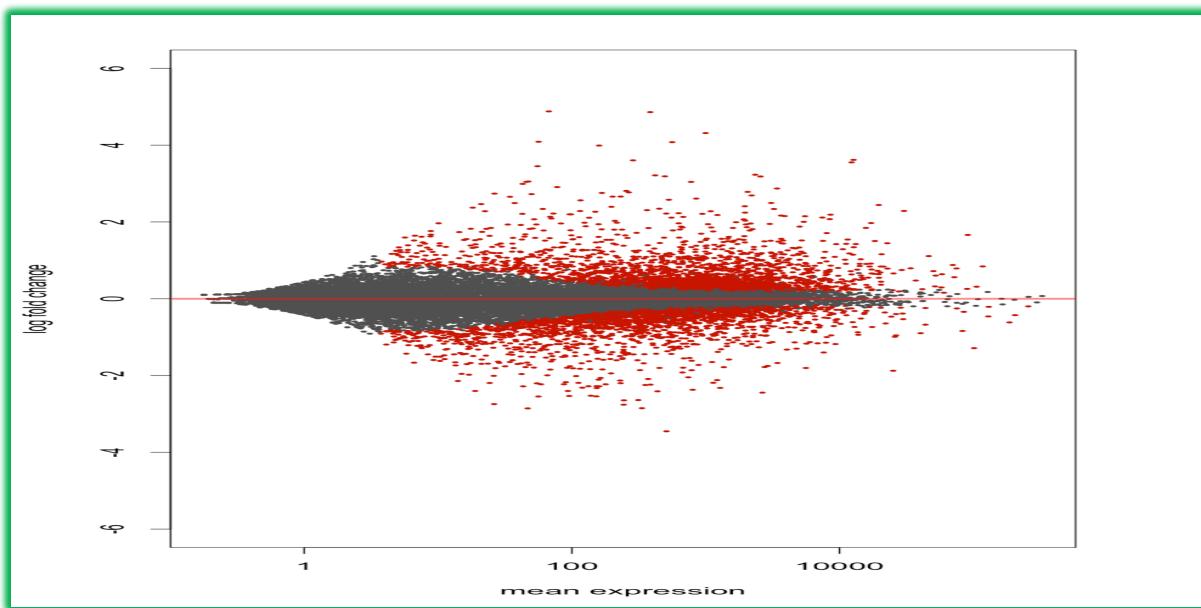
- 8) PCA plot using the rlog-transformed values. Each unique combination of Healthy and Schizophrenia sample is given its own colour.



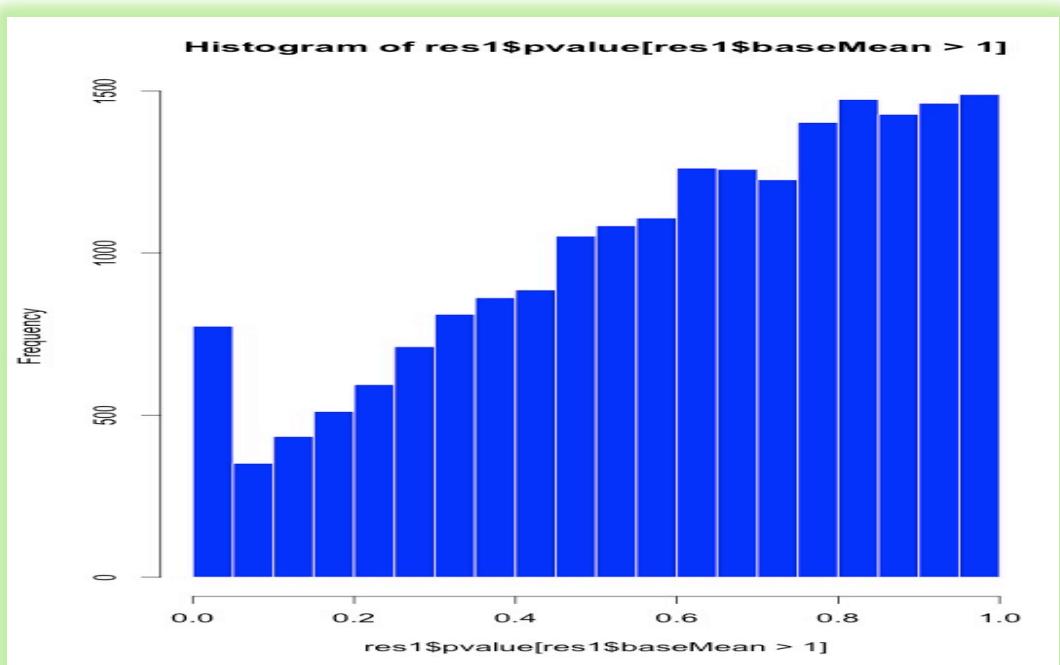
- 9) PCA plot (Plotting symbol for sample condition).



