

Bioinformatics for RNA sequencing analysis

Wenwen Huo, Postdoctoral Scholar, School of Medicine,
Isberg lab

Rebecca Batorsky, Senior Bioinformatics Specialist

Albert Tai, Research Assistant Professor of Immunology

June 2nd 2020

Using command line and R via OnDemand

Structure of Tufts HPC
cluster

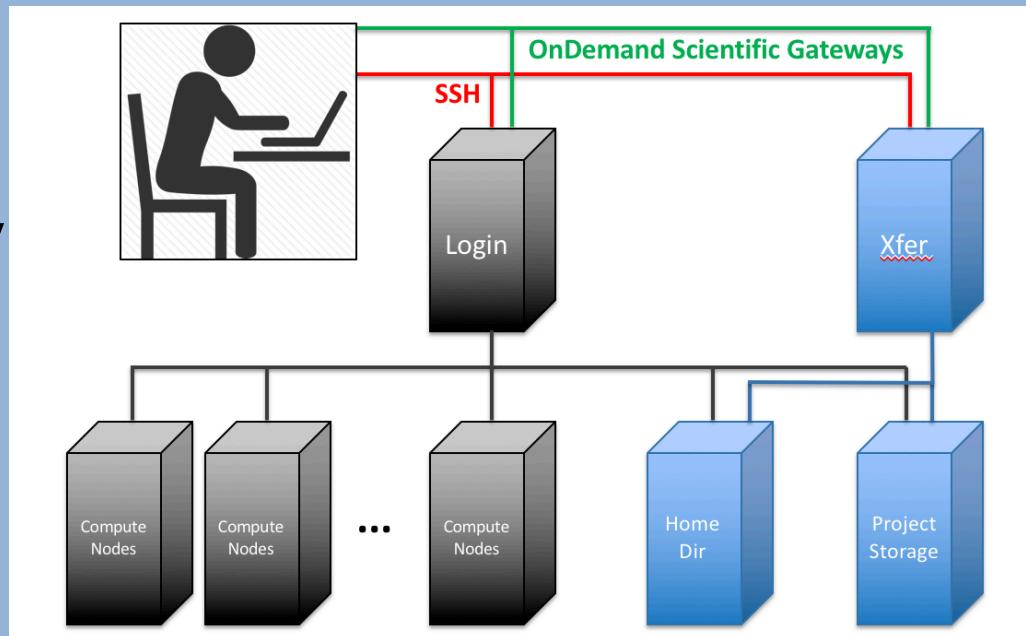
Command Line Basics

R Basics

Tufts High Performance Compute Cluster

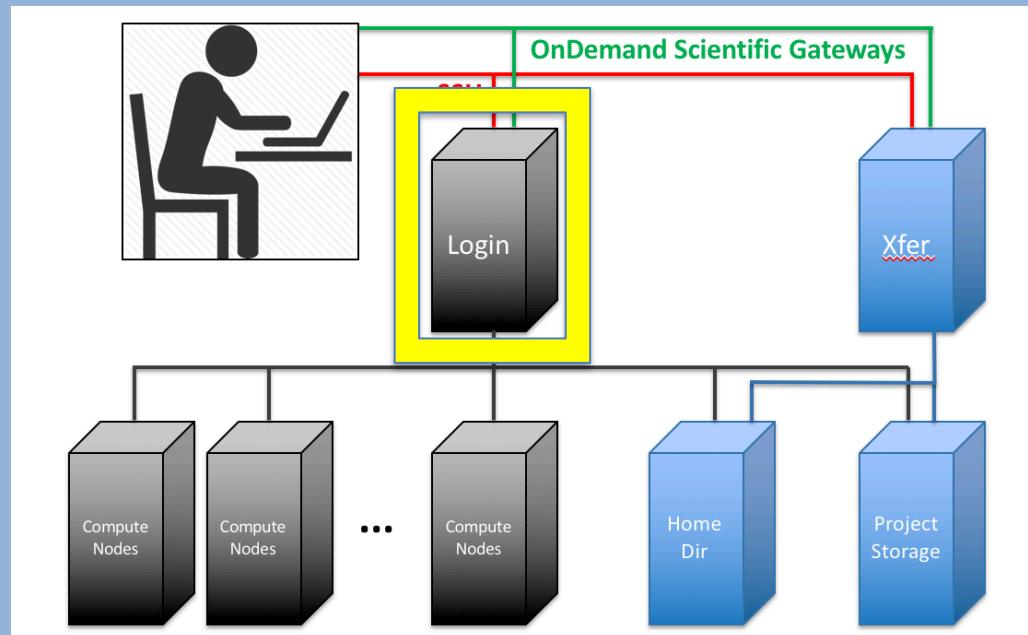
- **Cluster**

Cluster computing is the result of connecting many local computers (nodes) together via a high speed connection to provide a single shared resource.



Structure of Tufts HPC Cluster

- **Login Node**



- OnDemand login

➤ Go to:

<https://ondemand.cluster.tufts.edu>

➤ login with your UTLN and password

➤ Clusters/Tufts HPC (FastX11) Shell Access

Open OnDemand Files ▾ Jobs ▾ Clusters ▾ Interactive Apps ▾ Misc ▾ Help ▾ Logged in as whuo01 Log Out

OUTAGE NOTIFICATIONS and Alerts

➤ Tufts HPC Shell Access ←

➤ Tufts HPC FastX11 Shell Access

- Cluster Outage: Currently, there is no outage on the Tufts Cluster.
- Request Assistance: Email tts-research@tufts.edu for questions regarding Tufts High Performance Compute (HPC) cluster.
- Upload/Download: Via OnDemand web interface is limited to 976MB which will be increased in the future.

Structure of Tufts HPC Cluster

- **Login Node**

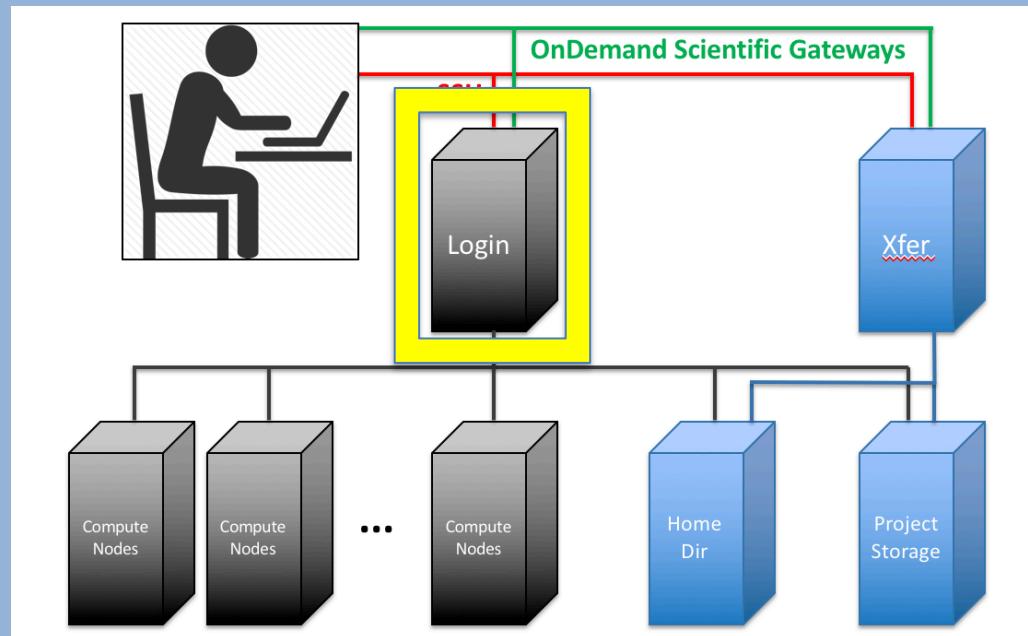
- OnDemand login

➤ Go to:

<https://ondemand.cluster.tufts.edu>

➤ login with your UTLN and password

➤ Clusters/Tufts HPC (FastX11) Shell Access



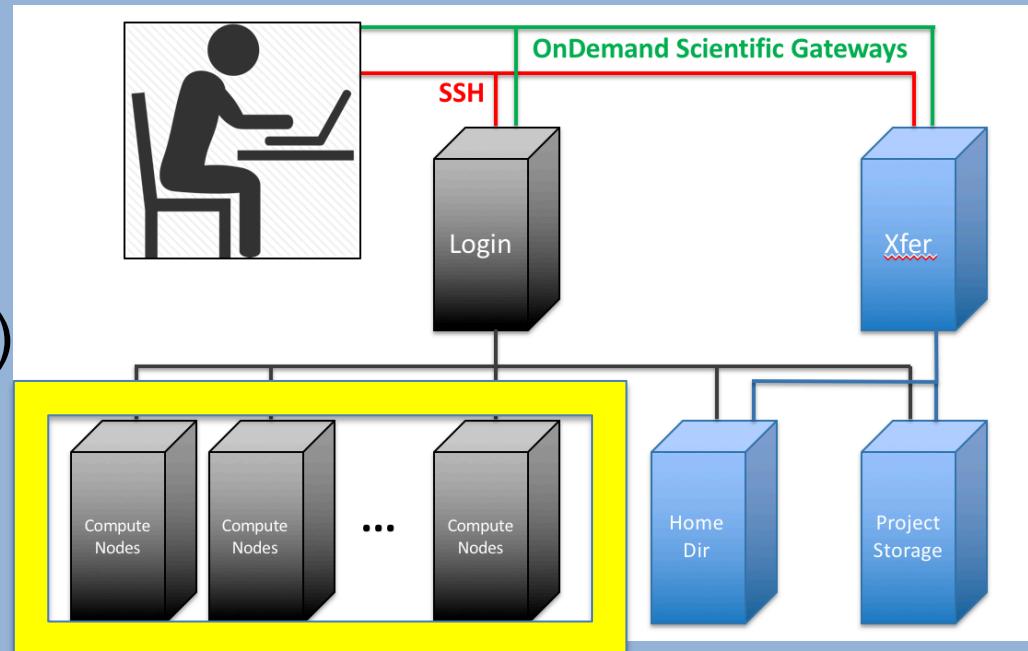
```
(base) [whuo01@login001 ~]$
```

Structure of Tufts HPC Cluster

- **Compute Nodes**

Different partitions:

- batch (3days)
- gpu (3days)
- interactive (4hours)
- largemem (7days)
- mpi (7days)
- m4 (7days)**
- long (21days)***



- Please **ALWAYS** run your programs/jobs on compute nodes!

time: hr:min:sec Memory requested per node
 ↑
 ↑
 ↑

nodes requested

```
(base) [whuo01@login001 ~]$ srun -t 3:00:00 --mem 16G -N 1 -n 4 -p preempt --reservation bioworkshop --pty bash
```

