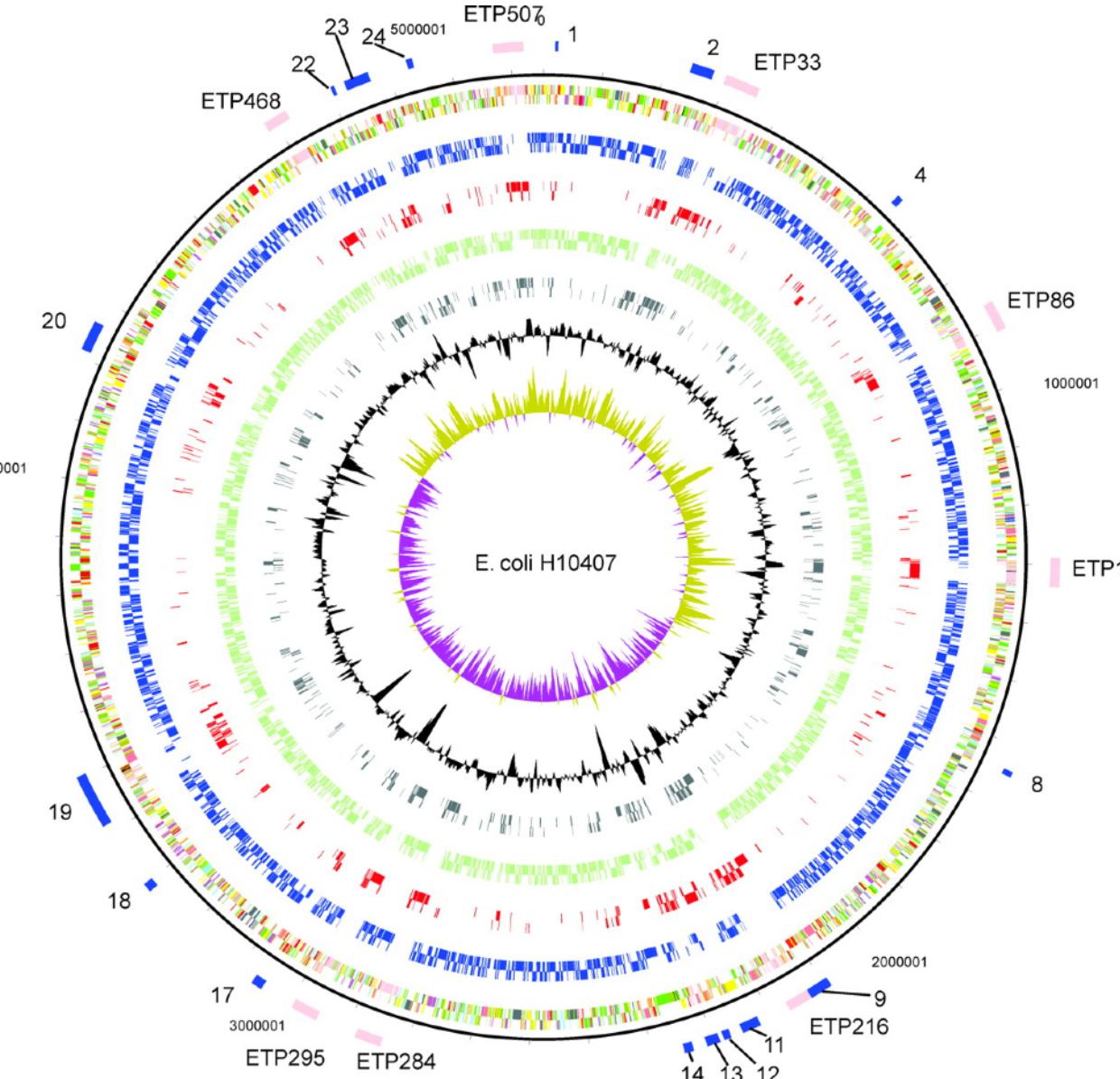
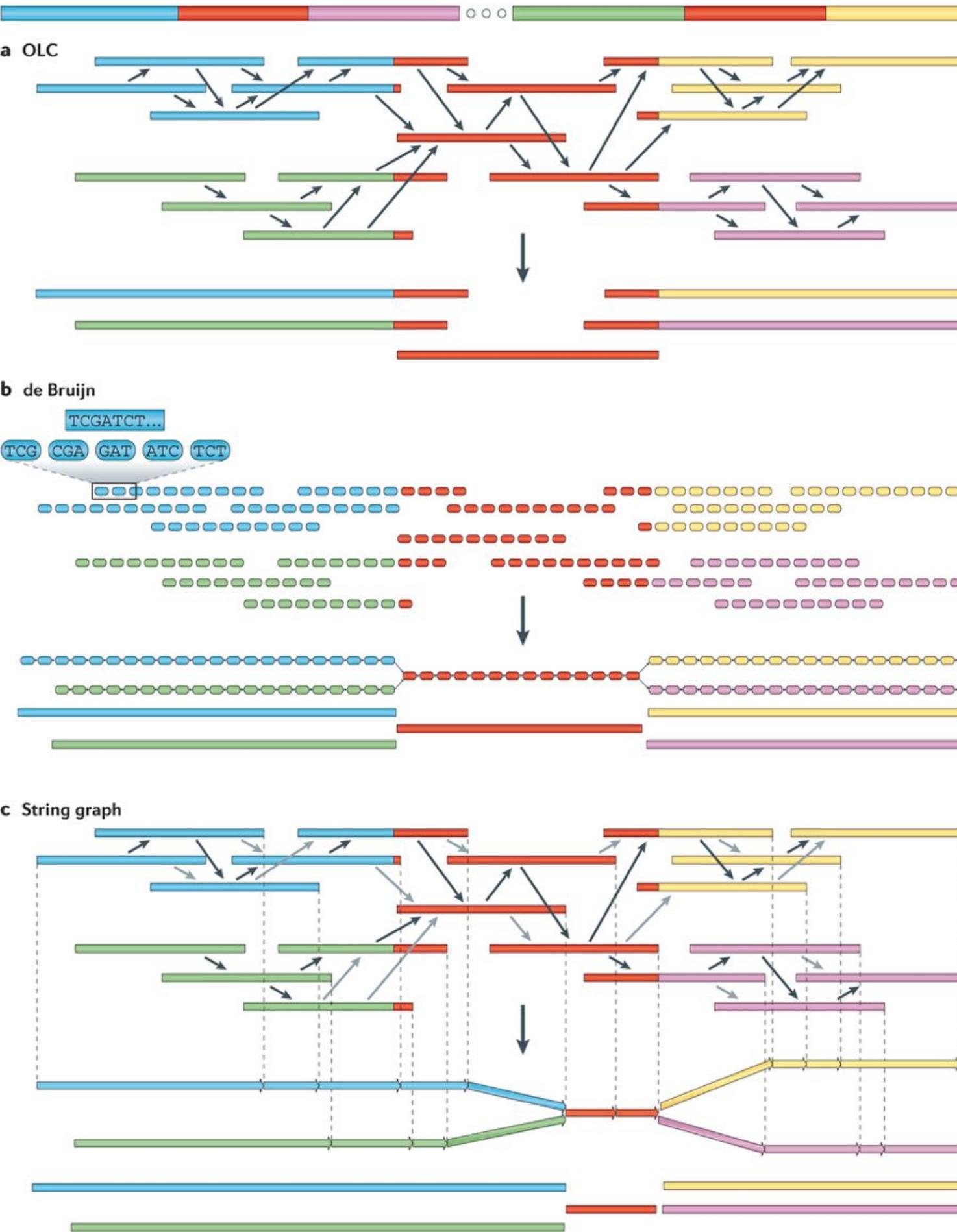


Genome assembly Practical



**HTS Workshop Genomics &
Transcriptomics
KAUST 2019**

Robert Lehmann
Octavio Salazar

Preparation

First Steps - Preparing the folder structure

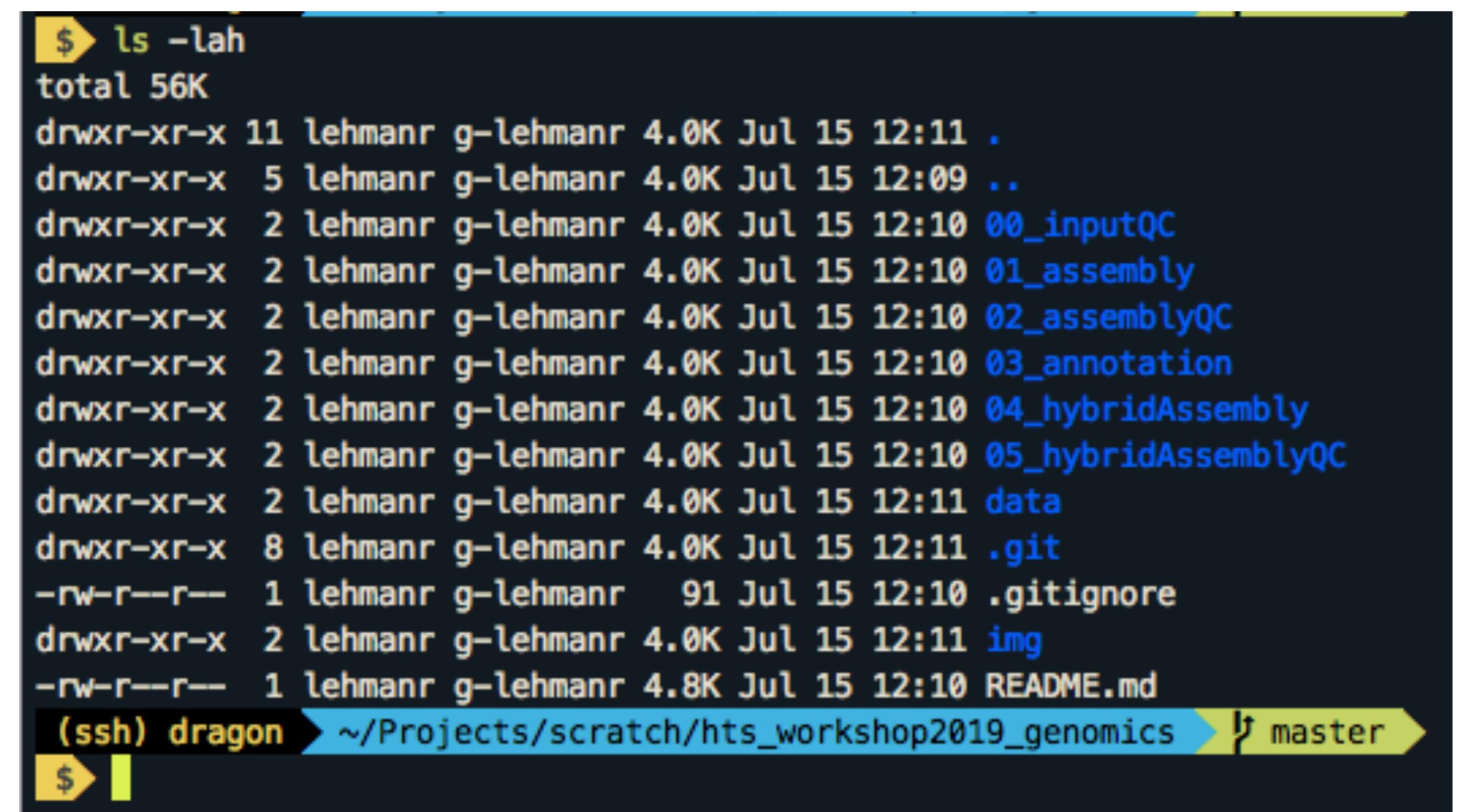
Download scripts and reference data for analyses:

In order to download the scripts and reference data for the analyses, we first need to login to the Ibex cluster:

```
$ ssh <username>@ilogin.ibex.kaust.edu.sa
```

Then, we download the prepared scripts from github:

```
$ mkdir <workshop dir>
$ cd <workshop dir>
$ git clone https://github.com/roblehmann/
hts_workshop2019_genomics.git
$ cd hts_workshop2019_genomics
```



The screenshot shows a terminal window with a black background and white text. At the top, it displays the command '\$ ls -lah'. Below this, it shows a detailed listing of files and directories in the current working directory. The listing includes:

- total 56K
- drwxr-xr-x 11 lehmanr g-lehmanr 4.0K Jul 15 12:11 .
- drwxr-xr-x 5 lehmanr g-lehmanr 4.0K Jul 15 12:09 ..
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:10 00_inputQC
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:10 01_assembly
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:10 02_assemblyQC
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:10 03_annotation
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:10 04_hybridAssembly
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:10 05_hybridAssemblyQC
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:11 data
- drwxr-xr-x 8 lehmanr g-lehmanr 4.0K Jul 15 12:11 .git
- rw-r--r-- 1 lehmanr g-lehmanr 91 Jul 15 12:10 .gitignore
- drwxr-xr-x 2 lehmanr g-lehmanr 4.0K Jul 15 12:11 img
- rw-r--r-- 1 lehmanr g-lehmanr 4.8K Jul 15 12:10 README.md

At the bottom of the terminal window, there is a blue status bar with the text '(ssh) dragon > ~/Projects/scratch/hts_workshop2019_genomics > p master'.

First Steps - Preparing the folder structure

Some of the tools we will use require additional data, such as a reference genome assembly and annotation for E. coli. Let's download these data before starting the assembly:

```
$ cd data && bash copyDataFromIbex.sh
```

The result should look like this

```
$ ls -laht
total 7.7G
-rw-r--r-- 1 lehmanr g-lehmannr 1.7M Jul 15 15:02 GCF_000005845.2_ASM584v2_protein.faa
drwxr-xr-x 3 lehmanr g-lehmannr 4.0K Jul 15 15:02 .
-rw-r--r-- 1 lehmanr g-lehmannr 2.8M Jul 15 15:02 GCF_000005845.2_ASM584v2_genomic.gff
-rw-r--r-- 1 lehmanr g-lehmannr 4.5M Jul 15 15:02 GCF_000005845.2_ASM584v2_genomic.fna
-rw-r--r-- 1 lehmanr g-lehmannr 1.7M Jul 15 15:02 taxo.k2d
-rw-r--r-- 1 lehmanr g-lehmannr 56 Jul 15 15:02 opts.k2d
-rw-r--r-- 1 lehmanr g-lehmannr 7.5G Jul 15 15:02 hash.k2d
-rwrxr-xr-x 1 lehmanr g-lehmannr 570 Jul 15 15:01 copyDataFromIbex.sh
drwxr-xr-x 5 lehmanr g-lehmannr 4.0K Jul 15 15:00 enterobacterales_odb9
drwxr-xr-x 11 lehmanr g-lehmannr 4.0K Jul 15 14:58 ..
-rw-r--r-- 1 lehmanr g-lehmannr 40M Jul 15 14:58 tardigrade_SRR2986339_subsampled_2.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 39M Jul 15 14:58 tardigrade_SRR2986339_subsampled_1.fq.gz
-rwxr-xr-x 1 lehmanr g-lehmannr 1.1K Jul 15 14:58 downloadData.sh
-rw-r--r-- 1 lehmanr g-lehmannr 31M Jul 15 14:58 SRR2627019.1m.80xR2.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 24M Jul 15 14:58 SRR2627019.1m.80xR1.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 20M Jul 15 14:58 SRR2627019.1m.40xR2.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 16M Jul 15 14:58 SRR2627019.1m.40xR1.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 7.6M Jul 15 14:58 SRR2627019.1m.20xR2.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 6.0M Jul 15 14:58 SRR2627019.1m.20xR1.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 3.9M Jul 15 14:58 SRR2627019.1m.10xR2.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 3.1M Jul 15 14:58 SRR2627019.1m.10xR1.fq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 4.4M Jul 15 14:58 SRR1284073.1m.5x.pacbio.fastq.gz
-rw-r--r-- 1 lehmanr g-lehmannr 1.7K Jul 15 14:58 README.md
```

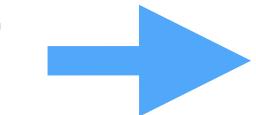
Quality control

Step 0 - Read data quality control

00_inputQC

```
$ cd 00_inputQC && ls -lah
total 0
drwxr-xr-x  2 lehmanr g-lehmanr 4.0K Jul 15 14:58 .
drwxr-xr-x 11 lehmanr g-lehmanr 4.0K Jul 15 14:58 ..
-rw xr-xr-x  1 lehmanr g-lehmanr 414 Jul 15 14:58 0-runFastQC.sh
-rw xr-xr-x  1 lehmanr g-lehmanr 1.1K Jul 15 14:58 1-runTrimming.sh
-rw xr-xr-x  1 lehmanr g-lehmanr 295 Jul 15 14:58 2-runMultiQC.sh
-rw xr-xr-x  1 lehmanr g-lehmanr 542 Jul 15 14:58 3-runKraken.sh
-rw xr-xr-x  1 lehmanr g-lehmanr 490 Jul 15 14:58 4-runKAT.sh
```

