

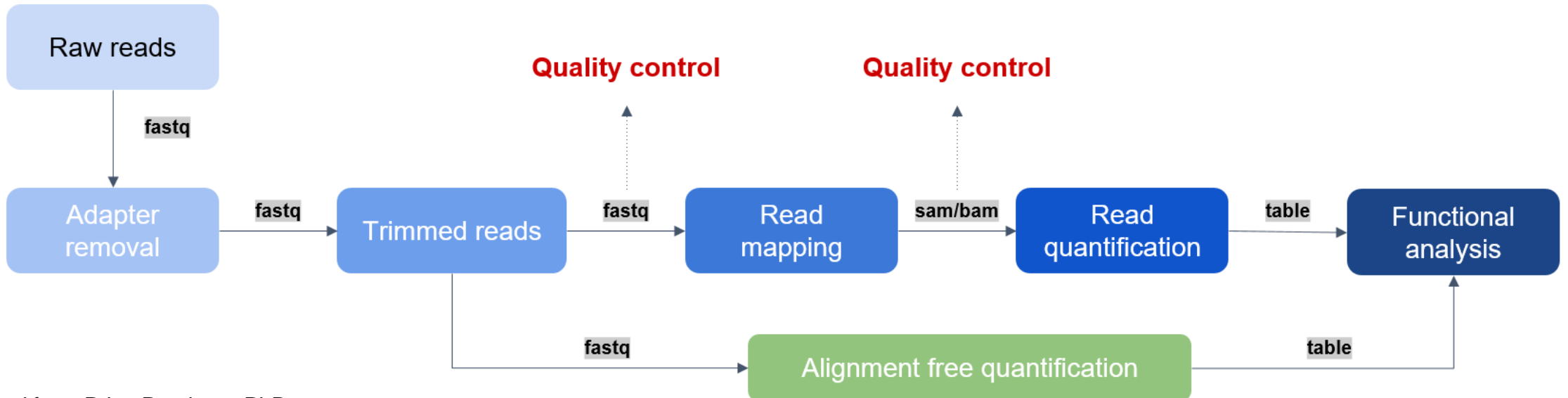
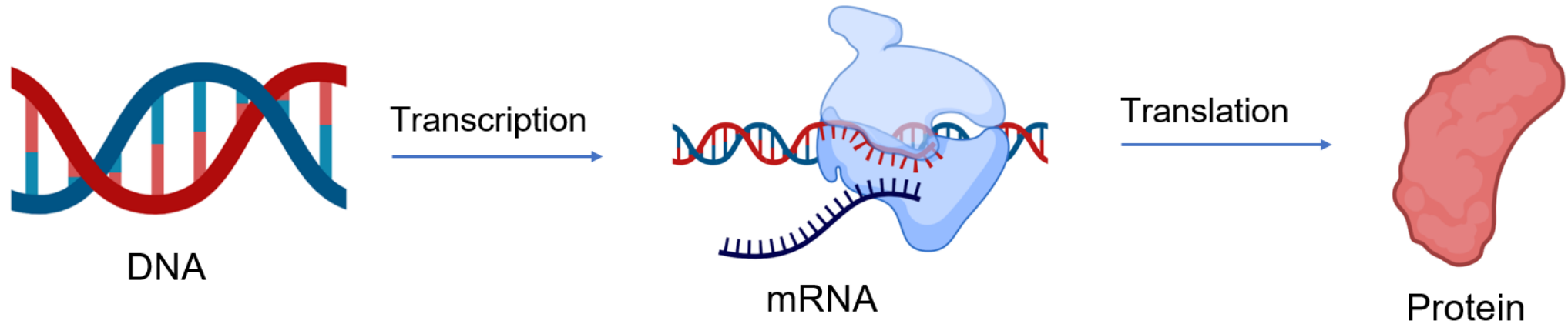
BGGN 239 – Week 2

Bulk and single-cell gene expression analysis

Ferhat Ay
ferhatay@lji.org

Associate Professor of Computational Biology, LJI
Department of Pediatrics & BISB PhD Program, UCSD

