



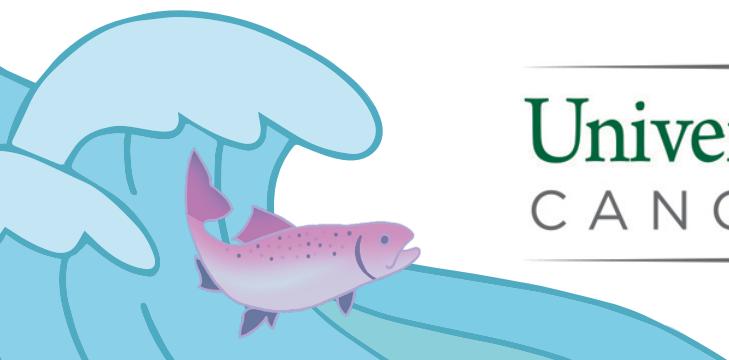
*A hands-on workshop for the experimental design
analysis and interpretation of RNA-seq data for
interdisciplinary research*

SACNAS NDiSTEM Conference 2023
Portland, OR



Seth Frietze, PhD

- Ph.D.; Harvard University
- Postdoctoral training; University of Southern California Medical School Norris Cancer Center
- Associate Professor
- Program co-leader UVM Cancer Center Cancer Cell Program



THE
University of Vermont
CANCER CENTER

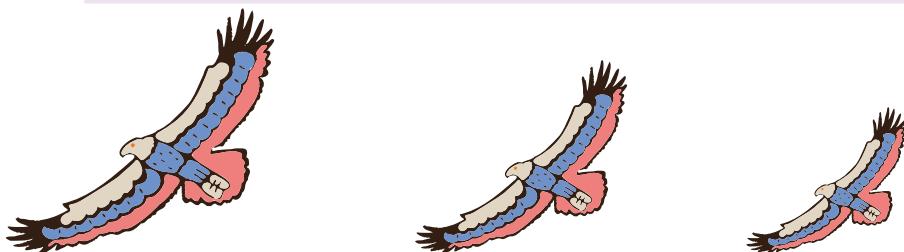


Princess Rodriguez, PhD

- MS in Immunology
- PhD, Cell and Molecular Biology
 - Started my bioinformatic journey in 2015 as a grad student in Seth's lab.
- Assistant Professor
 - Teach intro and advanced bioinformatics at UVM



THE
University of Vermont
CANCER CENTER



Noelle Gillis, PhD

- BS in Microbiology & Molecular Genetics
- PhD, Cell and Molecular Biology
 - Started my bioinformatic journey in 2016 as a grad student in the co-mentored by Seth Frietze.
- Currently a postdoc at U Minnesota Masonic Cancer Center.
- Working on projects related to steroid receptor transcriptional regulation in breast and ovarian cancer models.



Sophie Kogut

- BS in Biological Science
- Former research tech in Seth's lab
- PhD student in the Molecular and Cellular Biology Program at the University of Washington / Fred Hutchinson Cancer Center

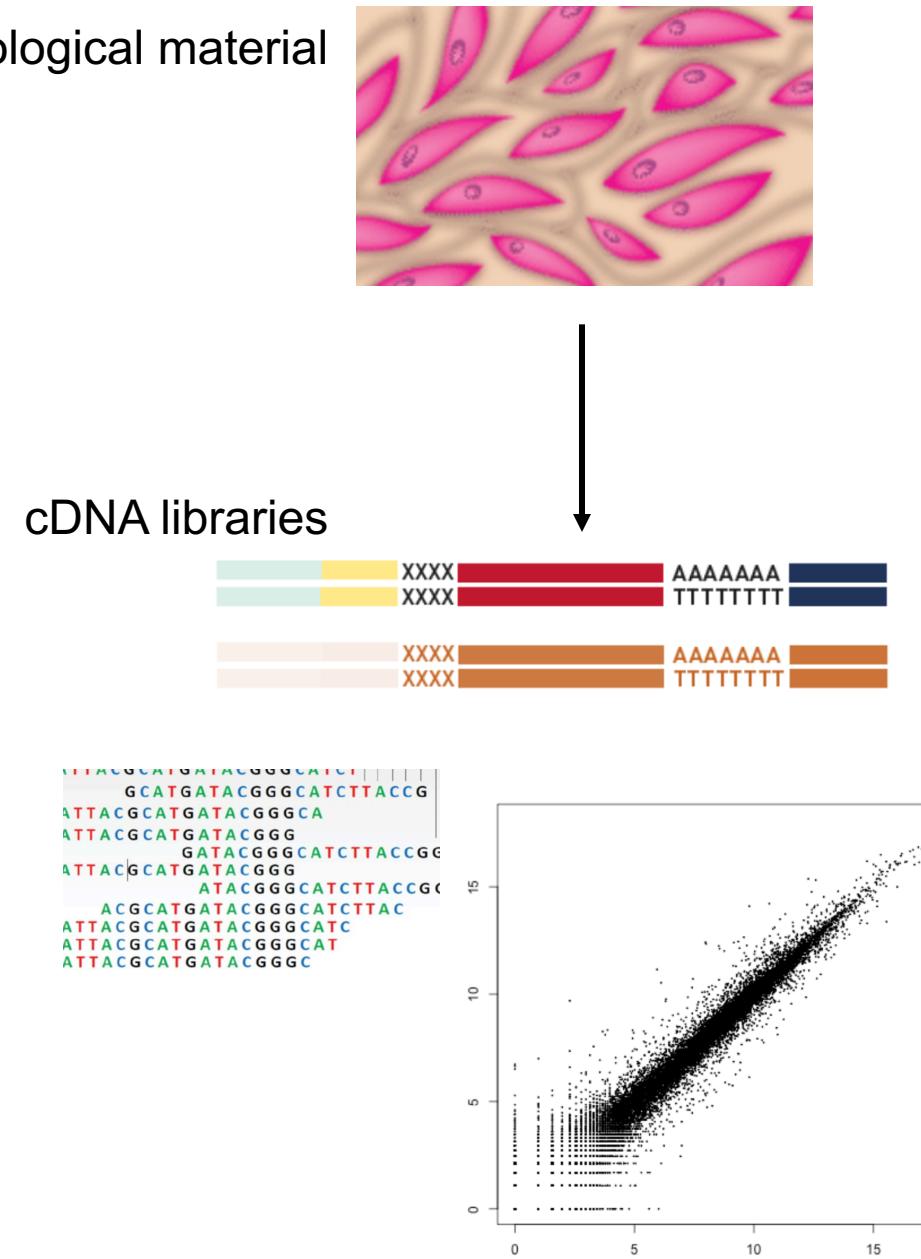
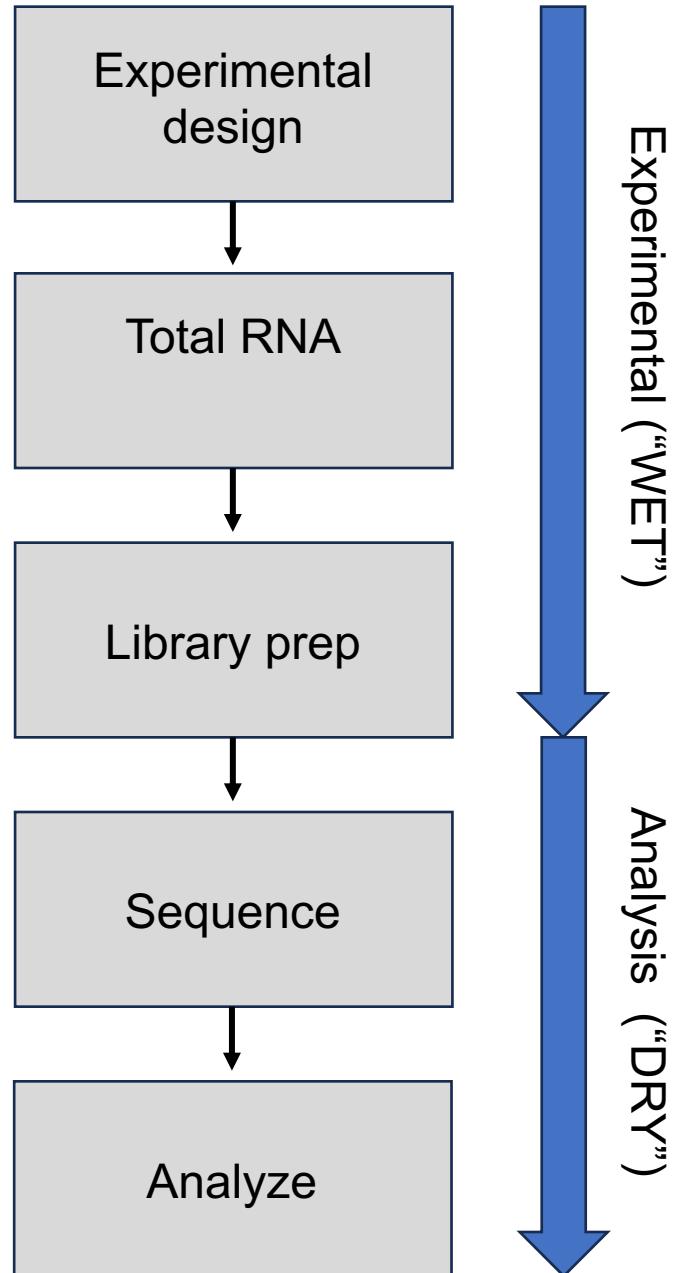
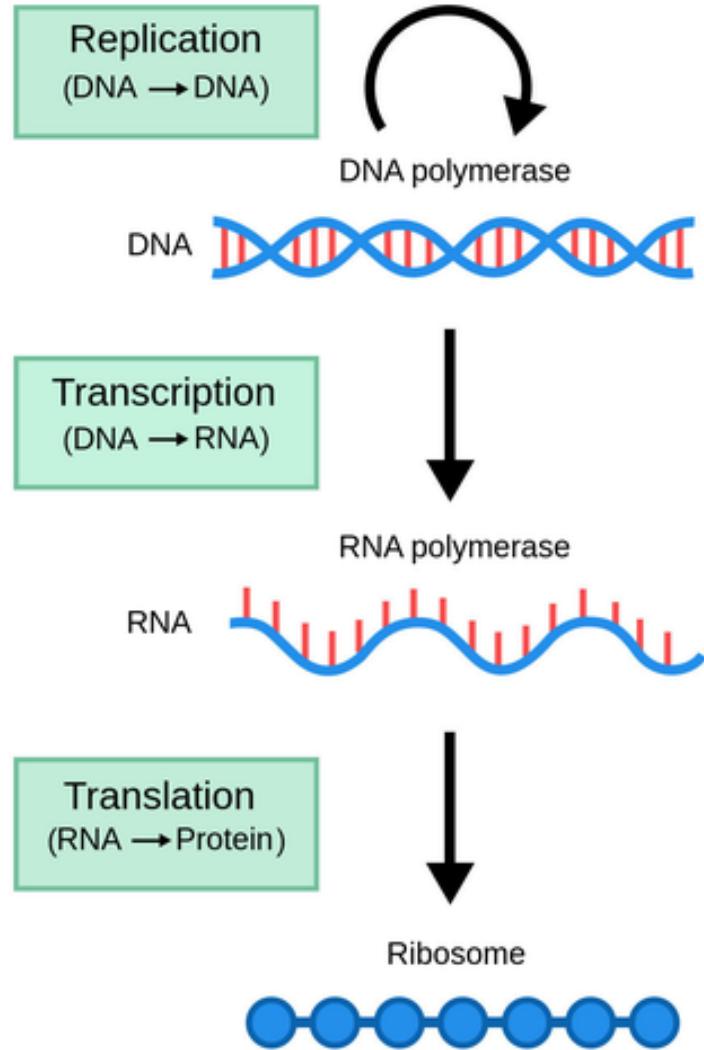