Job configuration part
for slurm



```
#!/bin/bash
#SBATCH --job-name=fastqc
#SBATCH --nodes=1
#SBATCH --ntasks=1
#SBATCH --cpus-per-task=2
#SBATCH --time=2:00:00
#SBATCH --mem=5000
#SBATCH --mail-type=end

# fastQC
# https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

set -vex

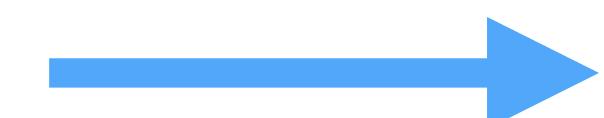
module load fastqc

outDir='output'
mkdir -p $outDir

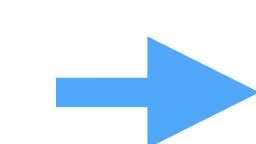
threads=2

for f in ../data/*fq.gz; do
    echo "fastQC'ing file"${f}
    fastqc $f -o ${outDir} -t ${threads}
done
```

Comment



Make fastqc available



Actual command

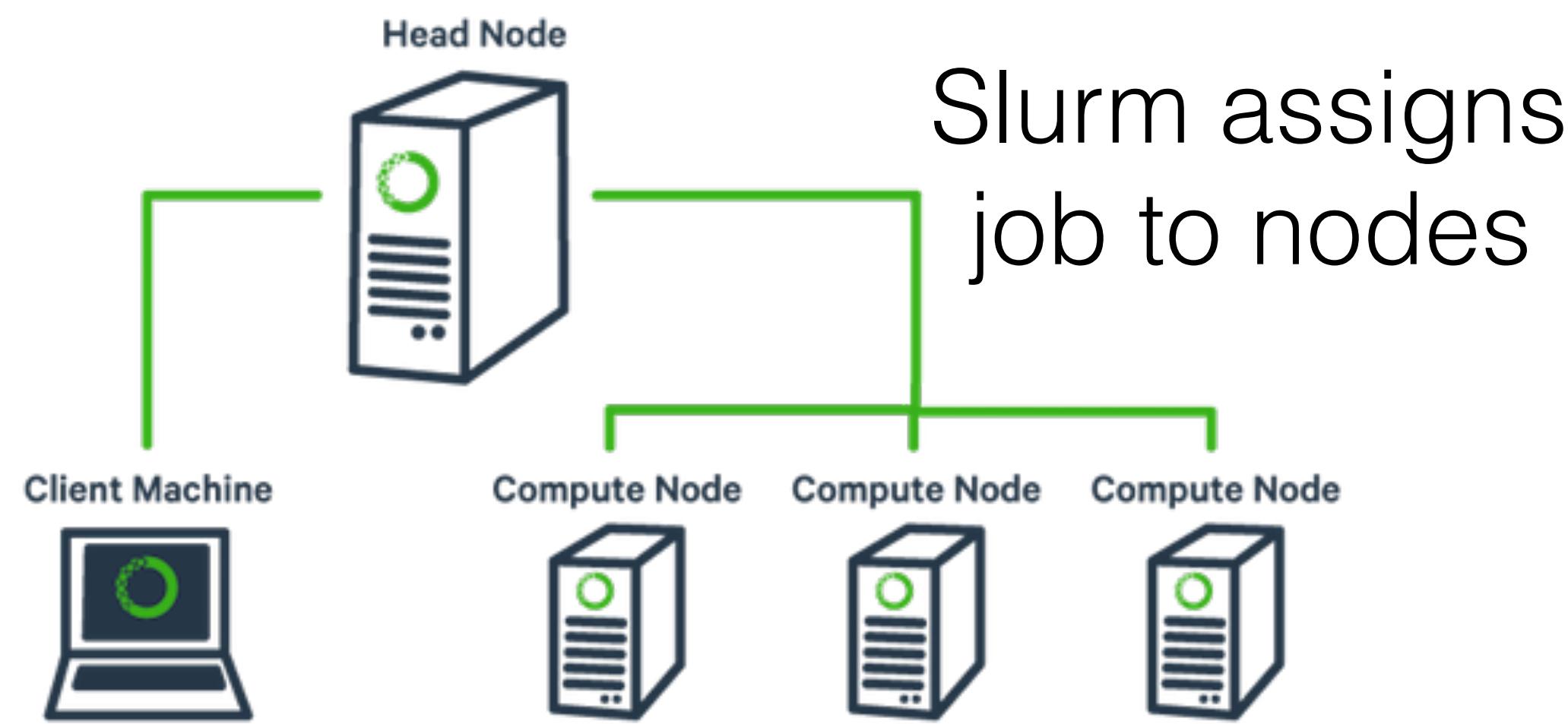


how to edit file:

\$ nano 0-runFastQC.sh

Slurm - Small refresher

00_inputQC



How to give a new job to Slurm?

sbatch <job_script.sh>
or

run <job_script.sh>

What is the cluster doing?

	JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST (REASON)
5457607_[338-2000]	group-sts	T-2-17	lenzia	PD		0:00	1	(Resources)
5458107	batch	BandR+UF	bordjiz	PD		0:00	4	(Priority)
	squeue tail -n 10							
5451125	batch	case_2	zhongp	R	14:55:06		1	ds505-17
5451101	batch	case_2	zhongp	R	14:55:07		1	ds506-25
5451103	batch	case_2	zhongp	R	14:55:07		1	ds503-13
5451104	batch	case_2	zhongp	R	14:55:07		1	ds505-29
5451105	batch	case_2	zhongp	R	14:55:07		1	ds506-05
5435734	batch	g09_m1_0	minenky	R	1-03:17:42		1	ds506-29
5292423	batch	hc_00_pm	minenky	R	3-16:38:16		1	dbn302-06-r
5292441	batch	hc_00_pm	minenky	R	3-16:38:16		1	dbn302-07-l
5292442	batch	hc_00_pm	minenky	R	3-16:38:16		1	dbn302-07-r
5458042	batch	bash	ramazan	R	3:38:55		1	cn603-07-r

What am I doing (on the cluster)?

```
$ squeue -u lehmanr
```

JOBid	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST(REASON)
-------	-----------	------	------	----	------	-------	------------------

Lets try to run FastQC

00_inputQC

Submit FastQC job

Check if its running

```
$ sbatch 0-runFastQC.sh
Submitted batch job 5459606
$ squeue -u lehmann
      JOBID PARTITION     NAME      USER ST      TIME  NODES NODELIST(REASON)
      5459606    batch   fastqc  lehmann R      0:44      1  cn605-07-1
$ ls -lah
total 256K
drwxr-xr-x  3 lehmann g-lehmann 4.0K Jul 15 15:36 .
drwxr-xr-x 11 lehmann g-lehmann 4.0K Jul 15 14:58 ..
-rw xr-xr-x  1 lehmann g-lehmann  414 Jul 15 14:58 0-runFastQC.sh
-rw xr-xr-x  1 lehmann g-lehmann  1.1K Jul 15 14:58 1-runTrimming.sh
-rw xr-xr-x  1 lehmann g-lehmann  295 Jul 15 14:58 2-runMultiQC.sh
-rw xr-xr-x  1 lehmann g-lehmann  542 Jul 15 14:58 3-runKraken.sh
-rw xr-xr-x  1 lehmann g-lehmann  490 Jul 15 14:58 4-runKAT.sh
drwxr-xr-x  2 lehmann g-lehmann 4.0K Jul 15 15:36 output
-rw-r--r--  1 lehmann g-lehmann 15K Jul 15 15:36 slurm-5459606.out
```

FastQC output

Slurm's log file

Email when job is done

00_inputQC

Check the content of the “output” folder

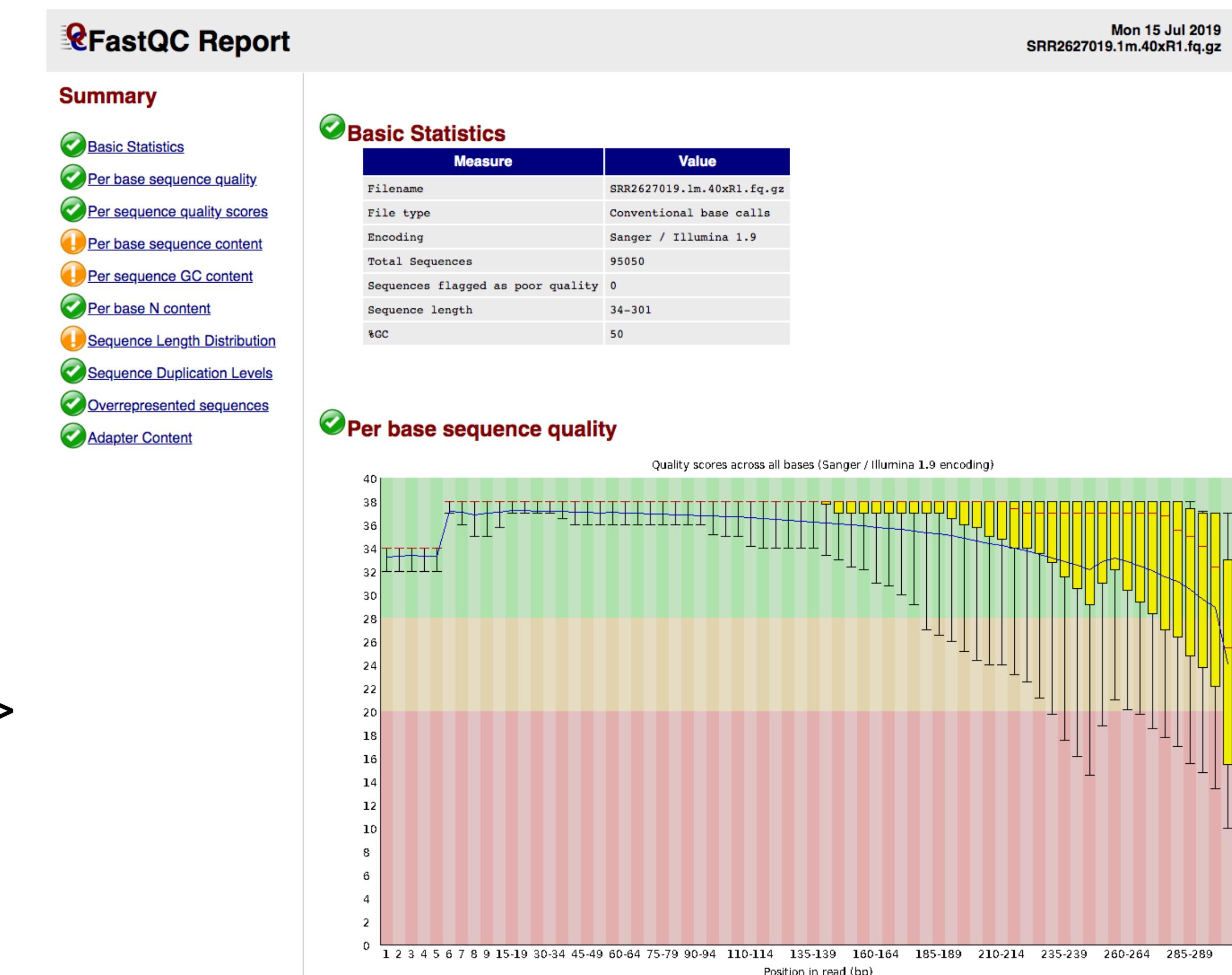
\$ ls -lah output

Download the html files to your local computer and open:

\$ scp <user>@illogin.ibex.kaust.edu.sa:<path to output folder>/SRR2627019.1m.40xR1_fastqc.html

\$ scp <user>@illogin.ibex.kaust.edu.sa:<path to output folder>/SRR2627019.1m.40xR2_fastqc.html

How to look at the results



This is very tedious for many samples ->
need to summarize reports...

00_inputQC

Trimming the reads

Read script 1-runTrimming.sh:

```
$ nano 1-runTrimming.sh
```

Submit trimming job for one coverage dataset to cluster:

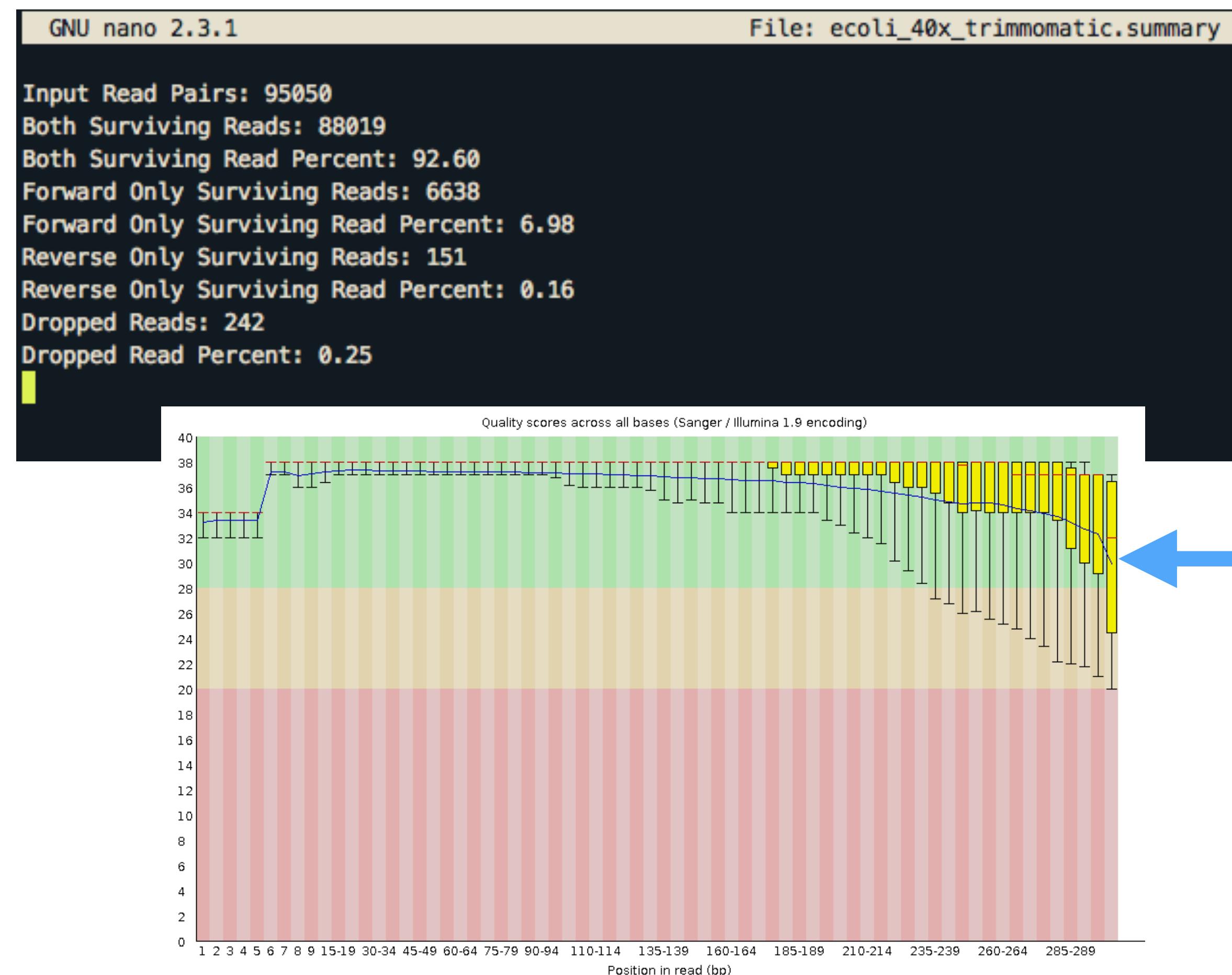
\$ sbatch 1-runTrimming.sh

Have a look at the trimming result:

```
$ nano output/ecoli_40x_trimmomatic.summary
```

```
$ scp <user>@illogin.ibex.kaust.edu.sa:<path to output  
folder>/SRR2627019.1m.40x_trimmed_R1_fastqc.html
```

Edit script to trim datasets for varying coverage and resubmit...



00_inputQC

Contamination check with Kraken2

Have a look at the script

```
$ nano 2-runKraken.sh
```

And submit as before:

```
$ batch 2-runKraken.sh
```

The output file is a flatfile, so we can check the result on IBEX using nano:

```
$ nano output/ecoli_40x_kraken2_report.txt
```

And modify parameters to run on the other input datasets as before...

