

# Generating & Assessing FASTQC Outputs

Dr. Princess Rodriguez

2026-02-09

# Learning Objectives:

- Describe the contents and format of a FASTQ file
- Create a quality report using FASTQC
- Evaluate the quality of your NGS data using FastQC

# Recap: Advantages of Batch Job Submissions

Batch job submission on an HPC (High-Performance Computing) system offers several advantages, particularly for computationally intensive tasks like bioinformatics, genomics, and large-scale data analysis.

## ① Efficient Resource Management:

- Jobs are queued and scheduled based on resource availability, ensuring optimal utilization of CPUs, memory, and GPUs.
- Users can specify resource requirements (e.g., nodes, cores, memory) to avoid wasting computational power.

# Recap: Advantages of Batch Job Submissions

## ② Scalability:

- HPC clusters handle jobs of varying sizes, from single-threaded processes to massively parallel workloads.
- Batch processing supports running multiple jobs concurrently, improving overall throughput.

# Recap: Advantages of Batch Job Submissions

## ③ Parallel Execution:

- Batch submission allows running thousands of jobs in parallel (e.g., processing multiple sequencing samples).

## ④ Job Monitoring:

- Provides insights into job status, resource usage, and debugging.

# Looking inside of `sra_fqdump.sh`

Purpose: Is to download FASTQ files from the SRA. FASTQ files to be downloaded from the GEO are listed in a text file as accession numbers (provided by user).

```
nano sra_fqdump.sh
```

```
#!/bin/bash
#SBATCH --partition=bluemoon
#SBATCH --nodes=1
#SBATCH --ntasks=2
#SBATCH --mem=50G
#SBATCH --time=30:00:00
#SBATCH --job-name=fastq
# %x=job-name %j=jobid
#SBATCH --output=%x_%j.out

#while there are lines in the list of SRRs file
while read p
do
#call the bash script that does the fastq dump, passing it the SRR number next $
sbatch inner_script.sh $p
done <list_of_SRRs.txt
```

## To submit a script, use the command:

```
sbatch your-script.sh
```

When you submit your job, Slurm will respond with the job ID. For example, where the job ID Slurm assigns is “123456,” Slurm will respond:

```
Submitted batch job 123456
```

## After submitting this script you will see .out files:

```
[pdrodrig@vacc-user1 GSE164713_Tcf1]$ ls
fastq_6008331.out  fastq_6426350.out      SRR13422709.fastq.gz
fastq_6366635.out  fastq_6426351.out      SRR13422710.fastq.gz
fastq_6366636.out  inner_script.sh        SRR13422711.fastq.gz
fastq_6366637.out  list_of_SRRs.txt       SRR13422712.fastq.gz
fastq_6366638.out  list.txt                SRR13422713.fastq.gz
fastq_6366639.out  sra_download.sh        SRR13423162.fastq.gz
fastq_6366640.out  sra_fqdump.sh         SRR13423163.fastq.gz
fastq_6366641.out  SRR13416485.fastq.gz  SRR13423164.fastq.gz
fastq_6426340.out  SRR13416486.fastq.gz  SRR13423165.fastq.gz
fastq_6426341.out  SRR13422702.fastq.gz  SRR13423166.fastq.gz
fastq_6426342.out  SRR13422703.fastq.gz  SRR13423167.fastq.gz
fastq_6426343.out  SRR13422704.fastq.gz  SRR17379677.fastq.gz
fastq_6426344.out  SRR13422705.fastq.gz  SRR17379678.fastq.gz
fastq_6426347.out  SRR13422706.fastq.gz  SRR17379679.fastq.gz
fastq_6426348.out  SRR13422707.fastq.gz  SRR17379680.fastq.gz
fastq_6426349.out  SRR13422708.fastq.gz
```