Module 1:

Overview of RNA-seq

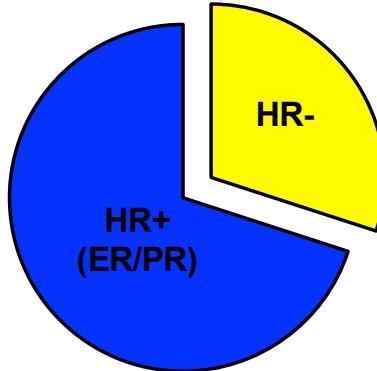
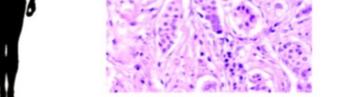


RNA-seq overview

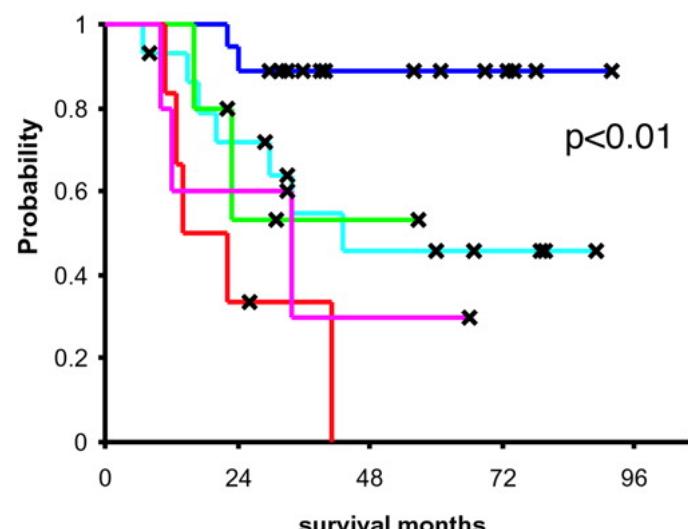


Breast cancer heterogeneity

Gene expression is diagnostic and prognostic



HR-
HR-/HER2-
(TNBC)

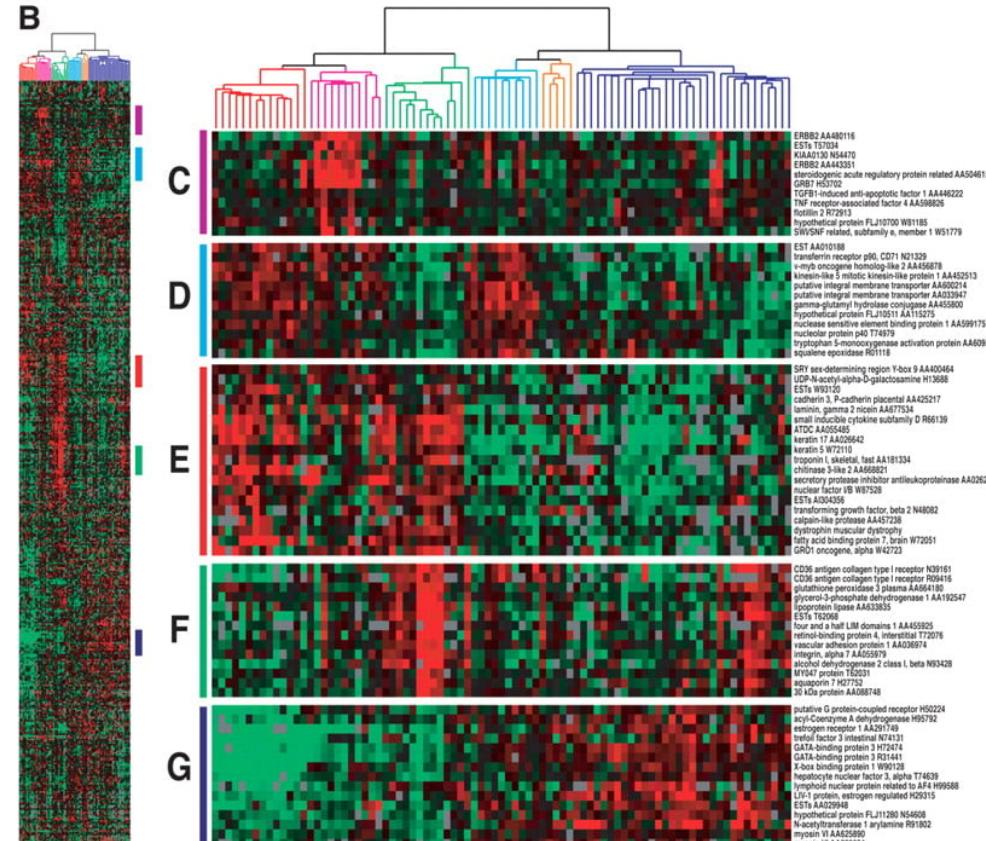
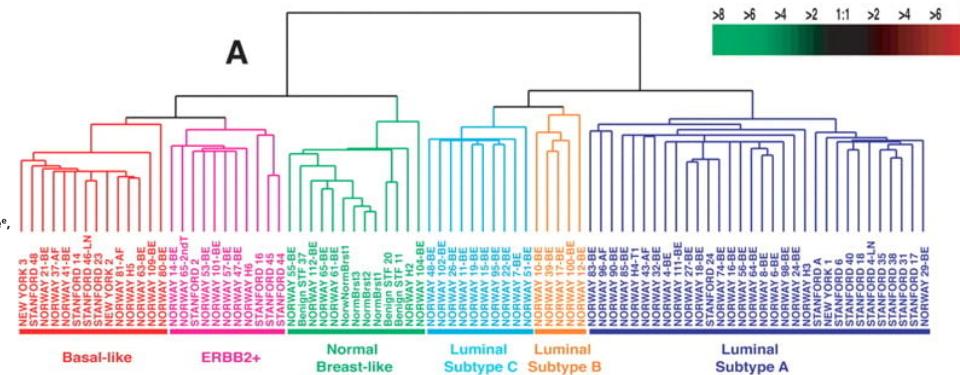


