



Cold  
Spring  
Harbor  
Laboratory

# Advanced Sequencing Technologies & Applications

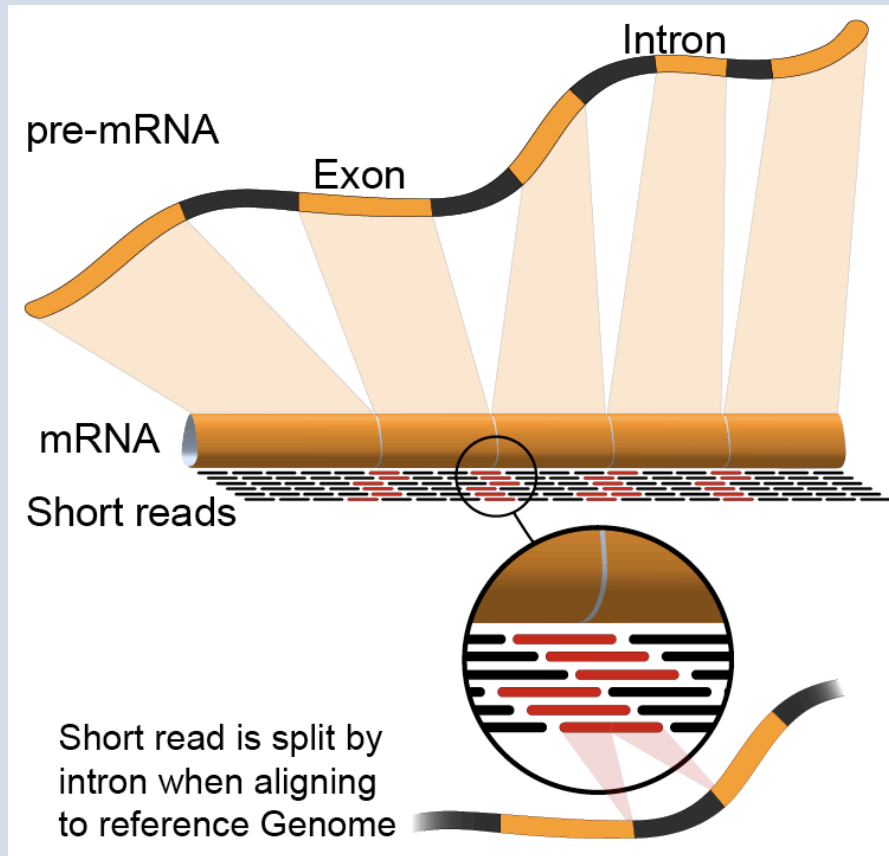
<http://meetings.cshl.edu/courses.html>



Cold  
Spring  
Harbor  
Laboratory

# RNA-Seq Module 3 Expression and Differential Expression (lecture)

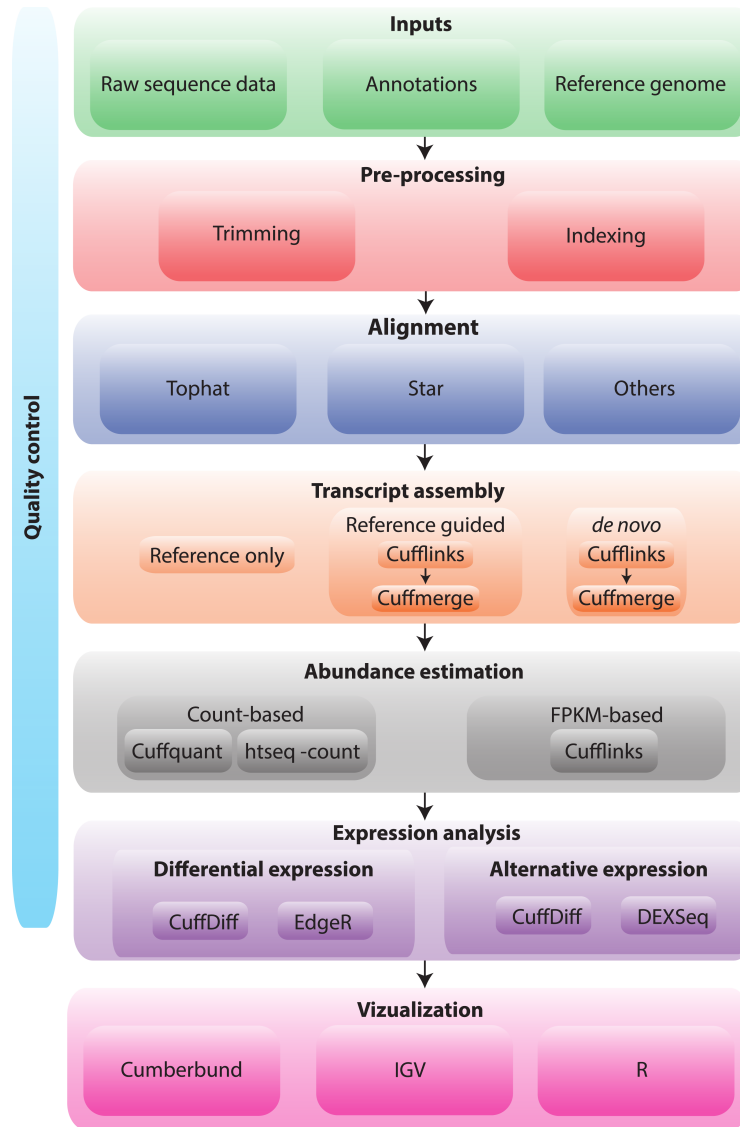
Malachi Griffith, Obi Griffith, Jason Walker, Alex Wagner  
Advanced Sequencing Technologies & Applications  
November 7 - 18, 2017



