



**ORT** OFFICE OF  
RESEARCH  
TRAINEES

Presents:

Translating Bioinformatics  
for Everyday Biology

# RNA-seq analysis

Musa Ahmed

Nov 15, 2019

# GOAL

- Basics of RNA-seq analysis
- Applications
- Challenges
- Practical
  - Alignment
  - DGE analysis

# What is RNA-seq

- RNA-seq works by sequencing every RNA molecule and profiling the expression of a particular gene by counting the number of time its transcripts have been sequenced.

# RNA sequencing

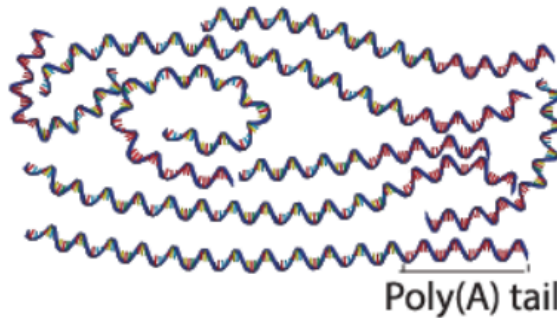
Samples of interest



Condition 1  
(e.g. tumor)

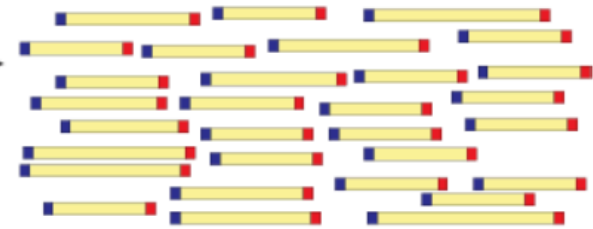
Condition 2  
(e.g. normal)

Isolate RNAs

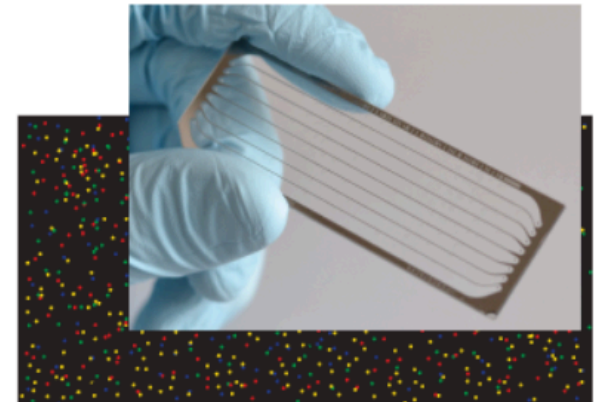


Poly(A) tail

Generate cDNA, Fragment, size select, add linkers

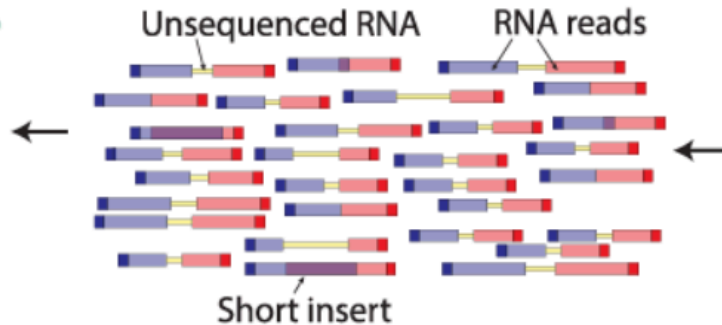
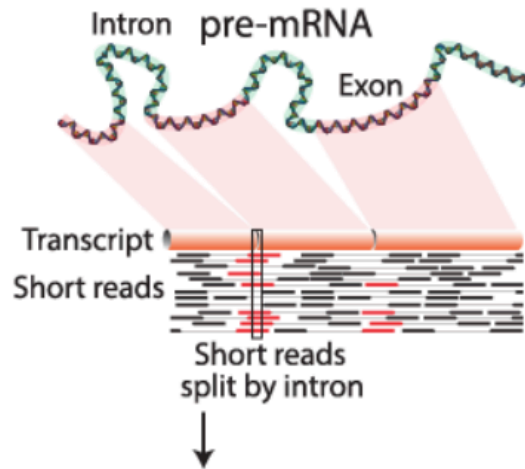


