# ANALYSIS OF RAW DATSETS AND DIFFERENTIAL EXPRESSION
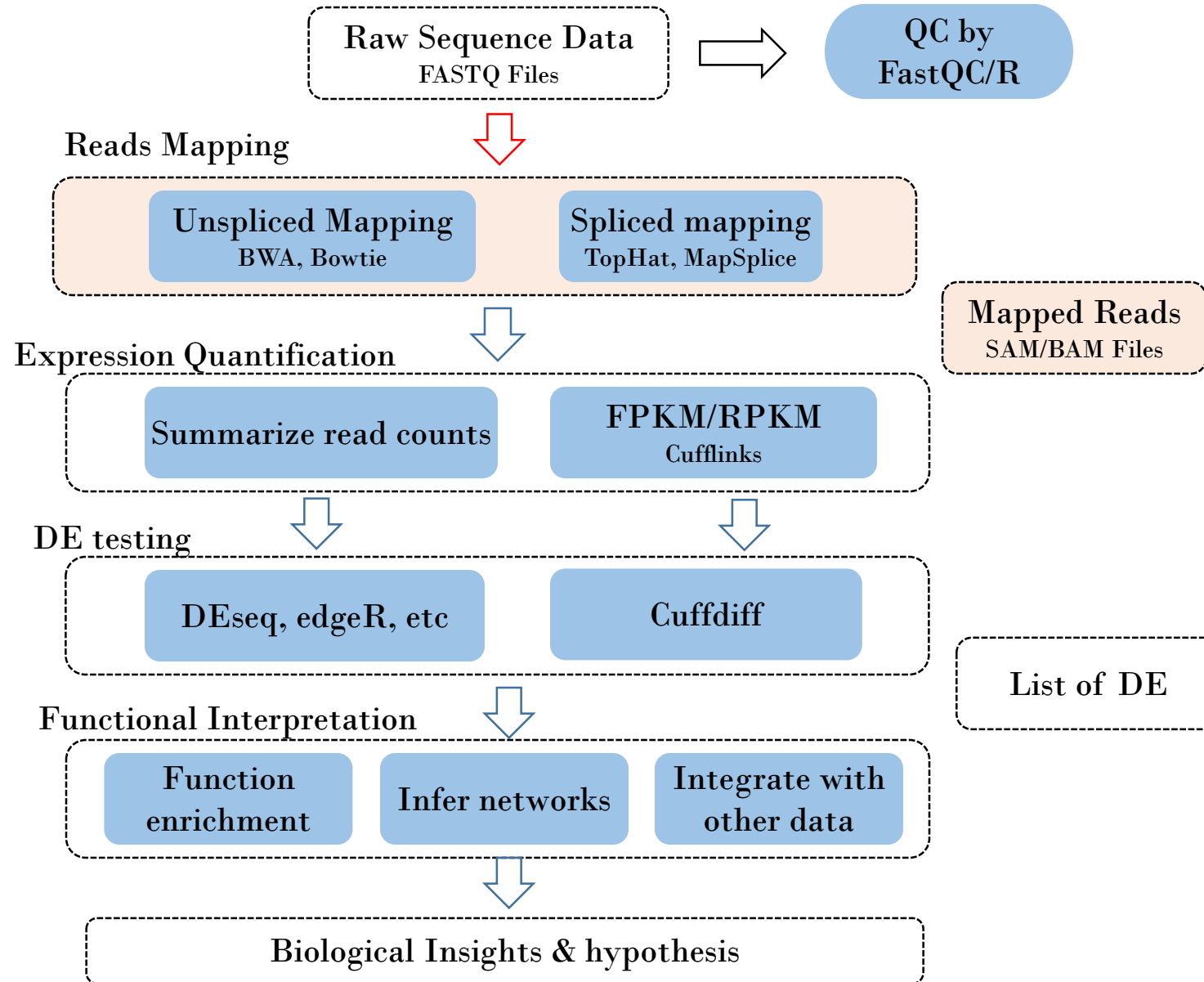
Presented by – Shikha Roy

Senior Research Fellow, ICGEB

# FASTQ format

The FASTQ format stores DNA sequence data as well as associated Phred quality scores of each base.