↑
pseudo terminal runs bash shell

of tasks

partition

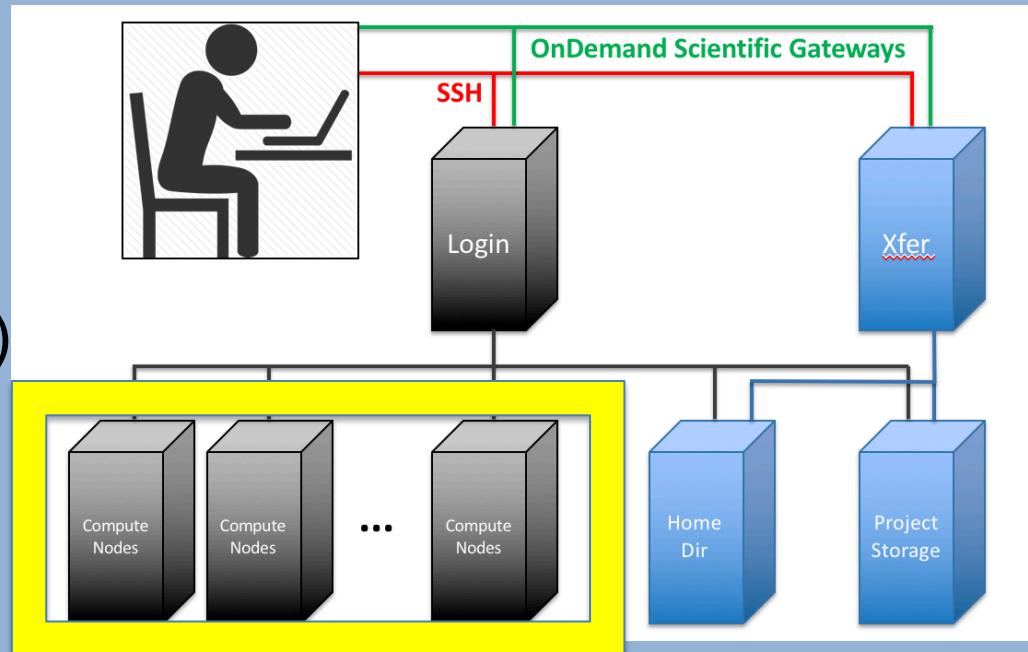
reserved spot

Structure of Tufts HPC Cluster

- **Compute Nodes**

Different partitions:

- batch (3days)
- gpu (3days)
- interactive (4hours)
- largemem (7days)
- mpi (7days)
- m4 (7days)**
- long (21days)***



- Please **ALWAYS** run your programs/jobs on compute nodes!

```
(base) [whuo01@login001 ~]$ srun -t 3:00:00 --mem 16G -N 1 -n 4 -p preempt --reservation bioworkshop --pty bash
(base) [whuo01@pcomp41 ~]$
```

Structure of Tufts HPC Cluster

- **Xfer Node**

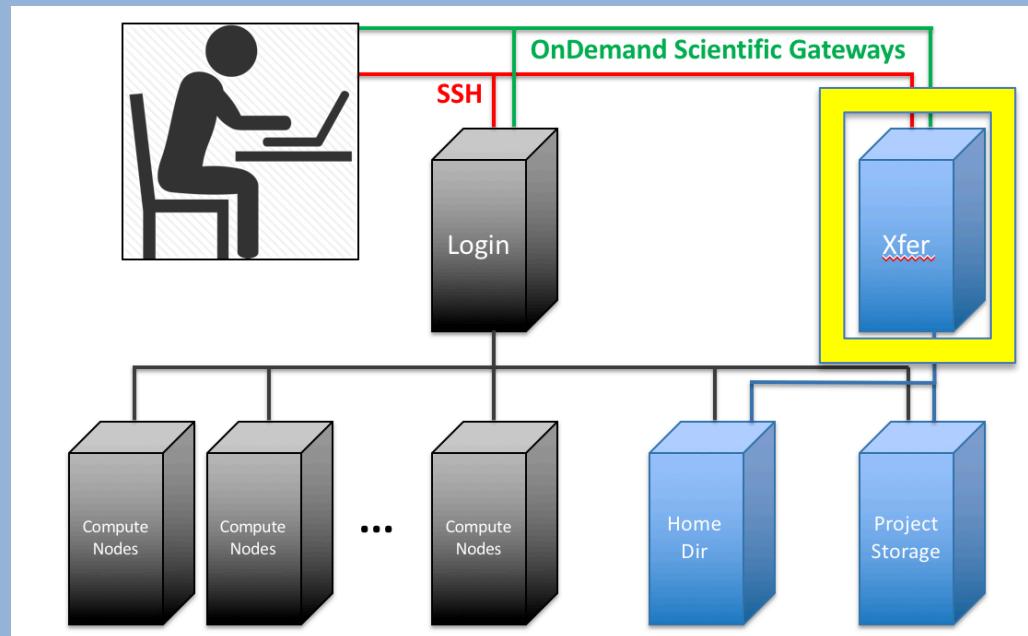
- data transfer node
- upload/download

- **OnDemand**

- Go to:

<https://ondemand.cluster.tufts.edu>

- Files (Upload/Download/View/Edit ...)



Open OnDemand Files ▾ Jobs ▾ Clusters ▾ Interactive Apps ▾ Misc ▾ ? Help ▾ Logged in as whuo01 Log Out

OUTAGE NOTIFICATION **SUPPORT REQUEST**

• **Cluster Outage:** Currently, there is no outage on the Tufts Cluster.

• **Request Assistance:** Email tts-research@tufts.edu for questions regarding Tufts High Performance Compute (HPC) cluster.

• **Upload/Download:** Via OnDemand web interface is limited to 976MB which will be increased in the future.

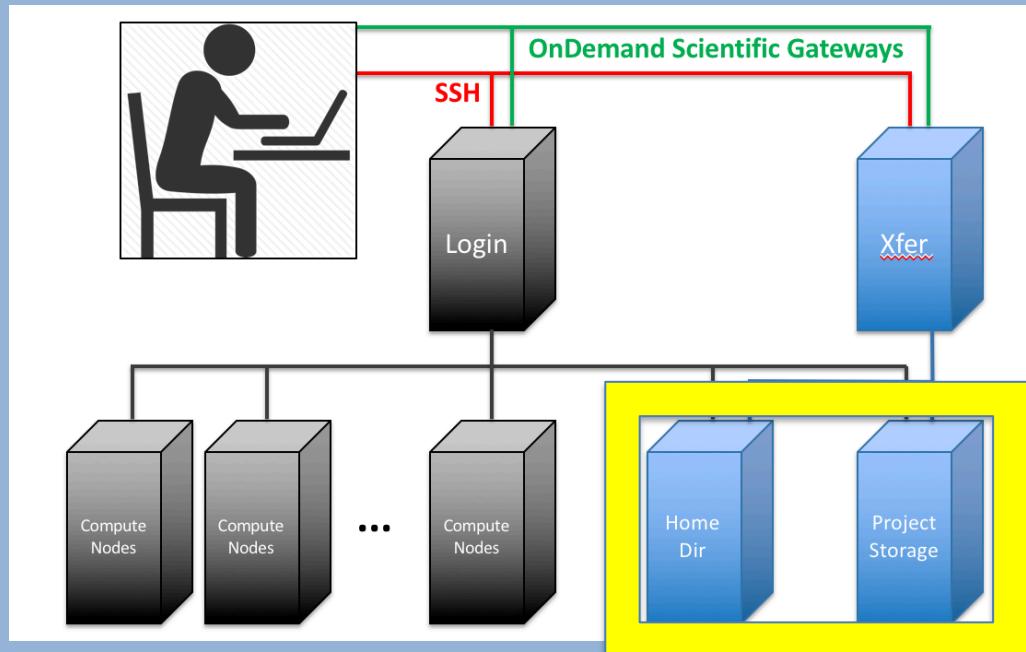
Structure of Tufts HPC Cluster

- **Home Directory**

- /cluster/home/utln
- 5GB (small!)
- Visible only to you

- **Project Storage**

- /cluster/tufts/lab_n
- 50GB+
- Collaboration friendly – visible to all lab members



OUTAGE NOTE

Home Directory Projects ← **SUPPORT REQUEST**

- **Cluster Outage:** Currently, there is no outage on the Tufts Cluster.
- **Request Assistance:** Email tts-research@tufts.edu for questions regarding Tufts High Performance Compute (HPC) cluster.
- **Upload/Download:** Via OnDemand web interface is limited to 976MB which will be increased in the future.

Structure of Tufts HPC Cluster

• Project Storage

/cluster/tufts/bio/tools/training/intro-to-rnaseq/users/

File Explorer v1.3.6

/cluster/tufts/

Go To... Open in Terminal New File New Dir Upload Show Dotfiles Show Owner/Mode

View Edit A-Z Rename Download Copy Paste (Un)Select All Delete

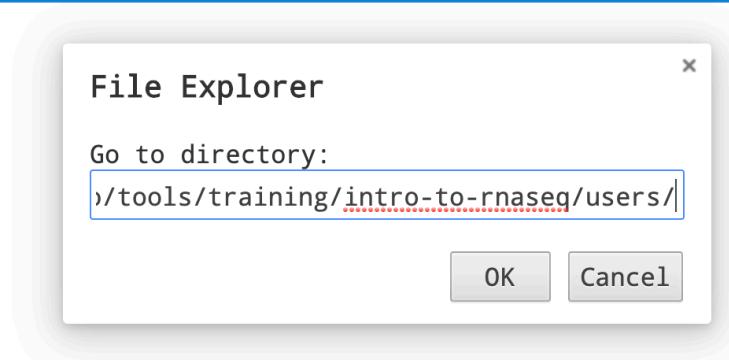
name size modified date

| name | size | modified date |
|---------------|-------|---------------|
| .. | <dir> | 10/15/2019 |
| abriolalab | <dir> | 05/14/2018 |
| acevedolab | <dir> | 05/26/2020 |
| aeronlab | <dir> | 05/28/2020 |
| aldrigelab | <dir> | 03/27/2018 |
| atlas16 | <dir> | 12/04/2019 |
| baiselab | <dir> | 05/15/2020 |
| beaucheminlab | <dir> | 03/21/2020 |
| berd | <dir> | 09/24/2018 |
| bio | <dir> | 05/04/2020 |
| blacklab | <dir> | 07/23/2019 |
| bohm | <dir> | 05/02/2020 |
| brovmanlab | <dir> | 02/21/2020 |
| cabcs | <dir> | 12/03/2019 |
| chatterjeelab | <dir> | 08/16/2018 |
| chiesalab | <dir> | 10/17/2019 |
| chinlab | <dir> | 10/31/2018 |
| chunglab | <dir> | 01/02/2020 |
| class | <dir> | 03/21/2020 |
| cochranlab | <dir> | 05/28/2020 |
| coffinlab | <dir> | 07/22/2019 |
| cohnlab | <dir> | 05/11/2020 |
| cordilab | <dir> | |

File Explorer

Go to directory: /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/

OK Cancel



Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Path: which folder you are in and how to change folder

pwd: print the current working directory

ls : list content in the folder

cd : change directory

cd .. : go back one directory

cd ../../ : go back two directories

```
(base) [whuo01@login001 ~]$ pwd  
/cluster/home/whuo01  
(base) [whuo01@login001 ~]$ cd /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/  
(base) [whuo01@login001 users]$ mkdir whuo01  
(base) [whuo01@login001 users]$ cd whuo01  
(base) [whuo01@login001 whuo01]$ cp /cluster/tufts/bio/tools/training/intro-to-rnaseq/intro-to-RNA-seq-May-2020.tar.gz .  
(base) [whuo01@login001 whuo01]$ tar -xvzf intro-to-RNA-seq-May-2020.tar.gz  
intro-to-RNA-seq/  
intro-to-RNA-seq/ERP004763_info.txt  
intro-to-RNA-seq/raw_data/ intro-to-RNA-seq/raw_data/sample_info.txt  
intro-to-RNA-seq/raw_data/WT/  
intro-to-RNA-seq/raw_data/WT/ERR458499.fastq.gz  
intro-to-RNA-seq/raw_data/WT/ERR458498.fastq.gz  
...  
...  
(base) [whuo01@login001 whuo01]$ cd intro-to-RNA-seq  
(base) [whuo01@login001 intro-to-RNA-seq]$ ls  
ERP004763_info.txt      raw_data      scripts  
(base) [whuo01@login001 intro-to-RNA-seq]$
```

Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Path: which folder you are in and how to change folder

pwd: get the current directory

ls : show content in the folder

ls -ltr : show content in list format with timestamp

cd : go to a specific directory

cd .. : go back one directory

tree : list all contents in a folder

```
(base) [whuo01@login001 ~]$ cd /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/whuo01/intro-to-RNA-seq  
(base) [whuo01@login001 intro-to-RNA-seq]$ tree
```

```
ERP004763_info.txt  
raw_data  
|   sample_info.txt  
SNF2  
|   ERR458500.fastq.gz  
|   ERR458501.fastq.gz  
|   ERR458502.fastq.gz  
|   ERR458503.fastq.gz  
|   ERR458504.fastq.gz  
|   ERR458505.fastq.gz  
|   ERR458506.fastq.gz  
WT  
|   ERR458493.fastq.gz  
|   ERR458494.fastq.gz  
|   ERR458495.fastq.gz  
|   ERR458496.fastq.gz  
|   ERR458497.fastq.gz  
|   ERR458498.fastq.gz  
|   ERR458499.fastq.gz  
scripts  
|   fastqc.sh  
|   featurecounts.sh  
|   featurecount_stat.R  
|   intro.R  
|   mapping_percentage.R  
|   star_align_individual.sh  
|   star_align_practice.sh
```

4 directories, 23 files

```
(base) [whuo01@login001 intro-to-RNA-seq]$
```

Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Path: which folder you are in and how to change folder

pwd: get the current directory

ls : show content in the folder

ls -ltr : show content in list format with timestamp

cd : go to a specific directory

cd .. : go back one directory

tree : list all contents in a folder

2. Run scripts:

sh test.sh # run a script named test.sh

sh ./test.sh # run test.sh in the current folder

sh ./subfolder/test.sh # run test.sh in the subfolder

(base) [whuo01@login001 intro-to-RNA-seq]\$ tree

```
ERP004763_info.txt
raw_data
|   sample_info.txt
SNF2
|   ERR458500.fastq.gz
|   ERR458501.fastq.gz
|   ERR458502.fastq.gz
|   ERR458503.fastq.gz
|   ERR458504.fastq.gz
|   ERR458505.fastq.gz
|   ERR458506.fastq.gz
WT
|   ERR458493.fastq.gz
|   ERR458494.fastq.gz
|   ERR458495.fastq.gz
|   ERR458496.fastq.gz
|   ERR458497.fastq.gz
|   ERR458498.fastq.gz
|   ERR458499.fastq.gz
scripts
|   fastqc.sh
|   featurecounts.sh
|   featurecount_stat.R
|   intro.R
|   mapping_percentage.R
|   star_align_individual.sh
|   star_align_practice.sh
```

4 directories, 23 files

(base) [whuo01@login001 intro-to-RNA-seq]\$ sh fastqc.sh
sh: fastqc.sh: No such file or directory

```
(base) [whuo01@login001 intro-to-RNA-seq]$ tree
```

```
ERP004763_info.txt
raw_data
└── sample_info.txt
SNF2
├── ERR458500.fastq.gz
├── ERR458501.fastq.gz
├── ERR458502.fastq.gz
├── ERR458503.fastq.gz
├── ERR458504.fastq.gz
├── ERR458505.fastq.gz
└── ERR458506.fastq.gz
WT
├── ERR458493.fastq.gz
├── ERR458494.fastq.gz
├── ERR458495.fastq.gz
├── ERR458496.fastq.gz
├── ERR458497.fastq.gz
├── ERR458498.fastq.gz
└── ERR458499.fastq.gz
scripts
├── fastqc.sh
├── featurecounts.sh
├── featurecount_stat.R
├── intro.R
├── mapping_percentage.R
├── star_align_individual.sh
└── star_align_practice.sh
```

4 directories, 23 files

```
(base) [whuo01@login001 intro-to-RNA-seq]$ sh fastqc.sh
```

sh: fastqc.sh: No such file or directory

```
(base) [whuo01@login001 intro-to-RNA-seq]$ sh scripts/fastqc.sh
```

Started analysis of ERR458493.fastq.gz

Approx 5% complete for ERR458493.fastq.gz

Approx 10% complete for ERR458493.fastq.gz

Approx 15% complete for ERR458493.fastq.gz

^C 20% complete for ERR458493.fastq.gz

Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Path: which folder you are in and how to change folder

pwd: get the current directory

ls : show content in the folder

ls -ltr : show content in list format with timestamp

cd : go to a specific directory

cd .. : go back one directory

2. Run scripts:

sh test.sh # run a script named test.sh

sh ./test.sh # run test.sh in the current folder

sh ./subfolder/test.sh # run test.sh in the subfolder

3. Edit scripts or other files using nano:

nano filename: enters the file named filename and creates one if filename does not exist

nano ./subfolder/filename: similar with above

```
(base) [whuo01@login001 intro-to-RNA-seq]$ cp ./scripts/star_align_practice.sh ./scripts/star_align_ERR458500.sh  
(base) [whuo01@login001 intro-to-RNA-seq]$ nano ./scripts/star_align_ERR458500.sh
```

GNU nano 2.0.9

File: scripts/star_align_ERR458500.sh

```
## Use STAR aligner to align fastq files
module load STAR/2.6.1d
mkdir -p STAR

## Fastq files to align, separated by commas for multiple lanes of a single sample
FASTQ="raw_data/WT/ERR458493.fastq.gz"

## Name the output file
OUT="WT_ERR458493"

## Defing reference genome directory
REF_DIR="/cluster/tufts/bio/data/genomes/Saccharomyces_cerevisiae/UCSC/sacCer3"

# execute STAR in the runMode "alignReads"
STAR --genomeDir ${REF_DIR}/Sequence/STAR \
--readFilesIn ${FASTQ} \
--readFilesCommand zcat \
--outFileNamePrefix STAR/${OUT}_ \
--outFilterMultimapNmax 1 \
--outSAMtype BAM SortedByCoordinate \
--runThreadN 4 \
--alignIntronMin 1 \
--alignIntronMax 2500 \
--sjdbGTFfile ${REF_DIR}/Annotation/Genes/sacCer3.gtf \
--sjdbOverhang 49
```

Navigate



Exit
Ctrl + x

[Read 28 lines]

| | | | | | |
|----------|-------------|--------------|--------------|---------------|-------------|
| Get Help | ^O WriteOut | ^R Read File | ^Y Prev Page | ^K Cut Text | ^C Cur Pos |
| Exit | ^J Justify | ^W Where Is | ^V Next Page | ^U UnCut Text | ^T To Spell |

Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Path: which folder you are in and how to change folder

pwd: get the current directory

ls : show content in the folder

ls -ltr : show content in list format with timestamp

cd : go to a specific directory

cd .. : go back one directory

2. Run scripts:

sh test.sh # run a script named test.sh

sh ./test.sh # run test.sh in the current folder

sh ./subfolder/test.sh # run test.sh in the subfolder

3. Edit scripts or other files using nano:

nano filename: enters the file named filename and creates one if filename does not exist

nano ./subfolder/filename: similar with above

4. Quick view a file:

head filename: show the top 10 lines of a file

cat filename: print all the content in filename to screen

cat filename | head: chaining commands using |

```
(base) [whuo01@login001 intro-to-RNA-seq]$ cp ./scripts/star_align_practice.sh ./scripts/star_align_ERR458500.sh  
(base) [whuo01@login001 intro-to-RNA-seq]$ nano ./scripts/star_align_ERR458500.sh  
(base) [whuo01@login001 intro-to-RNA-seq]$ cat ./scripts/star_align_ERR458500.sh
```

```
(base) [whuo01@login001 intro-to-RNA-seq]$ cp ./scripts/star_align_practice.sh ./scripts/star_align_ERR458500.sh
(base) [whuo01@login001 intro-to-RNA-seq]$ nano ./scripts/star_align_ERR458500.sh
(base) [whuo01@login001 intro-to-RNA-seq]$ cat ./scripts/star_align_ERR458500.sh
## Use STAR aligner to align fastq files
module load STAR/2.6.1d
mkdir -p STAR

## Fastq files to align, separated by commas for multiple lanes of a single sample
FASTQ="raw_data/SNF2/ERR458500.fastq.gz"

## Name the output file
OUT="SNF2_ERR458500"

## Defing reference genome directory
REF_DIR="/cluster/tufts/bio/data/genomes/Saccharomyces_cerevisiae/UCSC/sacCer3"

# execute STAR in the runMode "alignReads"
STAR --genomeDir ${REF_DIR}/Sequence/STAR \
--readFilesIn ${FASTQ} \
--readFilesCommand zcat \
--outFileNamePrefix STAR/${OUT}_ \
--outFilterMultimapNmax 1 \
--outSAMtype BAM SortedByCoordinate \
--runThreadN 4 \
--alignIntronMin 1 \
--alignIntronMax 2500 \
--sjdbGTFfile ${REF_DIR}/Annotation/Genes/sacCer3.gtf \
--sjdbOverhang 49
```

comments start with

Oriniginal: WT/ERR458493.fastq.gz

Oriniginal: WT_ERR458493

```
(base) [whuo01@login001 intro-to-RNA-seq]$ sh ./scripts/star_align_ERR458500.sh
```

Takes ~ 5 to 10 min to finish

Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Run R script without using Rstudio:

module load R/3.6.3

Rscript sometest.R # run script that's written in R

Rscript ./subfolder/sometest.R # run script that's

in subfolder

Using command line and R via OnDemand

Command Line
Basics

R Basics

1. Run R script without using Rstudio:

```
module load R/3.6.3
```

```
Rscript sometest.R # run script that's written in R
```

```
Rscript ./subfolder/sometest.R # run script that's  
in subfolder
```

2. Working inside Rstudio (version 3.5.0):

```
# a). know your working directory
```

```
getwd(): get the current working directory
```

```
setwd("/cluster/tufts/bio/tools/training/intro-to-  
rnaseq/users/whuo01/intro-to-RNA-seq/") : set  
working directory
```

```
# b). load library
```

```
.libPaths('/cluster/tufts/bio/tools/R_libs/3.5')
```

```
library(DESeq2)
```

```
# c). assign variable
```

```
meta <- read.table("raw_data/sample_info.txt")
```

The screenshot shows the RStudio interface with several labeled components:

- Scripts editor**: The main workspace where R code is written. It contains the following script:

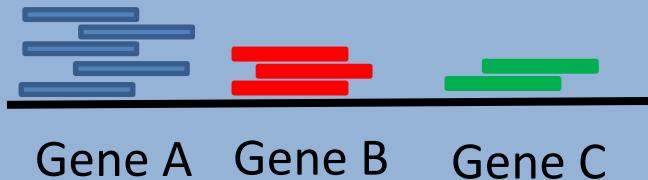
```
1 # get current directory
2 getwd()
3
4 # set to work directory
5 setwd("../")
6
7 # check current library path
8 .libPaths()
9
10 # load shared Tufts bio library path
11 .libPaths('/cluster/tufts/bio/tools/R_libs/3.5')
12
13 # load library
14 library(tidyverse)
15
16
```

