

RNAseq workshop Day1 (organized)

1. How to log onto server
2. download the data
3. FastQC the data

Part_1. log onto server and check the tools

login address: ssh -l ngscs33 aristotle.med.cornell.edu

password: hs!acgt@2001

type of aristotle: ssh scu-node02

find the tools stored in the 2017_rnaseq folder

```
[ngscs21@scu-node02 ~]$ ls /zenodotus/abc/store/courses/2017_rnaseq/
```

```
[ngscs21@scu-node02 ~]$ ls /zenodotus/abc/store/courses/2017_rnaseq/  
aln_tmp  rawReads_yeast_Gierlinski  referenceGenomes  software
```

create a symbolic link for the folder /zenodotus/abc/store/courses/2017_rnaseq and name the directory as 'mat'

```
[ngscs21@scu-node02 ~]$ ln -s /zenodotus/abc/store/courses/2017_rnaseq/ mat
```

ln means generates a link; -s means symbolic, not hard link

```
[ngscs21@scu-node02 ~]$ ls mat/
```

```
[ngscs21@scu-node02 ~]$ ls mat/  
aln_tmp  rawReads_yeast_Gierlinski  referenceGenomes  software
```

to log off the server

```
[ngscs21@scu-node02 ~]$ exit
```

Part_2. Download the data

2.1

go to ena <https://www.ebi.ac.uk/ena>

input the reference number for the data to be downloaded "ERP004763"

Search results for [ERP004763](#)

Show more data from EMBL-EBI

Read Run (672)	Run (672 results found) ERR458584 Illumina HiSeq 2000 sequencing View all 672 results
Study Study (1)	Study (1 results found) ERP004763 S. cerevisiae WT vs snf2 KO mutant RNA-seq data with 7 technical and 48 biological replicates (336 total) of each condition View all 1 results

Choose "Study (1 results found)"

[Navigation](#) [Read Files](#) [Portal](#) [Attributes](#) [Publications](#)

[Bulk Download Files](#) ⚠ (Please use Firefox to launch the bulk downloader app.)

Download: - of 672 results in [TEXT](#)

2.2

catch the TEXT link of the 672 ftp files, and name it as "samples_at_ENA.txt"

```
wget -O samples_at_ENA.txt "https://www.ebi.ac.uk/ena/data/warehouse/filereport?
accession=PRJEB5348&result=read_run&fields=study_accession,sample_accession,secondary_sample_acc
ession,experiment_accession,run_accession,tax_id,scientific_name,instrument_model,library_layout,fastq_ft
p,fastq_galaxy,submitted ftp,submitted_galaxy,sra ftp,sra_galaxy,cram_index ftp,cram_index_galaxy&down
load=txt"
```

Don't forget the quotation marks

Use wget to download a webpage; -O writes the downloaded documents to FILE

show the 11th column of the txt file, and only shows the first 10 items. Here it stores the filenames we need to download

```
[ngscls21@scu-node02 ~]$ cut -f11 samples_at_ENA.txt | head
```

```
[ngscls33@scu-node02 ~]$ cut -f11 samples_at_ENA.txt | head
```

```
fastq_galaxy
```

```
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458493/ERR458493.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458494/ERR458494.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458495/ERR458495.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458496/ERR458496.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458497/ERR458497.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458498/ERR458498.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458499/ERR458499.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458500/ERR458500.fastq.gz
ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458501/ERR458501.fastq.gz
```

2.3

catch the LINK that shows the correlation between ID name and sample name/info, and name it as "ERP004763_sample_mapping.tsv".

```
wget -O ERP004763_sample_mapping.tsv "http://dx.doi.org/10.6084/m9.figshare.1416210"
```

show the first 10 items of the code

```
head ERP004763_sample_mapping.tsv
```

```
[ngscls33@scu-node02 ~]$ head ERP004763_sample_mapping.tsv
```

RunAccession	Lane	Sample	BiolRep
ERR458493	1	WT	1
ERR458494	2	WT	1
ERR458495	3	WT	1
ERR458496	4	WT	1
ERR458497	5	WT	1
ERR458498	6	WT	1
ERR458499	7	WT	1
ERR458500	1	SNF2	1
ERR458501	2	SNF2	1

2.4 download the actual data for ERR458493 (one of the files)

check the files that need to be downloaded by finding replicate 1 samples:

```
[ngscls33@scu-node02 ~]$ awk '$4 == 1' ERP004763_sample_mapping.tsv | cut -f1
```

ERR458493
ERR458494
ERR458495
ERR458496
ERR458497
ERR458498
ERR458499
ERR458500
ERR458501
ERR458502
ERR458503
ERR458504
ERR458505
ERR458506

```
# method1. Download one by one with the ftp address
```

```
[ngsc1s33@scu-node02 ~]$ wget ftp.sra.ebi.ac.uk/vol1/fastq/ERR458/ERR458500/ERR458500.fastq.gz
```

#Awk breaks each line of input passed to it into fields. By default, a field is a string of consecutive characters delimited by whitespace, though there are options for changing this. Awk parses and operates on each separate field. This makes it ideal for handling structured text files -- especially tables -- data organized into consistent chunks, such as rows and columns.

Arbitrarily long lists of parameters cannot be passed to a command in some situations, so `xargs` breaks the list of arguments into sublists small enough to be acceptable.

```
# method2. Download one by one with the sample name
```

```
# method3. Generate a for loop to download selected files
```

```
[ngscls33@scu-node02 ~]$ for i in `seq 4 6`; do SAMPLE=ERR45850${i}; egrep ${SAMPLE}
samples_at_ENA.txt | cut -f11 | xargs wget; done
```

```
# the above code download the files with ERR458504 ~ ERR458506
```

```
# there is not gap between SAMPLE, = and ERR45850${i}
```

```
#method4. Generate a txt file to download all the files
```

????????????????????????????

?????????wget vs asw????????????????????

How FastQ file look like

```
[ngsc1s33@scu-node02 WT_rep1]$ zcat ERR458493.fastq.gz | head
```

APR 15 1965 4 51/2 PM '65 - 540 - 505771000 - 4 - 4404 - 4754 - 5000/4

```
# generate multiQC result
[ngscsls33@scu-node02 fastqc results]$ ~/mat/software/anaconda2/bin/multiqc WT rep1/ --dirs
```

```

[[ngscls33@scu-node02 fastqc_results]$ ~/mat/software/anaconda2/bin/multiqc WT_re
p1/ --dirs
[WARNING]      multiqc : MultiQC Version v1.2 now available!
[INFO  ]      multiqc : This is MultiQC v1.1
[INFO  ]      multiqc : Template      : default
[INFO  ]      multiqc : Prepending directory to sample names
[INFO  ]      multiqc : Searching 'WT_rep1/'
Searching 68 files.. [#####] 100%
[INFO  ]      fastqc : Found 10 reports
[INFO  ]      multiqc : Compressing plot data
[INFO  ]      multiqc : Report       : multiqc_report.html
[INFO  ]      multiqc : Data        : multiqc_data
[INFO  ]      multiqc : MultiQC complete
[[ngscls33@scu-node02 fastqc_results]$ ls
multiqc_data multiqc_report.html WT_rep1

```

3.2.3 send multiple QC results to email with multiQC

move to the folder where multiqc_report.html is stored (I moved them to WT_rep1 in advance)

```

[[ngscls33@scu-node02 WT_rep1]$ echo "here are the results of multiQC" | mailx -s "MultiQC
results" -a multiqc_report.html fta2001@med.cornell.edu

```