```
@EAS54_6_R1_2_1_413_324
CCCTTCTTGTCTTCAGCGTTTCTCC        ⬅  DNA read
+
;;3;;;;;;;;;;;7;;;;;;;88          ⬅  Base quality score
@EAS54_6_R1_2_1_540_792
TTGGCAGGCCAAGGCCGATGGATCA
+
;;;;;;;;;;;7;;;;;-;;;3;83
@EAS54_6_R1_2_1_443_348
GTTGCTTCTGGCGTGGGTGGGGGGG
+EAS54_6_R1_2_1_443_348
;;;;;;;;;;;9;7;;.7;393333
```
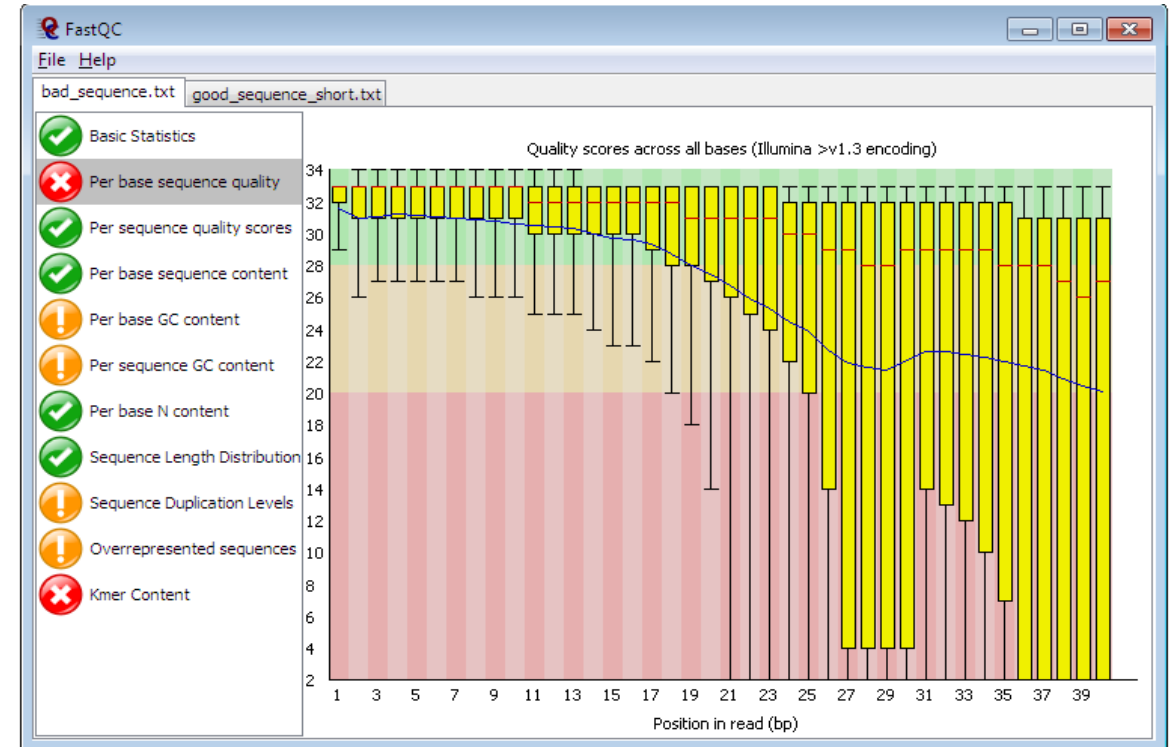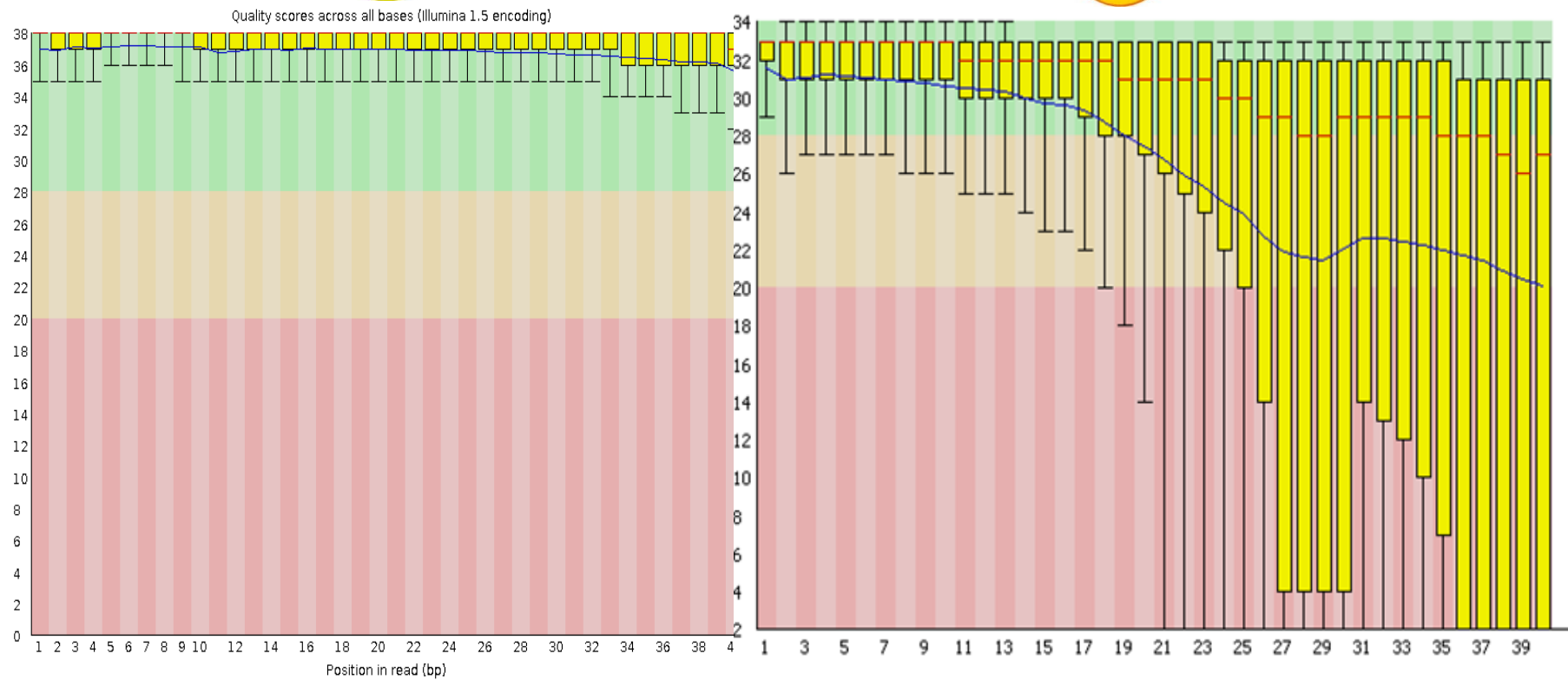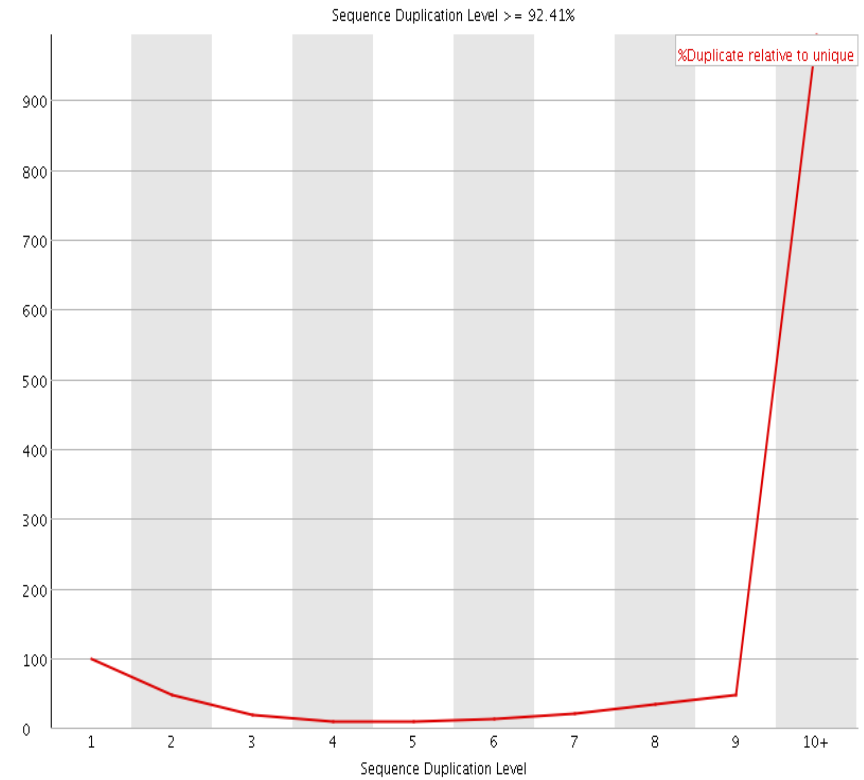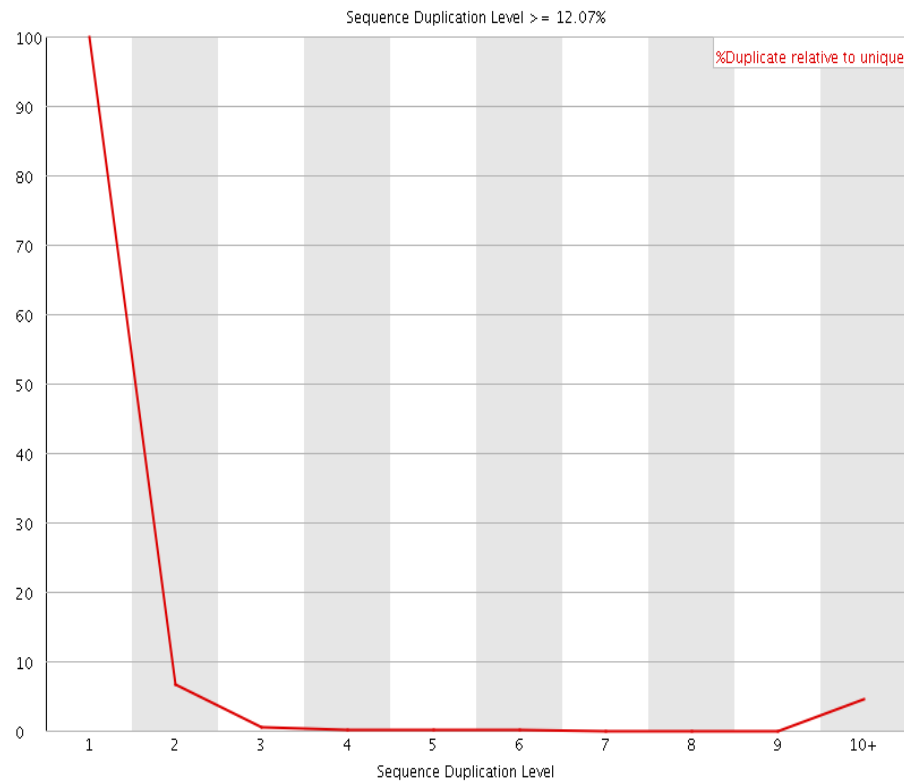
# FASTQC