- Environment**: Shows the Global Environment tab, which is currently empty.
- Console**: Shows the command `setwd("/cluster/tufts/bio/tools/training/intro-to-rnaseq/ users/YOUR_USERNAME/intro-to-RNA-seq/")` being run.
- files/plots/packages**: The Files pane, which lists the contents of the current directory:

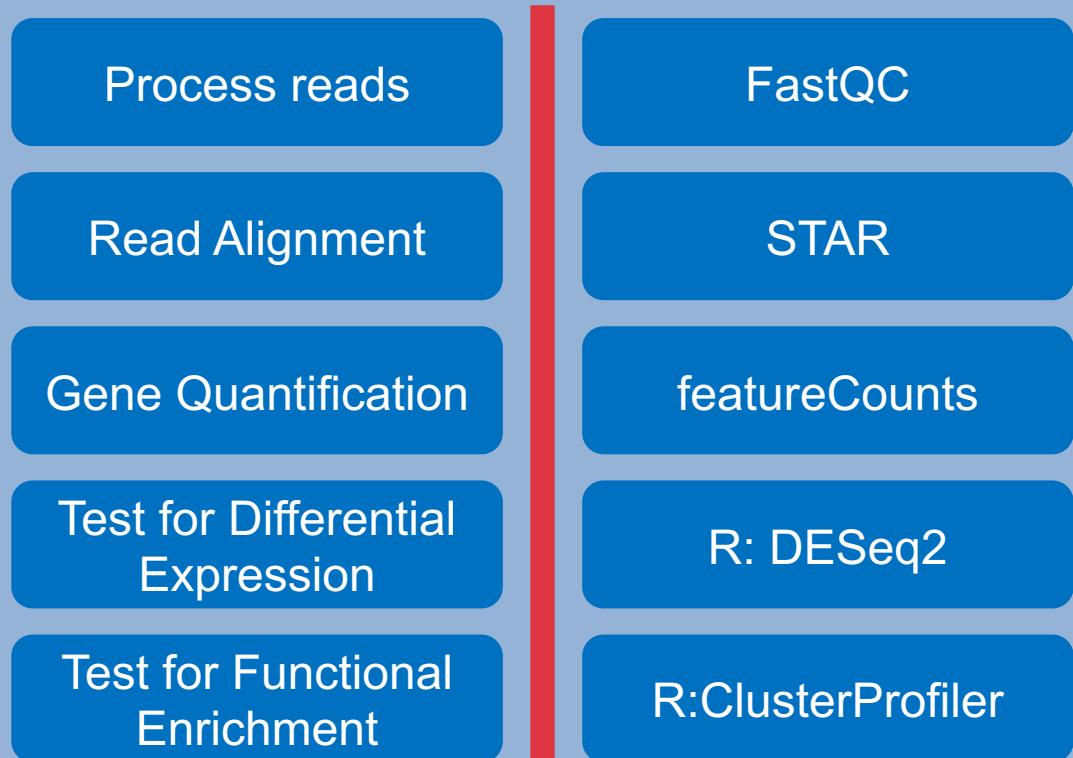
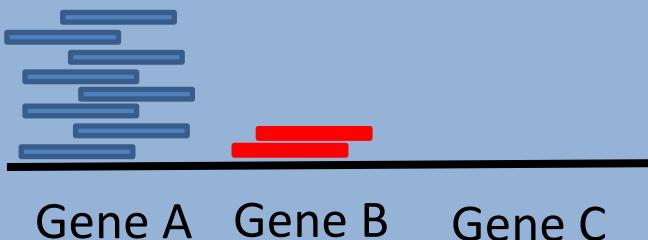
| Name | Size | Modified |
|----------------|---------|------------------------|
| .Rhistory | 13.8 KB | May 19, 2020, 11:09 PM |
| .gitignore | 192 B | Apr 26, 2020, 11:53 PM |
| .RData | 12.2 MB | Apr 11, 2020, 6:11 PM |
| install.sh | 213 B | Mar 7, 2019, 9:22 PM |
| tpp.cfg | 220 B | Jan 23, 2019, 2:49 PM |
| test.pl | 130 B | Apr 3, 2018, 5:04 PM |
| igv | | |
| ondemand | | |
| privatemodules | | |
| R | | |
| whuo1 | | |

Analysis pipeline

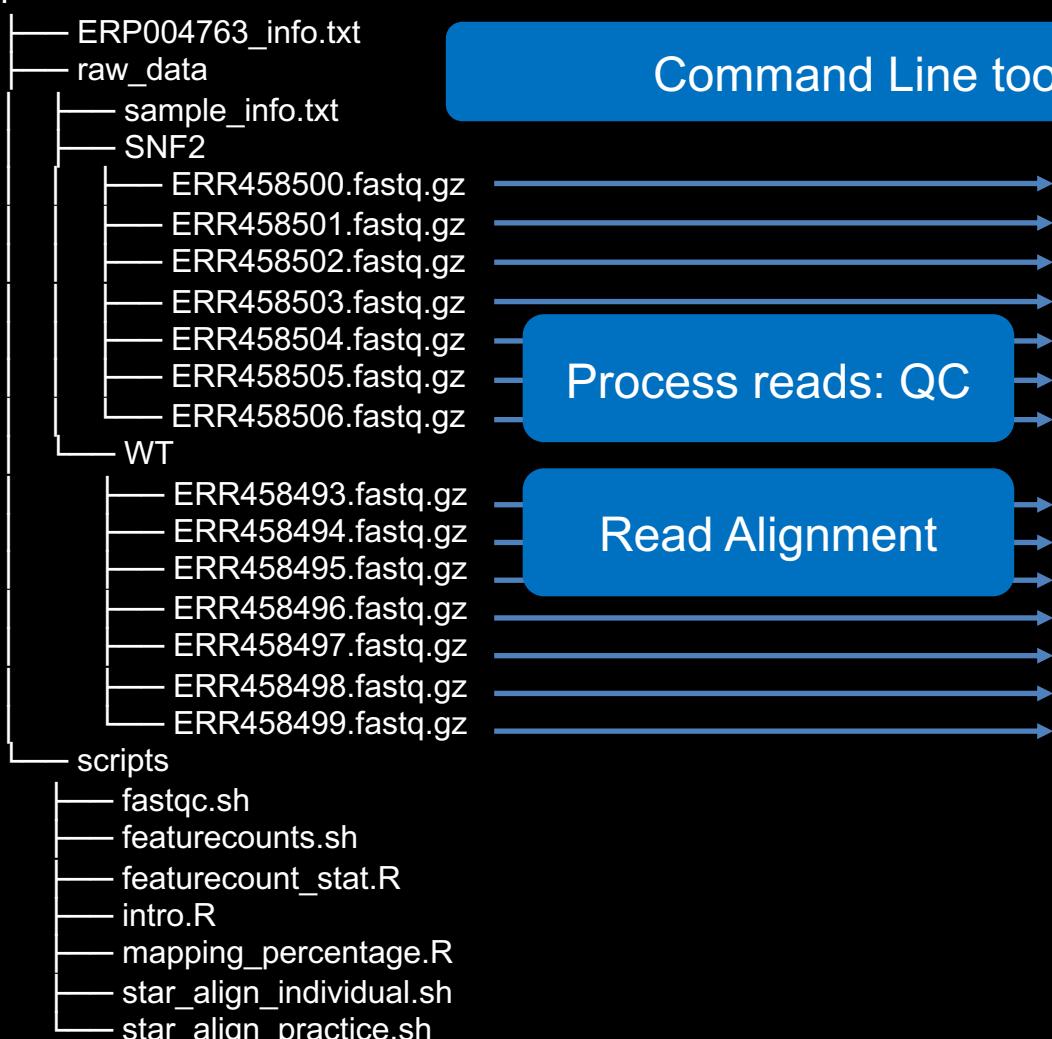
Yeast Wild Type Phenotype



Yeast Mutant Phenotype



```
(base) [whuo01@login001 ~]$ cd /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/whuo01/intro-to-RNA-seq  
(base) [whuo01@login001 intro-to-RNA-seq]$ tree
```



Command Line tools

R (Rstudio)

Test for Differential Expression

Test for Functional Enrichment

4 directories, 23 files

```
(base) [whuo01@login001 intro-to-RNA-seq]$
```

```
(base) [whuo01@omega001 intro-to-RNA-seq]$ module load fastqc  
(base) [whuo01@omega001 intro-to-RNA-seq]$ mkdir fastqc  
(base) [whuo01@omega001 intro-to-RNA-seq]$ fastqc raw_data/WT/*.fastq.gz -o fastqc --extract
```