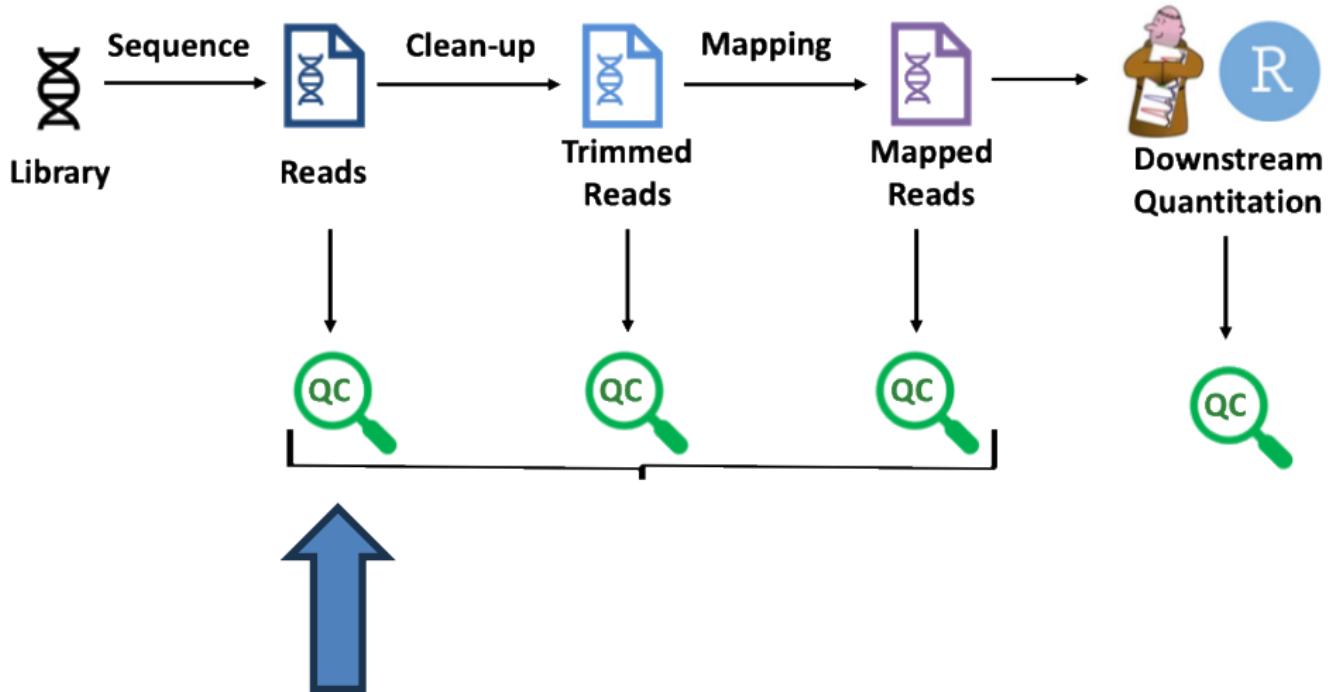
# STDOUT records the output of programs

Three data streams exist for all Linux programs:

- STDIN (Standard Input - a way to send data into the program)
- STDOUT (Standard Output - a way to send expected data out of the program)
- STDERR (Standard Error - a way to send errors or warnings out of the program)

# Next Step in Processing Sequencing Data

# Processing Sequencing Data



# What is the Point of QC?

Technical Problems ...

- Don't always cause pipelines to fail
- Don't prevent hits being generated
- These hits can look biologically real

Real biology...

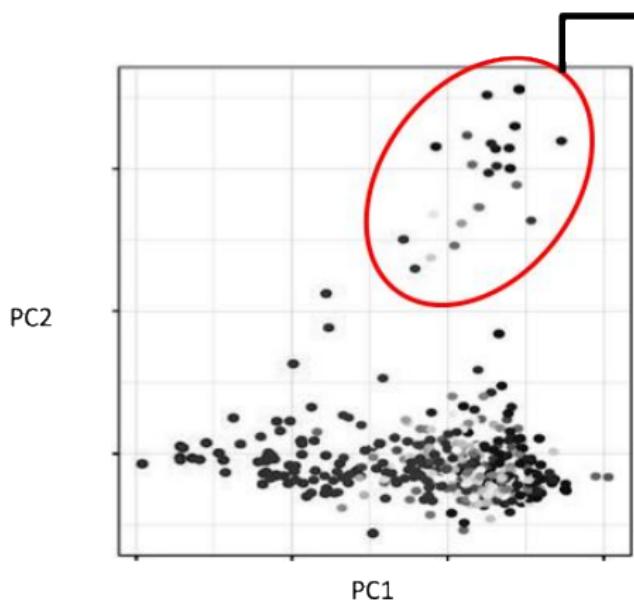
- Can cause unexpected, interesting behaviour of data

Set Pipelines can miss things....

QC Saves Time, Effort and Money!!