**How to find which taxons have many reads assigned?
-> awk**

Percentage Reads	Reads in clade rooted in taxon	Reads under this taxon	Taxon name
1.37	1302	1302	unclassified
98.63	93748	204	root
98.39	93524	2	cellular organisms
98.34	93469	196	Bacteria
98.12	93262	120	Proteobacteria
97.98	93133	126	Gammaproteobacteria
97.77	92930	739	Enterobacteriales
73.65	70008	33033	Enterobacteriaceae
34.99	33256	2497	Escherichia
30.60	29089	27680	Escherichia coli
0.17	158	158	Escherichia coli O15:H11
0.15	147	59	Escherichia coli K-12
0.07	69	0	Escherichia coli str. K-12 substr. M
0.07	69	69	Escherichia coli str. K-12 substr.
0.02	19	19	Escherichia coli str. K-12 substr. W

Contamination check with Kraken2

00_inputQC

```
$ awk '$1 > 1 {print}' ecoli_kraken2_report.txt | wc -l  
13
```

```
[lehmann@dbn503-35-r output]$ awk '$1 > 1' ecoli_kraken2_report.txt  
99.67 94735 298 R 1 root  
99.35 94436 2 R1 131567 cellular organisms  
99.30 94387 203 D 2 Bacteria  
99.08 94175 257 P 1224 Proteobacteria  
98.80 93909 69 C 1236 Gammaproteobacteria  
98.72 93831 917 O 91347 Enterobacterales  
74.22 70550 39261 F 543 Enterobacteriaceae  
31.88 30298 1603 G 561 Escherichia  
29.75 28282 27968 S 562 Escherichia coli  
23.52 22356 15 F 1903411 Yersiniaceae  
23.50 22337 205 G 629 Yersinia  
23.28 22124 16426 G1 1649845 Yersinia pseudotuberculosis complex  
5.76 5475 5410 S 632 Yersinia pestis
```

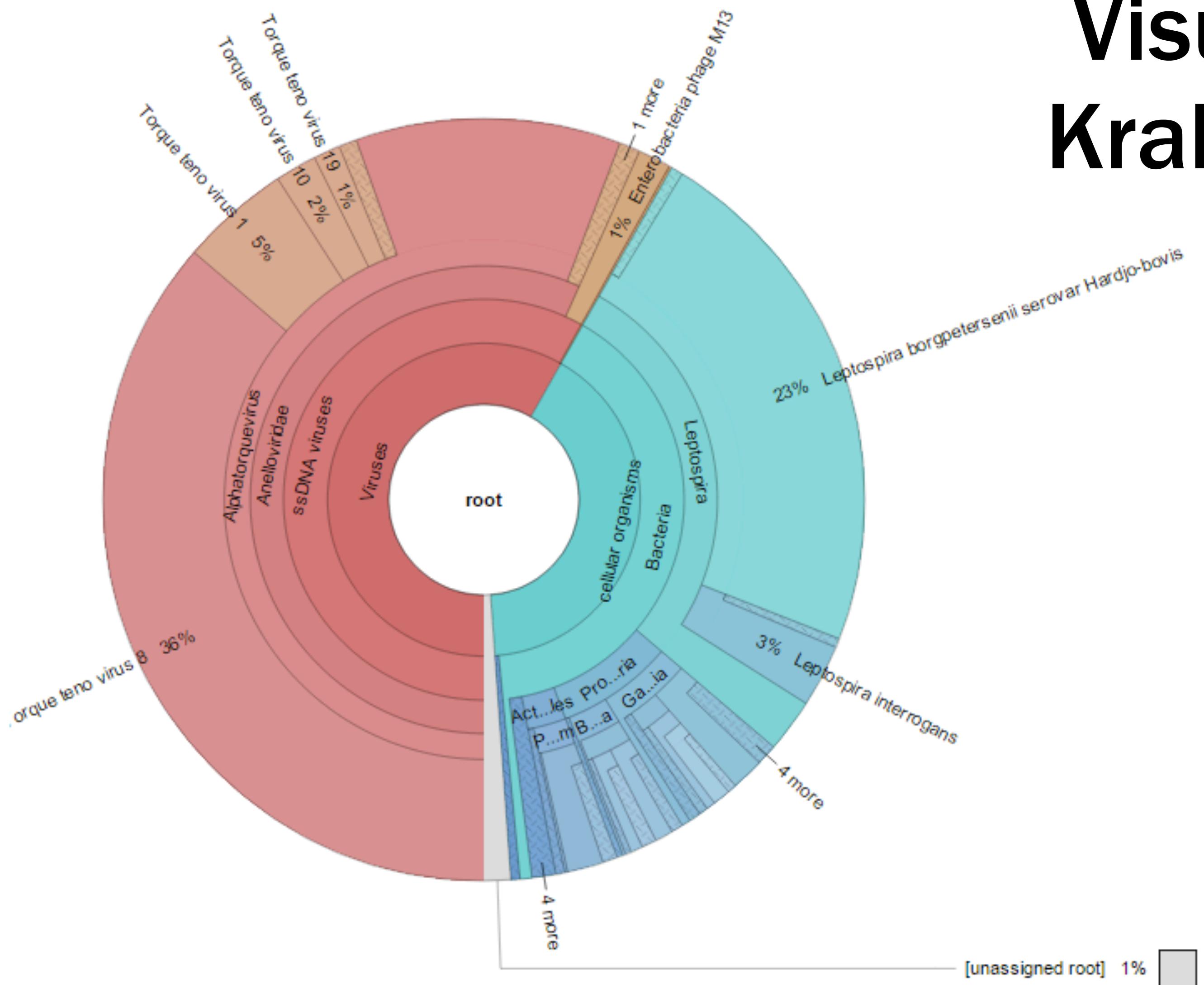


00_inputQC

Visualization of Kraken2 results

Krona

<https://github.com/marbl/Krona>



And also the GC vs. Coverage plots:

```
$ nano 3-runKAT.sh  
$ sbatch 3-runKAT.sh
```

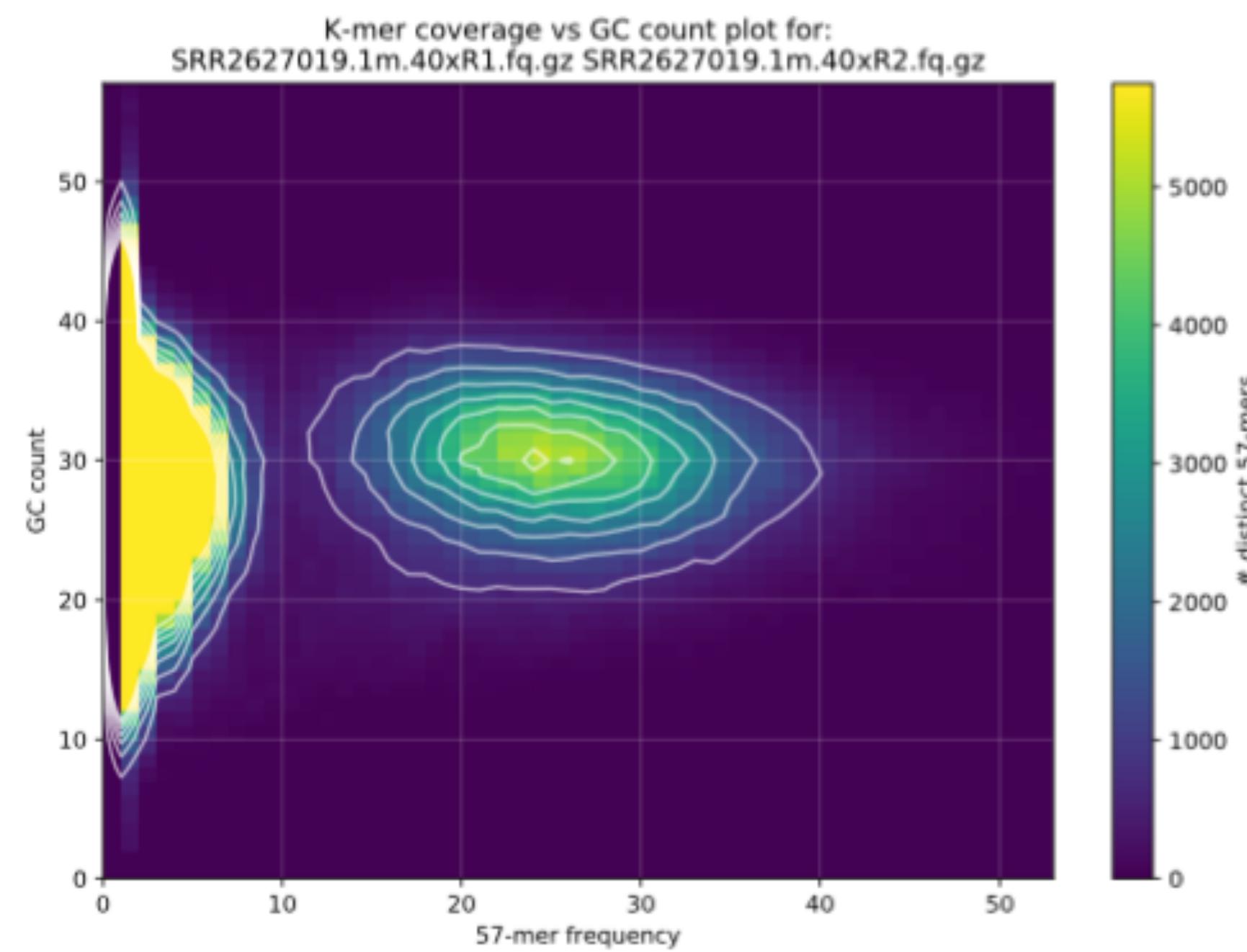
Change to dataset with different coverage, and re-run.
Then download all run results to your machine:

```
$ scp "<user>@ilogin.ibex.kaust.edu.sa:<path to  
output folder>/00_inputQC/output/*pdf" .
```

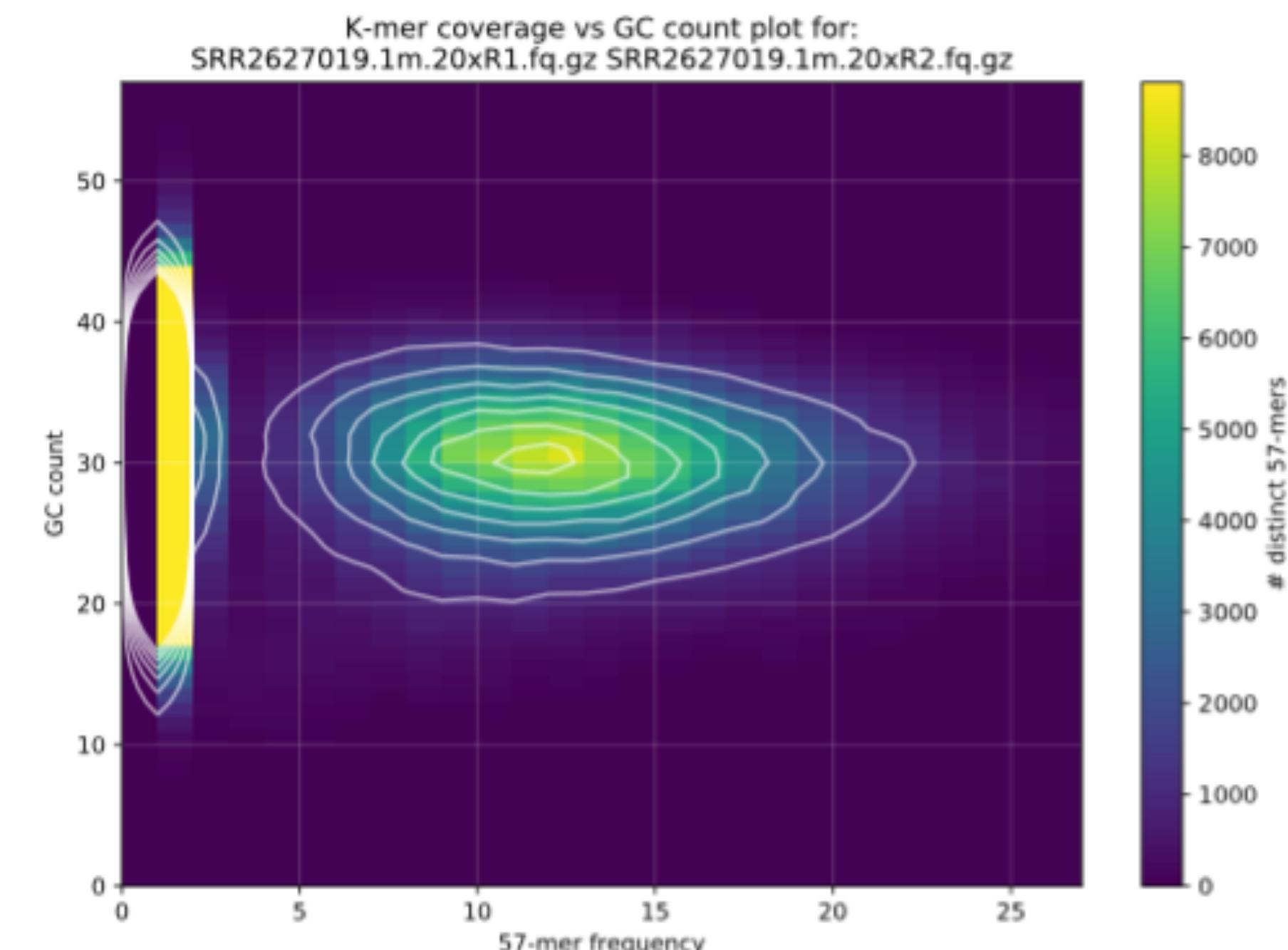
Contamination check with KAT

00_inputQC

40x dataset



20x dataset



Now we're stuck with many individual reports, how to compressed all that information?

MultiQC to the rescue

Summarize QC results

00_inputQC

Run MultiQC:

```
$ sbatch 4-runMultiQC.sh
```

Copy report to the local machine:

```
$ scp -r "<user>@idragon:/home/lehmanr/Projects/scratch/hts_workshop2019/00_inputQC/multiqc_**" .
```

And view the results:

```
$ open multiqc_report.html
```

Sample Name	% Dups	% GC	Length	M Seqs
SRR2627019.1m.10xR1	2.2%	51%	278 bp	0.0
SRR2627019.1m.10xR2	1.9%	52%	279 bp	0.0
SRR2627019.1m.10x_trimmed_R1	2.0%	51%	263 bp	0.0
SRR2627019.1m.10x_trimmed_R2	1.8%	51%	220 bp	0.0
SRR2627019.1m.20xR1	4.6%	51%	278 bp	0.0
SRR2627019.1m.20xR2	3.6%	52%	278 bp	0.0
SRR2627019.1m.20x_trimmed_R1	4.3%	51%	263 bp	0.0
SRR2627019.1m.20x_trimmed_R2	3.5%	51%	220 bp	0.0
SRR2627019.1m.40xR1	6.5%	50%	268 bp	0.1
SRR2627019.1m.40xR2	5.4%	51%	269 bp	0.1
SRR2627019.1m.40x_trimmed_R1	6.1%	50%	254 bp	0.1
SRR2627019.1m.40x_trimmed_R2	5.4%	50%	215 bp	0.1
SRR2627019.1m.80xR1	16.9%	51%	278 bp	0.1
SRR2627019.1m.80xR2	14.4%	52%	278 bp	0.1
SRR2627019.1m.80x_trimmed_R1	15.1%	51%	263 bp	0.1
SRR2627019.1m.80x_trimmed_R2	13.7%	51%	220 bp	0.1
SRR2986339_subsampled_1_trimmed_R1	0.5%	50%	97 bp	0.4
SRR2986339_subsampled_2_trimmed_R2	0.3%	50%	90 bp	0.4
tardigrade_SRR2986339_subsampled_1	0.6%	51%	101 bp	0.4
tardigrade_SRR2986339_subsampled_2	0.5%	51%	101 bp	0.4