Gene expression patterns of breast carcinomas
distinguish tumor subclasses with
clinical implications

Therese Sørlie^{a,b,c}, Charles M. Perou^{a,d}, Robert Tibshirani^b, Turid Aas^f, Stephanie Geisler^b, Hilde Johnsen^b, Trevor Hastie^b, Michael B. Eisen^b, Marc van de Rijn^c, Stefanie S. Jeffrey^c, Thor Thorsen^c, Hanne Quist^c, John C. Matees^c, Patrick O. Brown^d, David Botstein^e, Per Eystein Lonning^f, and Anne-Lise Børresen-Dale^{b,h}

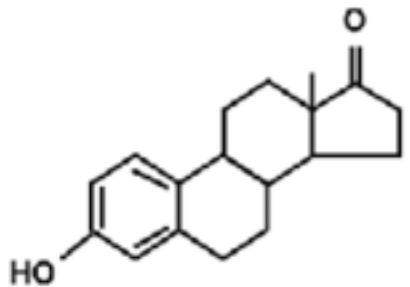
Departments of ^aGenetics and Surgery, The Norwegian Radium Hospital, Montebello, N-0310 Oslo, Norway, ^bDepartment of Genetics and Lineberger Comprehensive Cancer Center, University of North Carolina, Chapel Hill, NC 27599, ^cDepartments of ^cHealth Research and Policy and Statistics, Stanford, CA 94305, Departments of ^dMedicine (Section of Oncology), ^eSurgery, and ^fBiochemical Endocrinology, Haukeland University Hospital, N-5021 Bergen, Norway; and ^gLife Sciences Division, Lawrence Berkeley National Laboratories, and Department of Molecular and Cellular Biology, University of California, Berkeley, CA 94720