RNA-seq data flow



Reads and mapping them

One read

Identifier —● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50

Sequence —● TTGCCTGCCTATCATTTTAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT

'+' sign —● +

Quality scores —● hhhhhhhhhhhghhghhhhhfhhhhhfffffe'ee['X]b[d[ed'[Y[~Y

Another read

Identifier —● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50

Sequence —● GATTGTATGAAAGTATACAACTAAACTGCAGGTGGATCAGAGTAAGTC

'+' sign —● +

Quality scores —● hhhhgfhhcghghggfcffdhfehhhhcehdchhdhahehffffde'bVd

Aligned read

CGCATCAGT

Reference genome

...**ACTTGACGCATCAGTTGAAACGTA**..

chr1:105833

105843

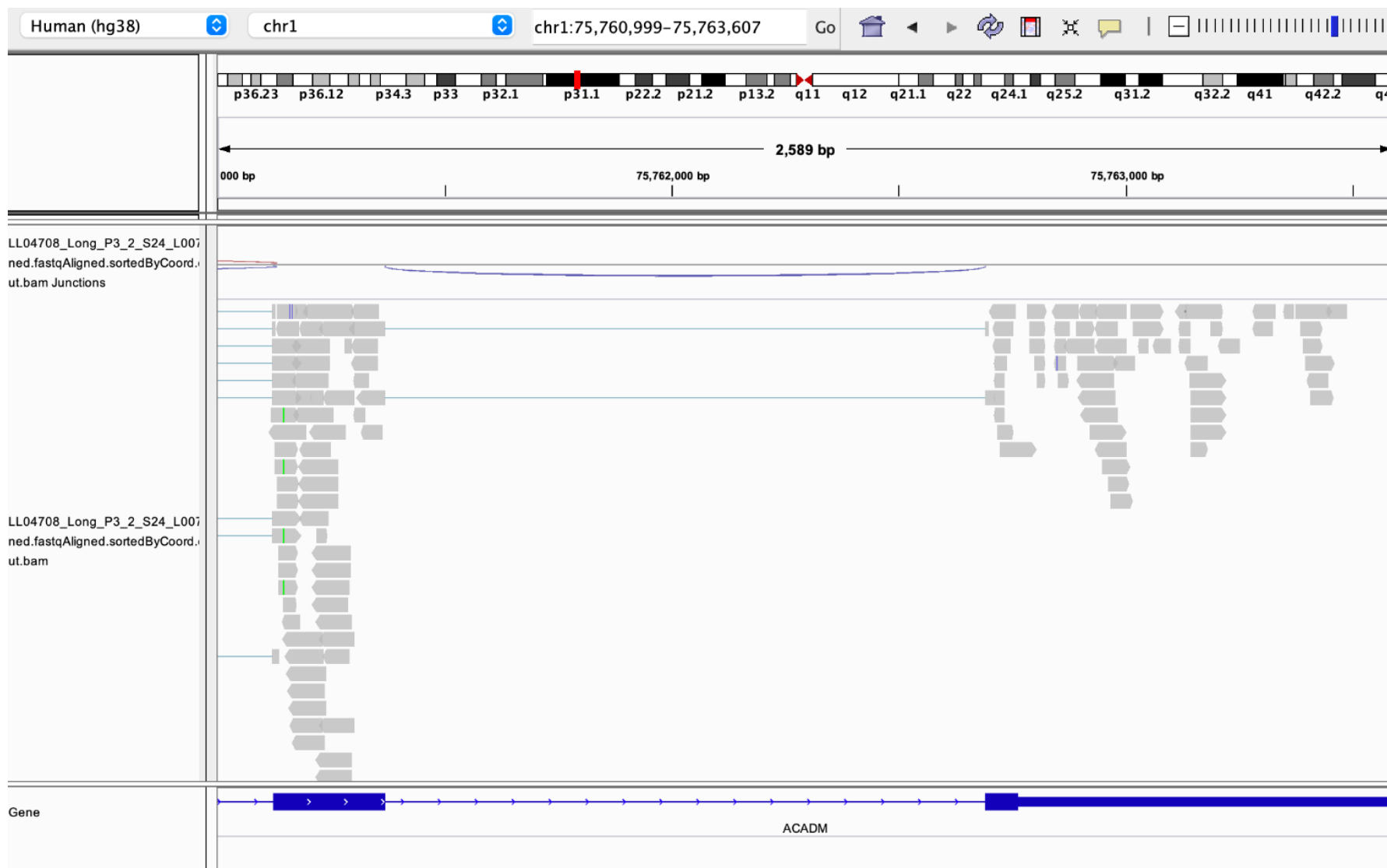
105853

Must match exactly?