- FastQC is a quality control application that allows users to perform numerous quality control checks on raw sequence data generated by high throughput sequencing pipelines such as Illumina and ABI SOLiD platforms in FASTQ format.

- It generates as output a comprehensive multi-page report on the composition and quality of reads in HTML format, with one page for each of the reads (e.g. Single End, Paired End: forward, Paired End: reverse). The modules included in the report are as follows:

➢ Per Base Sequence Quality

➢ Per Base Sequence Content

➢ Per Sequence GC content

➢ Per Base N Content

➢ Sequence Length Distribution

➢ Sequence Duplication Levels

➢ Adapter Content

# Per base sequence quality



Quality scores across all bases (Illumina 1.5 encoding)

# Duplication level

# Overrepresented Sequences

## Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---|---|---|
| GTGTCAGTCACTTCCAGCGGTCGTATGCCGTCTTCT | 2667259 | 7.236020826756234 | No Hit |
| TATCCCCGCCTGTCACGCGGGACGTGTCAGTCACTT | 703193 | 1.907695950497944 | No Hit |
| CTCGCTCCTCTCCTACTTGGATAACTGGTGTCAGTC | 352107 | 0.9552329133566171 | No Hit |
| TGTCAGTCACTTCCAGCGGTCGTATGCCGTCTTCTG | 351690 | 0.9541016318857297 | No Hit |
| CTCCTCTCCTACTTGGATAACTGGTGTCAGTCACTT | 247800 | 0.6722579100380558 | No Hit |
| CATCATATGGTGACCTCCCGGGTGTCAGTCACTTCC | 192614 | 0.5225435233416872 | No Hit |
| CATCAATATGGTGACCTCCCGGGTGTCAGTCACTTC | 192513 | 0.5222695199158848 | No Hit |
| CATCAATATGGTGACCTCCCGGAACGTGTCAGTCAC | 191604 | 0.5198034890836628 | No Hit |
| CATCAATATGGTGACCTCCCGGTGTCAGTCACTTCC | 163498 | 0.4435545753648186 | No Hit |
| CATCATATGGTGACCTCCCGGTGTCAGTCACTTCCA | 158547 | 0.43012298169008734 | No Hit |
| TATCCCCGCCTCACGCGGGACGTGTCAGTCACTTCC | 131347 | 0.3563319600878471 | No Hit |
| AAAACGTGTCAGTCACTTCCAGCGGTCGTATGCCGT | 127345 | 0.34547491345357634 | No Hit |
| CATGAGACTCTTAATCTCAGGTGTCAGTCACTTCCA | 109695 | 0.29759213656829914 | No Hit |

Adapter

# CUTADAPT

- Reads from small-RNA sequencing contain the 3' sequencing adapter because the read is longer than the molecule that is sequenced.