Contributed by David Botstein, July 17, 2001

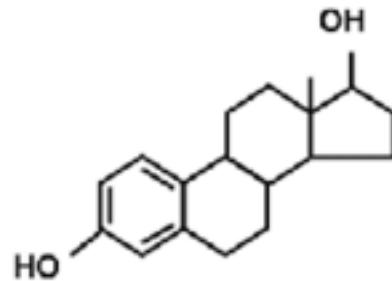


Endocrine therapies are based on estrogen ligands

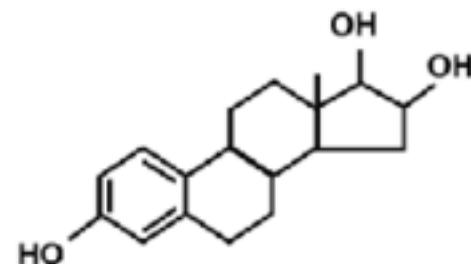
Estrogens



Estrone

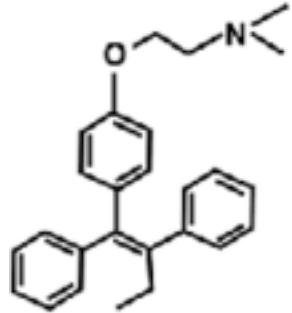


17 β -Estradiol

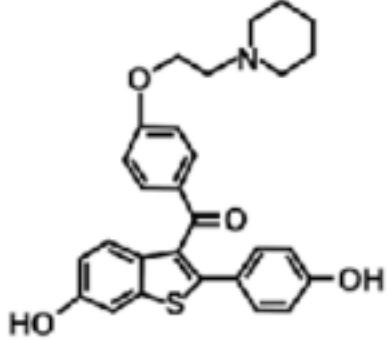


Estriol

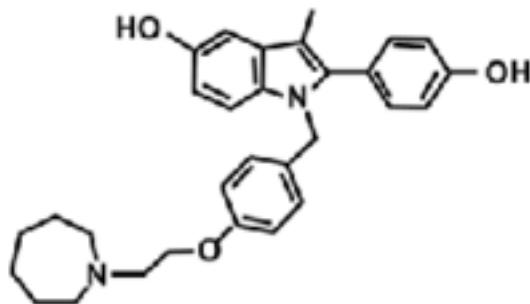
SERMs



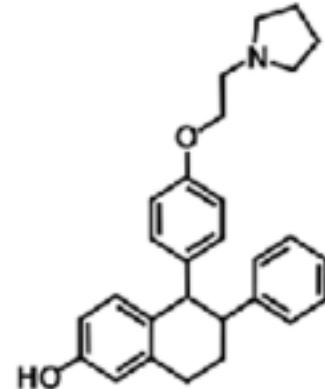
Tamoxifen



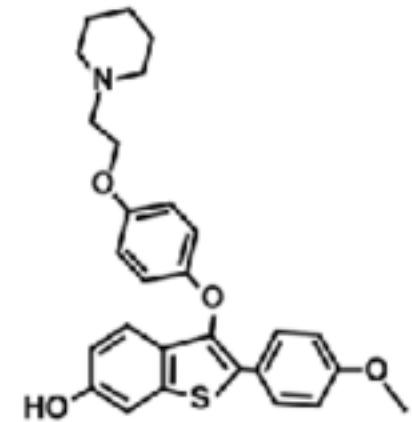
Raloxifene



Bazedoxifene



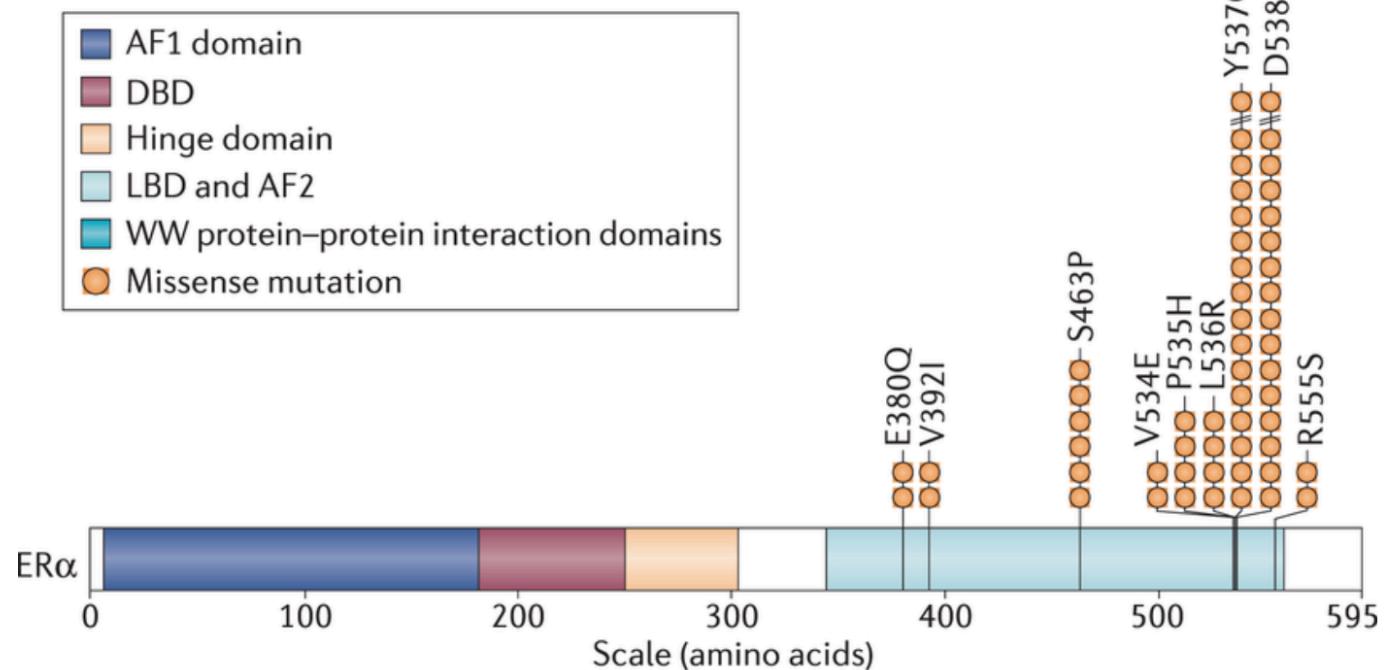
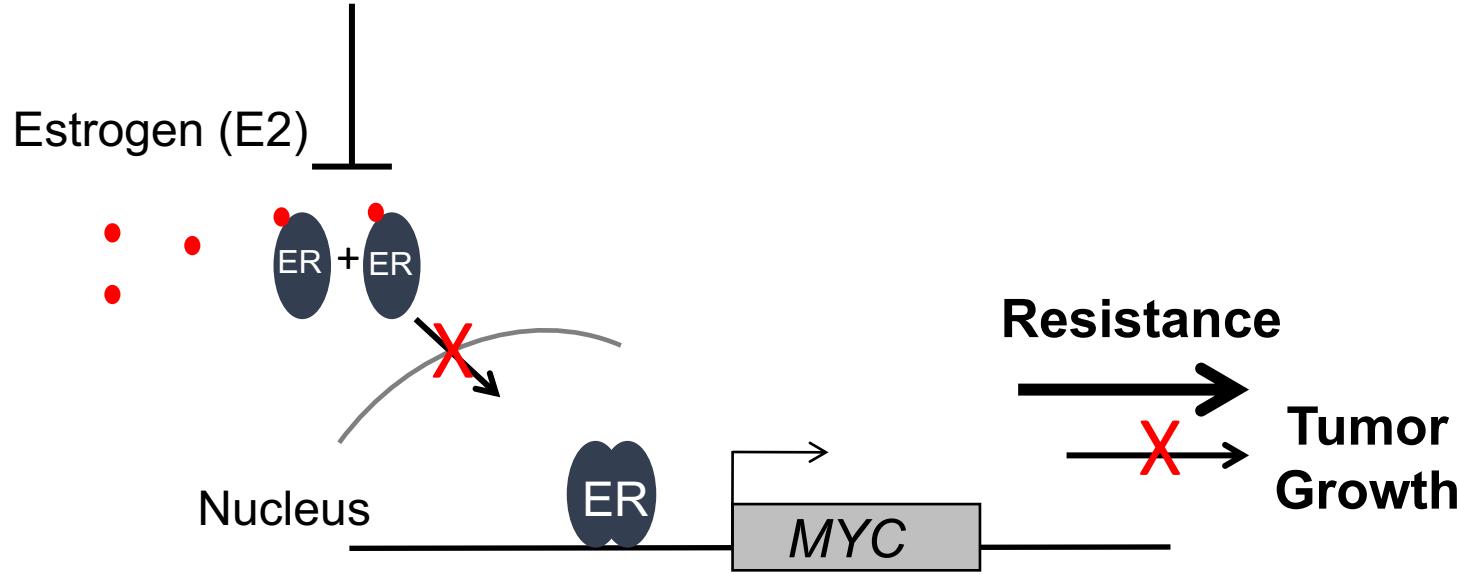
Lasofoxifene



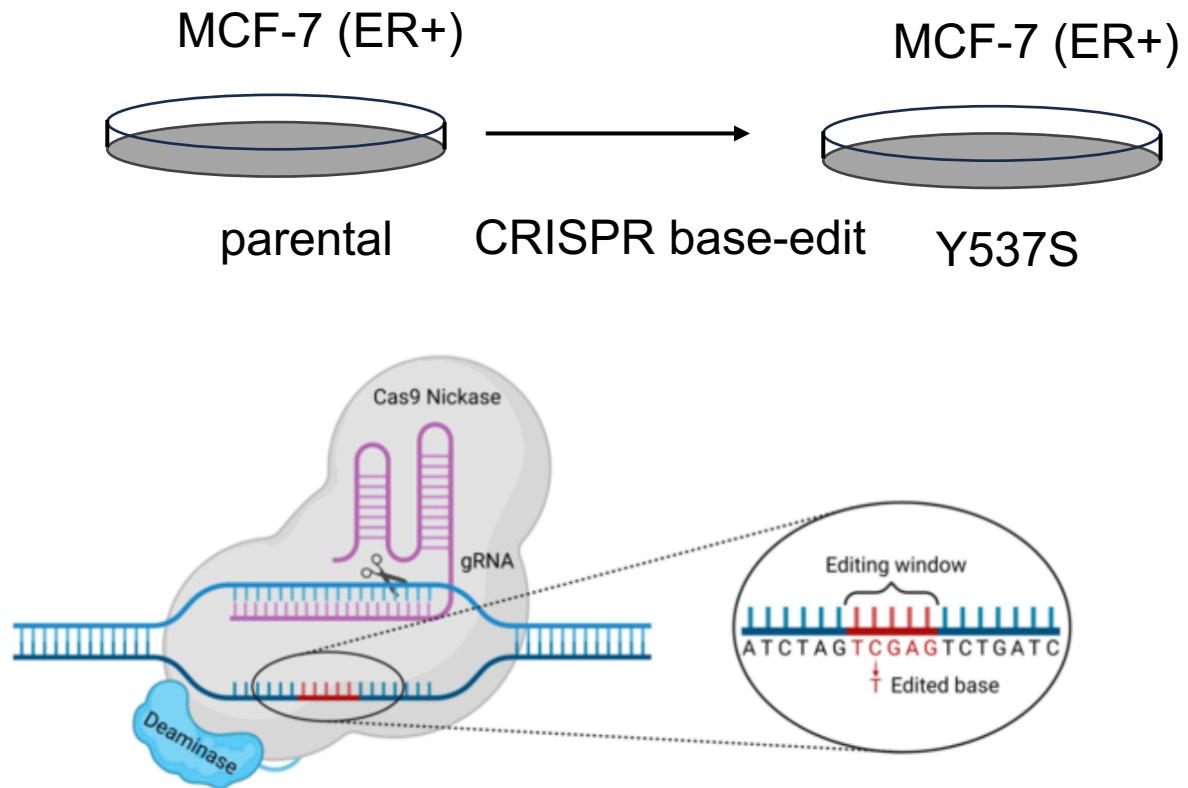
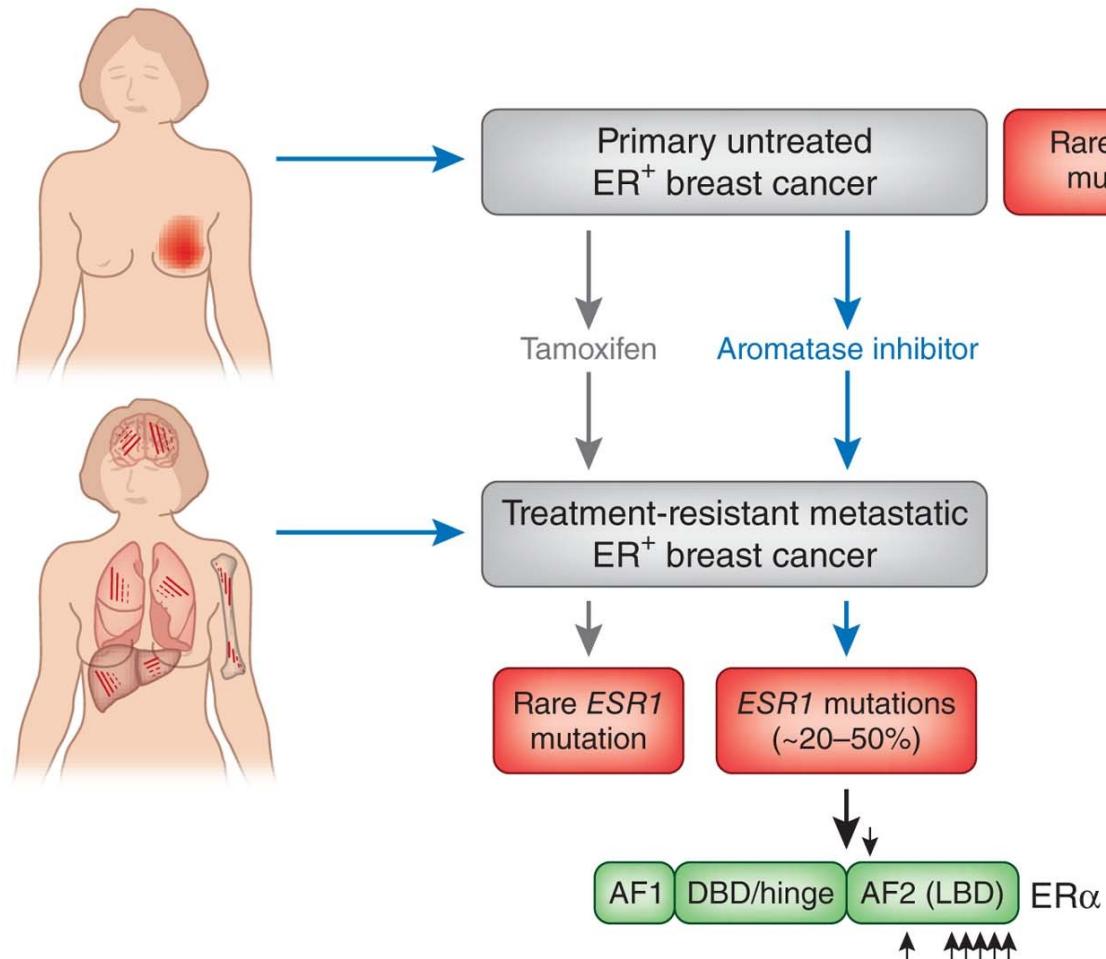
Arzoxifene