Read mapping – popular tools

- STAR, Bowtie, HISAT2, TopHat

Visualizing mapped reads



What to do with counts?

Issues:

- 1.
- 2.
- 3.

What to do with counts?

- Reads per million or Counts per million
- Does not account for transcript length
- OK to use for sequencing protocols where reads are generated irrespective of gene length

$$\text{RPM or CPM} = \frac{\text{Number of reads mapped to gene} \times 10^6}{\text{Total number of mapped reads}}$$

What to do with counts?

- RPKM: Reads Per Kilobase of transcript per Million mapped reads
- FPKM*: Fragments Per Kilobase of transcript per Million mapped reads
- FPKM (or RPKM) attempts to normalize for gene size and library depth

*Fragments can mean either individual reads (SE) or paired-end reads that map together (PE)

$$\text{RPKM} = \frac{\text{Number of reads mapped to gene} \times 10^3 \times 10^6}{\text{Total number of mapped reads} \times \text{gene length in bp}}$$

What to do with counts?

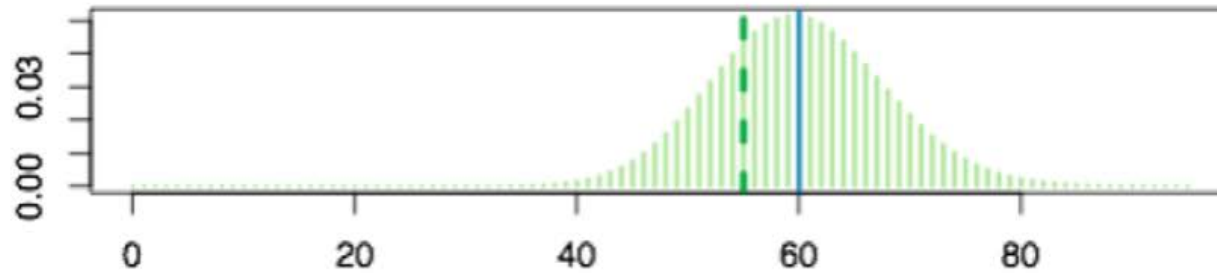
- TPM: Transcripts per million (Transcripts Per Kilobase Million)
- Another form of normalization for gene length and sequencing depth, but in a slightly different order

$$\text{TPM} = A \times \frac{1}{\sum(A)} \times 10^6$$

$$\text{Where } A = \frac{\text{total reads mapped to gene} \times 10^3}{\text{gene length in bp}}$$

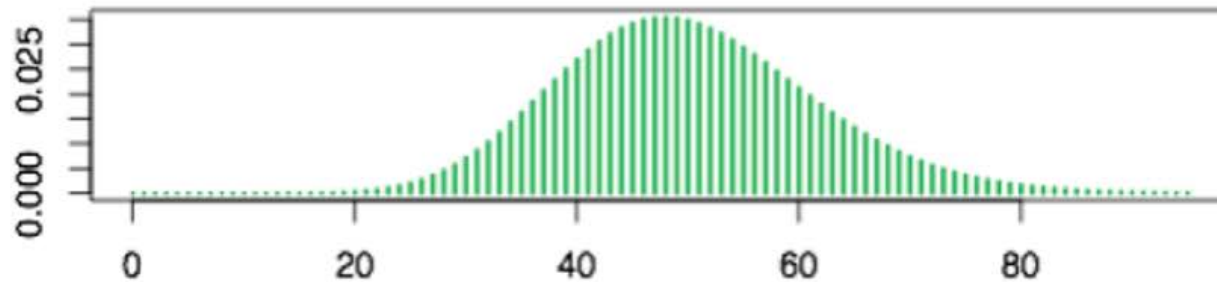
$$\text{TPM} = \frac{RPKM}{\sum(RPKM)} \times 10^6$$

Which distribution better captures count data?