10) Normalized counts for a single gene over Schizophrenia group.



11) An MA-plot of a test for large log<sub>2</sub> fold changes. The red points indicate genes for which the log<sub>2</sub> fold change was significantly higher than 1 or less than -1



---

12) Histogram of p values for genes with mean normalized count larger than 1.

## Second Pipeline

### A) Build the reference index

```
./bowtie2-build <reference_in> <bt2_index_base>
SULTAN_BIN -- bash - 164x45
...q_Analysis/SULTAN_BIN -- bash +
sultanaharbi@SULTANS-iMac:~/Desktop/RNA_seq_Analysis/SULTAN_BIN$ ./bowtie2-build /Users/sultanaharbi/Desktop/Arabidopsis_thaliana.TAIR10.Genome.fa/Arabidopsis_thaliana.TAIR10.31.dna_rm.genome.fa Arabidopsis_thaliana.TAIR10.31
Settings:
  Output files: "Arabidopsis_thaliana.TAIR10.31.*.bt2"
  Line rate: 6 (line is 64 bytes)
  Lines per side: 1 (side is 64 bytes)
  Offset rate: 4 (one in 16)
  FTable chars: 10
  Strings: unpacked
  Max bucket size: default
  Max bucket size, sqrt multiplier: default
  Max bucket size, len divisor: 4
  Difference-cover sample period: 1024
  Endianness: little
  Actual local endianness: little
  Sanity checking: disabled
  Assertions: disabled
  Random seed: 0
```

### B) Using Tophat (Version 2.0.10)

#### Mapping RNA-seq Reads to a reference

```
./tophat2 -o basename_Th2output/ -p 8 <index_genome> reads1 reads2
SULTAN_BIN -- bash - 148x32
...analysis/SULTAN_BIN -- bash +
sultanaharbi@SULTANS-iMac:~/Desktop/RNA_seq_Analysis/SULTAN_BIN$ ./tophat2 -p 32 -o Tophat2output1SRR671949_Root_polyA_KN03_Treatment_rep2 /Users/sultanaharbi/Desktop/Arabidopsis_thaliana_Eensembl_TAIR10/Arabidopsis_thaliana_Eensembl_TAIR10/Sequence/Bowtie2Index/genome /Users/sultanaharbi/Desktop/RNA_seq_Reads_SRA/SRR671949_Root_polyA_KN03_Treatment_rep2.fastq
[2016-07-01 06:25:11] Beginning TopHat run (v2.0.10)
[2016-07-01 06:25:11] Checking for Bowtie
  Bowtie version: 2.1.0.0
[2016-07-01 06:25:11] Checking for Samtools
  Samtools version: 0.1.7.0
[2016-07-01 06:25:11] Checking for Bowtie index files (genome)..
[2016-07-01 06:25:11] Checking for reference FASTA file
[2016-07-01 06:25:11] Generating SAM header for /Users/sultanaharbi/Desktop/Arabidopsis_thaliana_Eensembl_TAIR10/Arabidopsis_thaliana_Eensembl_TAIR10/Sequence/Bowtie2Index/genome
[2016-07-01 06:25:11] Preparing reads
  left reads: min_length=50, max_length=50, 5177274 kept reads (5143 discarded)
[2016-07-01 06:25:33] Mapping left_kept_reads to genome genome with Bowtie2
[2016-07-01 06:25:33] Mapping left_kept_reads to genome genome with Bowtie2 (1/2)
[2016-07-01 06:39:24] Mapping left_kept_reads_seg2 to genome genome with Bowtie2 (2/2)
[2016-07-01 06:40:32] Searching for junctions via segment mapping
  Coverage search algorithm is turned on, making this step very slow
    Please turn off TopHat along with the option (--no-coverage-search) if this step takes too much time or memory.
[2016-07-01 06:44:38] Retrieving splices
[2016-07-01 06:44:42] Indexing splices
[2016-07-01 06:44:56] Mapping left_kept_reads_seg1 to genome segment_juncs with Bowtie2 (1/2)
[2016-07-01 06:44:56] Mapping left_kept_reads_seg2 to genome segment_juncs with Bowtie2 (2/2)
[2016-07-01 06:46:08] Joining segment hits
[2016-07-01 06:46:56] Reporting output tracks
[2016-07-01 06:55:34] A summary of the alignment counts can be found in Tophat2output1SRR671949_Root_polyA_KN03_Treatment_rep2/align_summary.txt
[2016-07-01 06:55:34] Run complete! 00:30:42 elapsed
sultanaharbi@SULTANS-iMac:~/Desktop/RNA_seq_Analysis/SULTAN_BIN$
```