ESR1 mutations are common in endocrine therapy resistance (ETR)

Endocrine therapies

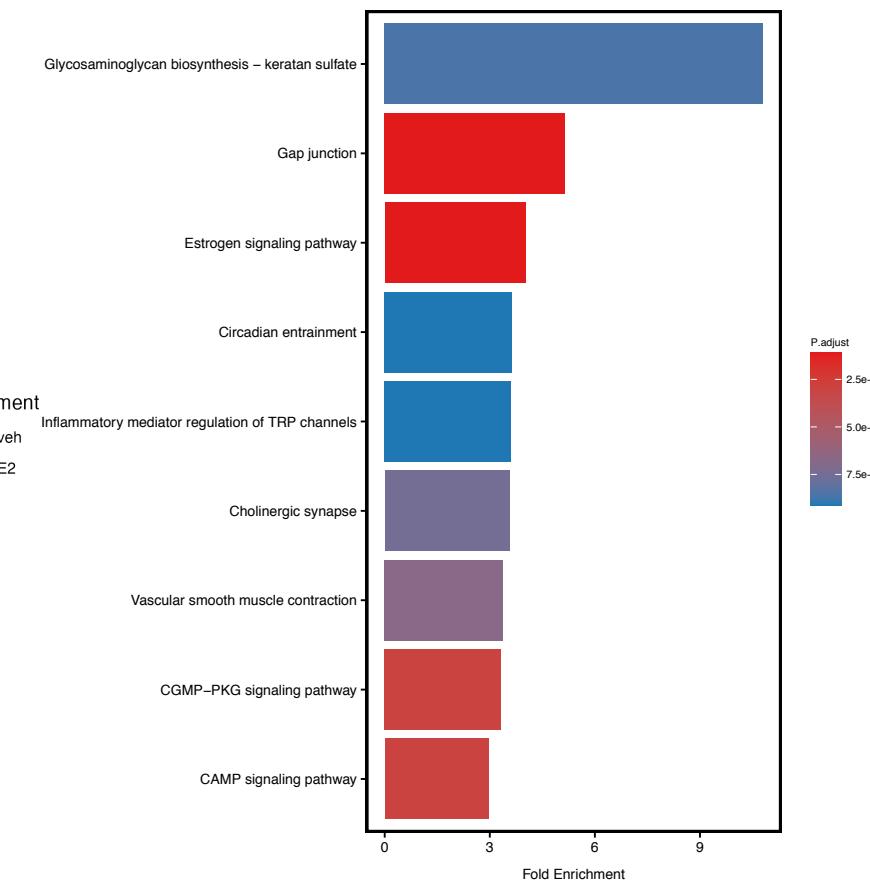
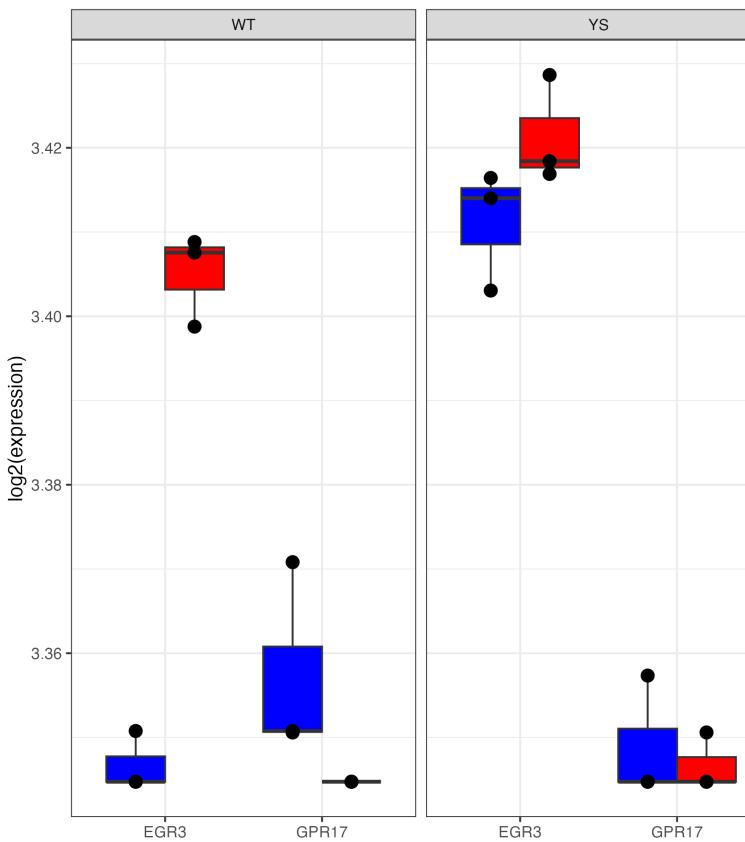
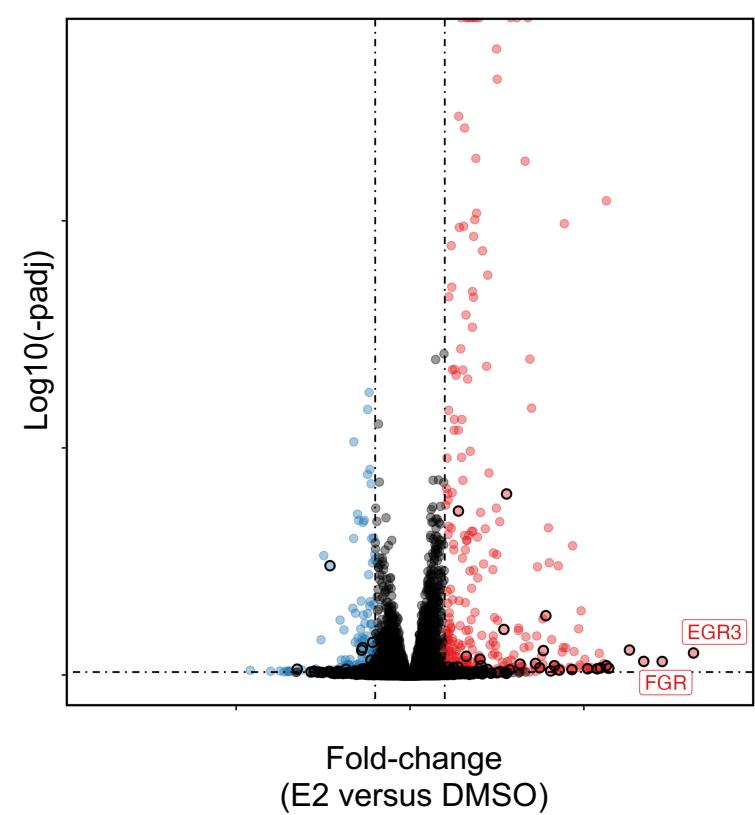
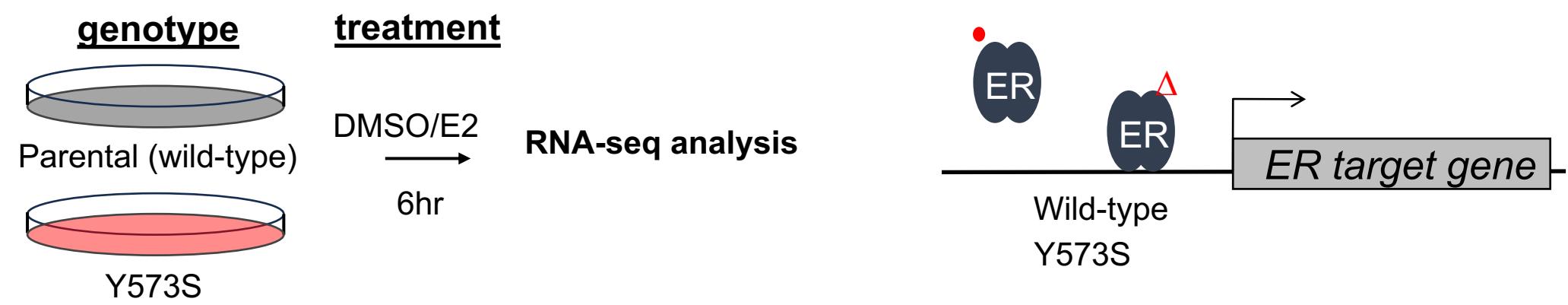