Poisson

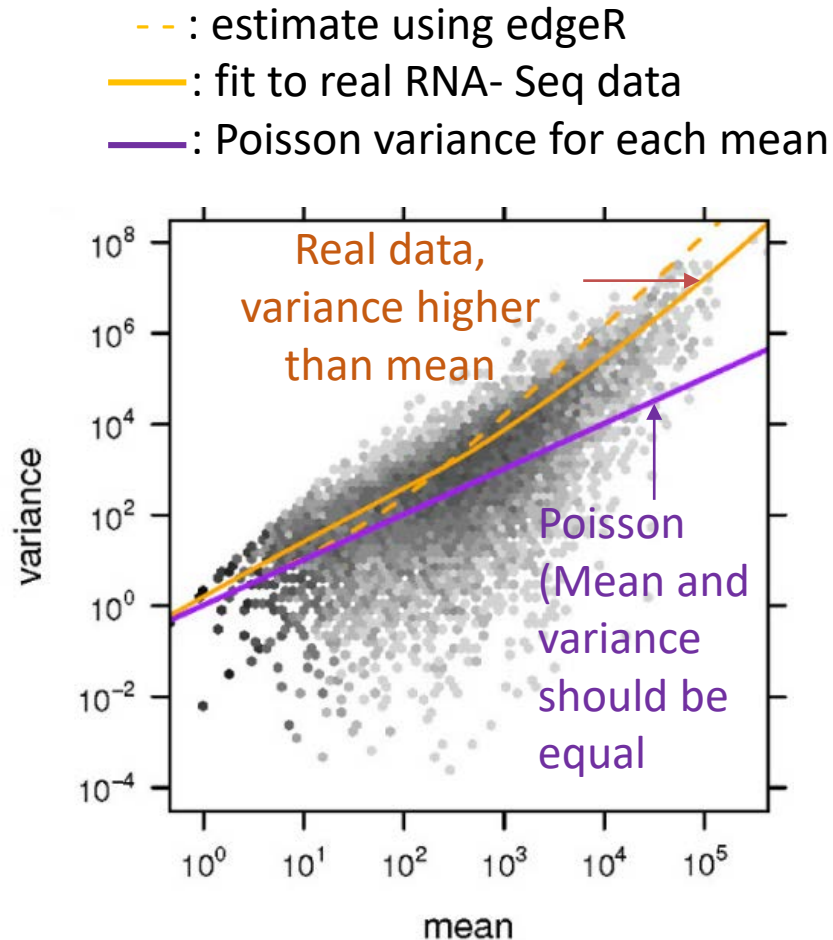
- Minimum variance of count data:
 $v = \mu$ (Poisson)