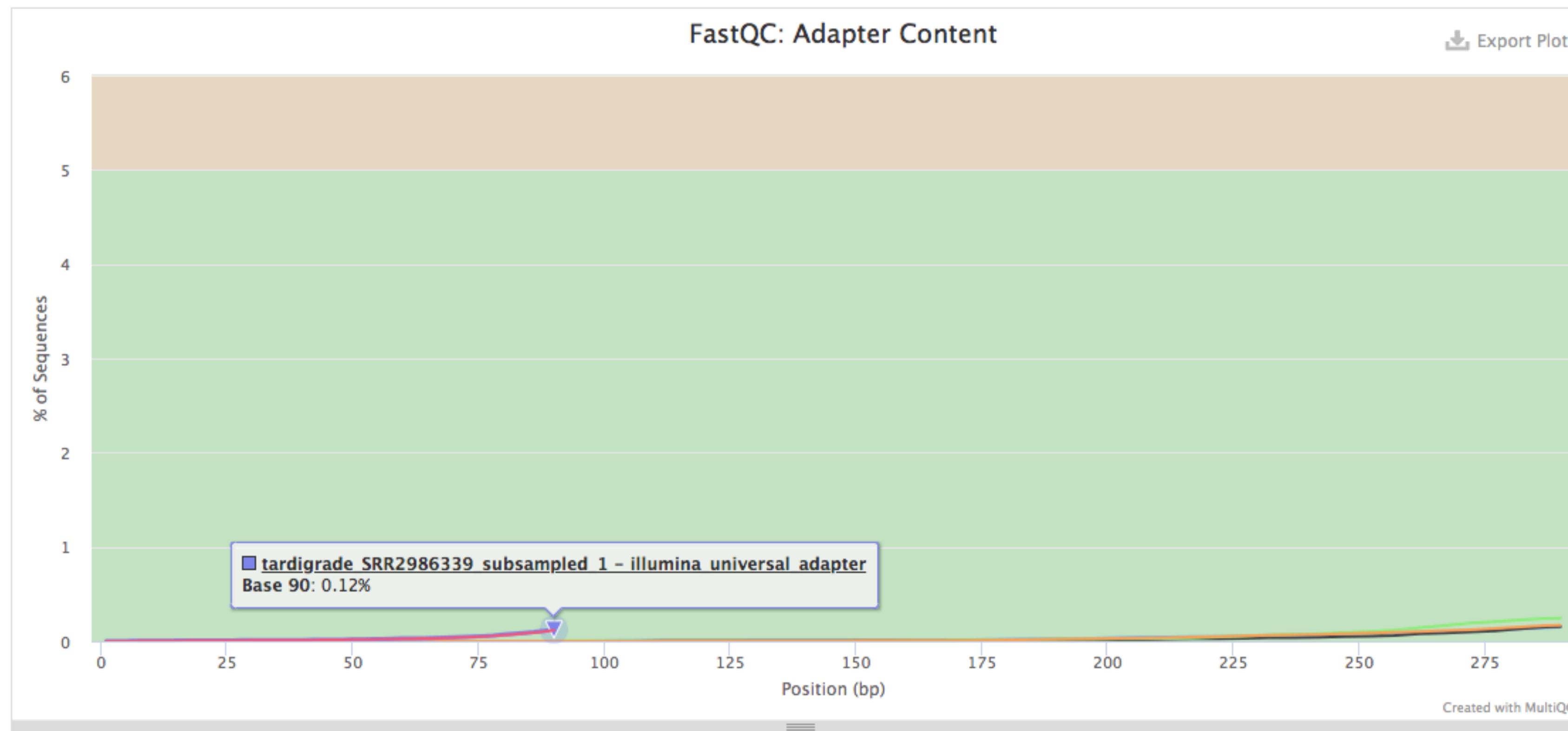
FastQC: Sequence Counts

Export Plot

00_inputQC

Adapter content

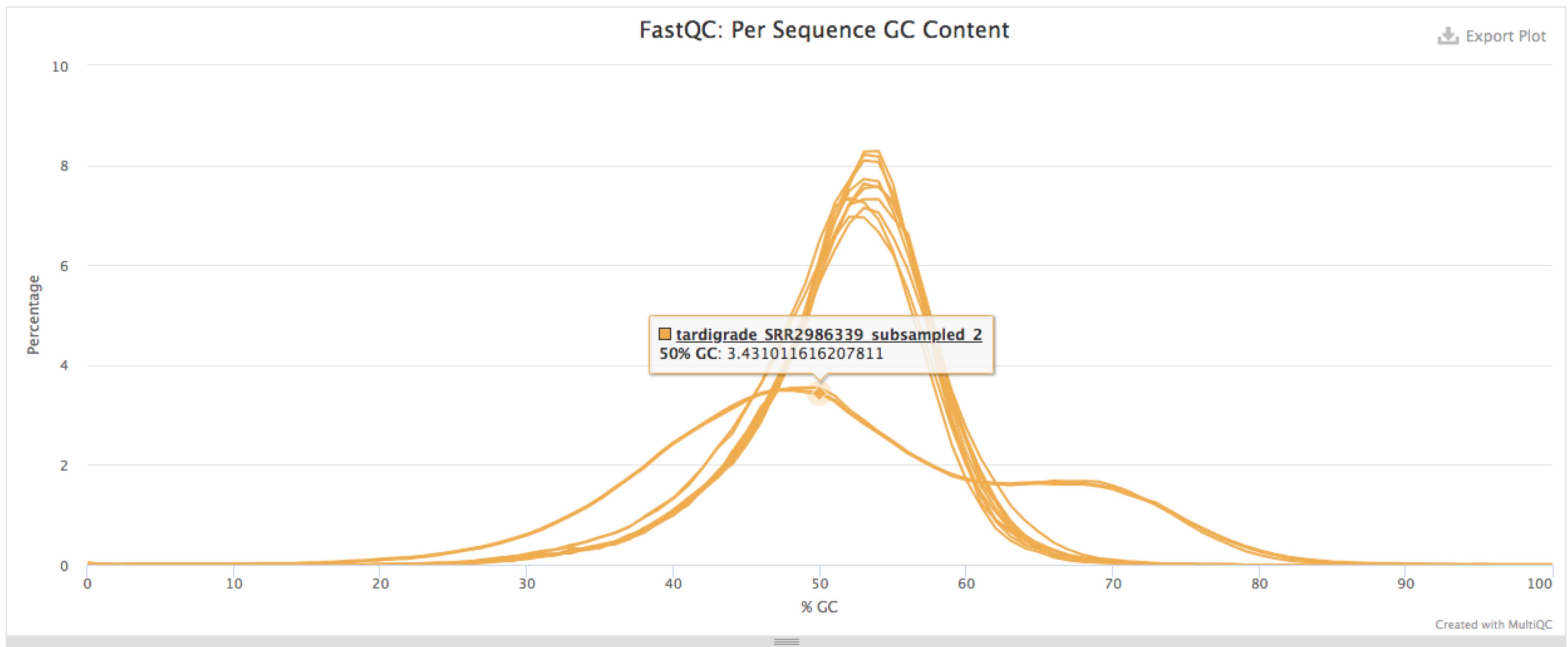
- Adapter content low
- removal by Illumina pipeline worked already well



00_inputQC

GC distribution

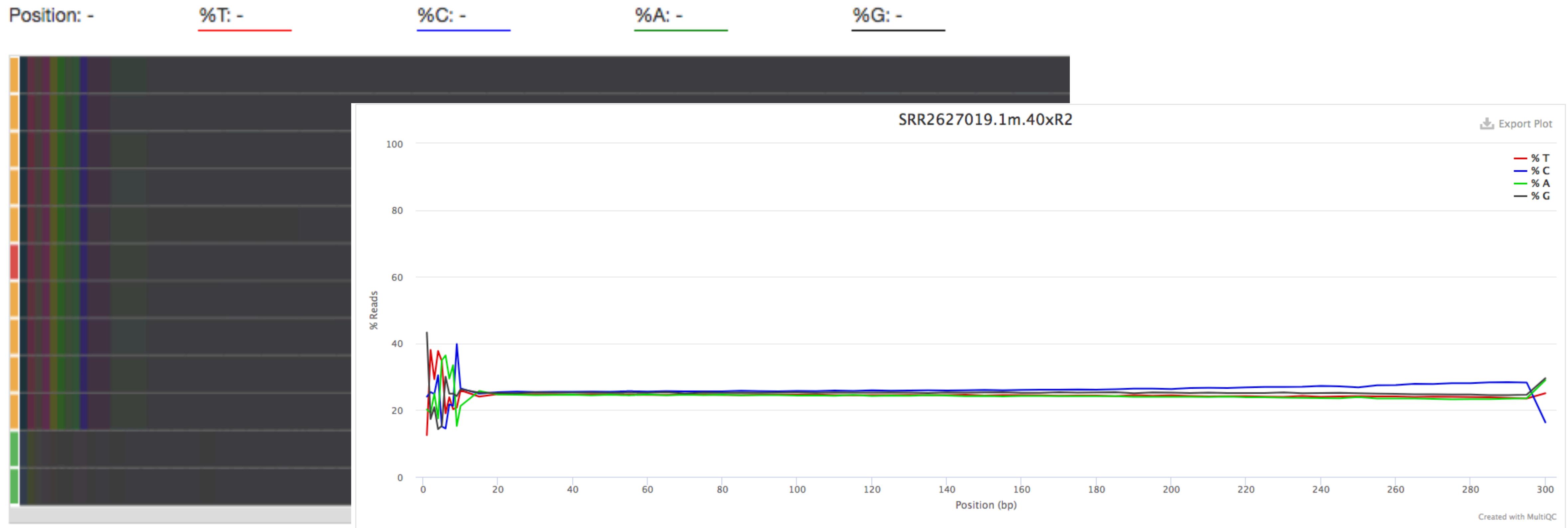
- Tardigrade bimodal GC content looks suspicious...
- E. coli looks ok



00_inputQC

- Some samples show bias to certain base at the very beginning
- virtually always the case, due to not-really random priming
- usually not a problem + cutting biased bases wouldn't help

Base Bias



Assembly

01_assembly

Assemble short reads

```
#!/bin/bash
#SBATCH --job-name=spades
#SBATCH --nodes=1
#SBATCH --ntasks=1
#SBATCH --cpus-per-task=5
#SBATCH --time=6:00:00
#SBATCH --mem=5000
#SBATCH --mail-type=end
#SBATCH --reservation=bioscience_ml

# SPAdes assembler
# https://www.liebertpub.com/doi/abs/10.1089/cmb.2012.0021

module load spades

set -vex

# run assembler on trimmed 40x ecoli data
spades.py -m30 -t5 -o ./ecoli_assembly-SRR2627019_40x \
-1 ../00_inputQC/output/SRR2627019.1m.40x_trimmed_R1.fq.gz \
-2 ../00_inputQC/output/SRR2627019.1m.40x_trimmed_R2.fq.gz \
> spades_ecoli_40x.log
```

01_assembly

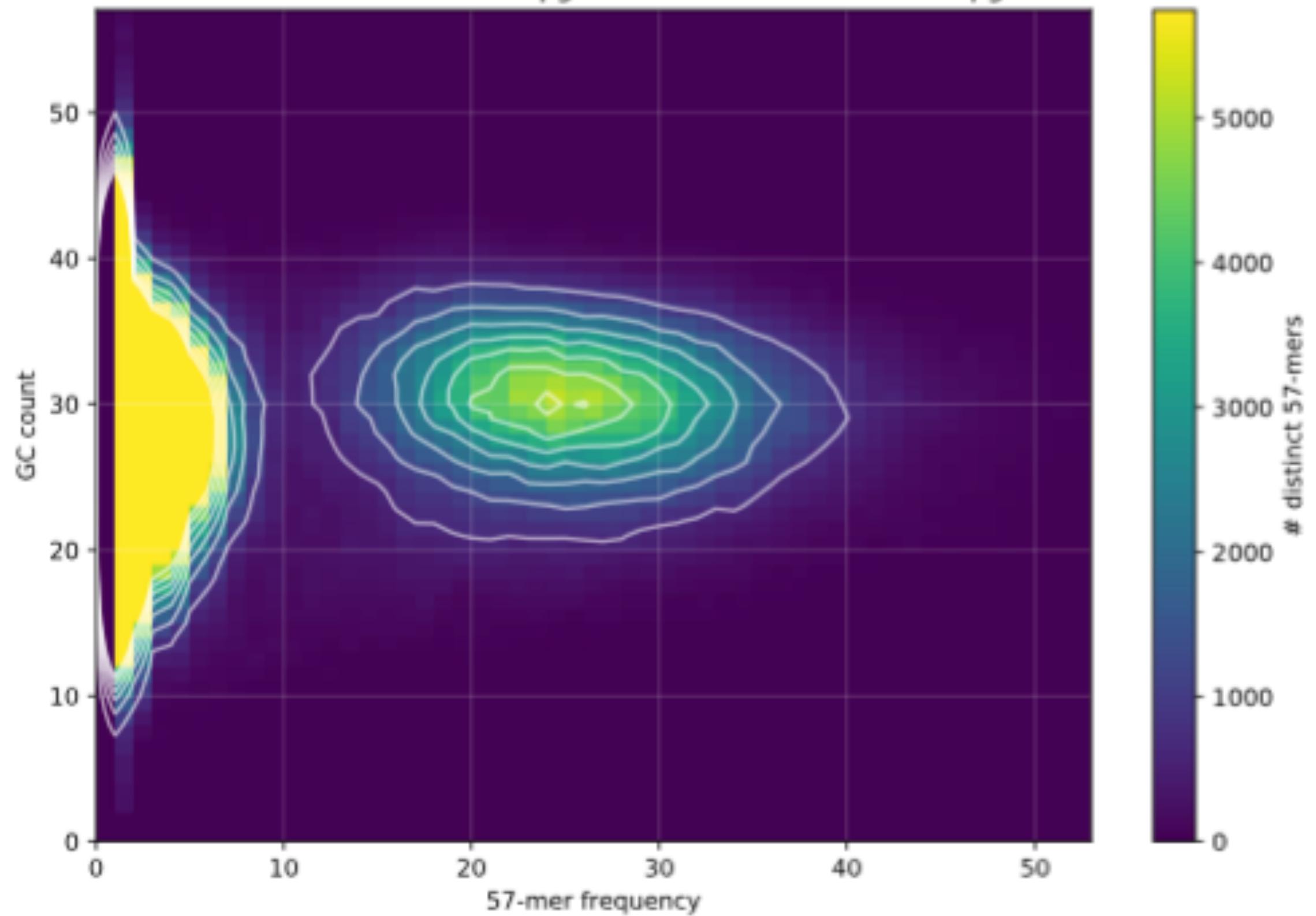
The result

- have a look at spades.log
did SPAdes terminate successfully?
 - where are the candidate assemblies for varying k?
- Good for reproducibility:
- params.txt lists parameters
 - spades.log gives detailed idea of what happened

```
$ ls -laht ecoli_assembly-SRR2627019_40x
total 20M
-rw-r--r-- 1 lehmannr g-lehmannr 200K Jul 16 17:20 spades.log
-rw-r--r-- 1 lehmannr g-lehmannr 532 Jul 16 17:20 warnings.log
drwxr-xr-x 2 lehmannr g-lehmannr 4.0K Jul 16 17:20 misc
drwxr-xr-x 11 lehmannr g-lehmannr 4.0K Jul 16 17:20 .
drwxr-xr-x 2 lehmannr g-lehmannr 4.0K Jul 16 17:20 tmp
-rw-r--r-- 1 lehmannr g-lehmannr 216K Jul 16 17:20 contigs.paths
-rw-r--r-- 1 lehmannr g-lehmannr 6.1M Jul 16 17:20 assembly_graph.fastq
-rw-r--r-- 1 lehmannr g-lehmannr 3.1M Jul 16 17:20 assembly_graph_with_scaffolds.gfa
-rw-r--r-- 1 lehmannr g-lehmannr 216K Jul 16 17:20 scaffolds.paths
-rw-r--r-- 1 lehmannr g-lehmannr 3.0M Jul 16 17:20 scaffolds.fasta
-rw-r--r-- 1 lehmannr g-lehmannr 3.0M Jul 16 17:20 contigs.fasta
-rw-r--r-- 1 lehmannr g-lehmannr 3.0M Jul 16 17:20 before_rr.fasta
drwxr-xr-x 4 lehmannr g-lehmannr 4.0K Jul 16 17:20 K127
drwxr-xr-x 3 lehmannr g-lehmannr 4.0K Jul 16 17:19 K99
drwxr-xr-x 3 lehmannr g-lehmannr 4.0K Jul 16 17:19 K77
drwxr-xr-x 3 lehmannr g-lehmannr 4.0K Jul 16 17:19 K55
drwxr-xr-x 3 lehmannr g-lehmannr 4.0K Jul 16 17:19 K33
drwxr-xr-x 3 lehmannr g-lehmannr 4.0K Jul 16 17:18 K21
-rw-r--r-- 1 lehmannr g-lehmannr 137 Jul 16 17:18 dataset.info
drwxr-xr-x 3 lehmannr g-lehmannr 4.0K Jul 16 17:18 corrected
-rw-r--r-- 1 lehmannr g-lehmannr 1.9K Jul 16 17:18 params.txt
-rw-r--r-- 1 lehmannr g-lehmannr 324 Jul 16 17:17 input_dataset.yaml
drwxr-xr-x 4 lehmannr g-lehmannr 4.0K Jul 16 17:17 ..
```

Now edit runSpades.sh and assemble the other datasets!