Sequence ends

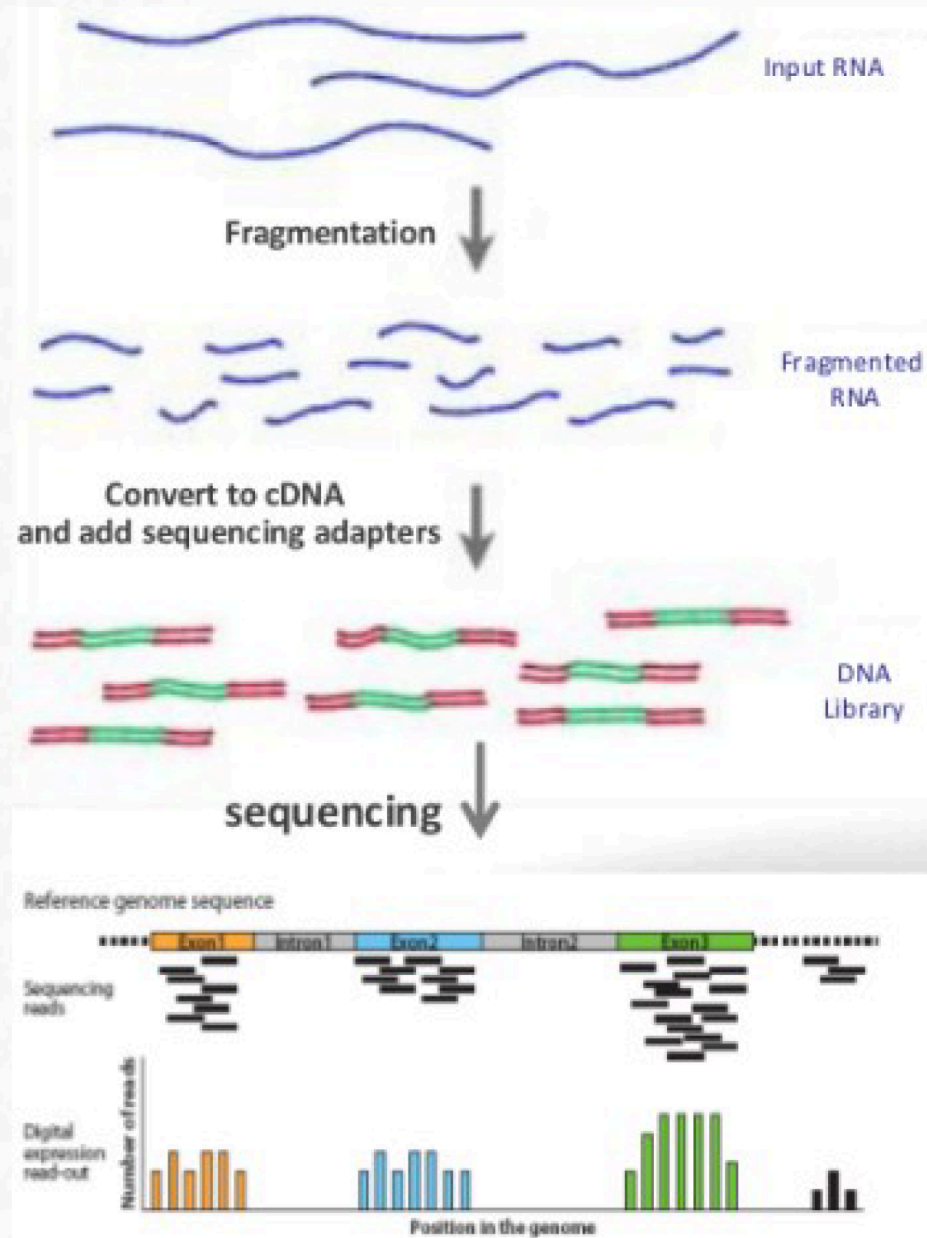
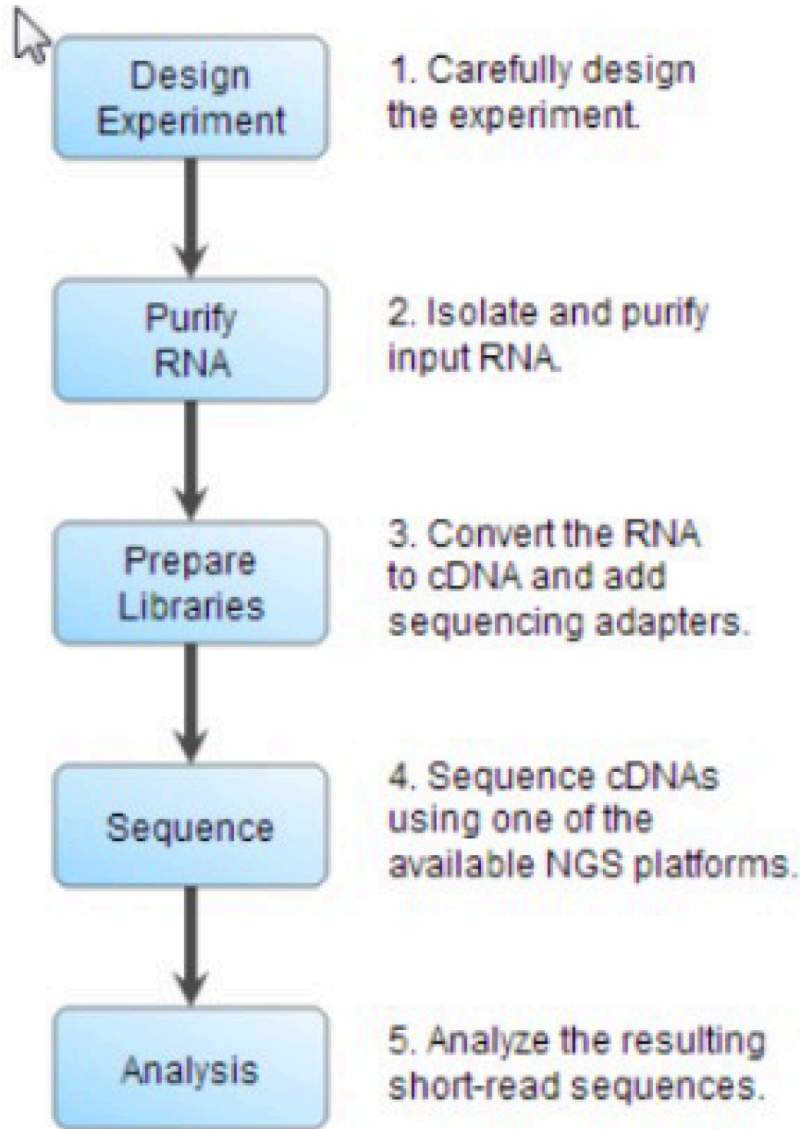


100s of millions of paired reads  
10s of billions bases of sequence

Map to genome, transcriptome, and predicted exon junctions



Downstream analysis

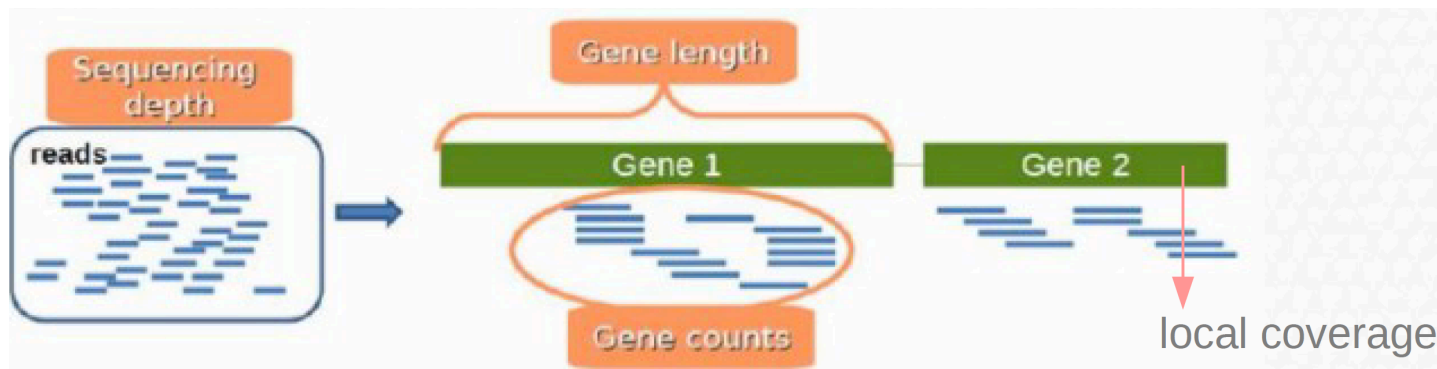


# Why RNA-seq?

- Gene expression
- Splicing/isoforms
- de-novo assembly of transcriptome
- Allele-specific expression

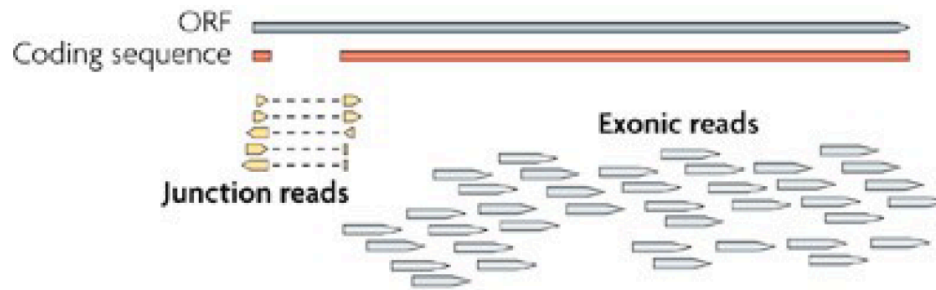
# Key concepts

- **Sequencing depth**
  - Total number of reads mapped to the genome. (Library size) Could also be applied to samples.
- **Coverage**
  - Number of reads mapped to a specific region (average of them if we are talking about the whole genome...) . Not typically used for RNA-seq
- **Gene length**
  - Number of bases that a gene has.