Cell models to study ETR: CRISPR engineering *ESR1* Y573S



CRISPR-base editing method

1. Find *ESR1* sequence on [Genome browser](#)
2. Protein sequence on [Uniprot](#)
3. Import sequence into [benchling](#), design gRNAs
4. Generate vectors, deliver and select stable cell lines
5. Characterize (RNA-seq)

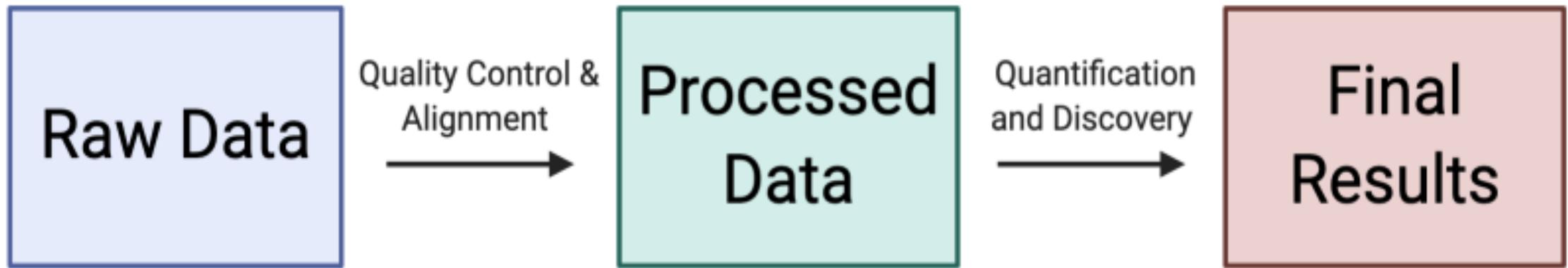


Module 2:

Post-Sequencing Processing Steps



Steps to get you from raw RNA-seq data to DEG's and beyond!



How I learned how to do this:

- Trial and Error.
- Spending lots of time reading the program manuals.
- Googling error messages.
- Asking for help when I get stuck.



Lots of this:



But eventually this:

Steps to get you from raw RNA-seq data to DEG's and beyond!

Phase 1: Processing

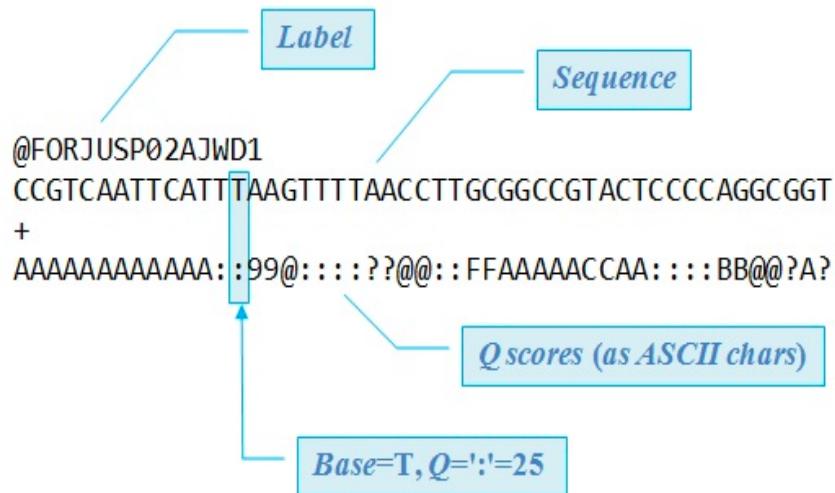
Processing Step	Tools
Quality control of raw sequences	FastQC, cutadapt, trimGalore
Alignment/Mapping to the genome	STAR
Generating gene counts	HT-seq, Salmon

Phase 2: Analysis

Analysis Step	Tools
Quality control of replicates	DESeq2, edgeR
Differential expression analysis	DESeq2, edgeR
Pathway analysis	GSEA, IPA, clusterProfiler

Quality Control of Raw Sequences

FASTQ File Format



Quality Control Steps

1. Generate FASTQC reports.
2. Check the quality of base-calls.
3. Check for “over-represented sequences.”
4. Trim low-quality bases.
5. Remove adapter sequences.

Example Commands:

```
fastqc 2D_WT_shGFP_E2_R50_001_S10_R1_001.fastq.gz

trim_galore --paired --Illumina 2D_WT_shGFP_E2_R50_001_S10_R2_001.fastq.gz
2D_WT_shGFP_E2_R50_001_S10_R2_001.fastq.gz
```

Alignment/Mapping to the Genome



Figure 4: Paired-End Sequencing and Alignment—Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in better alignment of reads, especially across difficult-to-sequence, repetitive regions of the genome.

Where do you get the reference genome files? https://www.ncbi.nlm.nih.gov/assembly/GCF_000001405.26/

Example Command:

```
STAR --runThreadN 12 \
--genomeDir /home/langeca/gilli431/software/STAR_hg38 \
--readFilesCommand zcat \
--readFilesIn 2D_WT_shGFP_E2_R50_001_S10_R1_001.fastq.gz
2D_WT_shGFP_E2_R50_001_S10_R2_001.fastq.gz \
--outFileNamePrefix 2D_WT_shGFP_E2_R50_001
```

STAR Alignment Summary:

Number of input reads	35202511
Average input read length	299
UNIQUE READS:	
Uniquely mapped reads number	33777023
Uniquely mapped reads %	95.95%
Average mapped length	298.68

Generating Gene Counts

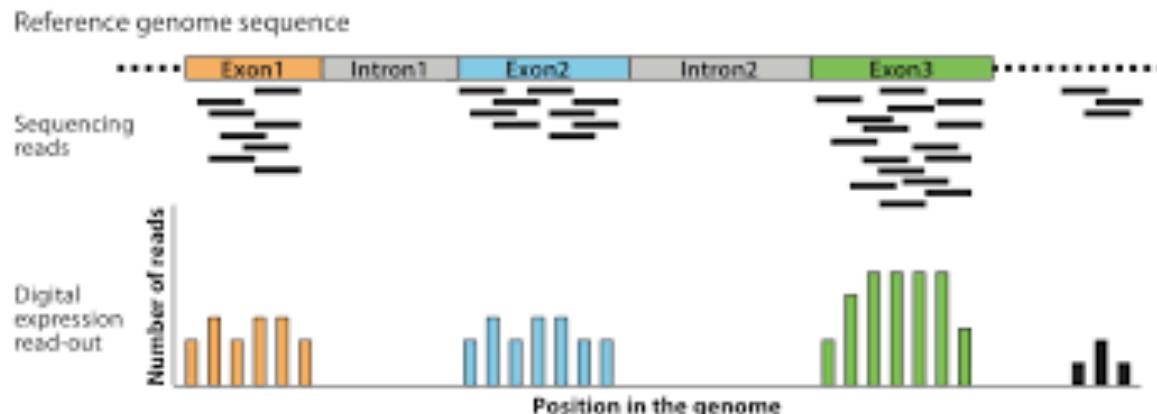


Example Command:

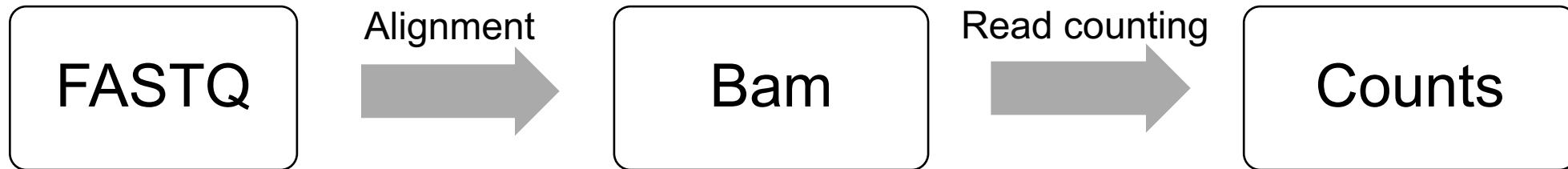
```
htseq-count -f bam -r name -s reverse -m union -i gene_id bam_files/${base}_sorted.bam  
~/software/gencode.v38.primary_assembly.annotation.gtf > count_files/${base}.counts
```

Gene Counting Steps

1. Check that bam files are sorted and indexed
2. Run HT-seq on bam files.
3. Save .cnts files for analysis.



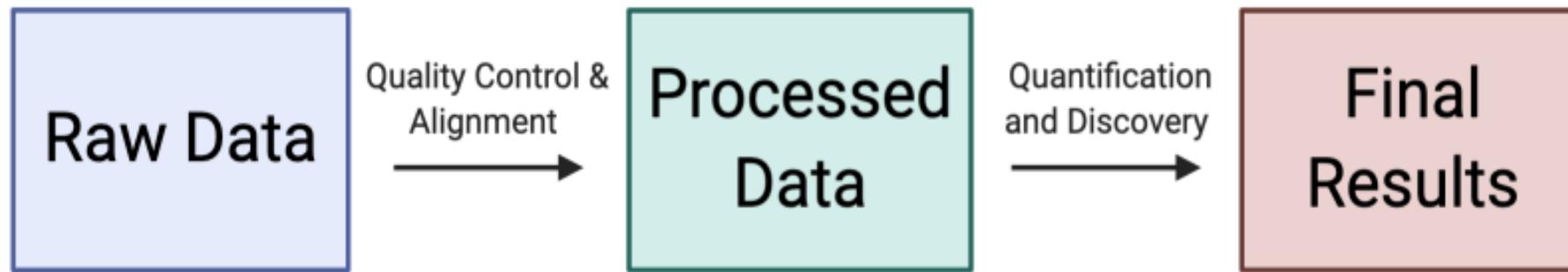
All these steps can be strung together into a “pipeline.”



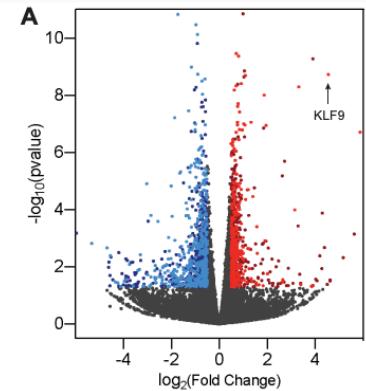
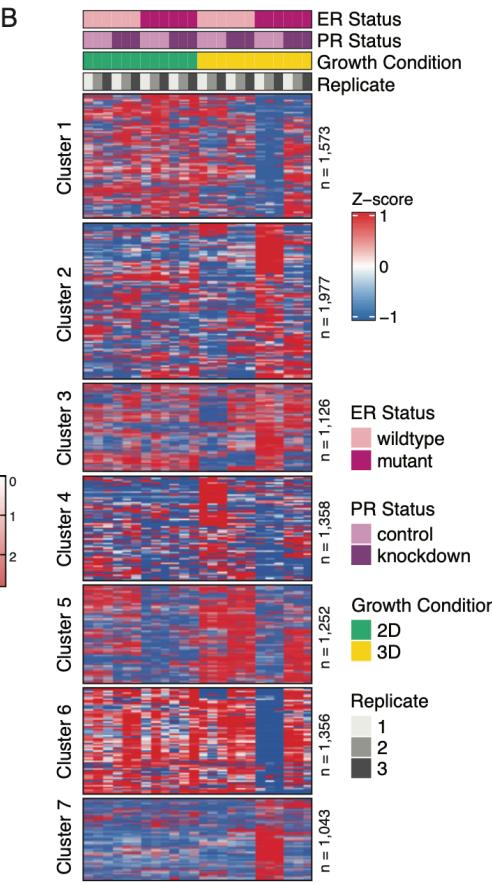
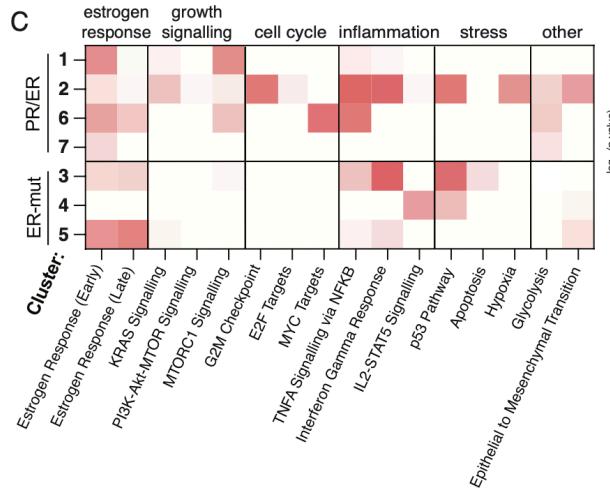
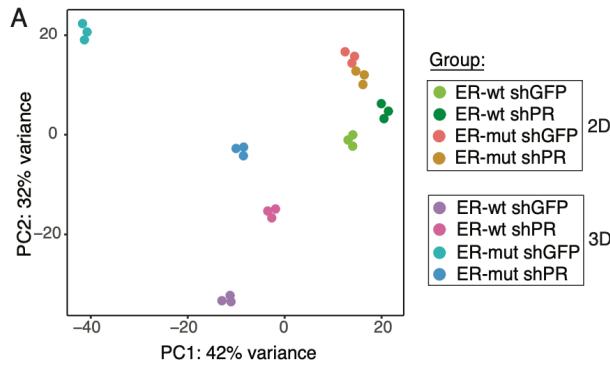
```
RNA-seq_pipeline_slurm.sh ● functional_RNA-seq_pipeline_slurm.sh ✘
16
17 # Load necessary modules – check what you have installed on your server vs. what is in your bash
18 # module load fastqc
19 # module load cutadapt
20 # module load star
21 # module load samtools
22 # module load htseq
23 # module load multiqc
24
25 # Create directory for FastQC reports
26 mkdir fastqc_reports
27
28 # Create directory for trimmed fastq files
29 mkdir trimmed_fastq
30
31 # Run FastQC on all input FASTQ files and perform adapter trimming with TrimGalore
32 for f in *_R1_001.fastq.gz; do #check the base name of your files and modify as needed
33 base=$(basename ${f} _R1_001.fastq.gz)
34 fastqc -o fastqc_reports ${base}_R1_001.fastq.gz ${base}_R2_001.fastq.gz; # Generate FastQC report
35 trim_galore --paired --illumina --fastqc --output_dir trimmed_fastq ${base}_R1_001.fastq.gz ${base}_R2_001.fastq.gz; # Trim
36 done
37
38 # Make a directory for the bam files
39 mkdir bam_files
40
41 # Run STAR on the trimmed FASTQ files
42 for f in trimmed_fastq/*_R1_001_val_1.fq.gz; do
43 base=$(basename ${f} _R1_001_val_1.fq.gz); # Extract basename of input file
44 r2="${base}''_R2_001_val_2.fq.gz"
45 STAR --runThreadN 12 \
46 --genomeDir /home/langeaca/gilli431/software/STAR_hg38 \
47 --readFilesCommand zcat \
48 --readFilesIn $f trimmed_fastq/$r2 \
49 --outFileNamePrefix bam_files/${base}_
50 --outSAMtype BAM \
51
52 # Sort and index the BAM files
53 samtools sort -o bam_files/${base}_sorted.bam bam_files/${base}_Aligned.sortedByCoord.out.bam
54 samtools index bam_files/${base}_sorted.bam
55
56 # Make a directory for the count files
57 mkdir count_files
```

1. Put the raw data for the experiment you are analyzing into a single folder (on MSI).
2. Modify the SLURM header in the pipeline script
3. Send the job out for analysis.
4. Wait 12-24 hours.
5. Check your QC metrics and retrieve the count files.