### C) Using cufflinks to assemble transcripts

```
./cufflinks -p 32 -o files to this directory <accepted bam file>
SULTAN_BIN -- bash - 145x32
...analysis/SULTAN_BIN -- bash +
sultanaharbi@SULTANS-iMac:~/Desktop/RNA_seq_Analysis/SULTAN_BIN$ ./cufflinks -p 32 -o cufflinkscoutput_SRR671949_Root_polyA_KN03_Treatment_rep2 /Users/sultanaharbi/Desktop/RNA_seq_Analysis/SULTAN_BIN/Tophat2output1SRR671949_Root_polyA_KN03_Treatment_rep2/accepted_hits.bam
Warning: Your version of Cufflinks is not up-to-date. It is recommended that you upgrade to Cufflinks v2.2.1 to benefit from the most recent features and bug fixes (http://cufflinks.cbcb.umd.edu).
[05:25:01] Inspecting reads and determining fragment length distribution.
> Processed 54756 loci. [*****] 100%
> Map Properties:
>   Normalized Map Mass: 1766793.00
>   Raw Map Mass: 1766793.00
>   Fragment Length Distribution: Truncated Gaussian (default)
>     Default Mean: 200
>     Default Std Dev: 80
[05:26:02] Assembling transcripts and estimating abundances.
> Processed 54781 loci. [*****] 100%
```



## Using STAR (version 020201)

STAR (Spliced Transcripts Alignment to a reference) is a relatively new aligner and faster than tophat

### 1) Generating genome indexes

```
./star --runMode genomeGenerate --genomeDir storedDirectory --genomeFastaFiles genome.fa --sjdbGTFfile annotation.gtf - sjdbOverhang 74 --runThreadN 8
sultanalharbi@SULTAN-IMac:~/Volumes/SULTAN2016/STAR-2.5.2a/bin/MacOSX_x86_64$ ./star --runThreadN 32 --runMode genomeGenerate --genomeDir /Volumes/SULTAN2016/Drosophila_melanogaster_Eensembl_BDGP5.25/Drosophila_melanogaster/Ensembl/BDGP5.25/Sequence/WholeGenomeFasta --genomeFastaFiles /Volumes/SULTAN2016/Drosophila_melanogaster_Eensembl_BDGP5.25/Drosophila_melanogaster/Ensembl/BDGP5.25/Sequence/WholeGenomeFasta/genome.fa --sjdbGTFfile /Volumes/SULTAN2016/Drosophila_melanogaster_Eensembl_BDGP5.25/Drosophila_melanogaster/Ensembl/BDGP5.25/Annotation/Archives/archive-2015-07-17-14-30-26/Genes/genes.gtf
Jun 26 17:28:35 ..... Started STAR run
Jun 26 17:28:38 ... Starting to generate Genome files
Jun 26 17:28:38 ... starting to sort Suffix Array. This may take a long time...
Jun 26 17:28:38 ... sorting Suffix Array chunks and saving them to disk...
Jun 26 17:29:34 ... loading chunks from disk, packing SA...
Jun 26 17:29:39 ... Finished generating suffix array
Jun 26 17:29:39 ... Generating Suffix Array index
Jun 26 17:30:09 ... Completed Suffix Array index
Jun 26 17:30:09 .... Processing annotations GTF
Jun 26 17:30:11 .... Inserting junctions into the genome indices
Jun 26 17:30:33 .... writing Genome to disk ...
Jun 26 17:30:34 ... writing Suffix Array to disk ...
Jun 26 17:30:47 ... writing SAindex to disk
Jun 26 17:31:05 .... Finished successfully
sultanalharbi@SULTAN-IMac:~/Volumes/SULTAN2016/STAR-2.5.2a/bin/MacOSX_x86_64$
```