K-mer coverage vs GC count plot for:
SRR2627019.1m.40xR1.fq.gz SRR2627019.1m.40xR2.fq.gz



Read the
Log Files!

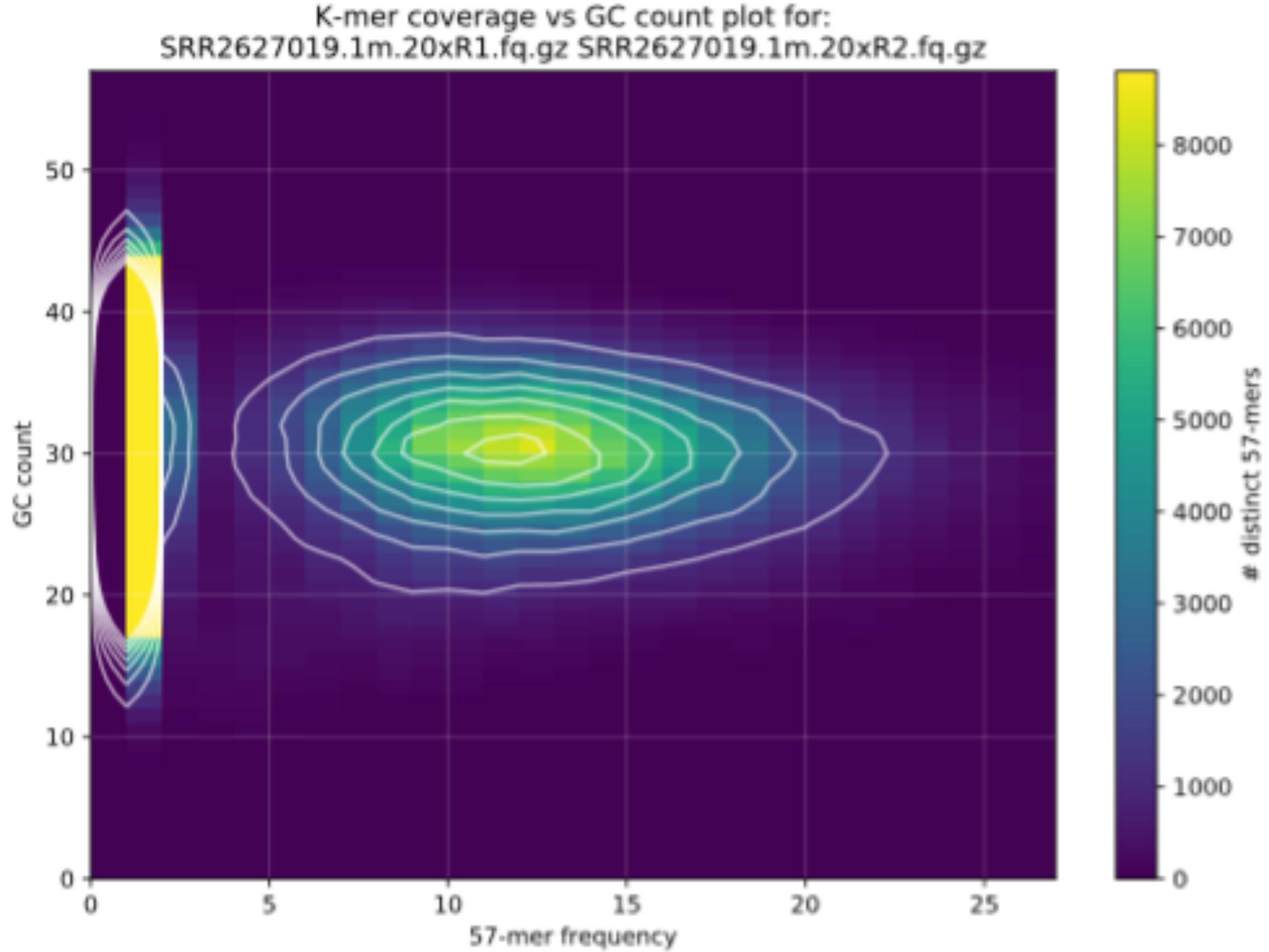
40x

01_assembly

===== SPAdes pipeline finished WITH WARNINGS!

== Error correction and assembling warnings:

```
* 0:00:09.658 112M / 3G  WARN  General          (kmer_coverage_model.cpp : 218)  Too many erroneous kmers, the estimates might be unreliable
* 0:00:12.967  96M / 3G  WARN  General          (simplification.cpp       : 578)  The determined erroneous connection coverage threshold may be determined
===== Warnings saved to /scratch/dragon/intel/lehmannr/bscratch/lehmannr/hts_workshop2019/01_assembly/ecoli_assembly-SRR2627019_40x/warnings.log
```



Read the
Log Files!

20x

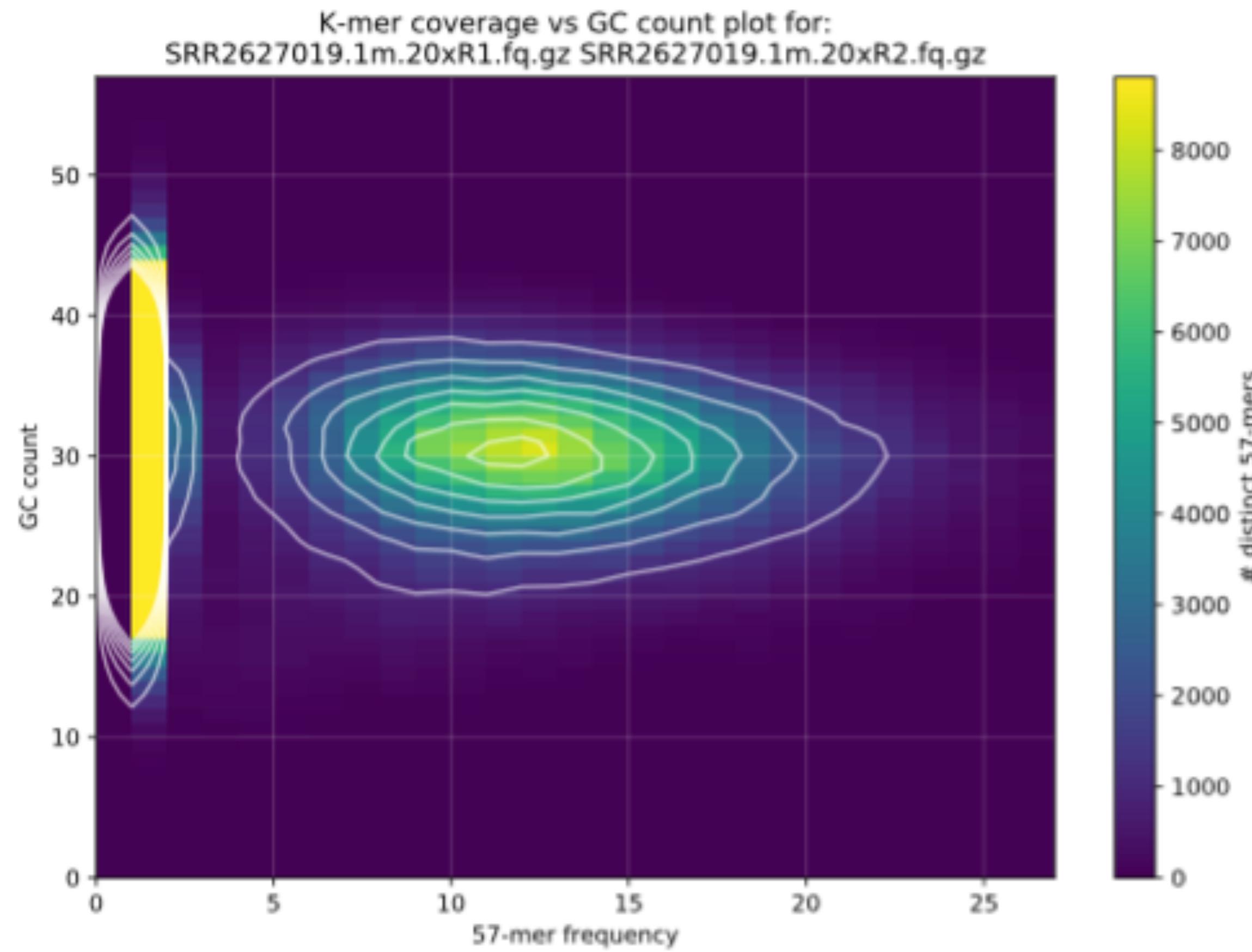
01_assembly

===== SPAdes pipeline finished WITH WARNINGS!

== Error correction and assembling warnings:

* 0:00:12.785 44M / 3G WARN General (kmer_coverage_model.cpp : 366) Failed to determine erroneous kmer threshold. Threshold set to: 0

===== Warnings saved to /scratch/dragon/intel/lehmanr/bscratch/lehmanr/hts_workshop2019/01_assembly/ecoli_assembly-SRR2627019_20x/warnings.log



**Read the
Log Files!**

01_assembly

**no error in 80x
assembly process**

===== SPAdes pipeline finished.

SPAdes log can be found here: /scratch/dragon/intel/lehmanr/bscratch/lehmanr/hts_workshop2019/01_assembly/ecoli_assembly-SRR2627019_80x/spades.log

How many pieces were assembled?

```
$ grep -c '>' K*/final_contigs.fasta  
K127/final_contigs.fasta:2668  
K21/final_contigs.fasta:6306  
K33/final_contigs.fasta:4915  
K55/final_contigs.fasta:3494  
K77/final_contigs.fasta:3223  
K99/final_contigs.fasta:2972
```

SPAdes chose the K=127 assembly!

```
$ grep -c '>' contigs.fasta  
2668
```

```
$ grep -c '>' scaffolds.fasta  
2666
```

And scaffolding closed 2 gaps

But isn't 2666 scaffolds on 1 Mb a bit much?

What about the other datasets?

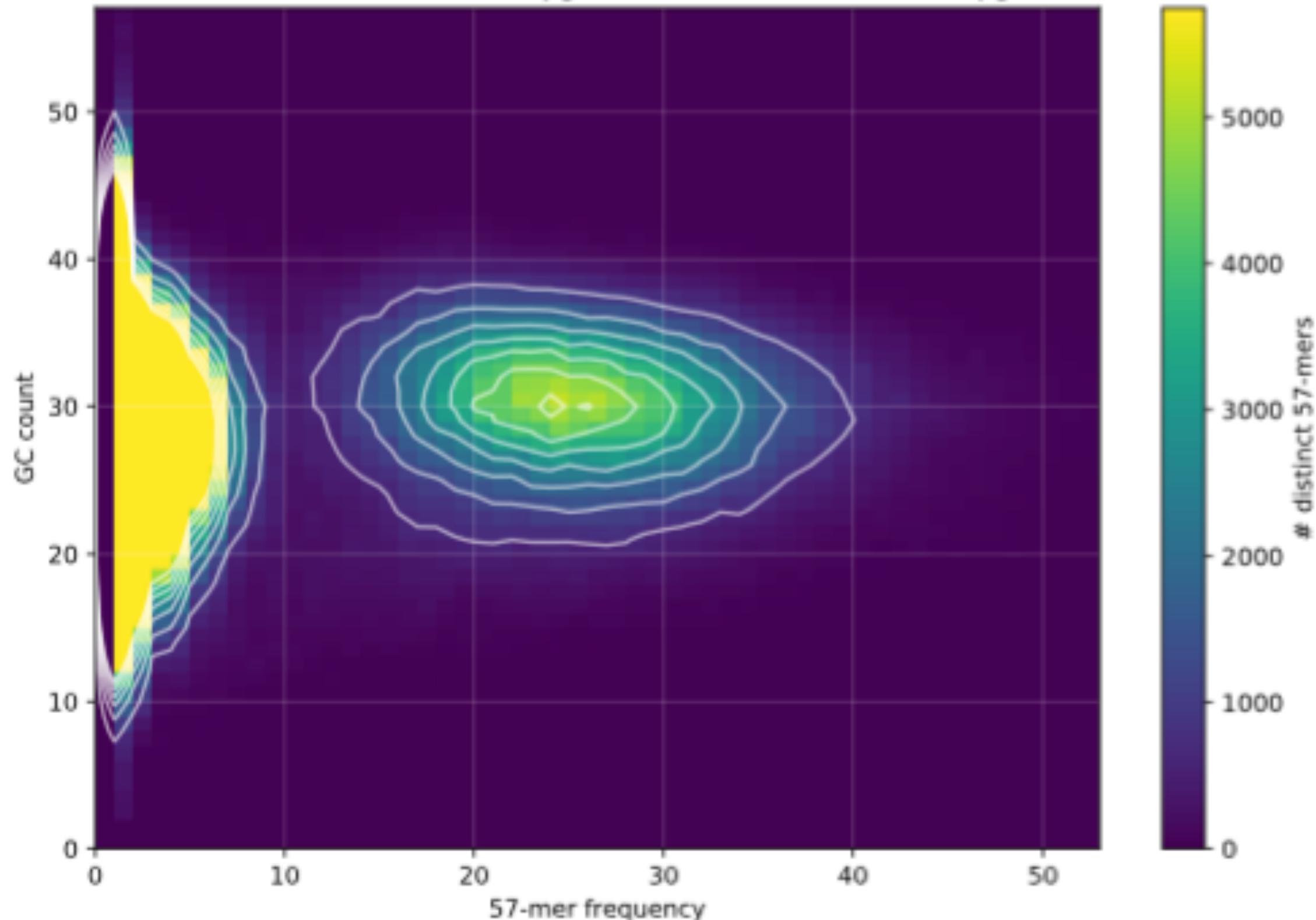
```
$ grep -c '>' K*/final_contigs.fasta
K127/final_contigs.fasta:23
K21/final_contigs.fasta:482
K33/final_contigs.fasta:278
K55/final_contigs.fasta:88
K77/final_contigs.fasta:79
K99/final_contigs.fasta:41
```

80x - only 23 contigs!
Due to higher coverage?

```
$ grep -c '>' ecoli_assembly-SRR2627019_*/contigs.fasta
ecoli_assembly-SRR2627019_10x/contigs.fasta:109
ecoli_assembly-SRR2627019_20x/contigs.fasta:30
ecoli_assembly-SRR2627019_40x/contigs.fasta:2668
ecoli_assembly-SRR2627019_80x/contigs.fasta:23
```

- Higher coverage does yield fewer contigs
- BUT: only 40x has much more contigs than all other?

K-mer coverage vs GC count plot for:
SRR2627019.1m.40xR1.fq.gz SRR2627019.1m.40xR2.fq.gz



Contamination assembled

01_assembly

The extra bits are likely Yersinia sequence that we detected earlier

[lehmanr@dbn503-35-r output]\$ awk '\$1 > 1' ecoli_kraken2_report.txt					
99.67	94735	298	R	1	root
99.35	94436	2	R1	131567	cellular organisms
99.30	94387	203	D	2	Bacteria
99.08	94175	257	P	1224	Proteobacteria
98.80	93909	69	C	1236	Gammaproteobacteria
98.72	93831	917	O	91347	Enterobacterales
74.22	70550	39261	F	543	Enterobacteriaceae
31.88	30298	1603	G	561	Escherichia
29.75	28282	27968	S	562	Escherichia coli
23.52	22356	15	F	1903411	Yersiniaceae
23.50	22337	205	G	629	Yersinia
23.28	22124	16426	G1	1649845	Yersinia pseudotuberculosis complex
5.76	5475	5410	S	632	Yersinia pestis



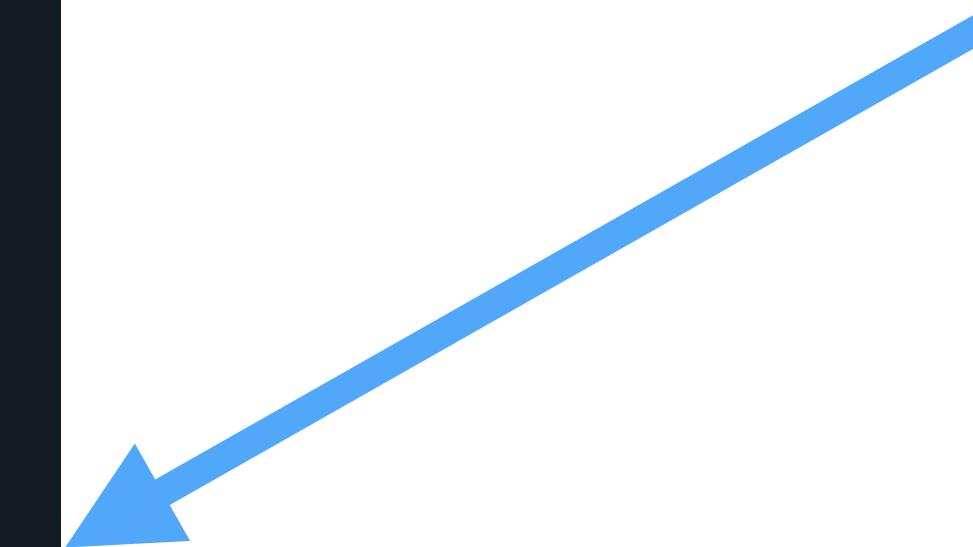
Are all assembled contigs trustworthy?