# Learning Objectives of Tutorial

- Generate gene/transcript expression estimates with StringTie
- Perform differential expression analysis with Ballgown
- Summarize and visualize results
  - Ballgown
  - Old school R methods

# RNA-seq Analysis Flow Chart



# 4-i. Generate expression estimates

- The alignment SAM/BAM files generated in the previous step will now be used by StringTie to calculate expression estimates
  - For all transcripts on the target chromosome
- For this step options ‘-G’ and ‘-e’ are used
  - ‘-e’ forces StringTie to calculate expression values for known transcripts
  - To discover novel transcripts with StringTie you should:
    - **Not use the ‘-e’ or ‘-G’ option. De novo transcript assembly and estimation will be performed. (we will try this in Module 4) OR ...**
    - Use the ‘-G’ option only. Known transcripts will be used as a ‘guide’, but novel transcripts will also be predicted.
- This step will generate one isoform and one gene expression file for each library
  - Expression values are reported as ‘FPKM’, or ‘**F**ragments **P**er **K**ilobase of exon per million fragments **M**apped’
  - Where each ‘fragment’ corresponds to a read-pair mapped to the genome

# 4-i. Generate expression estimates (Optional Alternatives)

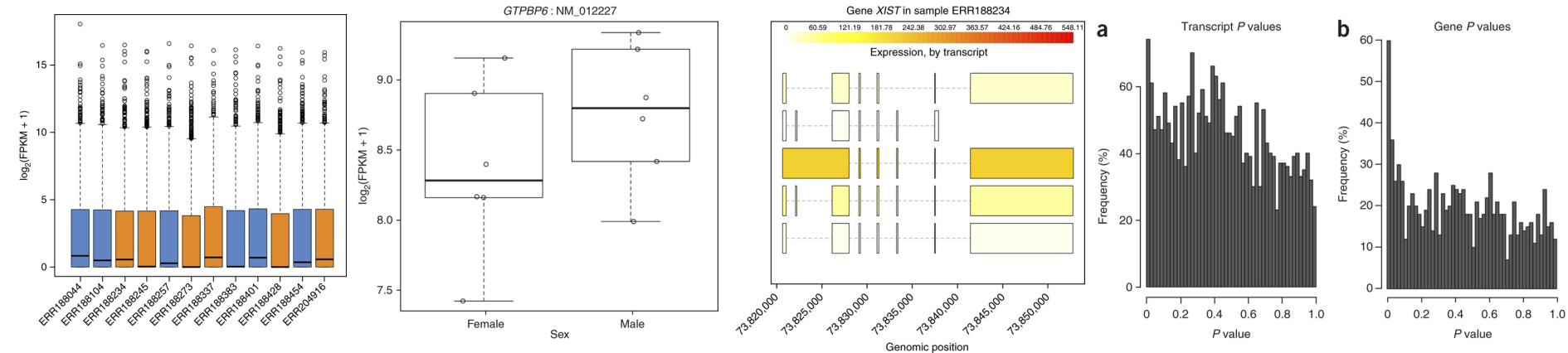
- Another alternative we will explore is a count-based method
  - We will use a program called htseq-count
    - Requires name-sorted SAM file
    - We will count at the gene level (transcript-level is also possible)
- In the end we will have two expression estimates for each sample
  - HISAT2/StringTie
  - HISAT2/Htseq-count

## 4-ii. Perform differential expression analysis

- In this step we will use Ballgown to:
  - Combine expression estimates from our 6 libraries into more convenient files
  - Combine expression estimates across replicates
  - Compare UHR vs. HBR and identify significantly differentially expressed genes and isoforms (transcripts)
- Note that these commands can get quite complicated when you have replicates
  - The positioning of spaces and commas, and grouping of libraries matters!
- Comparisons
  - Compare UHR vs. HBR using all replicates, for known (reference only mode) transcripts

# 4-iii. Summarize and visualize results

- In this step we will run the R package Ballgown to visualize our expression and differential expression results.
  - See online tutorial for details
  - <https://github.com/alyssafrazee/ballgown>
  - <http://bioconductor.org/packages/release/bioc/html/ballgown.html>