### 2) Mapping RNA-seq reads with a Reference

```
./star --genomeDir Directory to your STARGenomeIndex --runThreadN 8 --readFilesIn RNA-seqFASTQ --outFileNamePrefix <OutputPrefix>
```

```
sultanalharbi@SULTAN-IMac:~/Volumes/SULTAN2016/STAR-2.5.2a/bin/MacOSX_x86_64$ ./STAR --genomeDir /Volumes/SULTAN2016/Drosophila_melanogaster_Eensembl_BDGP5.25/Drosophila_melanogaster/Ensembl/BDGP5.25/Sequence/WholeGenomeFasta --runThreadN 24 --readFilesIn /Volumes/SULTAN2016/GSE32038_simulated_fastq_files/GSM794483_C1_R1_2.fq --outFileNamePrefix C1_R1_dataStar
Jun 26 17:47:30 ..... Started STAR run
Jun 26 17:47:30 ..... Loading genome
Jun 26 17:48:00 ..... Started mapping
Jun 26 17:50:11 ..... Finished successfully
sultanalharbi@SULTAN-IMac:~/Volumes/SULTAN2016/STAR-2.5.2a/bin/MacOSX_x86_64$
```

The output of STAR aligner is five different files:

- 1) Basename\_dataStarAligned.out.sam
- 2) Basename\_dataStarLog.final.out
- 3) Basename\_dataStarLog.out
- 4) Basename\_dataStarLog.progress.out
- 5) Basename\_dataStarSJ.out.tab

The Basename\_dataStar.final.out file provides useful mapping statistics

```

MacOSX_x86_64 — less /Volumes/SULTAN2016/STAR-2.5.2a/bin/MacOSX_x86_64/C1_R1_dataStarLog.final.out — 92x33
      Started job on | Jun 26 17:47:30
      Started mapping on | Jun 26 17:48:00
      Finished on | Jun 26 17:50:11
      Mapping speed, Million of reads per hour | 318.98

      Number of input reads | 11607353
      Average input read length | 150
      UNIQUE READS:
      Uniquely mapped reads number | 11516982
      Uniquely mapped reads % | 99.22%
      Average mapped length | 149.99
      Number of splices: Total | 3124770
      Number of splices: Annotated (sjdb) | 3124770
      Number of splices: GT/AG | 3099913
      Number of splices: GC/AG | 22959
      Number of splices: AT/AC | 1283
      Number of splices: Non-canonical | 615
      Mismatch rate per base, % | 0.00%
      Deletion rate per base | 0.00%
      Deletion average length | 0.00
      Insertion rate per base | 0.00%
      Insertion average length | 0.00

      MULTI-MAPPING READS:
      Number of reads mapped to multiple loci | 80311
      % of reads mapped to multiple loci | 0.69%
      Number of reads mapped to too many loci | 9795
      % of reads mapped to too many loci | 0.08%
      UNMAPPED READS:
      % of reads unmapped: too many mismatches | 0.00%
      % of reads unmapped: too short | 0.00%
      % of reads unmapped: other | 0.00%

      CHIMERIC READS:

```

Option	meaning	comment
--runMode	genomeGenerate	generate genome files
--genomeDir	path/to/your prefered/directory	string: path to the directory where genome files will be stored.
--genomeFastaFiles	The fasta files with genomic sequence	
--sjdbGTFfile	the GTF file with annotations	
--sjdbOverhang	int>0: length of the donor/acceptor sequence on each side of the junctions	ideally = (mate_length - 1)
--runThreadN	number of threads to run STAR	
--readFilesIn	RNA-seq reads (Fastq) file	
--outFileNamePrefix	output files name prefix	