# Key concepts

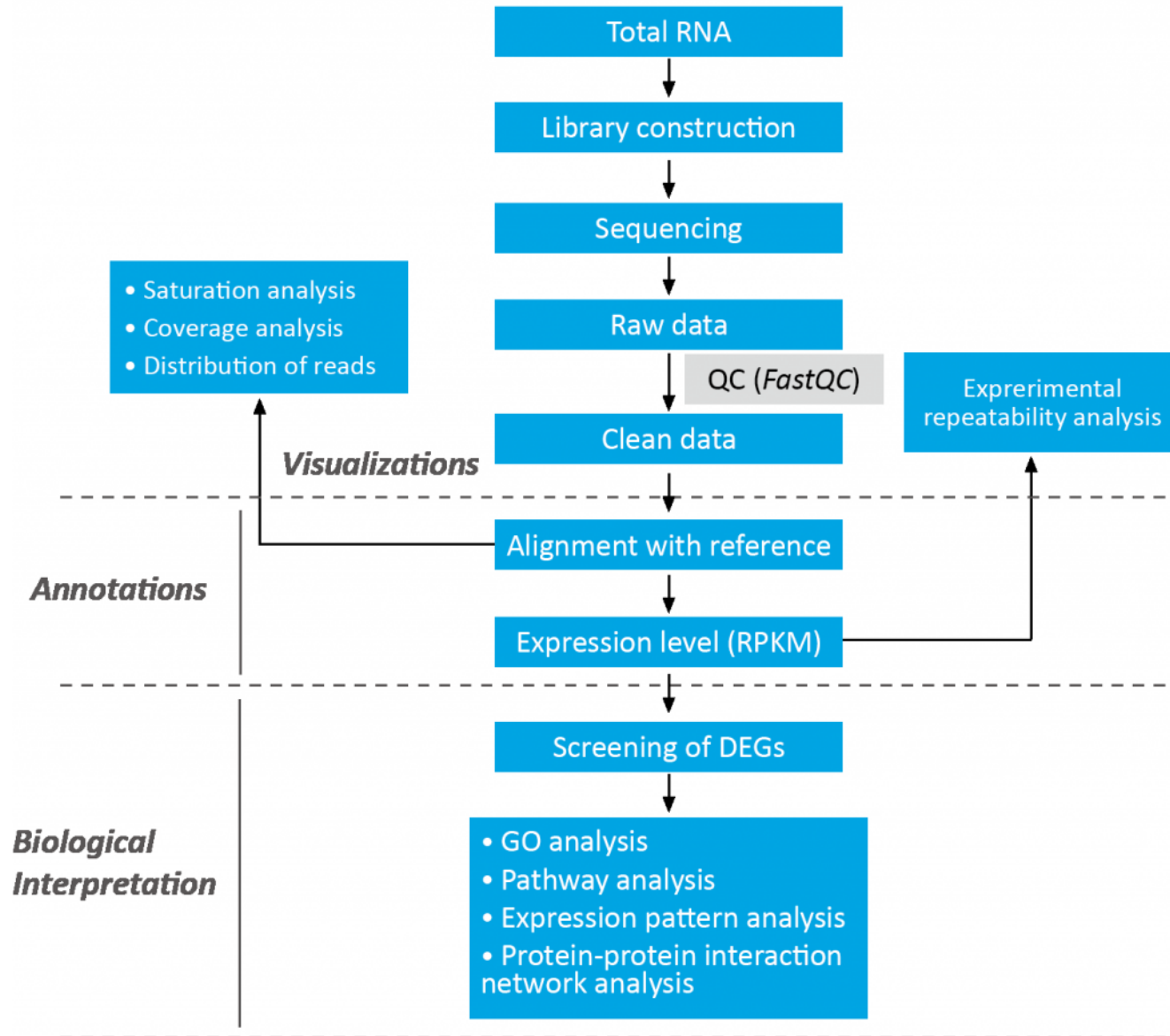
- **Exonic reads:** Reads within exons
- **Junction reads:** Reads spanning exon junctions