01_assembly

80x Assembly

- a first idea how the assembly looks like
- e.g. noise contigs like Node 23 (short and very high k-mer coverage compared to rest - repetitive / low complexity region?)

```
$ grep '>' scaffolds.fasta
>NODE_1_length_256701_cov_21.974647
>NODE_2_length_210887_cov_21.111340
>NODE_3_length_173798_cov_20.763576
>NODE_4_length_78753_cov_20.050085
>NODE_5_length_69295_cov_20.220275
>NODE_6_length_64076_cov_21.067022
>NODE_7_length_36688_cov_21.507262
>NODE_8_length_32471_cov_19.661792
>NODE_9_length_15646_cov_26.884335
>NODE_10_length_15264_cov_20.240602
>NODE_11_length_10776_cov_17.752371
>NODE_12_length_8831_cov_18.608915
>NODE_13_length_6805_cov_17.641360
>NODE_14_length_4267_cov_21.709903
>NODE_15_length_1345_cov_48.496716
>NODE_16_length_1255_cov_71.635638
>NODE_17_length_1195_cov_68.991573
>NODE_18_length_550_cov_6.846336
>NODE_19_length_258_cov_24.335878
>NODE_20_length_256_cov_41.186047
>NODE_21_length_192_cov_40.476923
>NODE_22_length_136_cov_49.444444
>NODE_23_length_128_cov_245.000000
```



Lets compare the candidate assemblies from varying K

Which K is best?

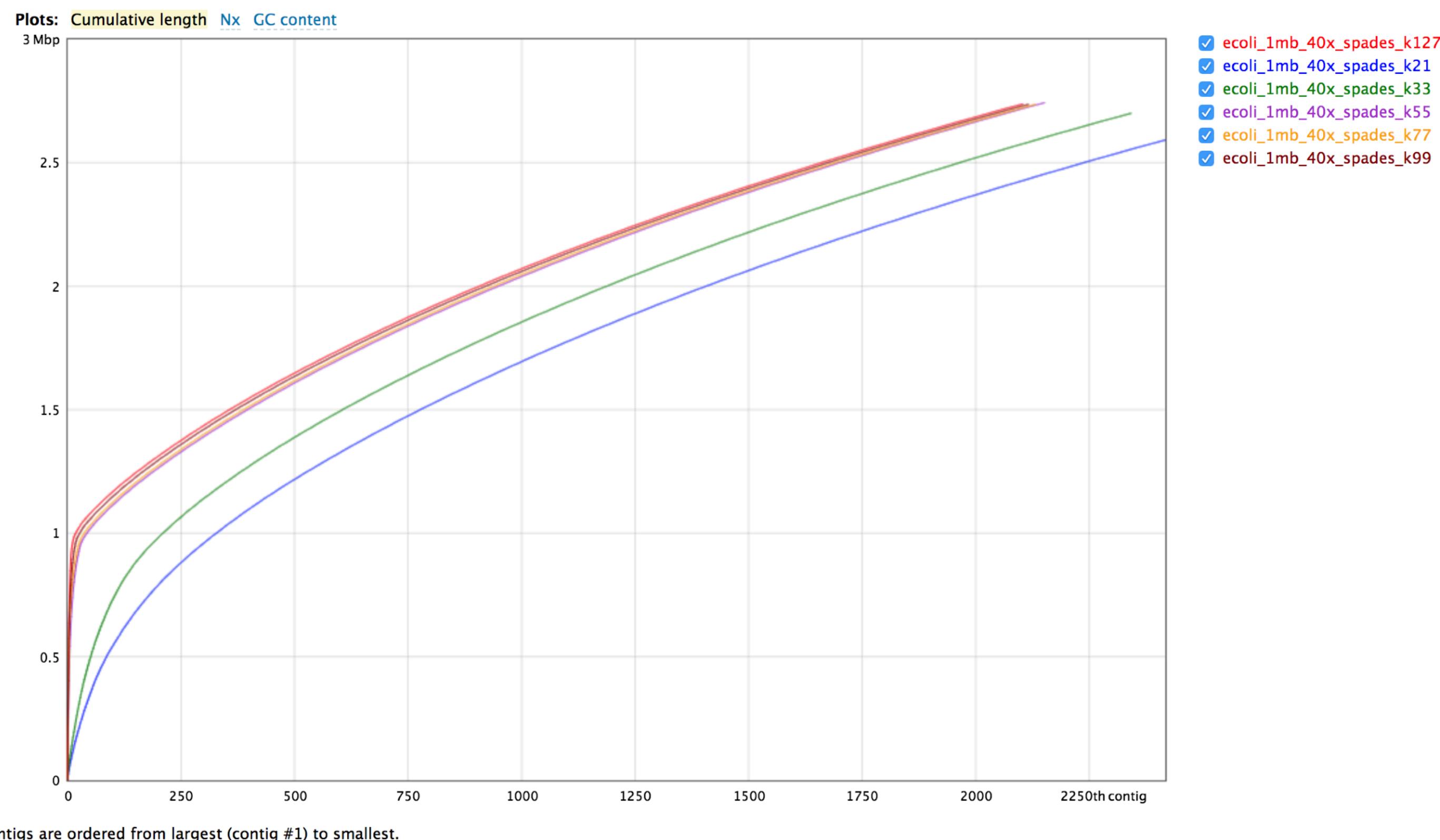
01_assembly

```
$ nano 2-runQuast_compK.sh  
$ sbatch 2-runQuast_compK.sh
```

And download the report to the local machine:

```
$ scp -r "<user>@ilogo.ibex.kaust.edu.sa:<path>/output/quast_results".  
$ open quast_results/latest/report.html
```

K=127 yields most continuous assembly



01_assembly

- here: contiguity = more completeness -> spades chooses right candidate
- final assembly has some fewer contigs due to scaffolding

Which K is best?

Worst Median Best Show heatmap

	ecoli_1mb_40x_spades_k127	ecoli_1mb_40x_spades_k21	ecoli_1mb_40x_spades_k33	ecoli_1mb_40x_spades_k55	ecoli_1mb_40x_spades_k77	ecoli_1mb_40x_spades_k99
# contigs	2105	2418	2343	2153	2132	2117
# contigs (>= 0 bp)	2668	6306	4915	3494	3223	2972
# contigs (>= 1000 bp)	469	660	625	500	494	480
# contigs (>= 5000 bp)	14	52	67	29	25	20
# contigs (>= 10000 bp)	11	5	21	21	18	16
# contigs (>= 25000 bp)	8	0	0	11	8	9
# contigs (>= 50000 bp)	6	0	0	5	4	4
Largest contig	252 996	12 664	20 787	159 849	196 920	204 292
Total length	2 738 532	2 592 217	2 700 462	2 743 708	2 736 906	2 736 802
Total length (>= 0 bp)	2 998 608	3 432 198	3 418 406	3 284 335	3 198 317	3 113 983
Total length (>= 1000 bp)	1 617 788	1 387 918	1 519 521	1 608 478	1 611 886	1 615 215
Total length (>= 5000 bp)	985 580	362 198	590 823	956 258	960 080	976 586
Total length (>= 10000 bp)	960 478	58 450	270 770	893 876	903 326	946 795
Total length (>= 25000 bp)	918 727	0	0	740 990	750 088	841 263
Total length (>= 50000 bp)	849 636	0	0	543 686	606 934	653 923
N50	1241	1065	1137	1214	1226	1241
N75	749	711	731	746	749	749
L50	245	572	467	285	274	257
L75	978	1327	1222	1025	1007	990
GC (%)	49.33	49.3	49.31	49.3	49.32	49.34
Mismatches						
# N's	0	0	0	0	0	0
# N's per 100 kbp	0	0	0	0	0	0

But which dataset yields best assembly?

Lets make a proper comparison of our assembly candidates with Quast

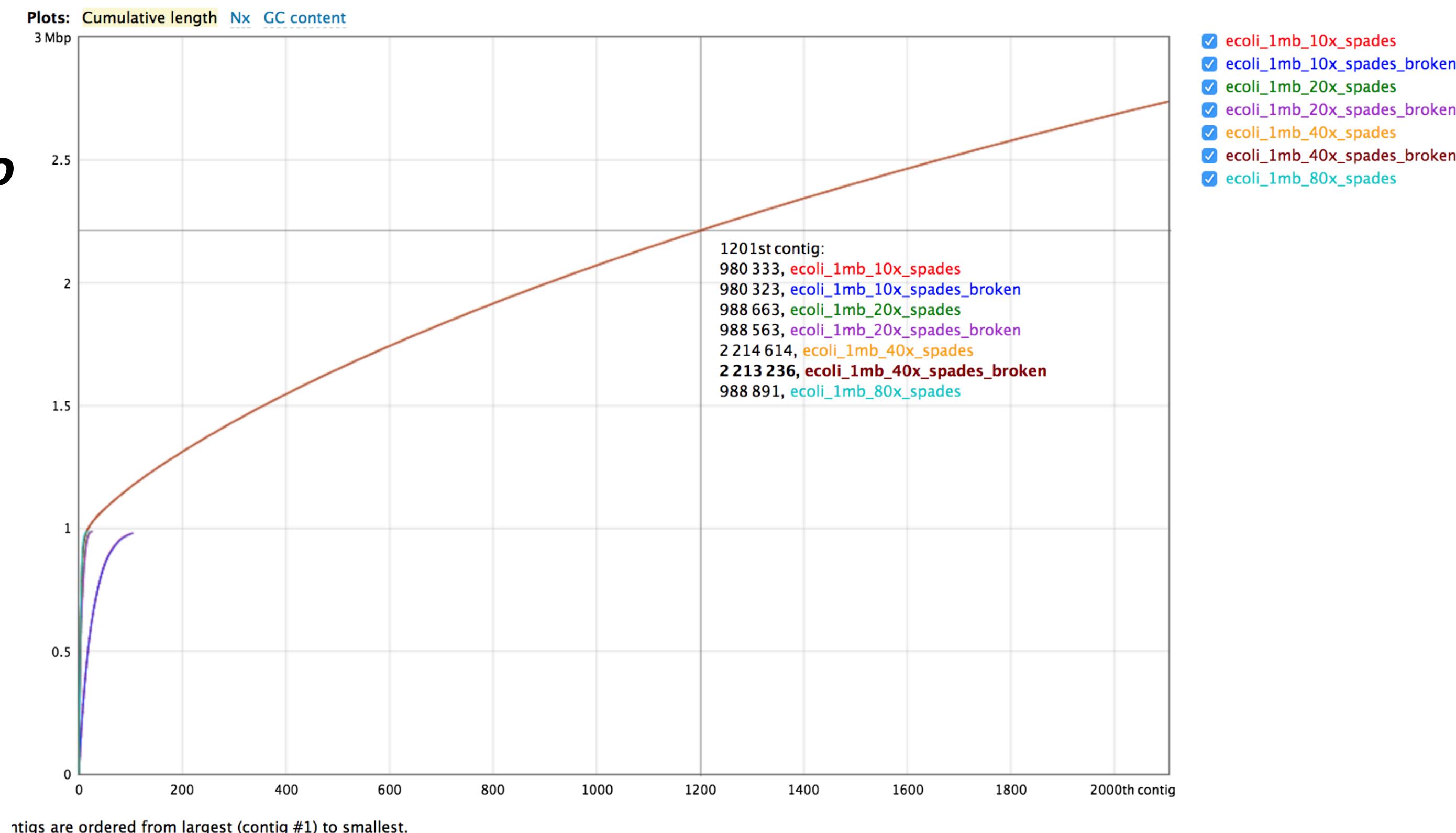
Read the script, it does something new:

```
$ nano 3-runQuast_compCoverage.sh  
$ sbatch 3-runQuast_compCoverage.sh
```

And download the report to the local machine:

```
$ scp -r "<user>@ilogo.ibex.kaust.edu.sa:<path to output>/quast_results".  
$ open quast_results/latest/report.html
```

80x assembly is most continuous!
BUT: we only look at how long fragments are...



**Assembly
Quality**

How to compare to a genome reference?

02_assemblyQC

Lets compare our assemblies from different K values with a reference.

```
$ nano 1-runQuastDifferentKmer.sh  
$ sbatch 1-runQuastDifferentKmer.sh
```

And download the report to the local machine:

```
$ scp -r "<user>@ilogin.ibex.kaust.edu.sa:<path to  
output>/quast_results".  
$ open quast_results/latest/report.html
```

Worst Median Best Show heatmap

	K127_final_contigs	K21_final_contigs	K33_final_contigs	K55_final_contigs	K77_final_contigs	K99_final_contigs
Genome statistics						
Genome fraction (%)	21.276	20.016	20.766	21.156	21.183	21.246
Duplication ratio	1.005	1.005	1.005	1.003	1.004	1.006
# genes	917 + 9 part	634 + 252 part	751 + 143 part	893 + 12 part	900 + 11 part	914 + 10 part
Largest alignment	252 996	12 664	20 787	159 849	196 920	204 292
Total aligned length	988 335	931 785	967 268	982 590	984 089	987 615
NG50	636	582	617	640	637	635
NGA50	-	-	-	-	-	-
LG50	1365	1915	1661	1406	1395	1380
Misassemblies						
# misassemblies	1	0	1	1	1	1
# relocations	1	0	1	1	1	1
# translocations	0	0	0	0	0	0
# inversions	0	0	0	0	0	0
# misassembled contigs	1	0	1	1	1	1
Misassembled contigs length	69 295	0	2749	11 138	11 182	27 947
# local misassemblies	1	1	1	3	3	2
# scaffold gap size misassemblies	0	0	0	0	0	0
# unaligned mis. contigs	3	1	1	1	1	3
Unaligned						
# fully unaligned contigs	2075	2048	2100	2103	2085	2076
Fully unaligned length	1 742 532	1 656 110	1 728 885	1 756 303	1 746 124	1 741 187
# partially unaligned contigs	5	3	3	3	5	5
Partially unaligned length	3962	2498	2510	2532	3962	3962
Mismatches						
# mismatches	38	32	36	45	49	55
# indels	1	0	0	0	0	1
Indels length	1	0	0	0	0	1
# mismatches per 100 kbp	3.85	3.44	3.73	4.58	4.98	5.58
# indels per 100 kbp	0.1	0	0	0	0	0.1
# indels (<= 5 bp)	1	0	0	0	0	1
# indels (> 5 bp)	0	0	0	0	0	0
# N's	0	0	0	0	0	0
# N's per 100 kbp	0	0	0	0	0	0
Statistics without reference						
# contigs	2105	2418	2343	2153	2132	2117
# contigs (>= 0 bp)	2668	6306	4915	3494	3223	2972
# contigs (>= 1000 bp)	469	660	625	500	494	480
# contigs (>= 5000 bp)	14	52	67	29	25	20
# contigs (>= 10000 bp)	11	5	21	21	18	16
# contigs (>= 25000 bp)	8	0	0	11	8	9
# contigs (>= 50000 bp)	6	0	0	5	4	4
Largest contig	252 996	12 664	20 787	159 849	196 920	204 292
Total length	2 738 532	2 592 217	2 700 462	2 743 708	2 736 906	2 736 802
Total length (>= 0 bp)	2 998 608	3 432 198	3 418 406	3 284 335	3 198 317	3 113 983
Total length (>= 1000 bp)	1 617 788	1 387 918	1 519 521	1 608 478	1 611 886	1 615 215
Total length (>= 5000 bp)	985 580	362 198	590 823	956 258	960 080	976 586
Total length (>= 10000 bp)	960 478	58 450	270 770	893 876	903 326	946 795
Total length (>= 25000 bp)	918 727	0	0	740 990	750 088	841 263
Total length (>= 50000 bp)	849 636	0	0	543 686	606 934	653 923
N50	1241	1065	1137	1214	1226	1241
N75	749	711	731	746	749	749
L50	245	572	467	285	274	257
L75	978	1327	1222	1025	1007	990
GC (%)	49.33	49.3	49.31	49.3	49.32	49.34
Similarity statistics						
# similar correct contigs	4	0	1	5	5	5
# similar misassembled blocks	0	0	0	0	0	0