Negative binomial

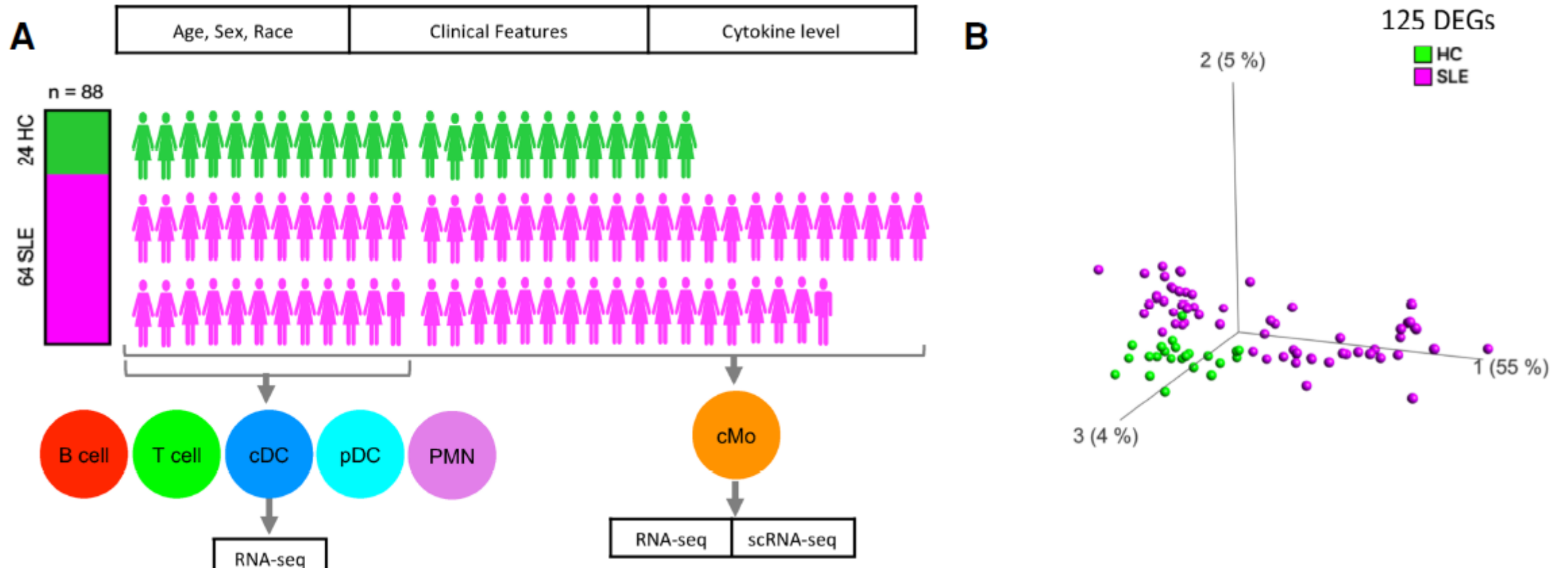
- Actual variance:
 $v = \mu + \alpha \mu^2$

Dispersion matters!