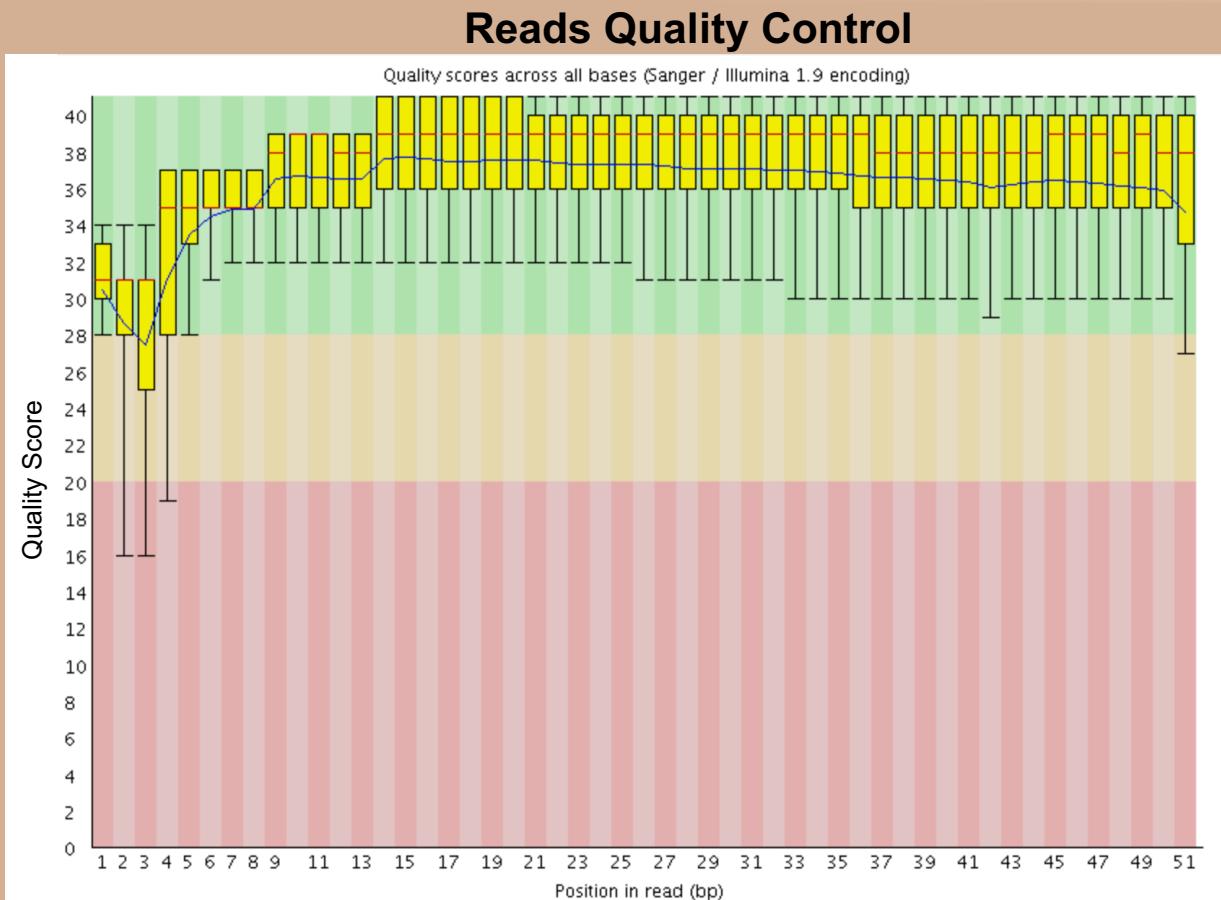
Process reads

Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment



```
(base) [whuo01@login001 ~]$ cd /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/whuo01/intro-to-RNA-seq  
(base) [whuo01@login001 intro-to-RNA-seq]$ tree
```

```
ERP004763_info.txt  
raw_data  
|   sample_info.txt  
SNF2  
|   ERR458500.fastq.gz  
|   ERR458501.fastq.gz  
|   ERR458502.fastq.gz  
|   ERR458503.fastq.gz  
|   ERR458504.fastq.gz  
|   ERR458505.fastq.gz  
|   ERR458506.fastq.gz  
WT  
|   ERR458493.fastq.gz  
|   ERR458494.fastq.gz  
|   ERR458495.fastq.gz  
|   ERR458496.fastq.gz  
|   ERR458497.fastq.gz  
|   ERR458498.fastq.gz  
|   ERR458499.fastq.gz  
scripts  
|   fastqc.sh  
|   featurecounts.sh  
|   featurecount_stat.R  
|   intro.R  
|   mapping_percentage.R  
|   star_align_individual.sh  
|   star_align_practice.sh
```

Command Line tools

R (Rstudio)

Test for Differential Expression

Test for Functional Enrichment

Process reads: QC

Read Alignment

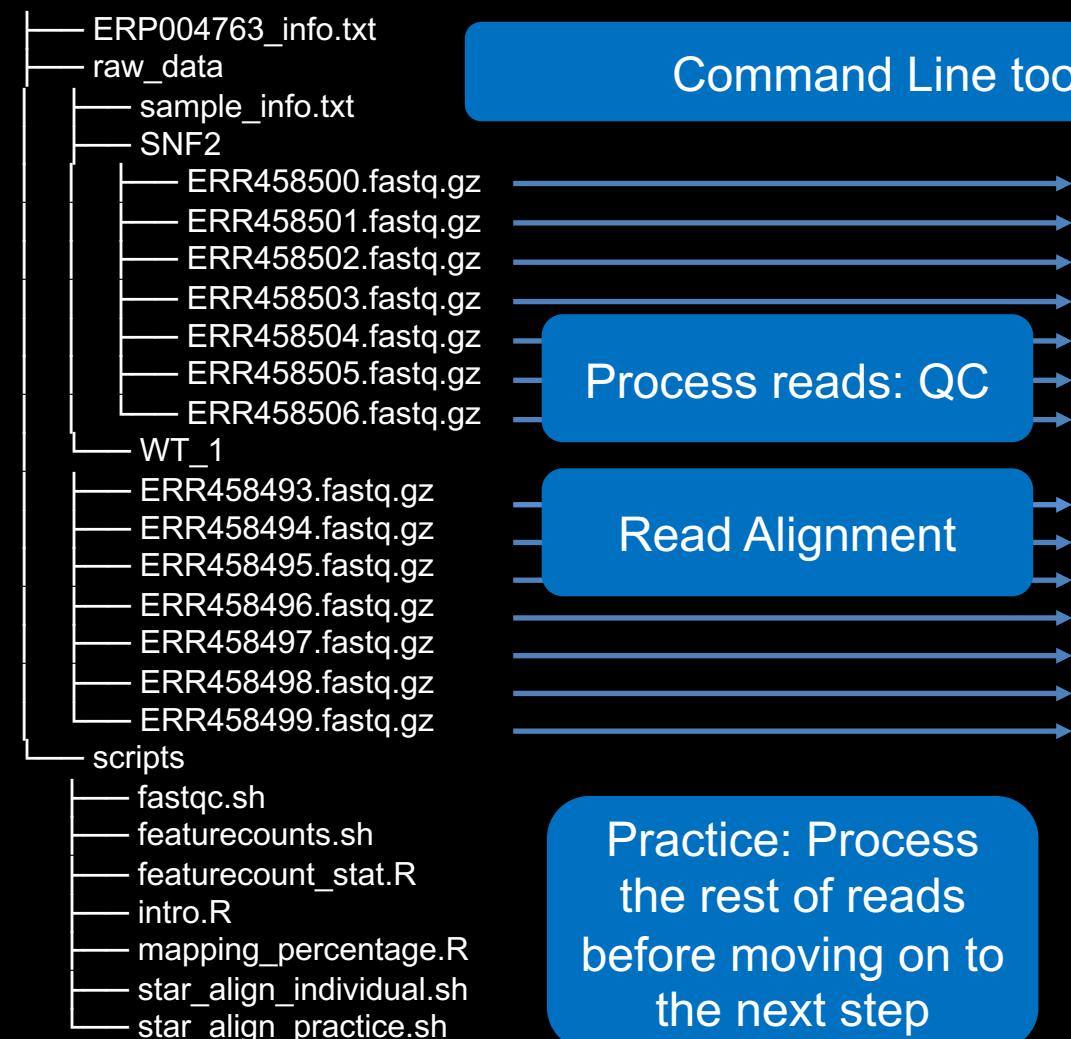
Gene Quantification using Feature Count

Practice: Process the rest of reads before moving on to the next step

4 directories, 23 files

```
(base) [whuo01@login001 intro-to-RNA-seq]$
```

```
(base) [whuo01@login001 ~]$ cd /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/whuo01/intro-to-RNA-seq  
(base) [whuo01@login001 intro-to-RNA-seq]$ tree
```



Command Line tools

R (Rstudio)

Process reads: QC

Read Alignment

Gene Quantification using Feature Count

Test for Differential Expression

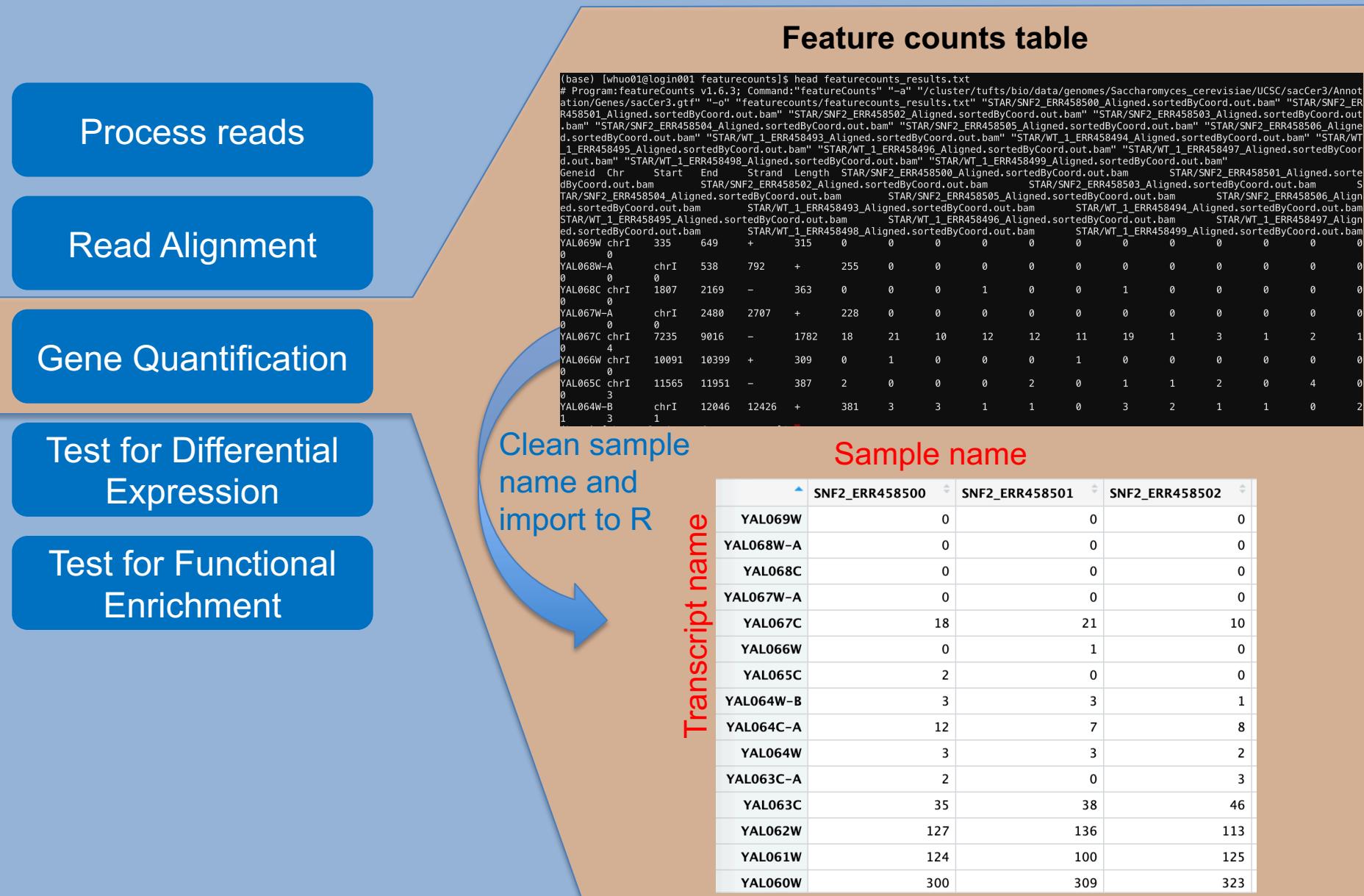
Test for Functional Enrichment

Practice: Process the rest of reads before moving on to the next step

4 directories, 23 files

```
(base) [whuo01@login001 intro-to-RNA-seq]$ sh ./scripts/star_align_individual.sh
```

```
(base) [whuo01@pcomp41 intro-to-RNA-seq]$ sh scripts/featurecounts.sh  
(base) [whuo01@pcomp41 intro-to-RNA-seq]$ head featurecounts/featurecounts_results.txt
```