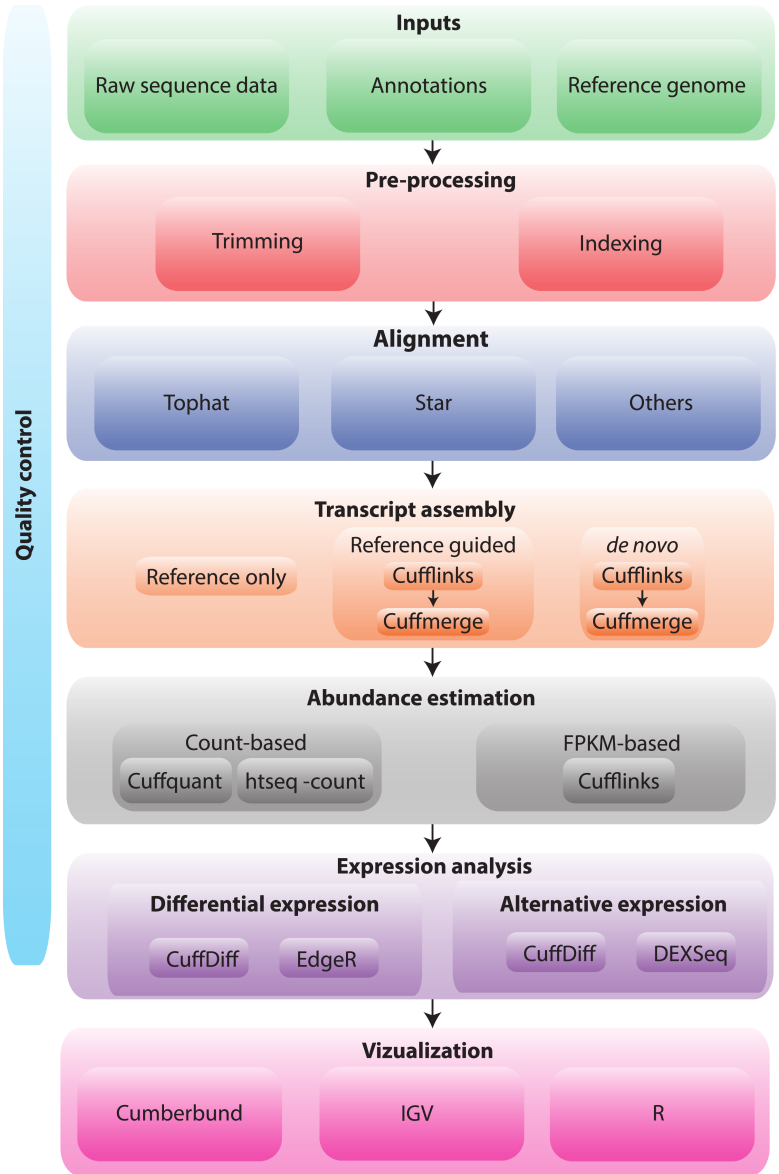




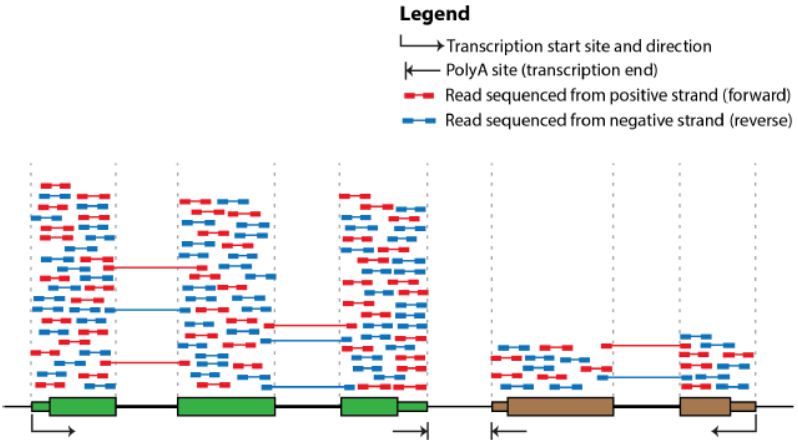
# BGI workflow

## Technique Workflow

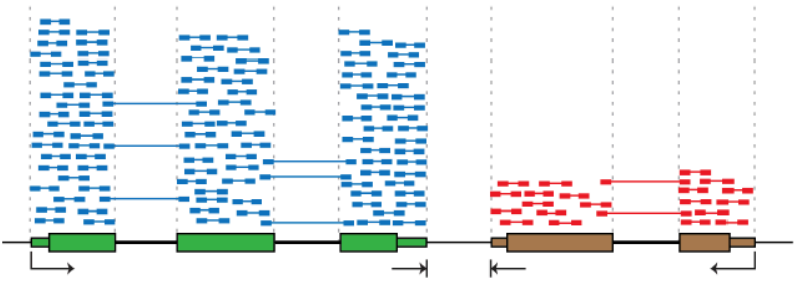




**A. Depiction of cDNA fragments from an unstranded library**

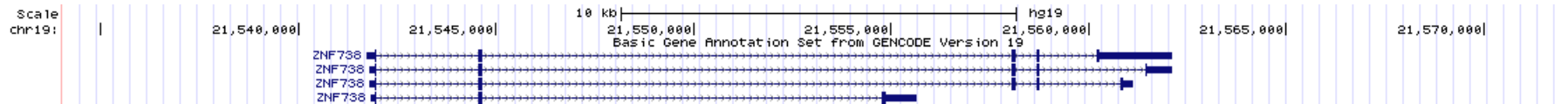


**B. Depiction of cDNA fragments from a stranded library**



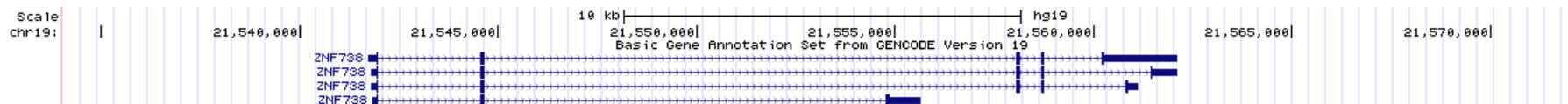
# Reference genomes by Genome Reference Consortium

- Assembled whole genome using a bunch of individual genome sequences
  - Human – Hg38 (GRCh38), hg19 (GRCh37), b37 etc
  - Mouse – mm10, mm9 etc



# Gene annotations

- Annotation of the whole genome (protein coding genes, noncoding RNAs etc)
  - RefSeq (good for general analysis)
  - GENCODE/ENSEMBL (good for noncoding genes)
  - miTranscriptome (good for noncoding genes, suitable for really deep sequencing)



# Annotations can be downloaded from UCSC genome browser

**Table Browser**

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Send data to [GenomeSpace](#) for use with diverse computational tools. Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