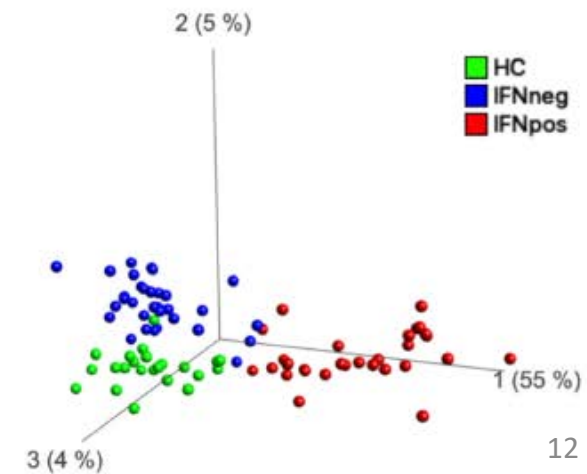
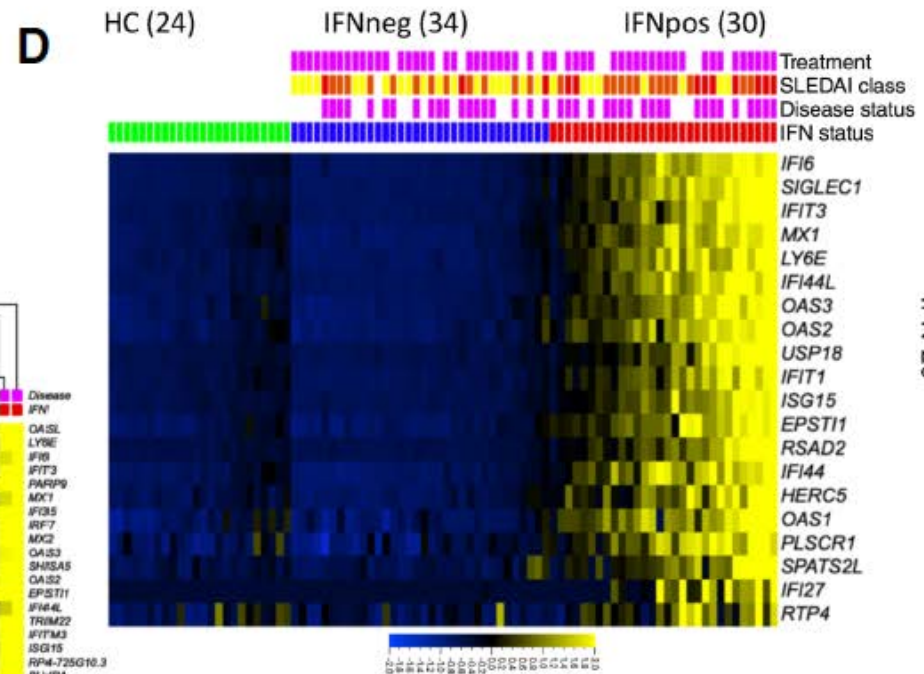
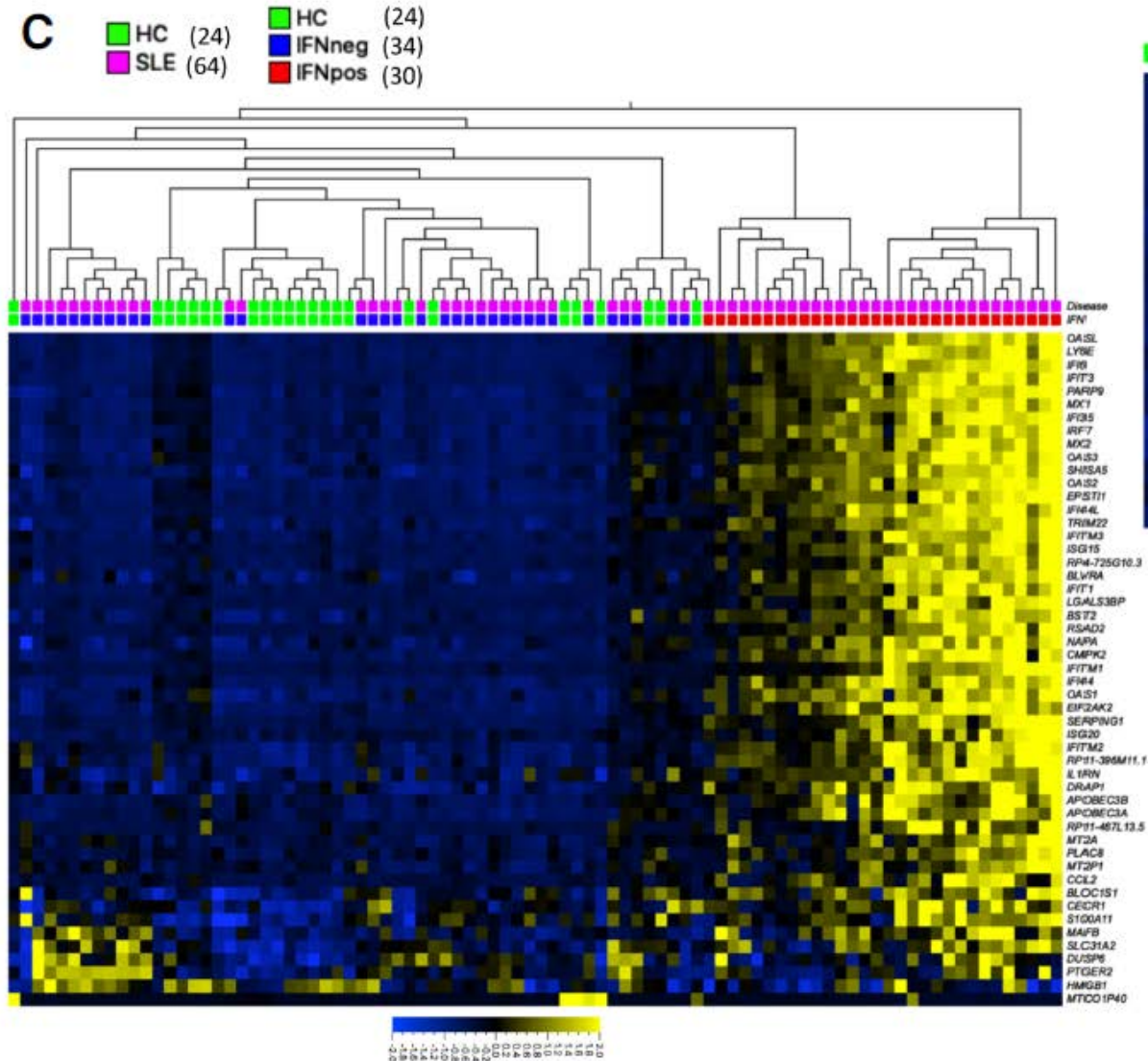