- Poly-A tails are useful for pulling out RNA from your sample, but often you don't want them to be in your reads.

- Cutadapt finds and removes adapter sequences, primers, poly-A tails and other types of unwanted sequence from your high-throughput sequencing reads.

- *sudo apt install cutadapt*

- *cutadapt –a adaptor sequence –o output.fastq input.fastq*

# BOWTIE

- Recent software tools allow the mapping (alignment) of millions or billions of short reads to a reference genome.

- For the human genome, this would take thousands of hours using BLAST.

- Indexing a genome can be explained similar to indexing a book. If you want to know on which page a certain word appears or a chapter begins, it is much more efficient/faster to look it up in a pre-built index than going through every page of the book until you found it.

*bowtie-build ~/hg38.fa hg38*

*bowtie -t hg38 -S ~/fastq/wt_H3K4me3_read1.fastq res.sam*

```
shikha@BIOINFO:~$ bowtie
bowtie            bowtie2-align-l     bowtie2-build       bowt
bowtie2           bowtie2-align-s     bowtie2-build-l     bowt
shikha@BIOINFO:~$ bowtie-build hg
hg19.gff              hg38.chrom.sizes   hg38.fa
shikha@BIOINFO:~$ bowtie-build hg38.fa  hg38
Settings:
  Output files: "hg38.*.ebwt"
  Line rate: 6 (line is 64 bytes)
  Lines per side: 1 (side is 64 bytes)
  Offset rate: 5 (one in 32)
  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
  Sizeofs: void*:8, int:4, long:8, size_t:8
Input files DNA, FASTA:
  hg38.fa
Reading reference sizes
```

```
shikha@BIOINFO:~$ bowtie -t ~/RNAseq/hg38 -S ~/fastq/wt_H3K4me3_read1.fastq res.sam
Time loading forward index: 00:00:08
Time loading mirror index: 00:00:08
Seeded quality full-index search: 00:00:07
# reads processed: 50000
# reads with at least one reported alignment: 1354 (2.71%)
# reads that failed to align: 48646 (97.29%)
Reported 1354 alignments
Time searching: 00:00:23
Overall time: 00:00:23
```

# Alignment to a reference genome: example of short-read alignment (Bowtie) results

# SAMTOOLS

✧ **SAM – Sequence Alignment/Map format**

    ✧ SAM file format stores alignment information

✧ **Plain text**

✧ **Specification**:
http://samtools.sourceforge.net/SAM1.pdf

✧ Contains quality information, meta data, alignment
information, sequence etc.

✧ **Files can be very large:** Many 100's of GB or more

✧ Normally converted into **BAM** to save space (and text
format is mostly useless for downstream analyses)

SAM is a common format having sequence reads and
their alignment to a reference genome.

BAM is the binary form of a SAM file.

SAMTools is a software package commonly used to
analyze SAM/BAM files.

*samtools view -bS -o res.bam res.sam*

# Formats : **BAM**

✧ **BAM – BGZF compressed SAM format**

 ✧ Compressed/binary version of SAM and is **not human readable.** Uses a specialized compression algorithm optimized for indexing and record retrieval (bgzip)

 ✧ Makes the alignment information easily accessible to downstream applications (large genome file not necessary)

 ✧ Unsorted, sorted by sequence name, **sorted by genome coordinates**

 ✧ May be accompanied by an index file (.bai) (only if coordinate sorted)

✧ **Files are typically very large:** ~ 1/5 of SAM, but still very large

# CUFFLINKS

- Cufflinks assembles transcripts, estimates their abundances, and tests for differential expression and regulation in RNA-Seq samples.

- It accepts aligned RNA-Seq reads and assembles the alignments into a parsimonious set of transcripts.

- Cufflinks then estimates the relative abundances of these transcripts based on how many reads support each one, taking into account biases in library preparation protocols.

- Output tracks of Cufflinks is the Assembled transcripts track, output tables of Cufflinks are Gene expression and Transcript expression tables.