Optional: Read alignment QC

Process reads

Read Alignment

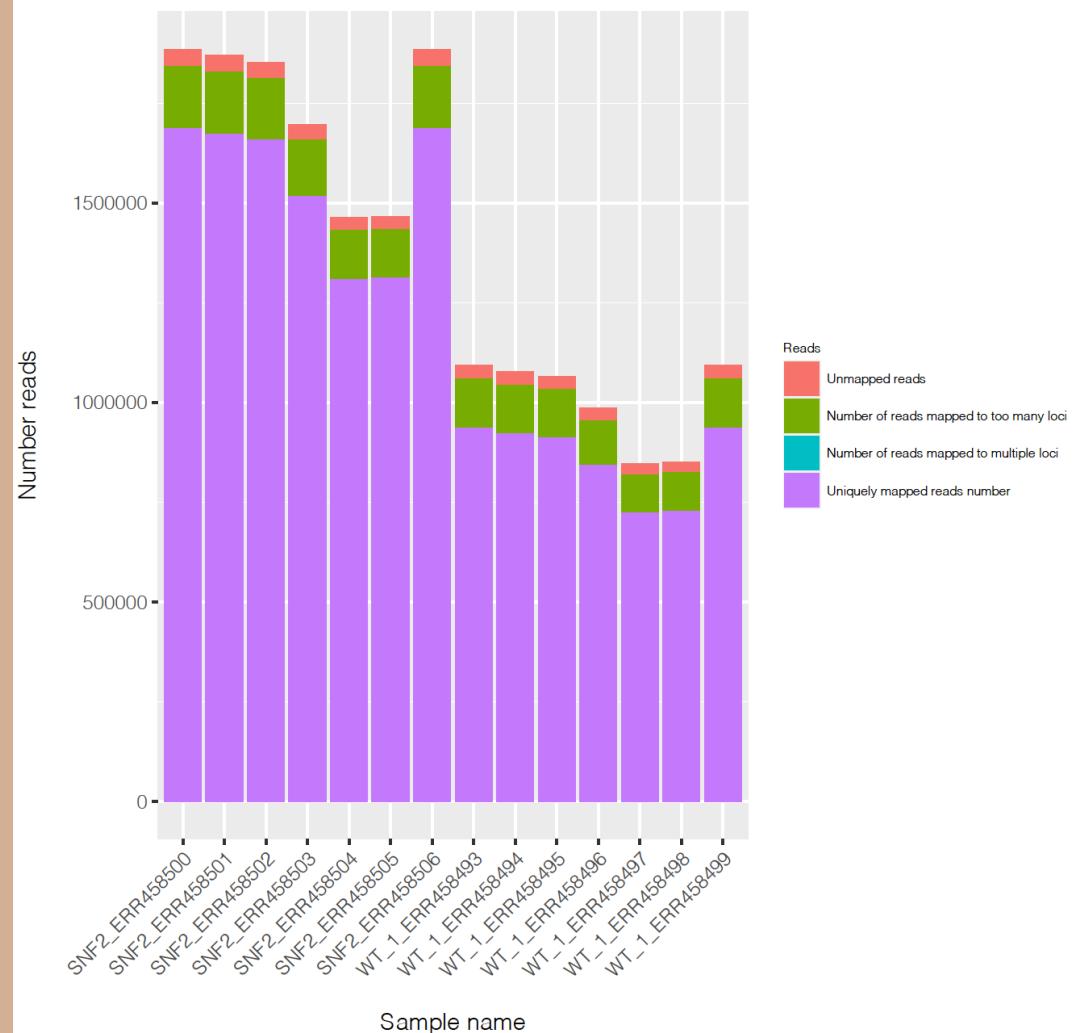
Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

Mapping Statistics

Mapping efficiency



(base) [whuo01@pcomp41 intro-to-RNA-seq]\$ Rscript scripts/mapping_percentage.R

Optional: Feature count QC

Process reads

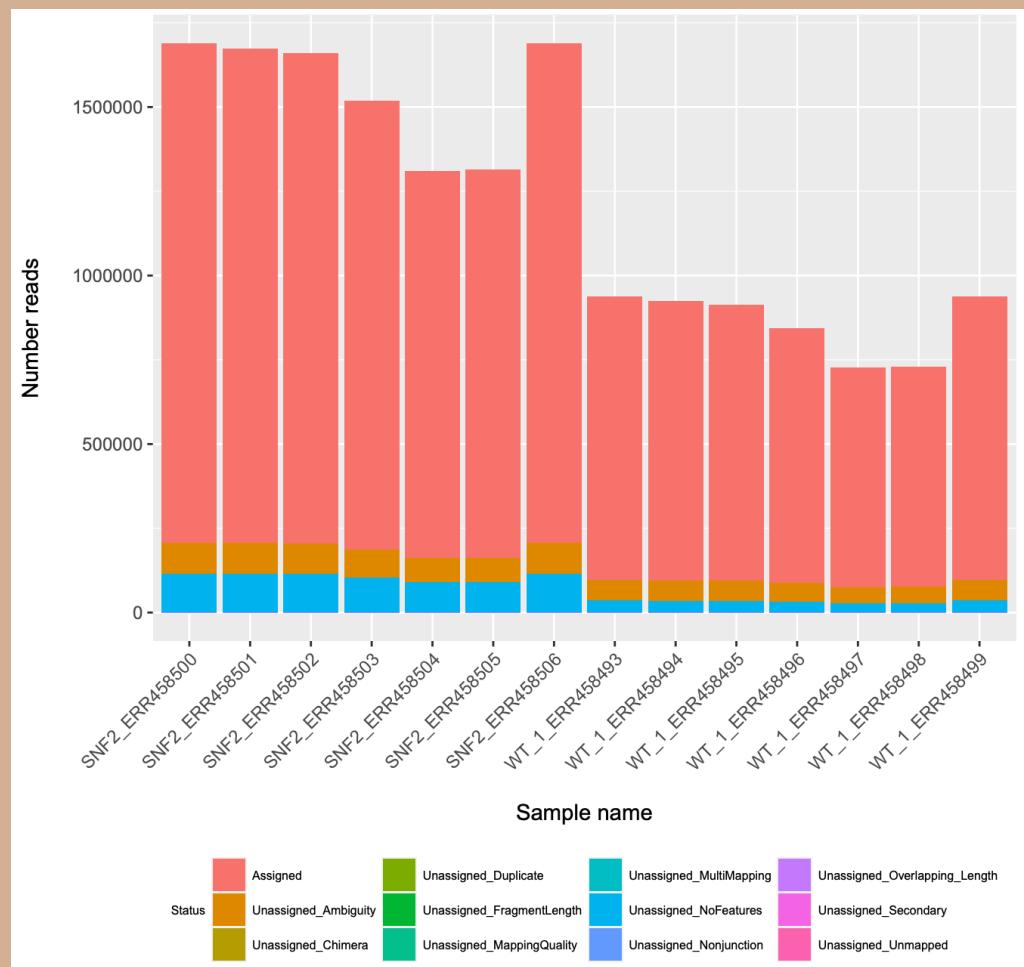
Read Alignment

Gene Quantification

Test for Differential Expression

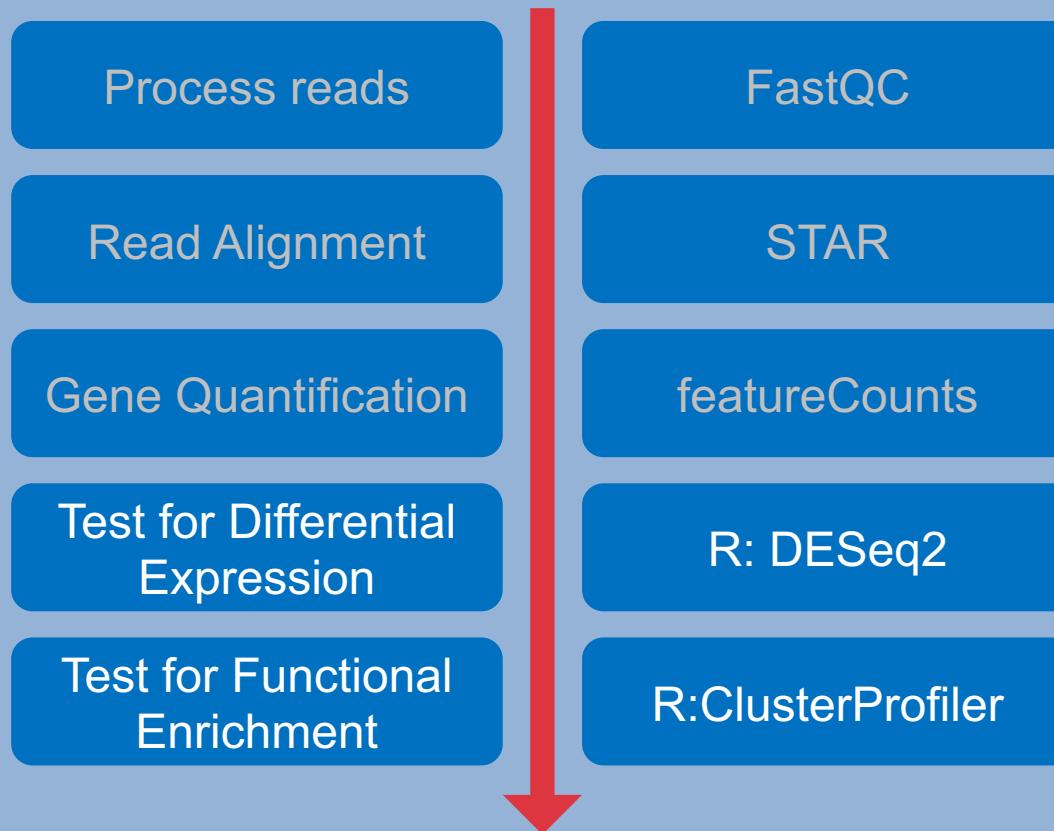
Test for Functional Enrichment

Feature counts statistics



(base) [whuo01@pcomp41 intro-to-RNA-seq]\$ Rscript scripts/featurecount_stat.R

Analysis pipeline



```

> # Input files: data, meta
> # DESeq2 transforms featurecount to log2FoldChange
> dds <- DESeqDataSetFromMatrix(countData = data, colData = meta, design = ~ condition)

```

Process reads

Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

Input data files

data

| Gene | Sample | | |
|-----------|----------------|----------------|----------------|
| | SNF2_ERR458500 | SNF2_ERR458501 | SNF2_ERR458502 |
| YAL069W | 0 | 0 | 0 |
| YAL068W-A | 0 | 0 | 0 |
| YAL068C | 0 | 0 | 0 |
| YAL067W-A | 0 | 0 | 0 |
| YAL067C | 18 | 21 | 10 |
| YAL066W | 0 | 1 | 0 |
| YAL065C | 2 | 0 | 0 |
| YAL064W-B | 3 | 3 | 1 |
| YAL064C-A | 12 | 7 | 8 |
| YAL064W | 3 | 3 | 2 |
| YAL063C-A | 2 | 0 | 3 |
| YAL063C | 35 | 38 | 46 |
| YAL062W | 127 | 136 | 113 |
| YAL061W | 124 | 100 | 125 |
| YAL060W | 300 | 309 | 323 |

meta

| Sample | | |
|----------------|------|----|
| | SNF2 | WT |
| SNF2_ERR458500 | SNF2 | |
| SNF2_ERR458501 | SNF2 | |
| SNF2_ERR458502 | SNF2 | |
| SNF2_ERR458503 | SNF2 | |
| SNF2_ERR458504 | SNF2 | |
| SNF2_ERR458505 | SNF2 | |
| SNF2_ERR458506 | SNF2 | |
| WT_ERR458493 | WT | |
| WT_ERR458494 | WT | |
| WT_ERR458495 | WT | |
| WT_ERR458496 | WT | |
| WT_ERR458497 | WT | |
| WT_ERR458498 | WT | |
| WT_ERR458499 | WT | |