**clade:** Mammal   
**genome:** Human   
**assembly:** Feb. 2009 (GRCh37/hg19)

**group:** Genes and Gene Predictions   
**track:** GENCODE Genes V19

**table:** Basic (wgEncodeGencodeBasicV19)

**region:**  genome  ENCODE Pilot regions  position chr21:33031597-33041570

**identifiers (names/accessions):**

**filter:**

**subtrack merge:**

**intersection:**

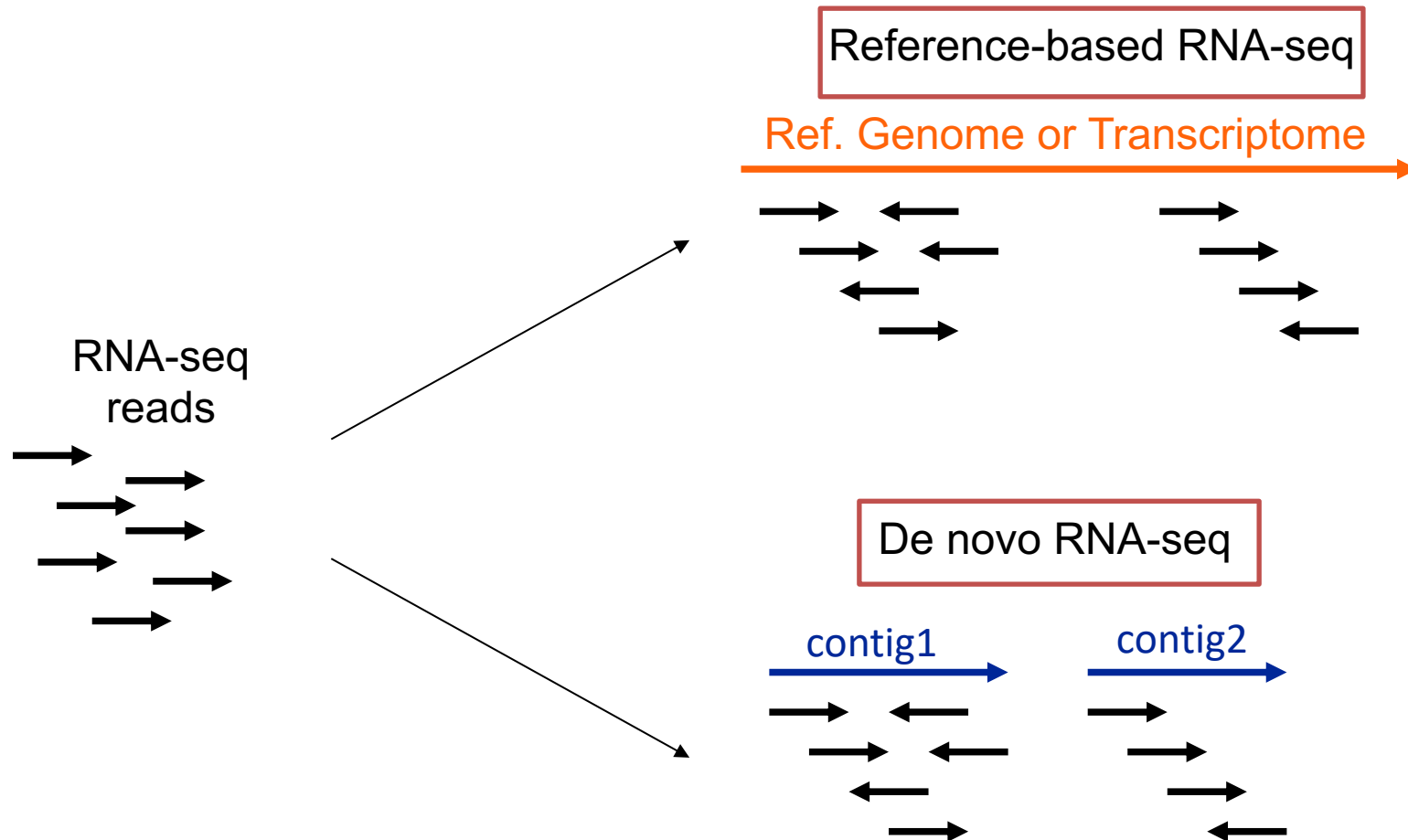
**correlation:**

**output format:** GTF - gene transfer format  Send output to  [Galaxy](#)  [GREAT](#)   
[GenomeSpace](#)

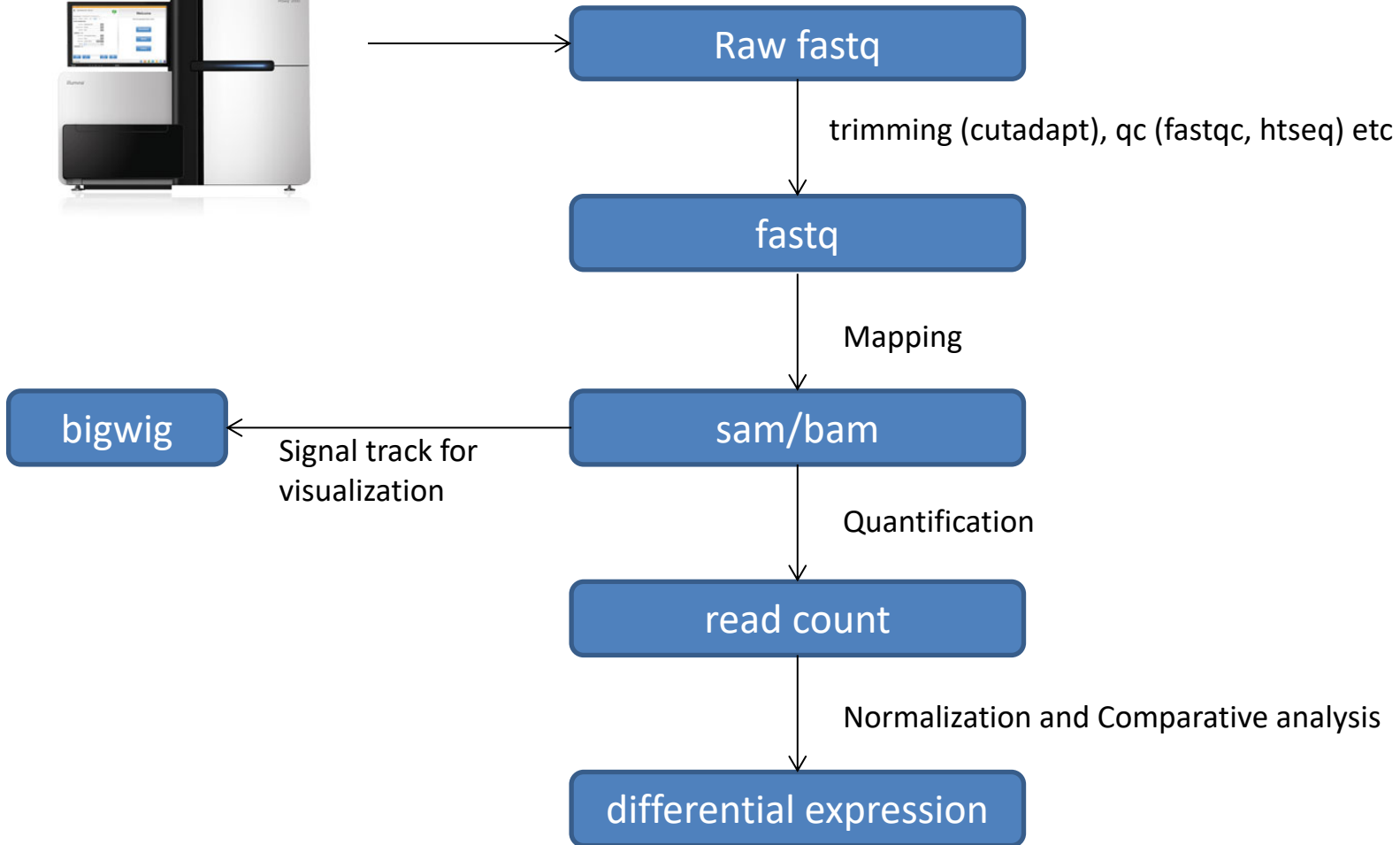
**output file:** genc19.gtf  (leave blank to keep output in browser)

**file type returned:**  plain text  gzip compressed

# RNA-seq alignment can be annotation-dependent or de novo



# RNA-seq analysis: overview



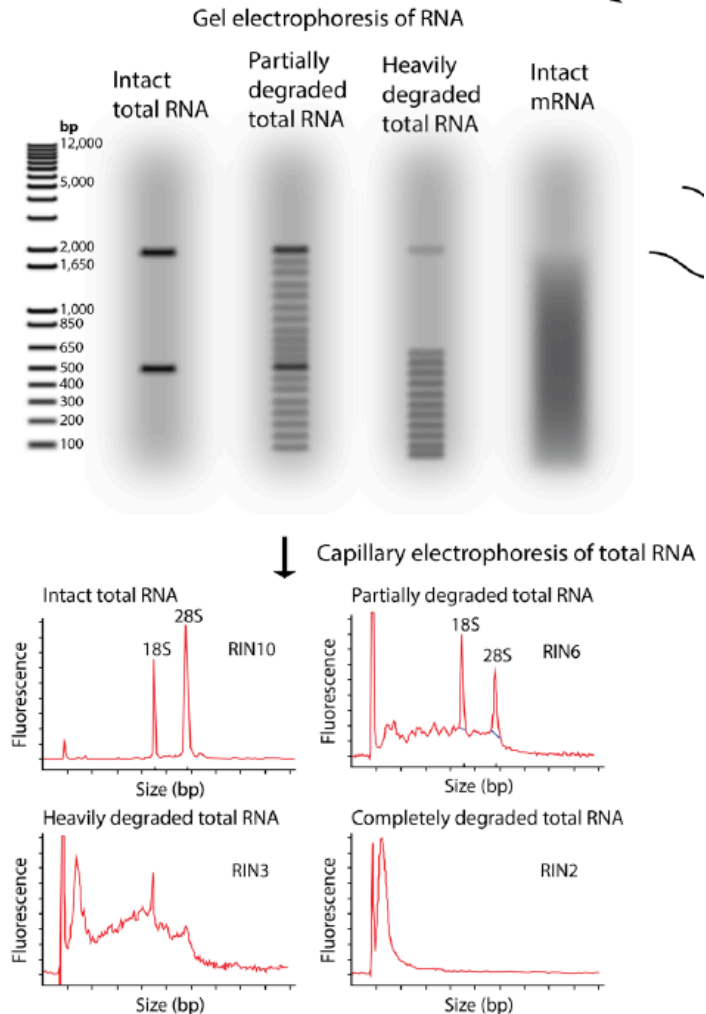
# RNA-seq questions during library preparation