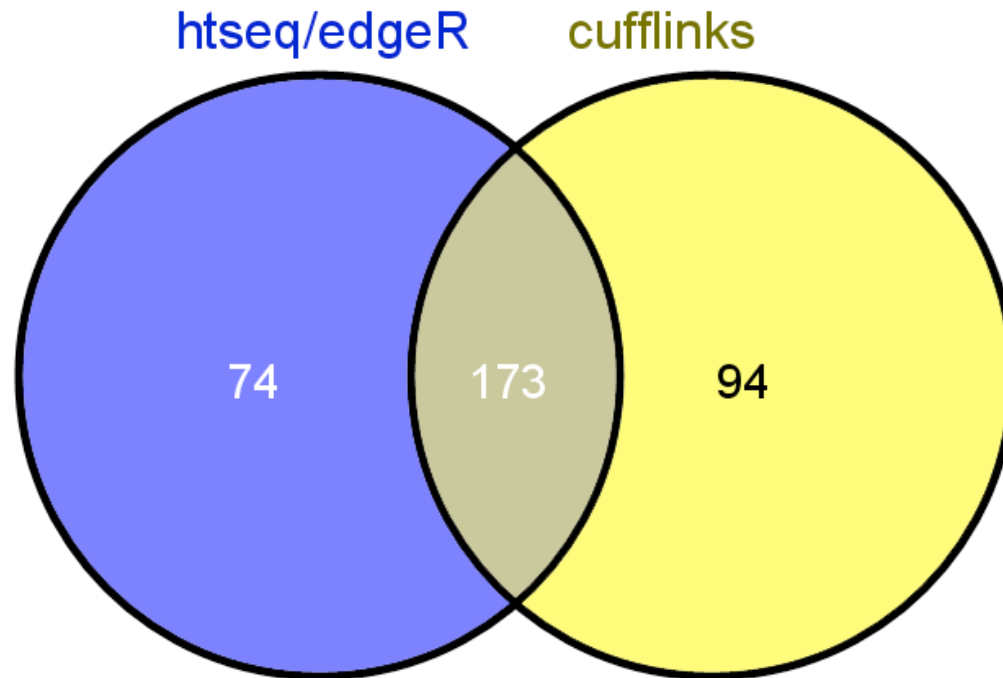


# Summarize and visualize results (optional)

- In this step we will use R to summarize and visualize the results of the previous steps
- Explanation of the R commands is provided in the online wiki
- Examples of the tasks performed:
  - Examine the expression estimates
    - How reproducible are the technical replicates?
    - How well do the different library construction methods correlate?
    - Visualize the differences between/among replicates, library prep methods and tumor versus normal
  - Examine the differential expression estimates
    - Visualize the expression estimates and highlight those genes that appear to be differentially expressed according to Ballgown
    - Generate a list of the top differentially expressed genes

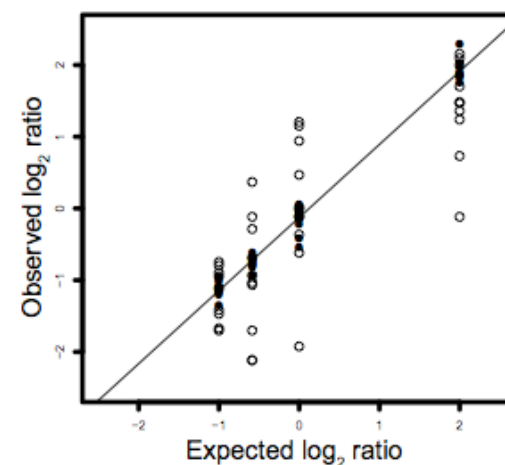
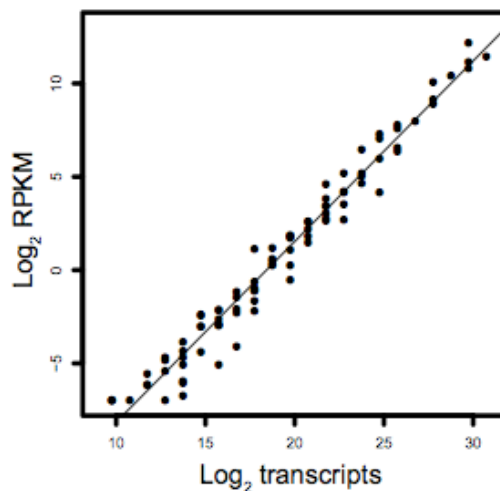
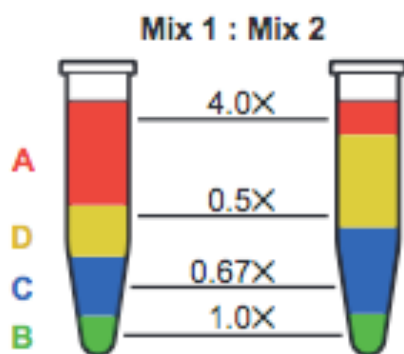
# Perform differential expression analysis with edgeR using htseq output (optional)

- Make use of raw counts generated by htseq-count
- Load into R and process with edgeR package
- Compare significantly differentially expressed genes from two methods



# Analysis of ERCC spike-in expression and differential expression (optional)

- [https://tools.lifetechnologies.com/content/sfs/manuals/cms\\_086340.pdf](https://tools.lifetechnologies.com/content/sfs/manuals/cms_086340.pdf)
- Lower Limit of Detection
- Dynamic Range (dose response)
- Fold-change response (DE)



We are on a Coffee Break &  
Networking Session