- Better to know asap what you're dealing with
- Want to be sure any follow-up work will be worth it

# What is the Point of QC? An Example...



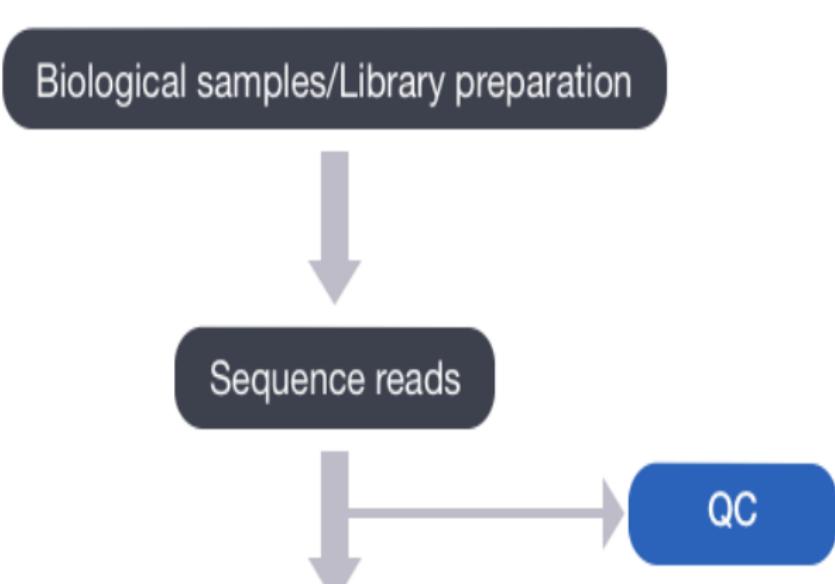
**PC2 Genes (85 total)**

- No clear biological theme
- No clear connection to system

**What is going on?**

# Quality Control of FASTQ files

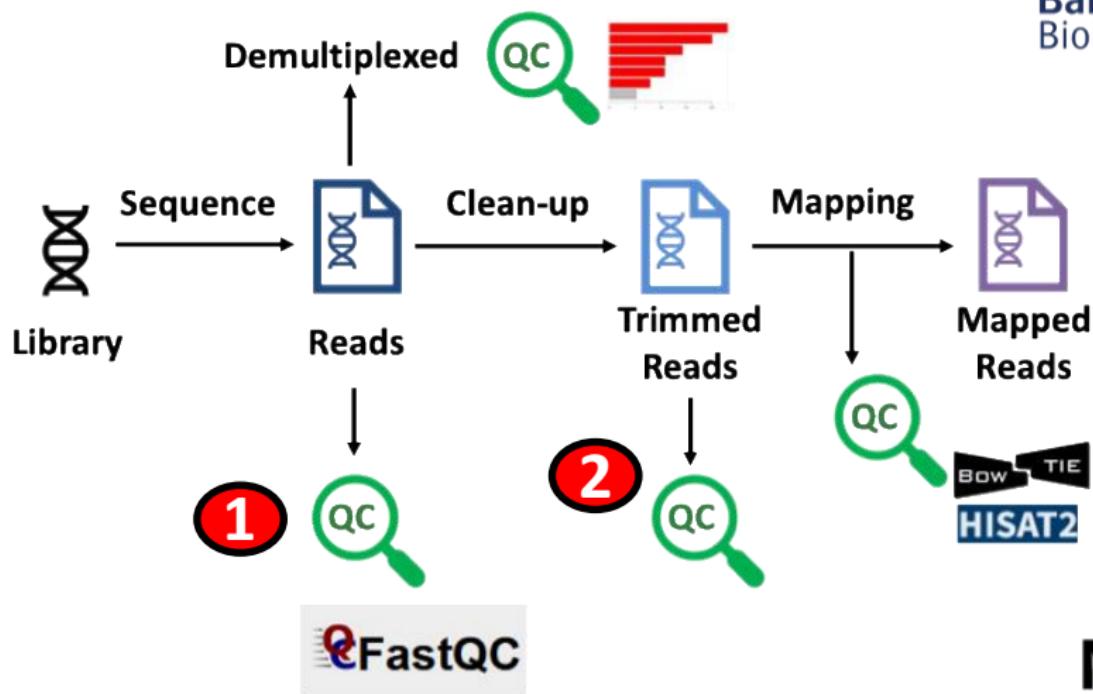
The first step in the bioinformatic pipeline is to assess the quality of the sequence reads retrieved from the sequencing facility.



## In a Genomics Workflow, QC helps us answer:

- 1) Did sequencing work as expected?
- 2) Are there systematic technical artifacts?
- 3) Do we need trimming, filtering, or adapter removal?
- 4) Are problems biological or technical?