Process reads

Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

```
> # DESeq2 transforms featurecount to log2FoldChange  
> dds <- DESeqDataSetFromMatrix(countData = data, colData = meta, design = ~ condition)  
> dds <- DESeq(dds)  
> # Tell DESeq2 which conditions you would like to compare  
> contrast <- c("condition", "SNF2", "WT") ←  
> res_unshrunken <- results(dds, contrast=contrast)  
> results <- lfcShrink(dds, contrast=contrast, res=res_unshrunken)  
> rld <- rlog(dds, blind=TRUE)  
> rld_counts <- assay(rld)  
> # output tables: results and rld_counts
```

output files

results: log2FoldChange with pvalue and p-adjusted (FDR)

| | baseMean | log2FoldChange | IfcSE | stat | pvalue | padj |
|-----------|--------------|----------------|------------|-------------|--------------|--------------|
| YAL069W | 0.000000e+00 | NA | NA | NA | NA | NA |
| YAL068W-A | 0.000000e+00 | NA | NA | NA | NA | NA |
| YAL068C | 9.021961e-02 | -0.043634441 | 0.35298047 | -0.12361413 | 9.016208e-01 | NA |
| YAL067W-A | 0.000000e+00 | NA | NA | NA | NA | NA |
| YAL067C | 5.977662e+00 | 1.632313960 | 0.37520919 | 4.12031328 | 3.783576e-05 | 1.234224e-04 |
| YAL066W | 9.812029e-02 | -0.046943453 | 0.36423350 | -0.12887995 | 8.974526e-01 | NA |
| YAL065C | 1.257050e+00 | -1.311117160 | 0.53025320 | -2.40606774 | 1.612527e-02 | 3.116017e-02 |
| YAL064W-B | 1.607497e+00 | -0.478036753 | 0.49638637 | -0.97070341 | 3.316960e-01 | 4.231887e-01 |

rld_counts: regularized log transformed counts. $y=\log_2(n + n_0)$

| | SNF2_ERR458500 | SNF2_ERR458501 | SNF2_ERR458502 | SNF2_ERR458503 | SNF2_ERR458504 | SNF2_ERR458505 |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|
| YAL069W | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.000000e+00 | |
| YAL068W-A | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.000000e+00 | |
| YAL068C | -2.05284510 | -2.05277066 | -2.052678613 | -2.03518019 | -2.050732e+00 | |
| YAL067W-A | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.000000e+00 | |
| YAL067C | 2.85518301 | 2.92415585 | 2.665424931 | 2.74363842 | 2.787779e+00 | |
| YAL066W | -2.05330217 | -2.03529421 | -2.053124206 | -2.05234800 | -2.051044e+00 | |
| YAL065C | 0.10262823 | 0.02271448 | 0.023084190 | 0.02600336 | 1.213956e-01 | |
| YAL064W-B | 0.58300398 | 0.58389721 | 0.509902895 | 0.51530992 | 4.804996e-01 | |
| YAL064C-A | 2.49883637 | 2.35760285 | 2.390449413 | 2.28097988 | 2.344918e+00 | |

Process reads

Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

```
> # DESeq2 transforms featurecount to log2FoldChange  
> dds <- DESeqDataSetFromMatrix(countData = data, colData = meta, design = ~ condition)  
> dds <- DESeq(dds)  
> # Tell DESeq2 which conditions you would like to compare  
> contrast <- c("condition", "SNF2", "WT") ←  
> res_unshrunken <- results(dds, contrast=contrast)  
> results <- lfcShrink(dds, contrast=contrast, res=res_unshrunken)  
> rld <- rlog(dds, blind=TRUE)  
> rld_counts <- assay(rld)  
> # output tables: results and rld_counts
```

output files

results: log2FoldChange with pvalue and p-adjusted (FDR)

| | baseMean | log2FoldChange | IfcSE | stat | pvalue | padj |
|-----------|--------------|----------------|------------|-------------|--------------|--------------|
| YAL069W | 0.000000e+00 | NA | NA | NA | NA | NA |
| YAL068W-A | 0.000000e+00 | NA | NA | NA | NA | NA |
| YAL068C | 9.021961e-02 | -0.043634441 | 0.35298047 | -0.12361413 | 9.016208e-01 | NA |
| YAL067W-A | 0.000000e+00 | NA | NA | NA | NA | NA |
| YAL067C | 5.977662e+00 | 1.632313960 | 0.37520919 | 4.12031328 | 3.783576e-05 | 1.234224e-04 |
| YAL066W | 9.812029e-02 | -0.046943453 | 0.36423350 | -0.12887995 | 8.974526e-01 | NA |
| YAL065C | 1.257050e+00 | -1.311117160 | 0.53025320 | -2.40606774 | 1.612527e-02 | 3.116017e-02 |
| YAL064W-B | 1.607497e+00 | -0.478036753 | 0.49638637 | -0.97070341 | 3.316960e-01 | 4.231887e-01 |

rld_counts: regularized log transformed counts. $y=\log_2(n + n_0)$

| | SNF2_ERR458500 | SNF2_ERR458501 | SNF2_ERR458502 | SNF2_ERR458503 | SNF2_ERR458504 | SNF2_ERR458505 |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|
| YAL069W | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.000000e+00 | |
| YAL068W-A | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.000000e+00 | |
| YAL068C | -2.05284510 | -2.05277066 | -2.052678613 | -2.03518019 | -2.050732e+00 | |
| YAL067W-A | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.000000e+00 | |
| YAL067C | 2.85518301 | 2.92415585 | 2.665424931 | 2.74363842 | 2.787779e+00 | |
| YAL066W | -2.05330217 | -2.03529421 | -2.053124206 | -2.05234800 | -2.051044e+00 | |
| YAL065C | 0.10262823 | 0.02271448 | 0.023084190 | 0.02600336 | 1.213956e-01 | |
| YAL064W-B | 0.58300398 | 0.58389721 | 0.509902895 | 0.51530992 | 4.804996e-01 | |
| YAL064C-A | 2.49883637 | 2.35760285 | 2.390449413 | 2.28097988 | 2.344918e+00 | |

Process reads

```
> # DESeq2 transforms featurecount to log2FoldChange  
> dds <- DESeqDataSetFromMatrix(countData = data, colData = meta, design = ~ condition)  
> dds <- DESeq(dds)  
> # Tell DESeq2 which conditions you would like to compare  
> contrast <- c("condition", "SNF2", "WT")  
> res_unshrunken <- results(dds, contrast=contrast)  
> results <- lfcShrink(dds, contrast=contrast, res=res_unshrunken)  
> rld <- rlog(dds, blind=TRUE)  
> rld_counts <- assay(rld)  
> # output tables: results and rld_counts  
> # plot PCA  
> plotPCA(rld, intgroup="condition") + geom_text(aes(label=name))
```

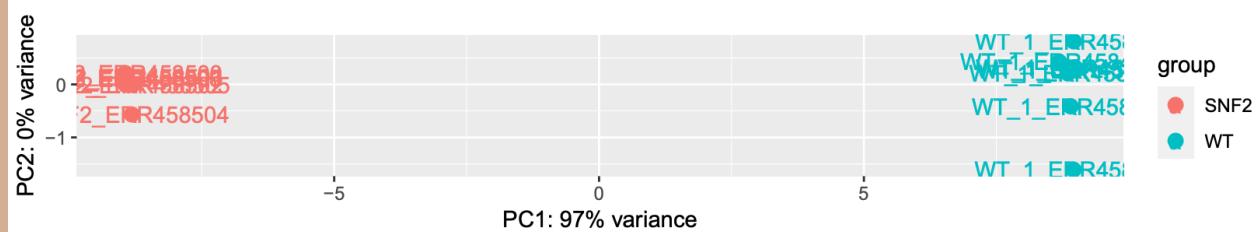
Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

How well are the replicates?



Process reads

Read Alignment

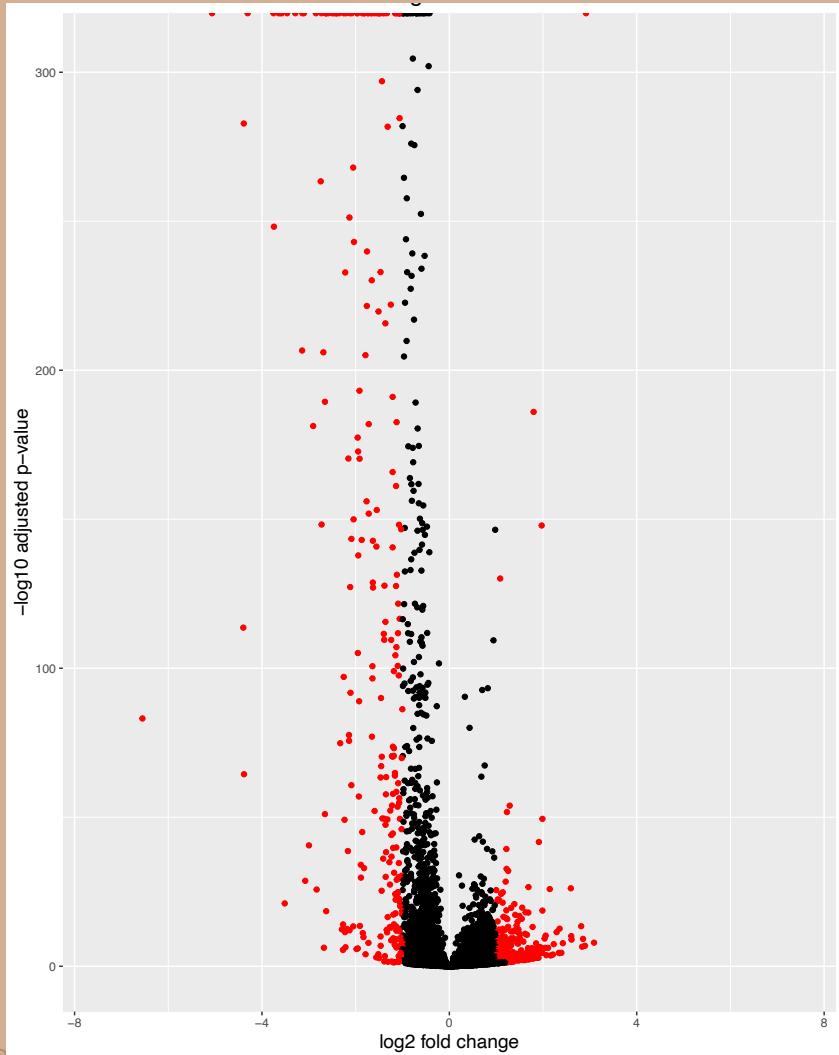
Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

```
> # output table: results  
> # load necessary library ggplot2  
> library(ggplot2)  
> # red dots: threshold of padj<0.05 and absolute value of log2foldchange >=1  
> res_table <- results %>% data.frame() %>% rownames_to_column(var="gene") %>% as_tibble()  
> res_table <- res_table %>% mutate(threshold_OE = padj < 0.05 & abs(log2FoldChange) >= 1)  
> # make volcano plot  
> ggplot(res_table) + ...
```

Volcano plot



```

> # output file: results rld_counts
> # select a subset of genes to plot heatmap
> rld_counts_sig <- rld_counts[gene_list, ]
> # plot heatmap
> pheatmap(rld_counts_sig, cluster_rows = T, show_rownames = T, annotation = meta,
border_color = NA, fontsize = 10, scale = "row", fontsize_row = 8, height = 20)