SPAdes candidate Assemblies

- here: contiguity = more completeness -> spades chooses right candidate
- final assembly has some fewer contigs due to scaffolding

02_assemblyQC

Add all candidate assemblies from different datasets

```
$ nano 2-runQuast_differentCoverage.sh  
$ sbatch 2-runQuast_differentCoverage.sh
```

02_assemblyQC

Candidate vs. reference

Interesting key figures:

- N50
- total length
- genome fraction*
- # ref. genes covered*
- duplication ratio*

* if ref. Available

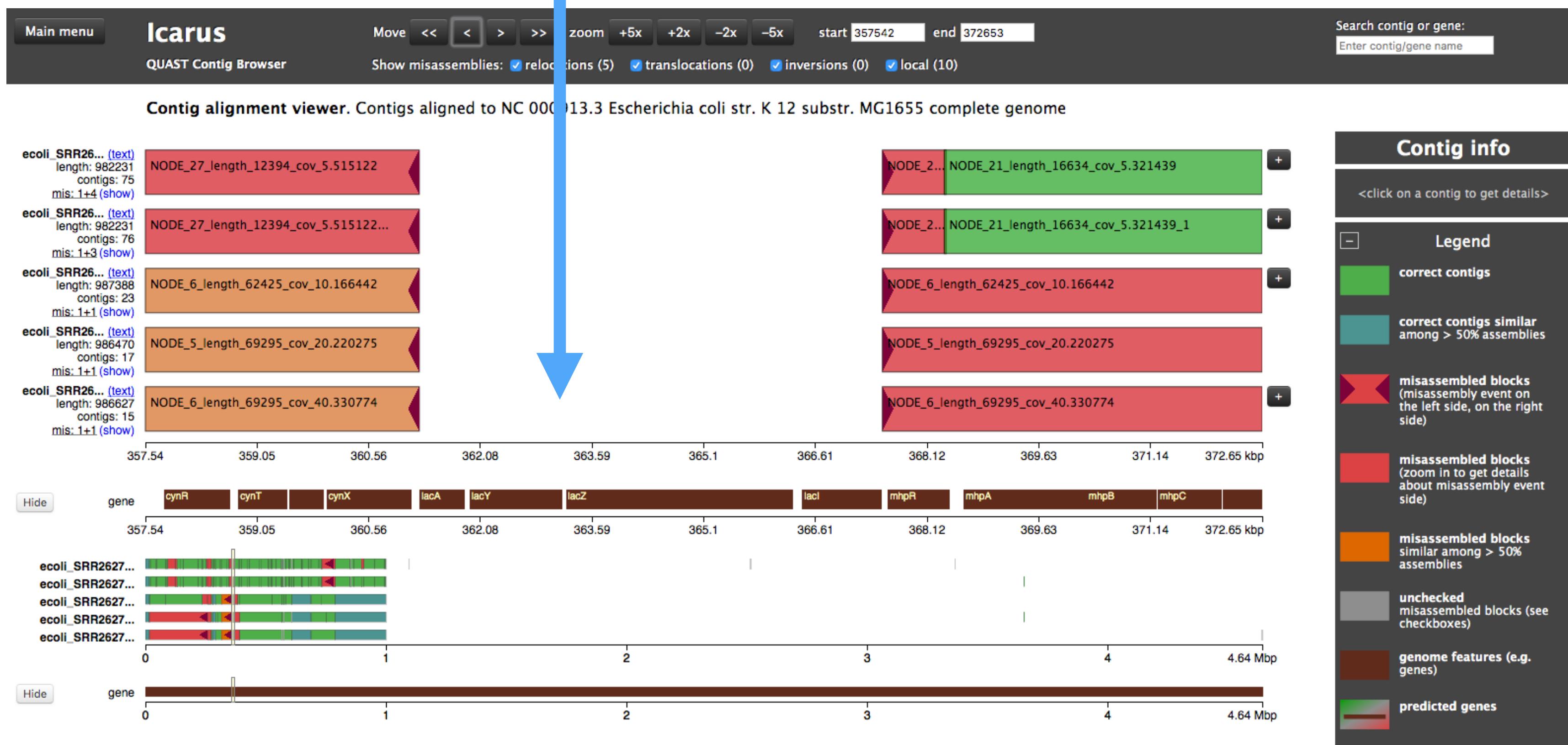
Genome statistics	ecoli_SRR2627019_scaffolds_10x	ecoli_SRR2627019_scaffolds_10...	ecoli_SRR2627019_scaffolds_20x	ecoli_SRR2627019_scaffolds_40x	ecoli_SRR2627019_scaffolds_80x
Genome fraction (%)	21.146	21.14	21.25	21.244	21.242
Duplication ratio	1.002	1.003	1.002	1.001	1.003
# genes	871 + 44 part	871 + 44 part	912 + 8 part	913 + 5 part	914 + 4 part
Largest alignment	60 888	60 888	210 281	256 701	256 701
Total aligned length	982 231	982 231	987 388	986 470	986 627
NGA50	-	-	-	-	-
Misassemblies					
# misassemblies	1	1	1	1	1
Misassembled contigs length	12 394	12 394	62 425	69 295	69 295
Matches					
# mismatches per 100 kbp	24.25	24.26	3.85	1.12	1.12
# indels per 100 kbp	0.1	0.1	0	0	0
# N's per 100 kbp	1.02	0	0	0	0
Statistics without reference					
# contigs	76	77	24	18	15
Largest contig	60 888	60 888	210 453	256 701	256 701
Total length	984 319	984 309	988 799	988 603	988 908
Total length (>= 1000 bp)	980 111	979 264	987 668	988 053	988 908
Total length (>= 10000 bp)	828 190	817 978	953 722	964 355	974 041
Total length (>= 50000 bp)	175 489	175 489	763 963	853 510	890 822

[Extended report](#)

How to compare our candidate assemblies?

02_assemblyQC

Lac negative E. coli strain sequenced

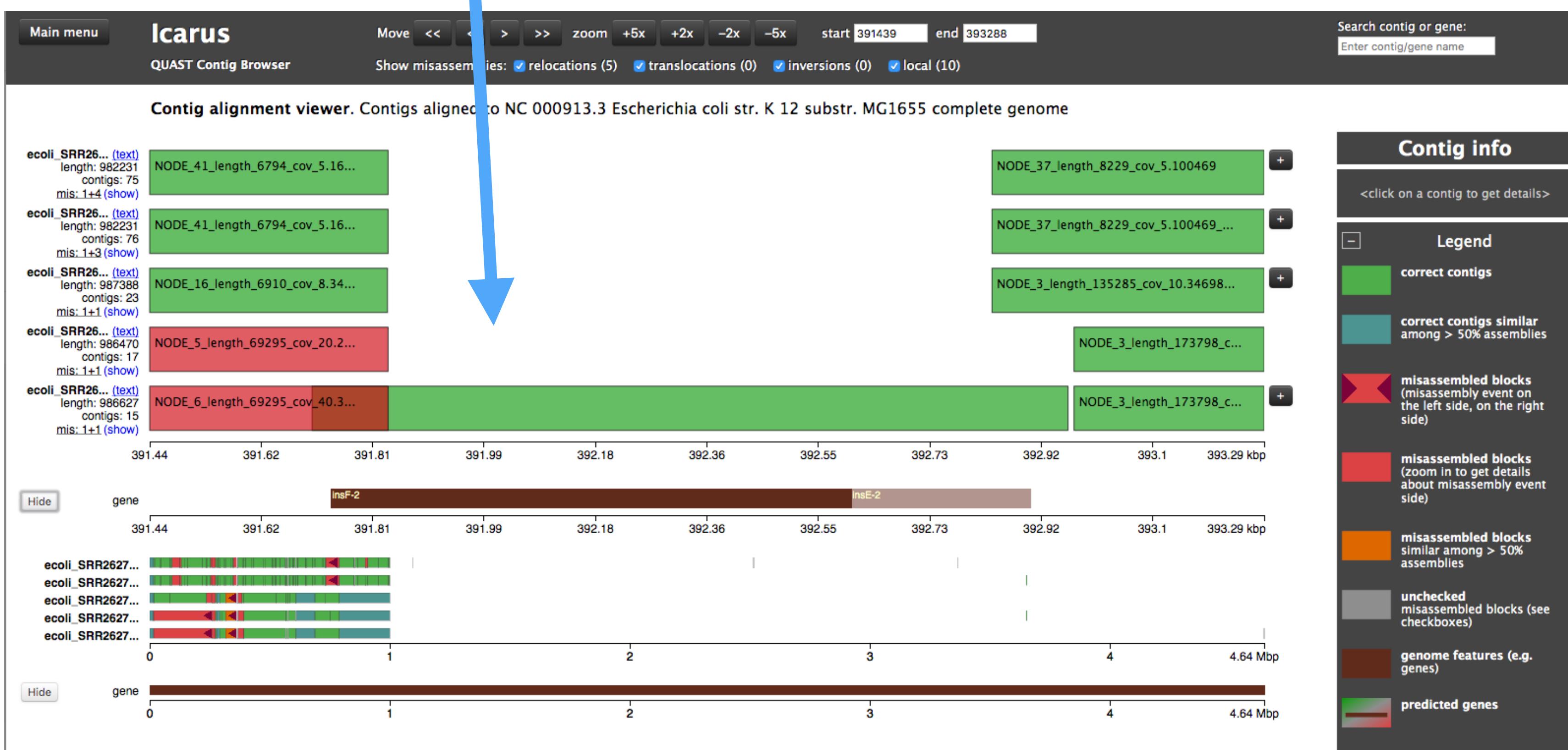


02_assemblyQC

How to compare our candidate assemblies?

Transposon genes

- repetitive sequences hard to assemble
- possibly polymorphic



Gene Annotation + QC

Let's annotate our best assembly

Since we know we are assembling E. coli, we can help the genome annotation pipeline (PROKKA) by supplying the reference protein sequences so that it knows what to look for.

```
$ nano 1-runProkka.sh  
$ sbatch 1-runProkka.sh
```

Count number of genes in gff:

```
awk '$3 == "gene"' <annotation.gff> | wc -l
```

Gene Annotation

Run annotation for all coverage levels...

BUSCO completeness

Use BUSCO to test how many “typical” bacterial genes we find in the reference annotation and in our new annotation (Keep in mind we only assembled the first 1 Mb)

```
$ nano 2-runBUSCO.sh  
$ sbatch 2-runBUSCO.sh
```

Count number of genes in gff:

```
awk '$3 == "gene"' <annotation.gff> | wc -l
```

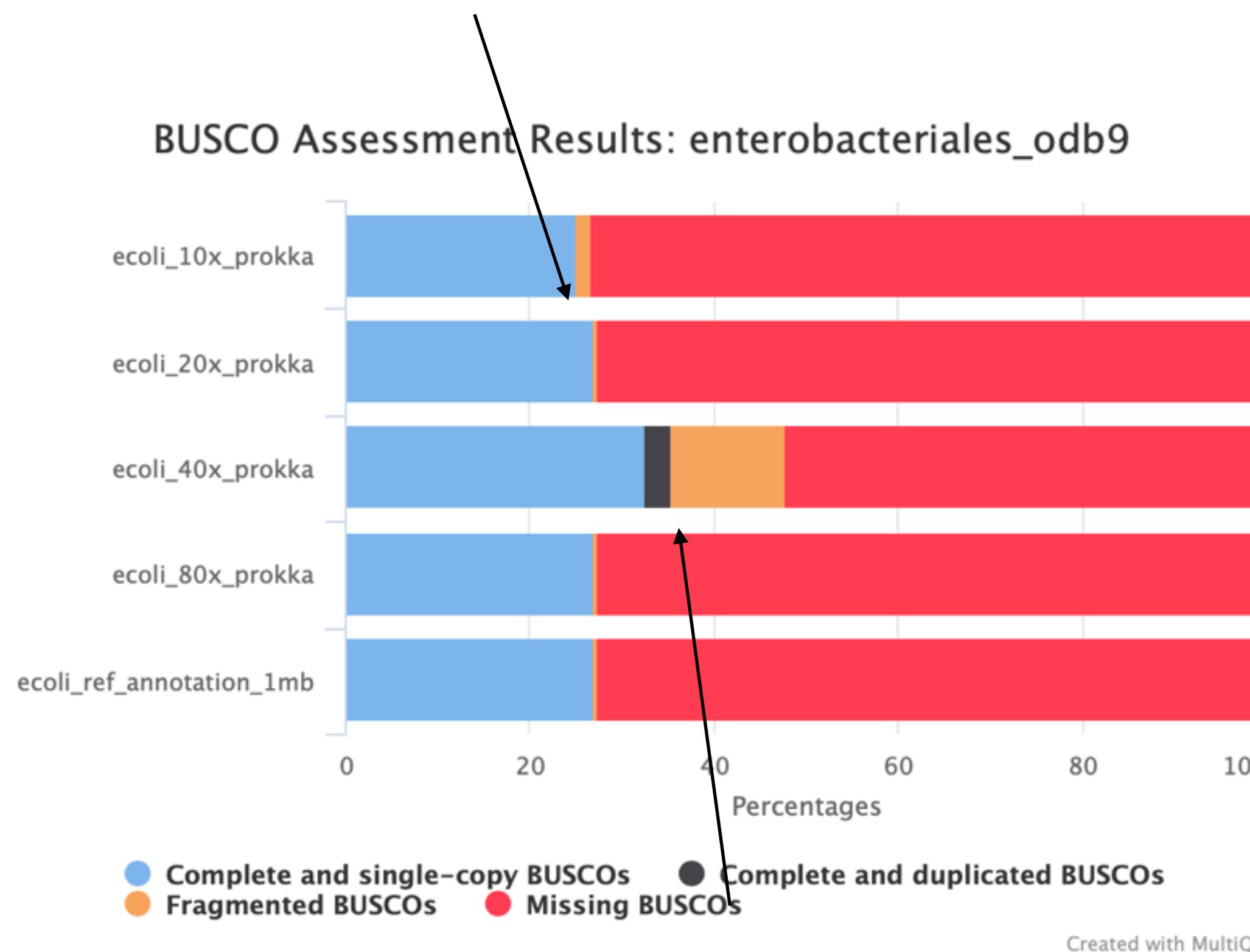
Run annotation for all coverage levels. And summarize again with MultiQC.

```
$ sbatch 3-runMultiQC.sh
```

Download report and check the result.

Coverage increases completeness

BUSCO completeness



Contamination creates
“artificial” completeness

Hybrid Assembly

Let's have a look at the Pacbio reads

FastQC is the tool for the job

```
$ nano 0-runFastQC.sh  
$ sbatch 0-runFastQC.sh
```

How many reads do we have, how long are they, how are they different to short reads?

Hybrid Assembly

Hybrid Assembly

We can use again SPAdes

```
$ nano 1-runSpades.sh  
$ sbatch 1-runSpades.sh
```

Check if SPAdes had any errors / warnings

Which K did it select?