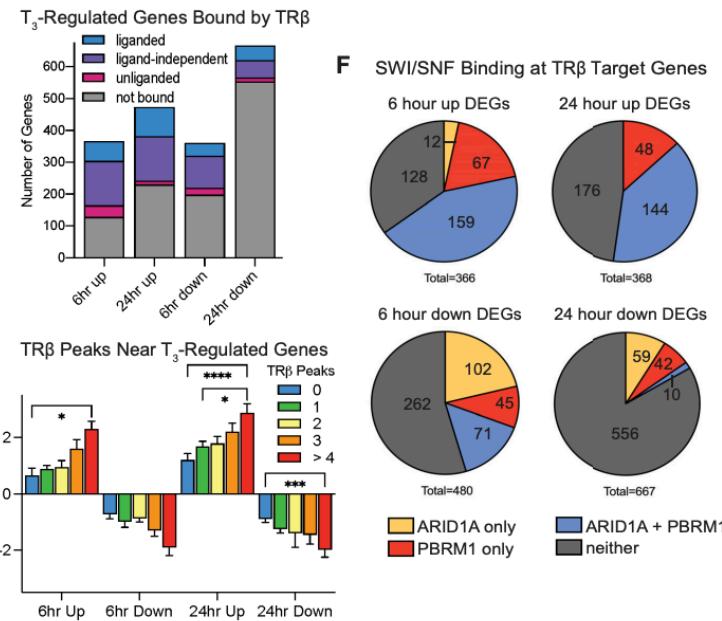
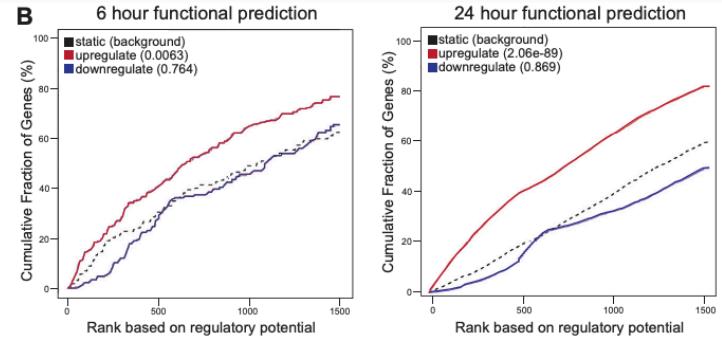
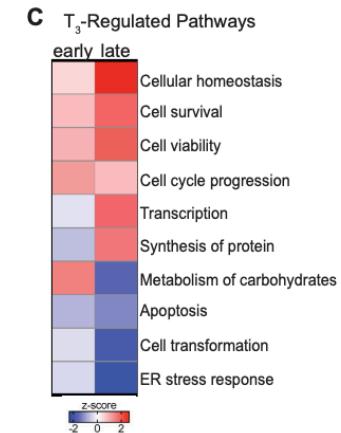
What can be done with the processed data?



What can be done with the processed data?



● 6hr up; n=366 ● 5hr down; n=368
 ● 24hr up; n=480 ● 24hr down; n=667
 ● unchanged; n=22,846



Module 3: Hands-On Workshop

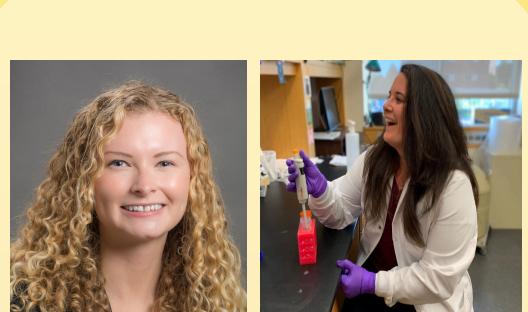


Group A

Is this your first time analyzing an RNA-Seq dataset?



Participants without a laptop or are new to RNA-Seq are encouraged to join Group A



Group B

Do you have previous experience with R/Rstudio?

Participants with a laptop are encouraged to join Group B



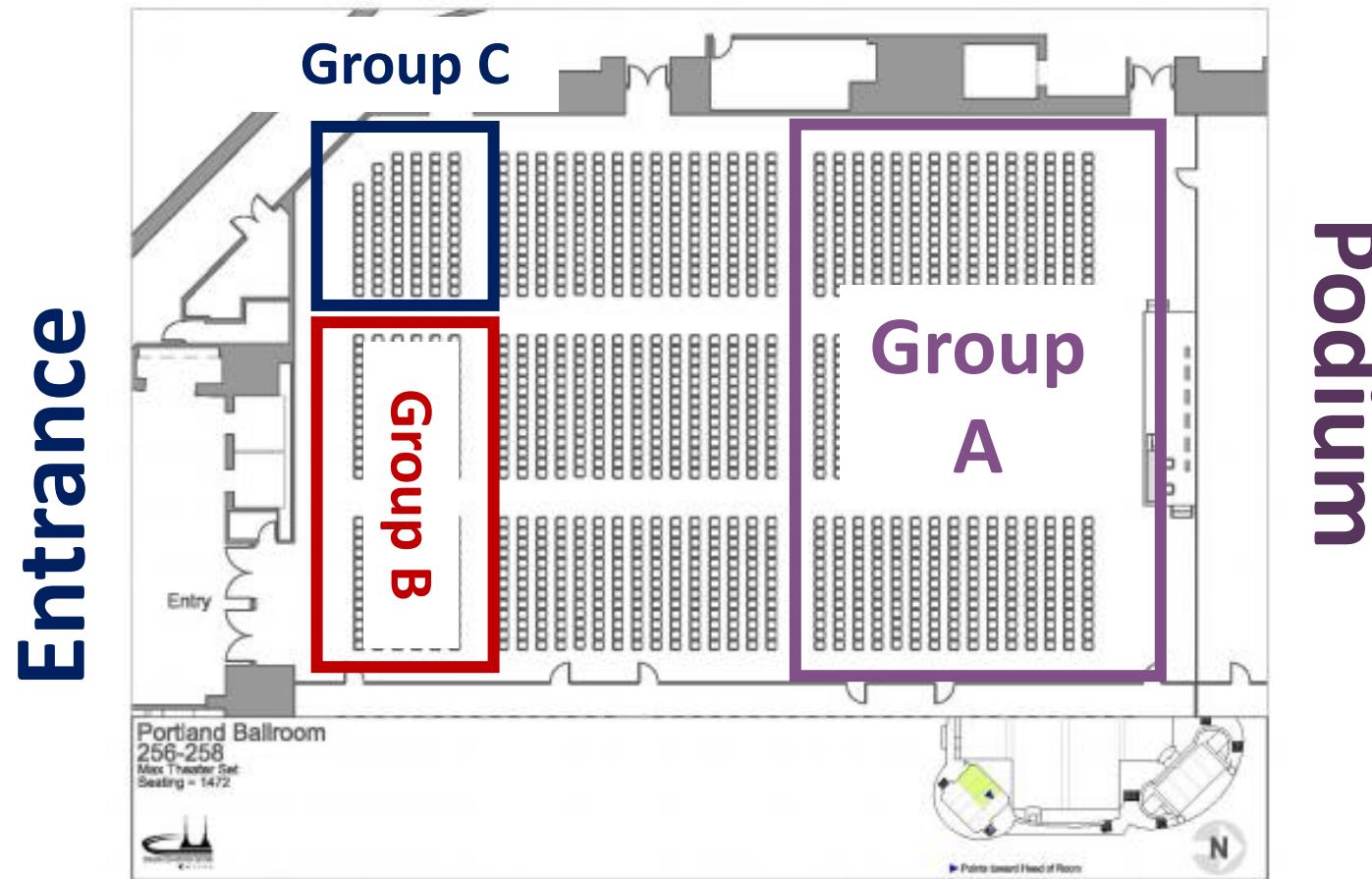
Group C

Do you have questions about an RNA-Seq dataset you generated or would like to analyze?

Informal discussion about NGS



Hands-On Workshop



Entrance

Podium

Please take a few minutes to join a group!