```

Process reads

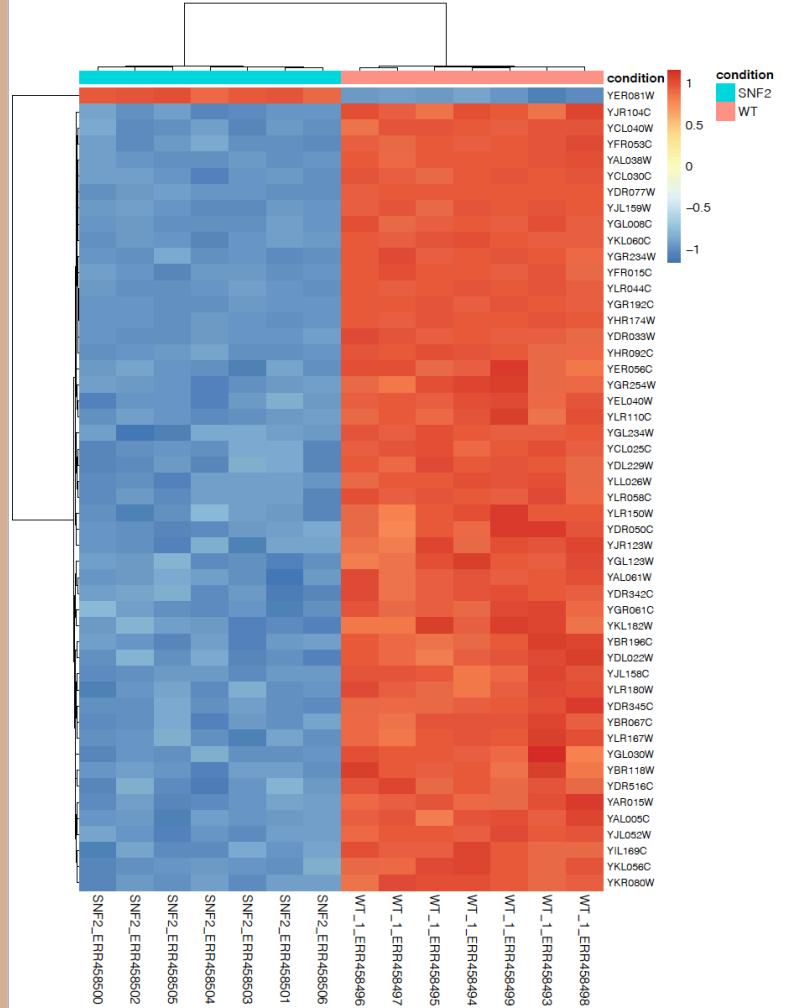
Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

Heatmap and clustering analysis



Process reads

Read Alignment

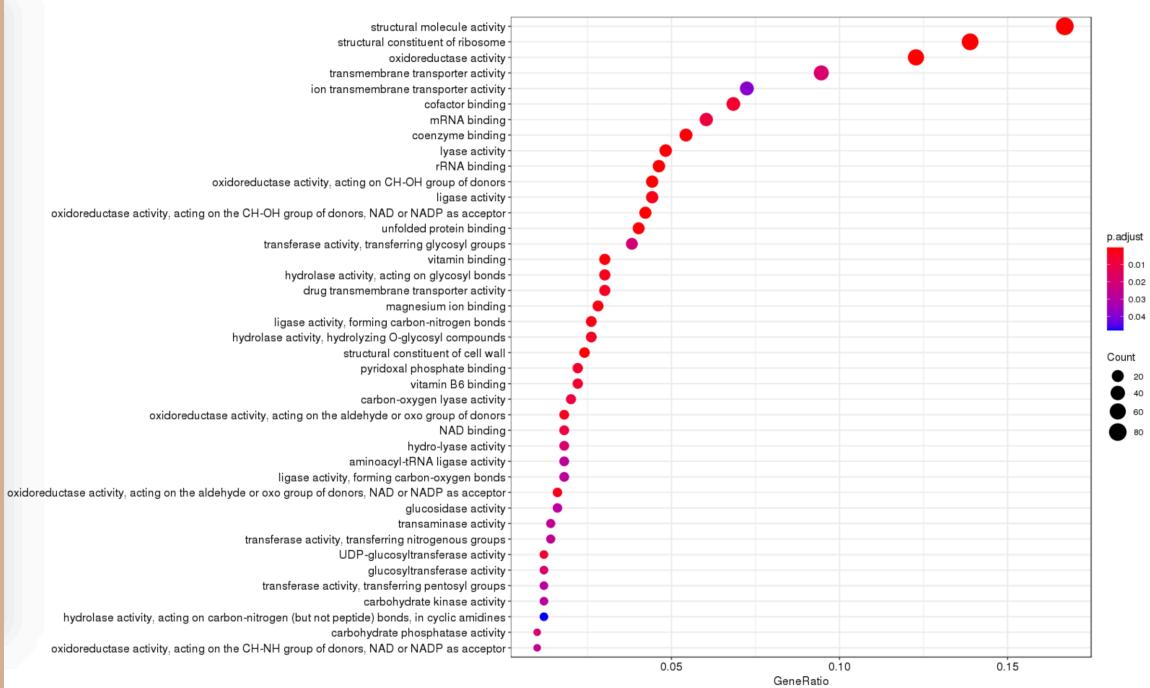
Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

```
> # output file: results and rld_counts -> extract genes of interest  
> # load necessary libraries  
> library(clusterProfiler)  
> library(org.Sc.sgd.db)  
> # Run GO enrichment for a list of genes: gene_list  
> ego <- enrichGO(gene = gene_list, keyType = "ENSEMBL", OrgDb = org.Sc.sgd.db)  
> cluster_summary <- data.frame(ego)  
> dotplot(ego, showCategory=50)
```

dotplot: Enrichment analysis using clusterProfiler



Process reads

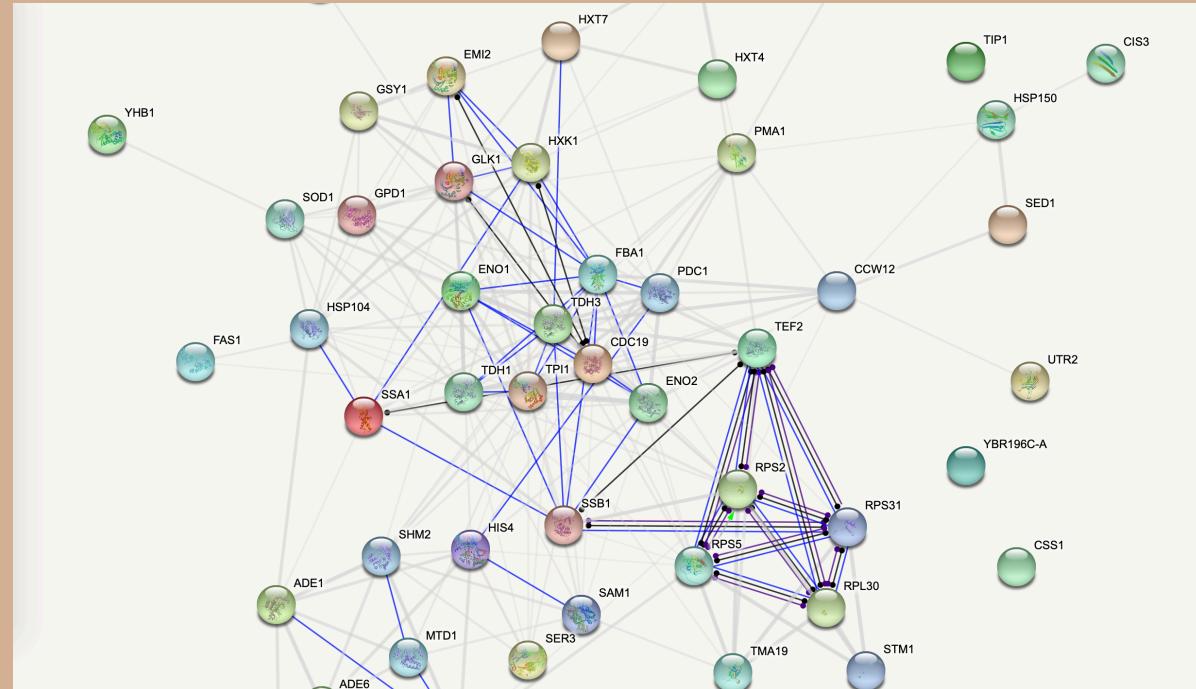
Read Alignment

Gene Quantification

Test for Differential Expression

Test for Functional Enrichment

Enrichment analysis using web server: String-db.org



Molecular Function (GO)

| GO-term | description | count in gene set | false discovery rate |
|------------|-------------------------------------|-------------------|----------------------|
| GO:0005199 | structural constituent of cell wall | 5 of 17 | 0.00011 |
| GO:0036094 | small molecule binding | 21 of 931 | 0.00021 |
| GO:0051287 | NAD binding | 5 of 30 | 0.00030 |
| GO:0003824 | catalytic activity | 33 of 2212 | 0.00030 |
| GO:0050662 | coenzyme binding | 8 of 139 | 0.00051 |

Cellular Component (GO)

| GO-term | description | count in gene set | false discovery rate |
|------------|-----------------------|-------------------|----------------------|
| GO:0071944 | cell periphery | 25 of 790 | 6.91e-09 |
| GO:0005829 | cytosol | 26 of 845 | 6.91e-09 |
| GO:0005886 | plasma membrane | 19 of 528 | 2.38e-07 |
| GO:0009277 | fungal-type cell wall | 10 of 112 | 5.55e-07 |
| GO:0005576 | extracellular region | 9 of 120 | 7.16e-06 |

What's next?

Next step: follow self-guided tutorial and get familiar with the pipeline.

What to do with your own data?

– one simply solution: move your reads.fastq files into raw_data folder and apply the same pipeline.

Questions?

– Ask questions and seek practice buddies on Piazza

Advanced reading material

An 1.5 day workshop for RNA sequencing analysis

https://hbctraining.github.io/DGE_workshop/

DESeq2 package description

<http://bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html>

Alternative analysis tools: Geneious Prime

<https://www.geneious.com/tutorials/expression-analysis/>

<https://www.youtube.com/watch?v=dwG-4i85CWw&t=16s>

Initiate shell access on compute node

```
srun -t 3:00:00 --mem 16G -N 1 -n 4 -p preempt --reservation bioworkshop --pty bash  
cd /cluster/tufts/bio/tools/training/intro-to-rnaseq/users/  
mkdir YOUR_USERNAME  
cd YOUR_USERNAME
```

Setup Rstudio on OnDemand

Number of cores: 1

Amount of Memory: 32 Gb

R version: 3.5.0

Partition: Preempt

Reservation for class, training, workshop: Bioinformatics workshop

Number of hours
3

Number of cores
1

Number of cores.

Amount of memory
32GB

Amount of memory (in GB).

Partition
Preempt

Partition: Default - Uses the default batch partition. Preempt - Uses extra contributed capacity. Your job will be canceled if contrib owner needs resources.

Reservation for class, training, workshop
Bioinformatics Workshop

If you don't know about specific reservation, select default.

R version
3.5.0

Only use version of R that works for you. Changing versions may require you to reinstall packages. Only change if necessary.