Hybrid Assembly QC

Let's compare all candidate assemblies

```
$ nano 1-runQuast.sh  
$ sbatch 1-runQuast.sh
```

Hybrid Assembly

05_hybridAssemblyQC

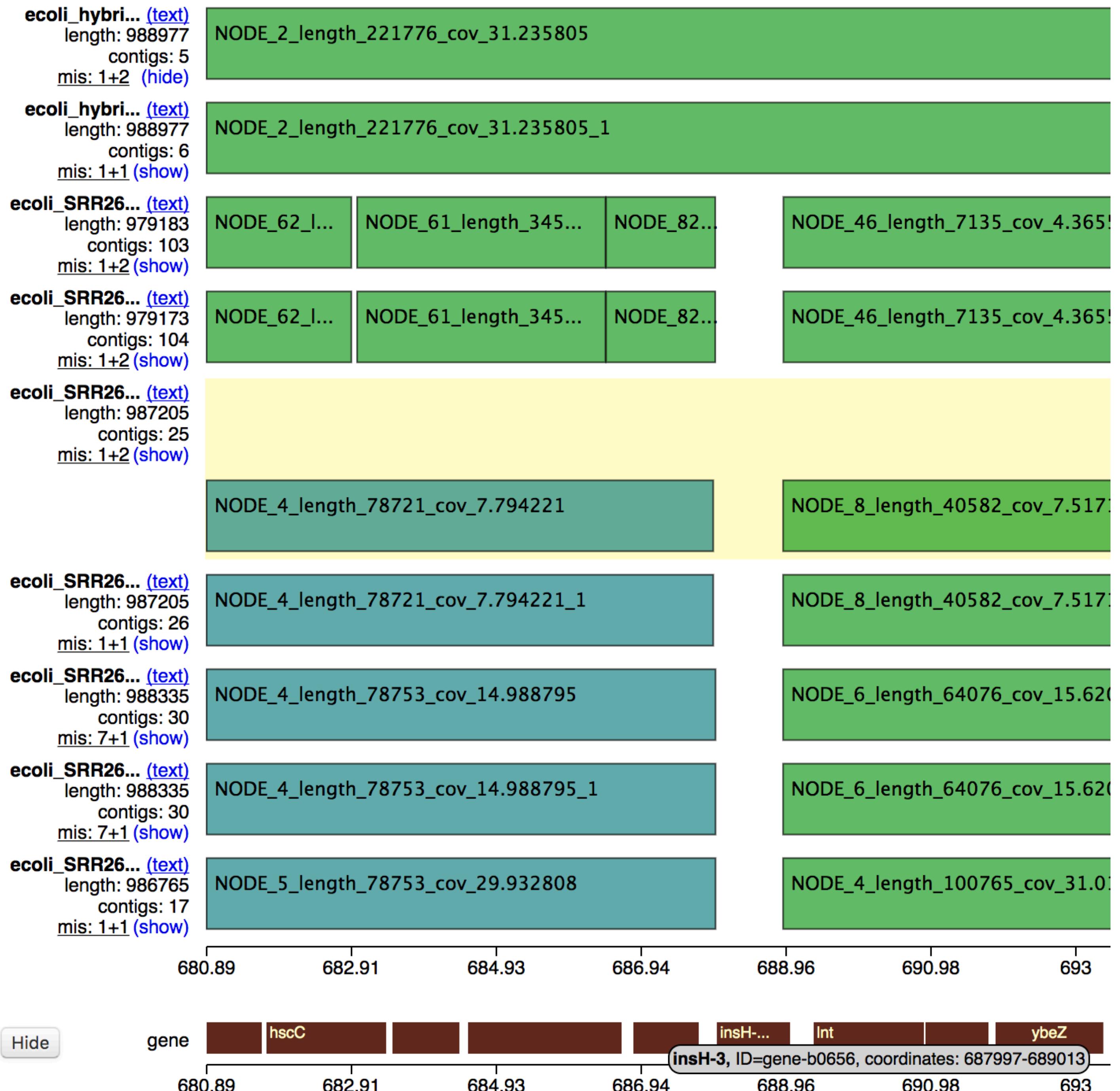
Table 1: QUAST report on assembly and annotation

Assembly	10x	20x	40x	80x	80x+5x
# contigs	104	26	2103	18	6
Largest contig	39,006	210,408	252,996	252,996	368,407
GC (%)	51.44	51.4	49.33	51.39	51.41
N50	20,693	135,285	1,244	173,798	221,776
Misassembled contigs length	12,345	62,652	69,295	69,295	368,407
# genes	842 + 73 part	912 + 7 part	917 + 9 part	912 + 7 part	920 + 2 part

- long reads boost contiguity + completeness
- would be even more pronounced with lower coverage short read data

05_hybridAssemblyQC

Hybrid Assembly

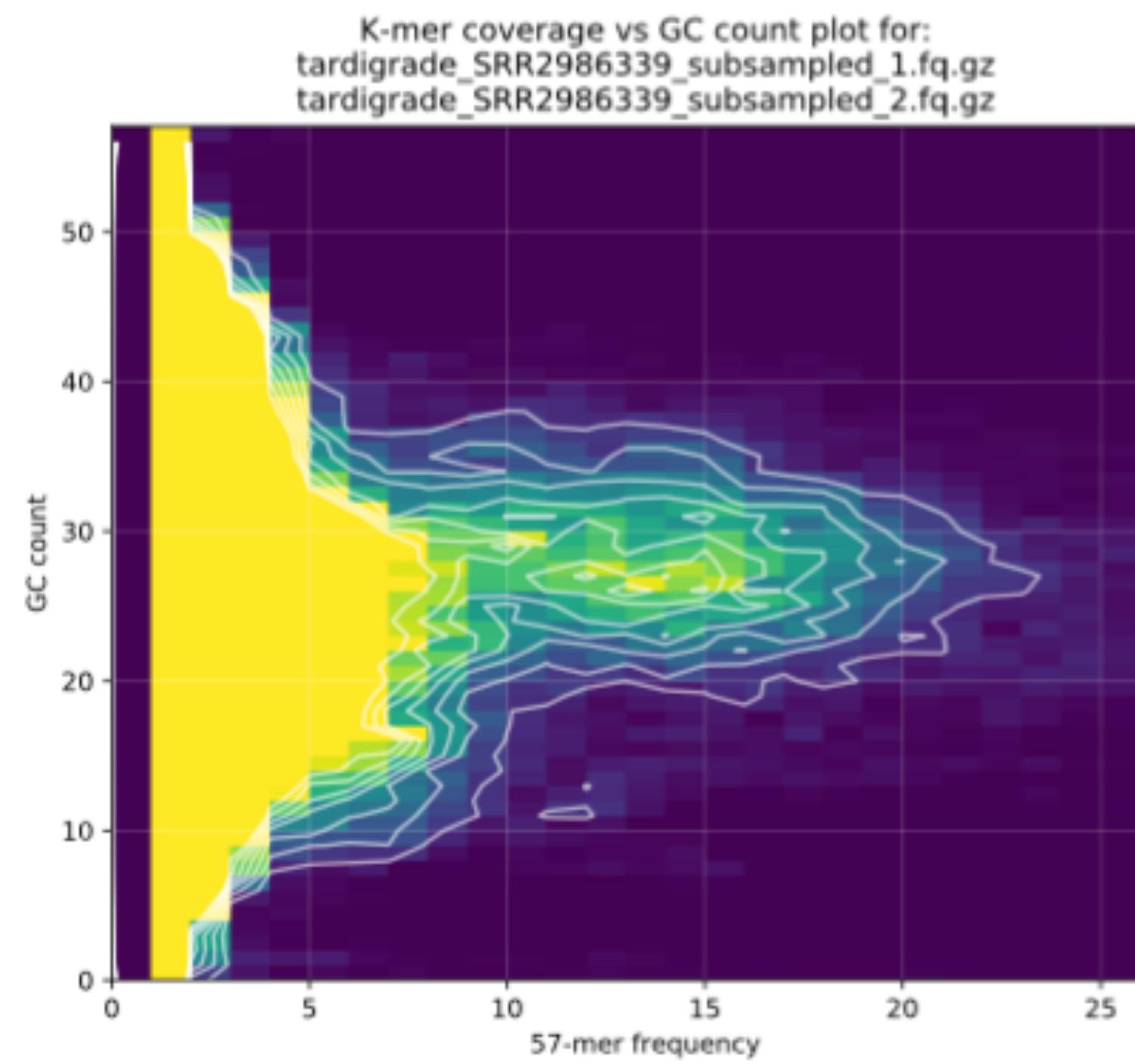


insH-3 transposase gene only
assembled in hybrid assembly

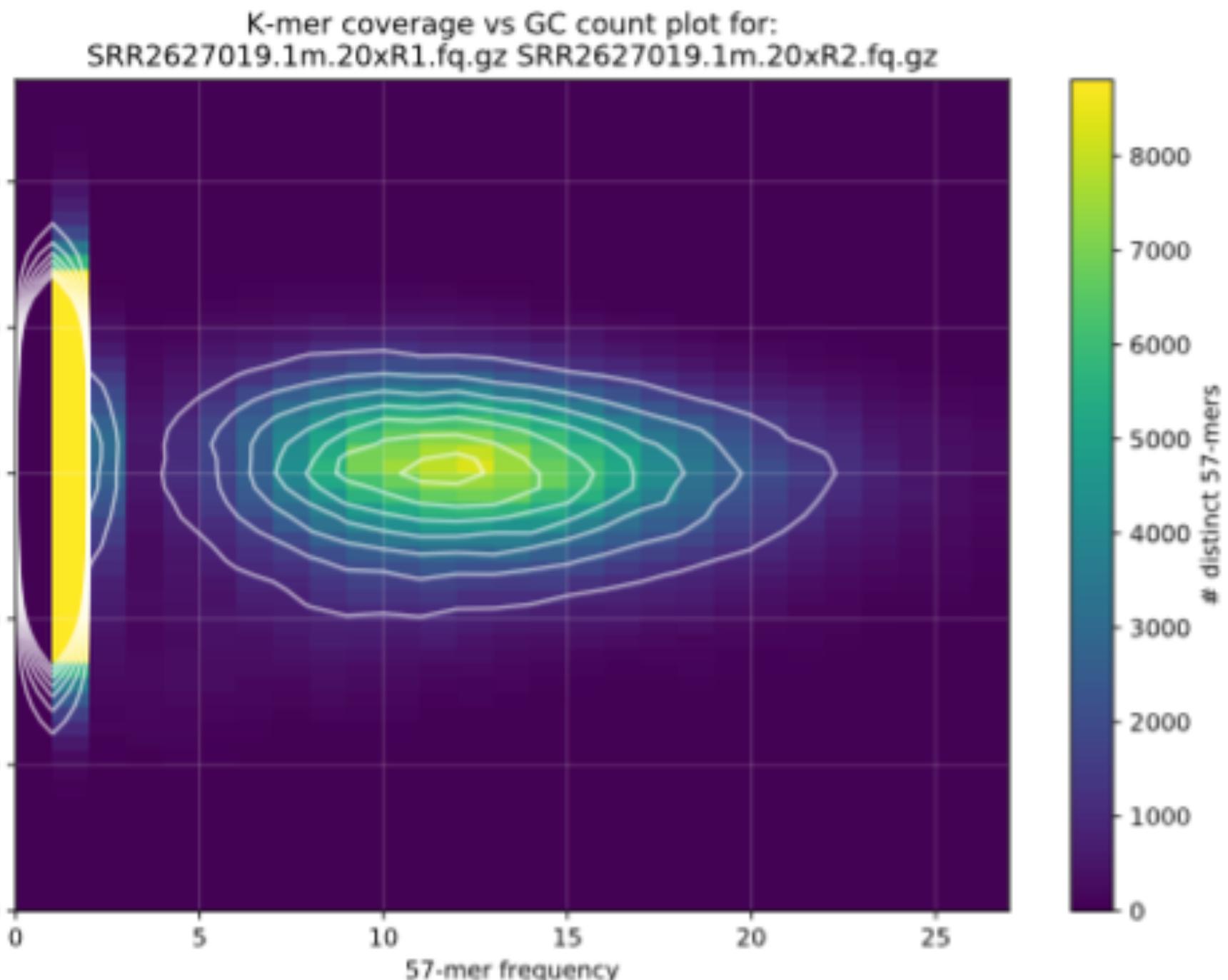
Tardigrade Input data QC

Contamination?

Tardigrade

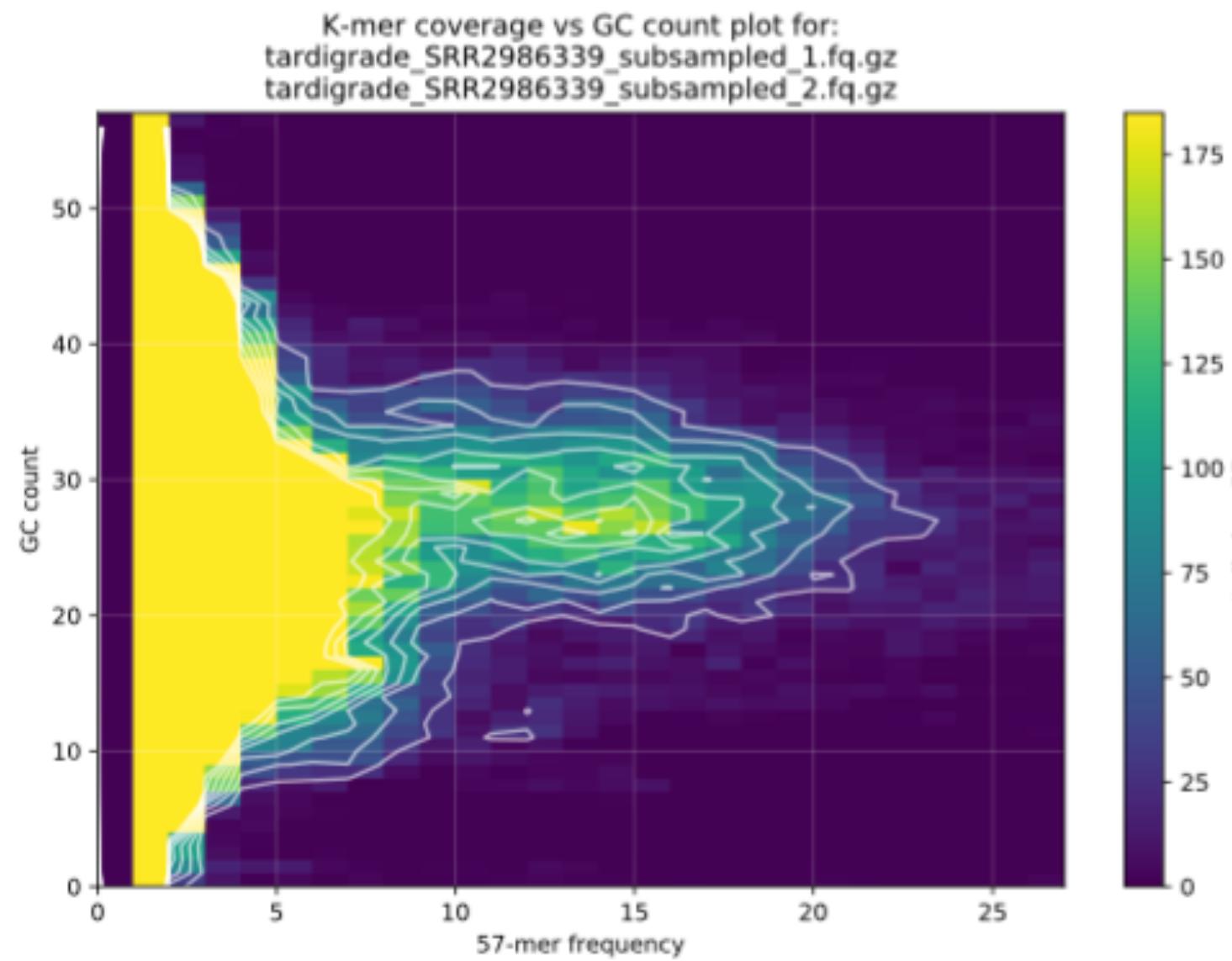


E. coli
20x



Tardigrade Input data QC

Contamination?



- 90 % unclassified - ok since not in database
- 6 % bacteria?
- 3.45 % human?
- no horizontal gene transfer...

```
[lehmanr@dbn503-33-r output]$ awk '$1 > 1' tardigrade_kraken2_report.txt
```

90.19	394434	394434	U	0	unclassified
9.81	42880	1066	R	1	root
9.55	41742	209	R1	131567	cellular organisms
6.01	26299	1182	D	2	Bacteria
4.23	18485	1053	P	1224	Proteobacteria
2.06	8989	686	C	28211	Alphaproteobacteria
1.07	4675	426	O	356	Rhizobiales
1.57	6844	189	C	28216	Betaproteobacteria
1.46	6374	585	O	80840	Burkholderiales
3.45	15085	0	D	2759	Eukaryota
3.45	15085	0	D1	33154	Opisthokonta
3.45	15085	0	K	33208	Metazoa
3.45	15085	0	K1	6072	Eumetazoa
3.45	15085	0	K2	33213	Bilateria
3.45	15085	0	K3	33511	Deuterostomia
3.45	15085	0	P	7711	Chordata
3.45	15085	0	P1	89593	Craniata
3.45	15085	0	P2	7742	Vertebrata
3.45	15085	0	P3	7776	Gnathostomata
3.45	15085	0	P4	117570	Teleostomi
3.45	15085	0	P5	117571	Euteleostomi
3.45	15085	0	P6	8287	Sarcopterygii
3.45	15085	0	P7	1338369	Dipnotetrapodomorpha
3.45	15085	0	P8	32523	Tetrapoda
3.45	15085	0	P9	32524	Amniota
3.45	15085	0	C	40674	Mammalia
3.45	15085	0	C1	32525	Theria
3.45	15085	0	C2	9347	Eutheria
3.45	15085	0	C3	1437010	Boreoeutheria
3.45	15085	0	C4	314146	Euarchontoglires
3.45	15085	0	O	9443	Primates
3.45	15085	0	O1	376913	Haplorrhini
3.45	15085	0	O2	314293	Simiiformes
3.45	15085	0	O3	9526	Catarrhini
3.45	15085	0	O4	314295	Hominidea
3.45	15085	0	F	9604	Hominidae
3.45	15085	0	F1	207598	Homininae
3.45	15085	0	G	9605	Homo
3.45	15085	15085	S	9606	Homo sapiens