*cufflinks testA.bam -g Homo_sapiens.GRCh37.63.gtf/data da -o cuff_res*

# Metrics for quantifying gene expression levels

- RPKM
  - **R**eads **P**er **K**ilobase per **M**illion mapped reads
  - Normalize relative to sequencing depth and gene length

- FPKM
  - Similar to RPKM but count **DNA fragments** instead of reads
  - Used in paired end RNA-Seq experiments to avoid bias

- TPM
  - **T**ranscripts **P**er **M**illion
  - Normalize for gene length, then normalize by sequencing depth

# DIFFERENTIAL EXPRESSION USING RSTUDIO

```
rawCountTable <- read.table("countData.txt", header=TRUE, sep="\t", row.names=1)
sampleInfo <- read.table("design.csv", header=TRUE, sep=",", row.names=1)
```

|               | Cond.WT.Rep.1 | Cond.WT.Rep.2 | Cond.WT.Rep.3 | Cond.Mt.Rep.1 |
|---------------|---------------|---------------|---------------|---------------|
| Solyc00g005000.2.1 | 0 | 0 | 0 | 0 |
| Solyc00g005020.1.1 | 0 | 0 | 0 | 0 |
| Solyc00g005040.2.1 | 0 | 0 | 0 | 0 |
| Solyc00g005050.2.1 | 306 | 502 | 468 | 369 |
| Solyc00g005060.1.1 | 0 | 0 | 0 | 0 |
| Solyc00g005070.1.1 | 0 | 0 | 0 | 0 |

|               | Cond.Mt.Rep.2 | Cond.Mt.Rep.3 |
|---------------|---------------|---------------|
| Solyc00g005000.2.1 | 0 | 0 |
| Solyc00g005020.1.1 | 0 | 0 |
| Solyc00g005040.2.1 | 0 | 0 |
| Solyc00g005050.2.1 | 366 | 294 |
| Solyc00g005060.1.1 | 0 | 0 |
| Solyc00g005070.1.1 | 0 | 0 |

| files | condition |
|-------|-----------|
| Cond.WT.Rep.1 | WT |
| Cond.WT.Rep.2 | WT |
| ... | ... |
| Cond.Mt.Rep.1 | M |
| ... | ... |

Save this file under the name design.csv (csv format) inside your working directory. In my case, this file is separated by commas, as in the following picture:

```
file,condition
Cond.WT.Rep.1,WT
Cond.WT.Rep.2,WT
Cond.WT.Rep.3,WT
Cond.Mt.Rep.1,M
Cond.Mt.Rep.2,M
Cond.Mt.Rep.3,M
```

Create a DGEList data object
dgeFull <- DGEList(rawCountTable,
group=sampleInfo$condition)
pseudoCounts <- log2(dgeFull$counts+1)
boxplot(pseudoCounts, col="gray", las=3)

estimate the normalization factors
dgeFull <- calcNormFactors(dgeFull, method="TMM")
eff.lib.size <-
dgeFull$samples$lib.size*dgeFull$samples$norm.factors
normCounts <- cpm(dgeFull)
pseudoNormCounts <- log2(normCounts + 1)
boxplot(pseudoNormCounts, col="gray", las=3)

# Differential Gene Expression overview

④ Set up to do differential gene expression (DGE)

*Identify read counts associated with genes*

   a. Do you want to obtain raw read counts or normalized read counts? This will depend on the statistical analysis you wish to perform downstream