- Construction strategies
  - total RNA or polyA+ RNA?
  - Ribo minus?
  - Stranded or unstranded?
  - Read length?
  - Single end vs paired end?
  - size selection – microRNA?
  - Depth?
- RNA quantity
- RNA quality
  - RNA is fragile and easily degraded
  - Low quality material can bias the data
- Replicates

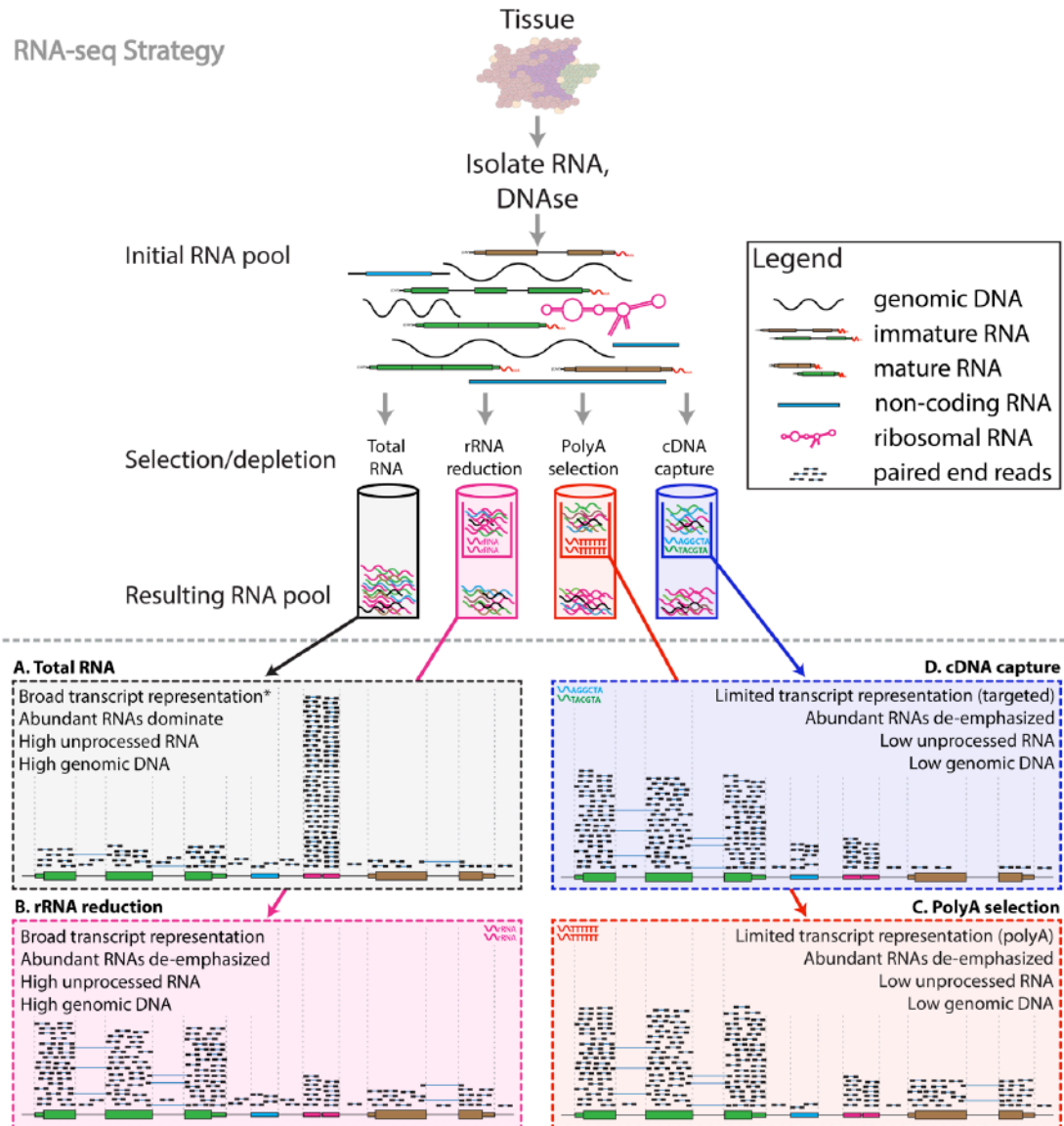


# Agilent

- [https://github.com/griffithlab/rnaseq\\_tutorial/wiki/Resources/Agilent\\_Trace\\_Examples.pdf](https://github.com/griffithlab/rnaseq_tutorial/wiki/Resources/Agilent_Trace_Examples.pdf)
- ‘RIN’ = RNA integrity number
  - 0 (bad) to 10 (good)



# Strategies

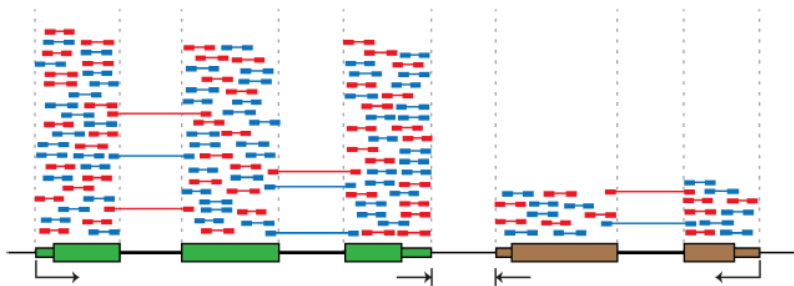


# Stranded vs unstranded

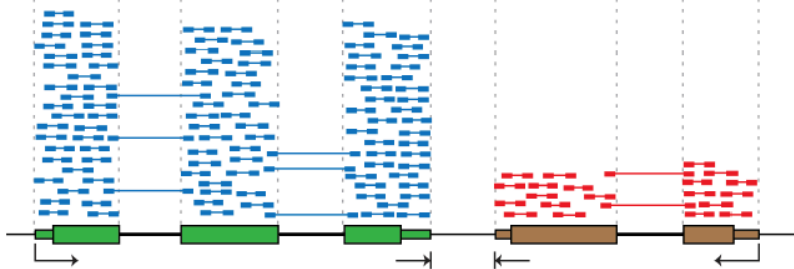
## A. Depiction of cDNA fragments from an unstranded library

### Legend

- ↳ Transcription start site and direction
- └ PolyA site (transcription end)
- Read sequenced from positive strand (forward)
- Read sequenced from negative strand (reverse)