- α : “dispersion” $\alpha = (\mu - v) / \mu^2$
(squared coefficient of variation of extra-Poisson variability)
- Dispersion is a measure of the spread or variability in the data
- Biological Data is often ‘overdispersed’. With increasing mean the variance grows disproportionately
- Negative binomial model can account for this overdispersion

SLE paper



SLE paper



Homework #1

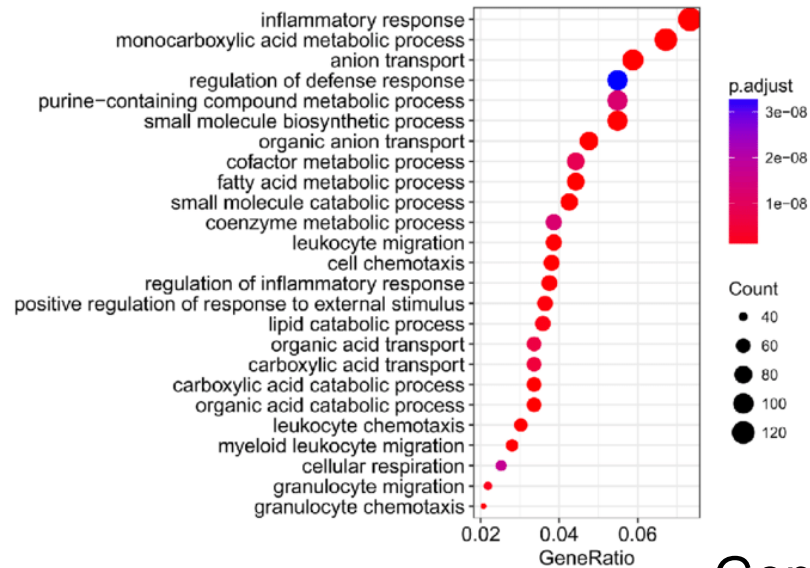
SLE mini project

- Do the same exercise (SLE.qmd) for another cell type and answer the questions below.
 - How many genes were differentially expressed at adjusted p-value cutoff of 5%? how many up and down-regulated?
 - How many genes remain when you filter with log2FC greater than 1 versus absolute log2FC greater than 1?
 - Write the resulting short list of genes ($p_{adj} < 0.05$ and $\log_2FC > 1$) in a csv file.
 - From that short list, select one gene and write 3-4 sentences about how this gene in this specific cell type may be relevant to SLE. Ask Google and ChatGPT for help if you like.