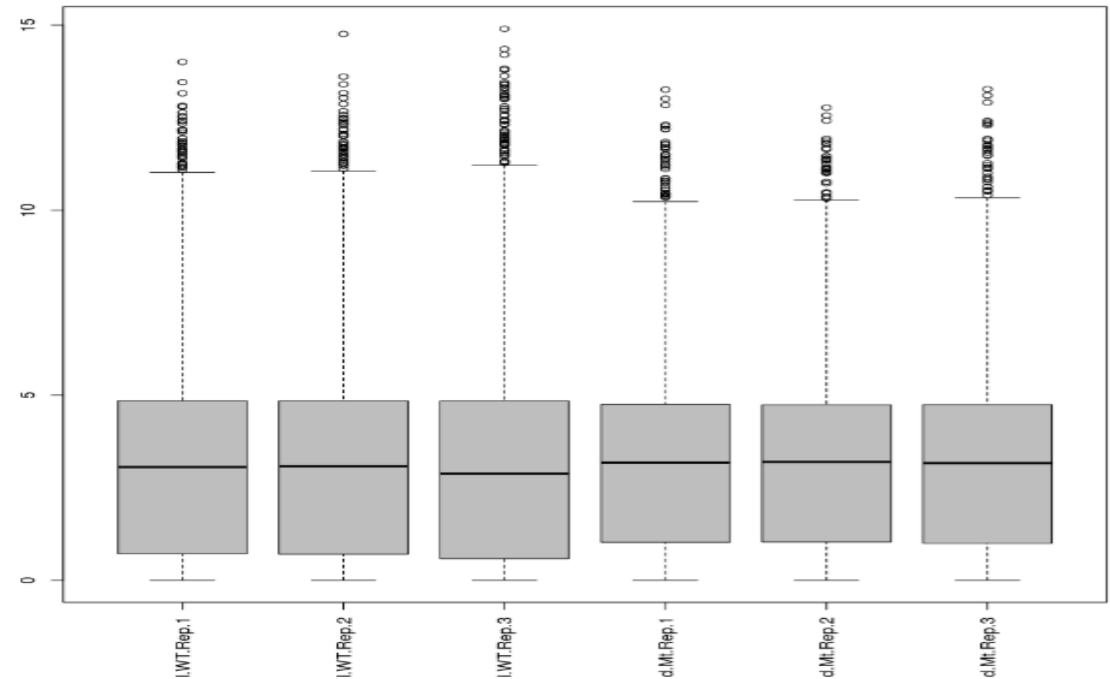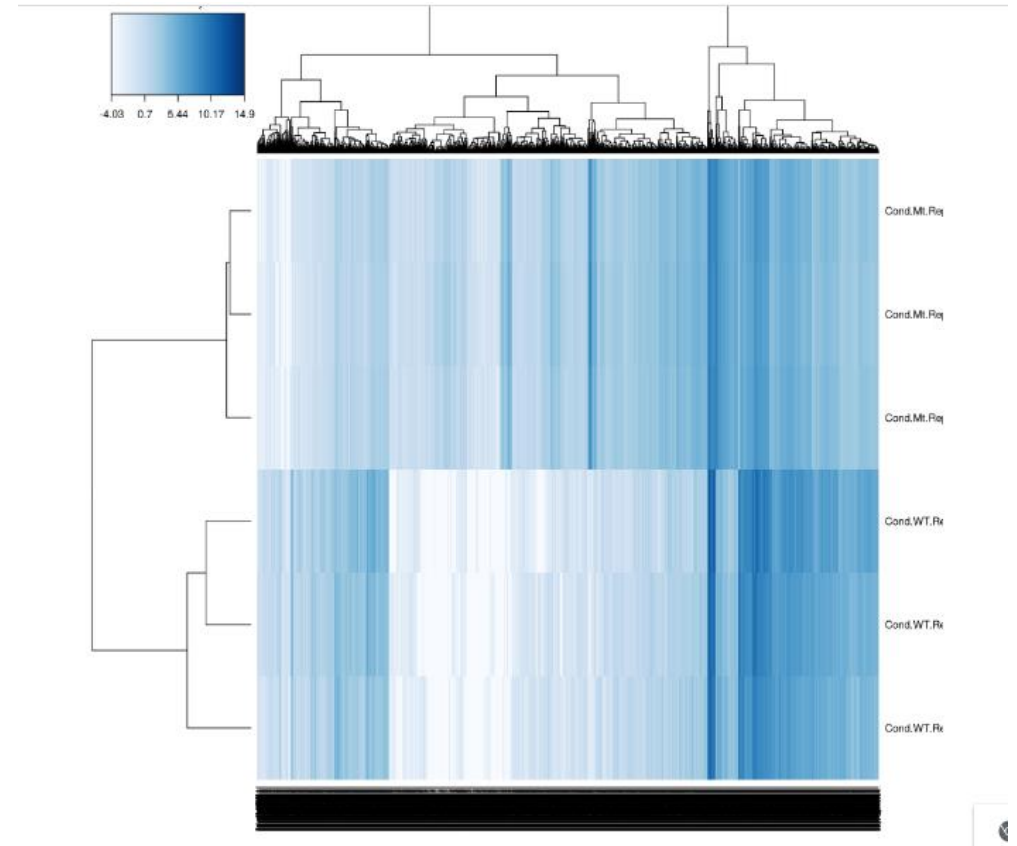- ✧ htseq & feature-counts return raw read counts
  - ✧ Required for R programs like DESeq & EdgeR
- ✧ Ballgown & Cufflinks return FPKM normalized counts for each gene