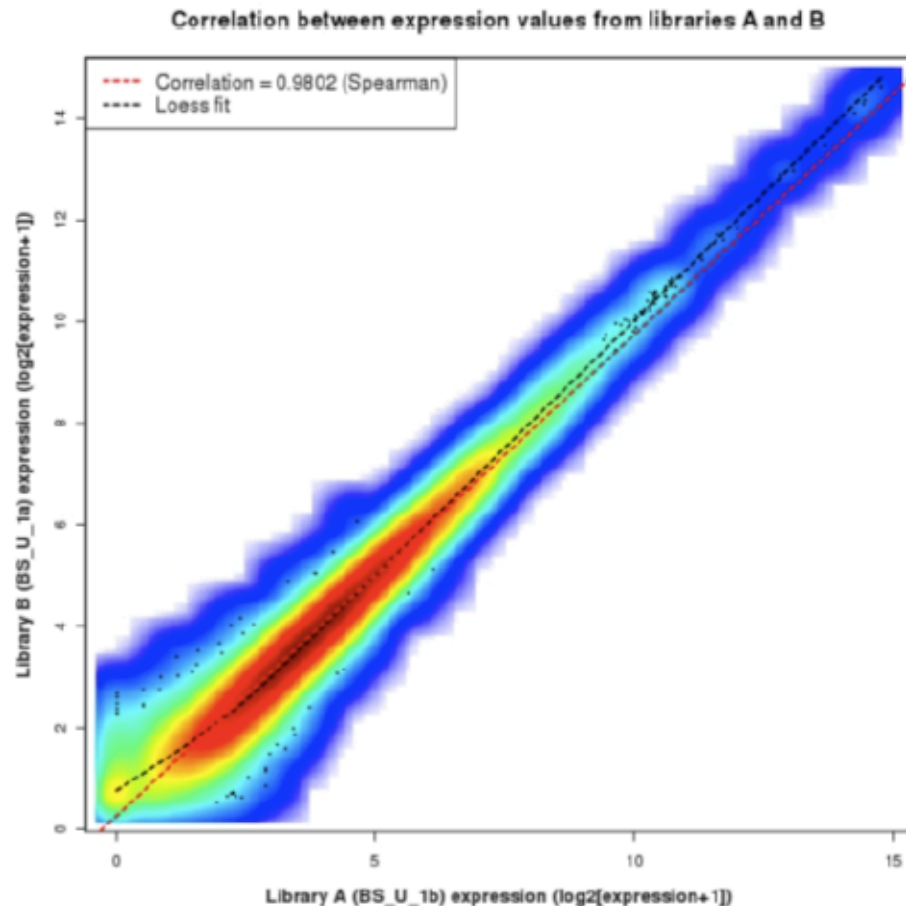
## B. Depiction of cDNA fragments from a stranded library



Library Kit	Stranded	5' to 3' IGV	TopHat (--library-type parameter)	HTSeq (--stranded/-s)	Picard (STRAND_SPECIFICITY option of CollectRnaSeqMetrics)
TruSeq Strand Specific Total RNA	Yes	F2R1	fr-firststrand	reverse	SECOND_READ_TRANSCRIPTION_STRAND
NuGEN Encore	Yes	F1R2	fr-secondstrand	yes	FIRST_READ_TRANSCRIPTION_STRAND
NuGEN OvationV2	No	F2R1 or F1R2	fr-unstranded	no	NONE

# Replicates

- Technical Replicate
  - Multiple instances of sequence generation
    - Flow Cells, Lanes, Indexes
- Biological Replicate
  - Multiple isolations of cells showing the same phenotype, stage or other experimental condition
  - Some example concerns/challenges:
    - Environmental Factors, Growth Conditions, Time
  - Correlation Coefficient 0.92-0.98



# RNA-seq questions during mapping

- Reference genome version – the latest version may have compatibility issues with other analysis
- Annotation – refseq or gencode or ENSEMBL
- Want junction read or not
- Remove duplicates? (No!)
- How many mismatches to allow?

# RNA-seq questions during quantification

- Keep reads mapping to multiple loci?
- Keep reads overlapping multiple genes?

# Alignment/Mapping tools

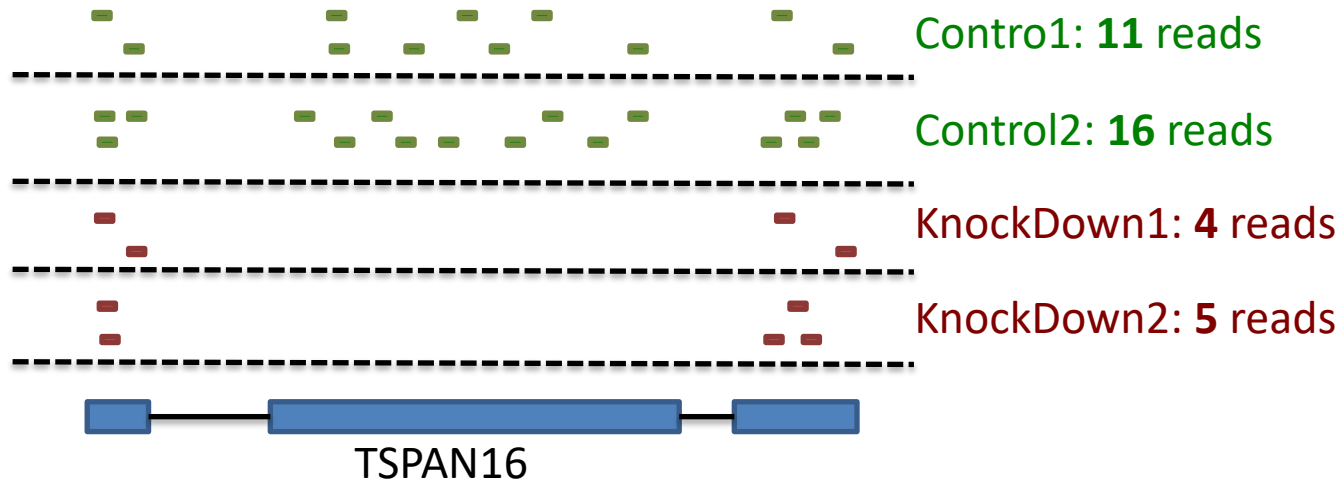
- TopHat2 (widely used, slow)
- STAR (super fast, very popular now, very demanding)
- RSEM (getting popular, used by ENCODE)
- Rsubread (works in R, not very popular)
- Sailfish (good for isoforms, less popular)

# Quantification

- TopHat2 → Cufflinks2 (provides FPKM) → Cuffdiff (DGE analysis)
- STAR → HTSeq-count, featureCount (provides raw count) → DESeq2, EdgeR (DGE analysis and normalized count)
- RSEM → RSEM (provides expected count, TPM and FPKM) → EBSeq
- Rsubread → featureCount → DESeq2, EdgeR
- Sailfish → Sailfish (provides raw count, TPM) → DESeq2, EdgeR



# Summarized RNA-seq



	Control1	Control2	KnockDown1	KnockDown2
TSPAN6	11	16	4	5
TNMD	1	0	0	0
DPM1	435	743	836	739
SCYL3	203	218	416	352
C1orf112	216	643	714	704
FGR	2365	5011	2828	2294
CFH	6	1	4	0
FUCA2	380	865	431	523
...	...	...	...	...
NFYA	888	827	1674	1580