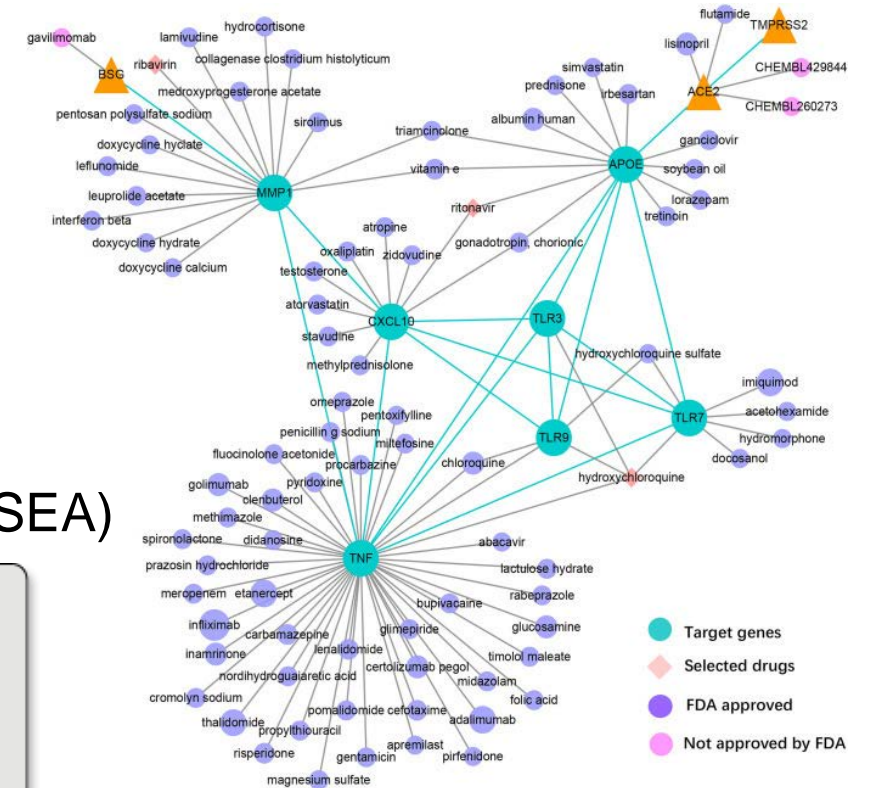
Questions?

Functional analysis of gene sets

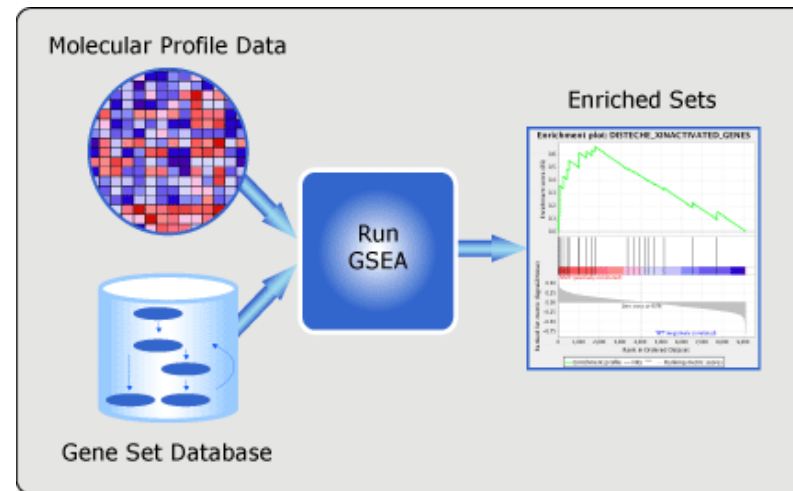
GO Term Enrichment Analysis



Gene Co-expression Networks



Gene Set Enrichment Analysis (GSEA)



Homework #2

GO Term enrichment

- In the `enrichGO` function, try setting `universe = names(sig_genes)` instead of `universe = names(all_genes_list)`. What happened? How many terms are statistically significant now?
- In the `enrichGO` function, set `ont = "CC"` rather than `ont = "BP"`. What did this do? Do you believe BP or CC will be more relevant for most use-cases?
- Go to the NYU link and select one other visualization you like to use. Add a code chunk that generates this visualization.
<https://learn.gencore.bio.nyu.edu/rna-seq-analysis/over-representation-analysis/>

Thank You!

- Paramita Dutta - LJI
- Priya Pantham - UCSD
- Barry Grant - UCSD

Resources

- <https://learn.gencore.bio.nyu.edu/rna-seq-analysis/over-representation-analysis/>
- <https://allisonhorst.com/r-packages-functions>
- <http://bioconductor.org/packages/devel/bioc/vignettes/DESeq2/inst/doc/DESeq2.html>
- <http://yulab-smu.top/biomedical-knowledge-mining-book/enrichment-overview.html>