*dgeFull <- DGEList(dgeFull$counts [apply(dgeFull$counts, 1, sum) != 0, ],group=dgeFull$samples$group)*
*dgeFull <- estimateCommonDisp(dgeFull)*
*dgeFull <- estimateTagwiseDisp(dgeFull)*
*dgeTest <- exactTest(dgeFull)*
*remove low expressed genes*
*filtData <- HTSFilter(dgeFull)filteredData*
*dgeTestFilt <- exactTest(filtData)*
*resFilt <- topTags(dgeTestFilt, n=nrow(dgeTest$table))*
*sigReg <- resFilt$table(resFilt$table$FDR<0.01,]*
*sigReg <- resFilt$table[order(sigReg$logFC),]*

```
plotSmear(dgeTestFilt, de.tags = rownames(resFilt$table)[which(resFilt$table$FDR<0.01)])
```



```
selY <- y[rownames(resFilt$table)[resFilt$table$FDR<0.01 & abs(resFilt$table$logFC)>1.5],]

cimColor <- colorRampPalette(rev(brewer.pal(9, "Blues")))(255)[255:1] finalHM <- cim(t(selY),
 color=cimColor, symkey=FALSE)
```

# Tools for analyzing differentially expressed genes

- Gene Ontology (GO) terms enrichment:
  - topGO (https://bioconductor.org/packages/release/bioc/html/topGO.html)
  - goSTAG (https://bioconductor.org/packages/release/bioc/html/goSTAG.html)
  - DAVID (https://david.ncifcrf.gov/)
- Pathway analysis:
  - GAGE (http://bioconductor.org/packages/release/bioc/html/gage.html)
  - Reactome (http://www.reactome.org/)
- Sample walkthrough:
  - From reads to genes to pathways: differential expression analysis of RNA-Seq experiments using Rsubread and the edgeR quasi-likelihood pipeline
    - https://www.bioconductor.org/help/workflows/RnaSeqGeneEdgeRQL/

# GENE ONTOLOGY ENRICHMENT USING DAVID



- ➢ Functional Annotation Tool
  - Gene Ontology
  - Protein interaction
  - Protein domain
  - Pathway
  - Disease
- ➢ Gene ID Conversion
- ➢ Gene Functional Classification

# RT (Related Term)

Any given gene is associating with a set of annotation terms. If genes share similar set of those terms, they are most likely involved in similar biological mechanisms. The algorithm adopts kappa statistics to quantitatively measure the degree of the agreement how genes share the similar annotation terms. Kappa result ranges from 0 to 1. The higher the value of Kappa, the stronger the agreement.

Any a biological process/term coming from all functional categories listed in DAVID.

## Functional Related Terms

**Options**

Similarity Score(Kappa)>= 0.3          Overlap>= 2

Rerun using options

6622 term(s) were searched. 143 term(s) passed the filter.          Download File

Similarity Score:   ■ Very High (0.75-1)   ■ High (0.5-0.75)   ■ Moderate (0.25-0.5)   □ Low (<0.25)

| # | Category | Term | Kappa |
|---|---|---|---|
| 1 | BIOCARTA | Cytokine Network | 1.00 |
| 2 | KEGG_PATHWAY | Allograft rejection | 0.86 |
| 3 | BIOCARTA | Selective expression of chemokine receptors during T-cell polarization | 0.86 |
| 4 | BIOCARTA | Cytokines and Inflammatory Response | 0.86 |
| 5 | SP_PIR_KEYWORDS | lymphokine | 0.80 |
| 6 | BIOCARTA | Th1/Th2 Differentiation | 0.80 |
| 7 | BIOCARTA | IL 5 Signaling Pathway | 0.80 |
| 8 | GOTERM_BP_FAT | regulation of activated T cell proliferation | 0.80 |
| 9 | GOTERM_BP_FAT | positive regulation of activated T cell proliferation | 0.80 |
| 10 | INTERPRO | Four-helical cytokine, core | 0.67 |
| 11 | KEGG_PATHWAY | Asthma | 0.67 |
| 12 | KEGG_PATHWAY | Intestinal immune network for IgA production | 0.67 |
| 13 | BIOCARTA | GATA3 participate in activating the Th2 cytokine genes expression | 0.67 |
| 14 | GOTERM_BP_FAT | positive regulation of gene-specific transcription | 0.67 |
| 15 | SP_PIR_KEYWORDS | T-cell | 0.57 |
| 16 | BIOCARTA | Regulation of hematopoiesis by cytokines | 0.57 |
| 17 | GOTERM_BP_FAT | positive regulation of peptidyl-tyrosine phosphorylation | 0.57 |
| 18 | GOTERM_BP_FAT | regulation of gene-specific transcription | 0.57 |
| 19 | UP_SEQ_FEATURE | chain:Interleukin-5 | 0.50 |
| 20 | UP_SEQ_FEATURE | chain:Interleukin-4 | 0.50 |
| 21 | UP_SEQ_FEATURE | chain:Interleukin-12 subunit beta | 0.50 |