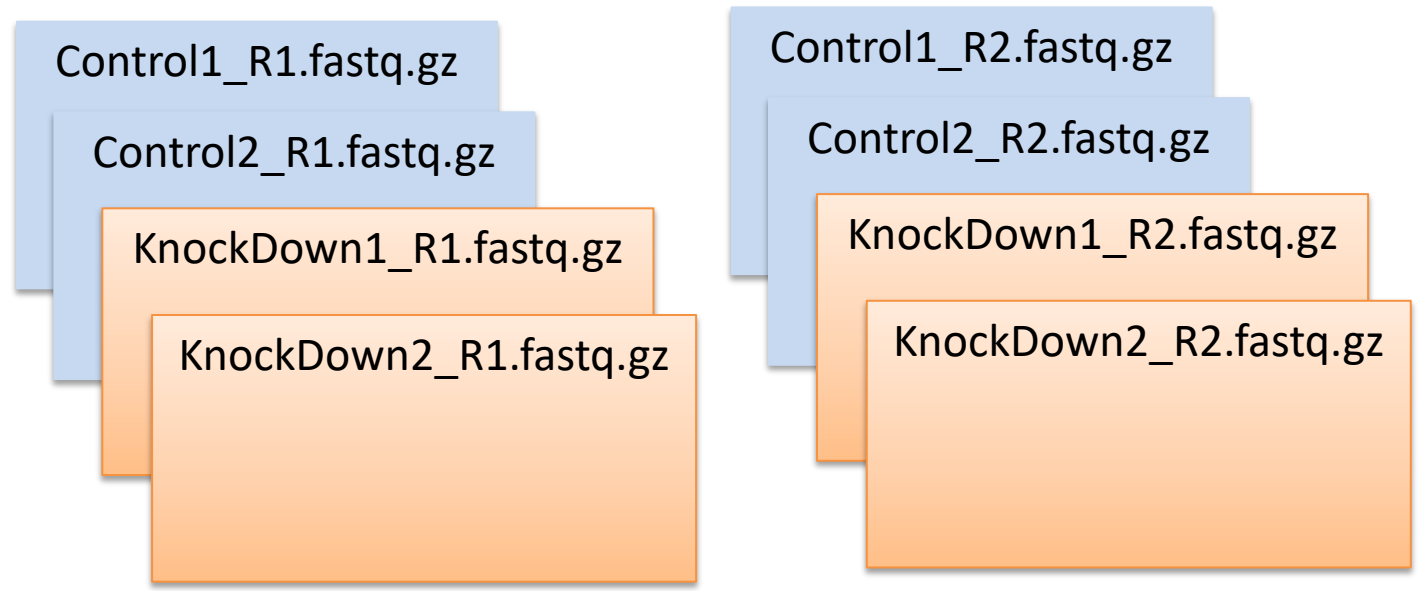
# Key concepts

- **Expression units:** There are several expression units available – RPKM/FPKM, CPM, TPM, Normalized expression
- **Fasta file:** Sequence storing file (can be opened in TextWrangler (unix) or Notepad++ (windows))
  - Format:

```
>sequence1  
ATCGTGCTGATGCGTGACG
```
- **Fastq file:** Sequence storing file with quality score, what you get from the sequencing centre
- **Bed file:** Standard file format for storing genomic coordinates
  - Format:

```
chr1 1033453 10443542 locus1  
chr3 4442235 45235256 locus2
```

# First file: fastq



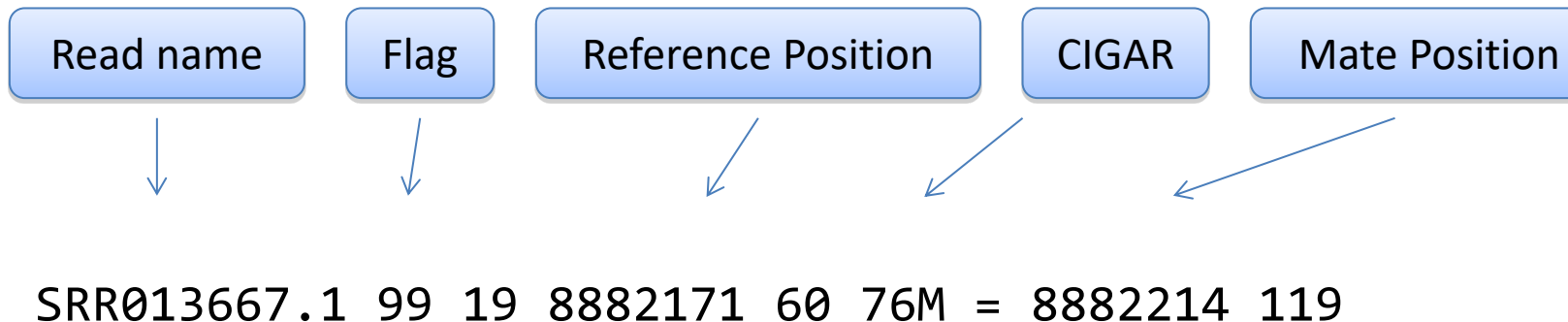
~ 10Gb each sample

```
@ERR127302.1 HWI-EAS350_0441:1:1:1055:4898#0/1
GGCTCATCTTGAAC TGGGTGGCGACCGTCCCTGGCCCCTTCTTGACACCCA
+
4=B@D99BDDDDDDDD:DD?B<<=?>6B#####
```

# From fastq to sam/bam



- Used to store alignments
- SAM = text, BAM = binary, CRAM=compressed binary



# Trivia time!

- What is the first step after getting the fastq file?
  - a) Contact bioinformatician    b) Alignment then QC
  - c) QC then alignment            d) Quantification and plotting
- Should we always have replicates?
  - a) Yes    b) No
- From fastq we make bam files. What do bam files contain?
  - a) Mapped reads    b) Treasure map    c) Raw reads            d) Nothing useful really
- Can I use FPKM in DESeq2?
  - a) Yes    b) No
- Which one of below is a major bottleneck in gene expression analysis?
  - a) My confidence                    b) High performance computers
  - c) Repeats in the genome          d) ribosomal RNAs
- What data should I use to generate an expression boxplot for *MYC* for 4 samples processed together?
  - a) raw read count    b) FPKM    c) normalized read count    d) TPM
- HTSeq-count or featureCount requires \_\_\_\_\_
  - a) fastq files    b) latest computer    c) bam files    d) annotation file

## Trivia time!

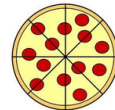
- What does GTF file contain?
  - a) Mapping information      b) Genome annotation
  - c) Quality information      d) Loci information
- What does bed file contain?
  - a) Mapping information      b) Genome annotation
  - c) Quality information      d) Loci information
- You should consider before deciding on sequencing depth -
  - a) single end or paired end      b) research purpose
  - c) RNA quantity      d) Read length
- I want to study alternative splicing. Factors in order of their importance for my study -
  - a) selection, paired-end, depth, read length
  - b) Read length, selection, paired-end, depth
  - c) Depth, read length, paired-end, selection
  - d) Read length, depth, paired-end, selection

# Resources

- <https://www.ncbi.nlm.nih.gov/pubmed/26248053>
- <http://www.bioconductor.org/help/workflows/rnaseqGene>
- <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4728800/>
- <https://usegalaxy.org/u/jeremy/p/galaxy-rna-seq-analysis-exercise>

# Are there any particular topics you would like to be discussed?

- “I'm sure it will be covered by quality control checks that determine if your RNAseq data is poor or good enough for further use.”
  - RNA integrity
  - fastQC
  - Check known genes
  - Compare replicates (correlation, PCA etc)
  - Visualize read distribution in IGV
- “TCGA”
  - cBioPortal
  - Xena browser (<https://tcga.xenahubs.net>)
- “N/A But I don't need pizza :)”



**NEVER  
TRUST  
SOMEONE WHO  
DOESN'T  
EAT PIZZA**