# QC Programs in Data Processing



# FastQC

- FastQC provides a simple way to do some quality checks on raw sequence data coming from high throughput sequencing pipelines.
- It provides a modular set of analyses, which you can use to obtain an impression of whether your data has any problems that you should be aware of before moving on to the next analysis.

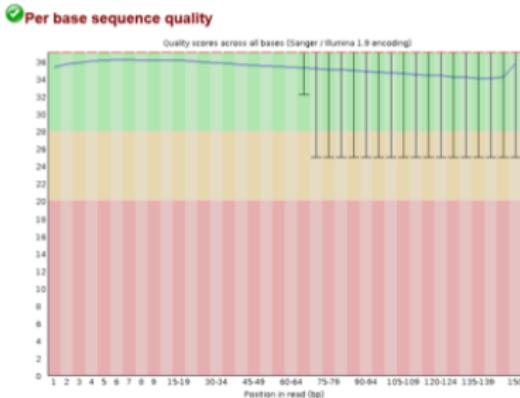
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

# FastQC

## FastQC Report

### Summary

- ✓ Basic Statistics
- ✓ Per base sequence quality
- ⚠ Per file sequence quality
- ✓ Per sequence quality scores
- ✗ Per base sequence content
- ⚠ Per sequence GC content
- ✓ Per base N content
- ⚠ Sequence Length Distribution
- ✓ Sequence Duplication Levels
- ✓ Overrepresented sequences
- ✓ Adapter Content



**fastqc seqfile1 seqfileN**

**fastqc \*.fastq.gz**

- Reads raw fastq file(s)
- Performs multiple checks
  - Pass/warn/fail
  - Compares to genomic library
- Generates a HTML Report

<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>



# Interpreting the HTML report

- Within the report, a summary of all of the modules is given on the left-hand side.
- Do not take the yellow “**WARNING**”s and red “**FAIL**”s as ***this sample is not usable***; they should be interpreted as flags!



# Context is Key for QC

Expectations

Observations

Critical Analysis

Understanding & Confidence



QC should be about what you expect and what you see



# Input: Raw FASTQ files

Similar to FASTA, the FASTQ file begins with a header line. The difference is that the FASTQ header is denoted by a @ character.

Sequence 1

```
@HWUSI-EAS611:34:6669YAXX:1:1:5069:1159 1:N:0:  
TCGATAATACCGTTTTTCCGTTGATGTTGATACCATT  
+  
IIHIIHIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
```

Sequence 2

```
@HWUSI-EAS611:34:6669YAXX:1:1:5243:1158 1:N:0:  
TATCTGTAGATTCACAGACTCAAATGTAATATGCAGAG  
+  
DF=DBD<BBFGGGGGGBD@GGD4@CA3CGG>DDD:D,B
```

Sequence 3

```
@HWUSI-EAS611:34:6669YAXX:1:1:5266:1162 1:N:0:  
GGAGGAAGTATCACTCCTGCCTGCCTCCTGGGGCCT  
+  
:GBGGGGGGGGGDGGDEDGGDGGGGDHHDHGHGBGG:GG
```

# Illumina Header Sections (Line 1)

@HWUSI-EAS611:34:6669YAXXX:5:1:5069:1159 1:N:0:

- Starts with @ (required by fastq spec)
- Instrument ID (HWUSI-EAS611)
- Run number (34)
- Flowcell ID (6669YAXXX)
- Lane (5)
- Tile (1)
- X-position (5069)
- Y-position (1159)
- [space]
- Read number (1)
- Was filtered (Y/N) (N) - You wouldn't normally see the Ys
- Control number (0 = no control)
- Sample number (only if demultiplexed using Illumina's software)

# A Single FASTQ Entry

1. @HWUSI-EAS611:34:6669YAXX:1:1:5266:1162 1:N:0:  
2. GGAGGAAGTATCACTCCTGCCTGCCTCCTGGGCCT  
3. +  
4. :GBGGGGGGGGDGGDEDGGDGDDHHDGHGBGG:GG

1. Header - starts with @
2. Base calls (can include N or IUPAC codes)
3. Mid-line - starts with + usually empty
4. Quality scores (= Phred Scores)

# Phred Scores (Line 4)

Base Calls	GGAGGAAAGTATCACTTCCTGCCTGCCTCCTCTGGGCCT
Phred Scores	 :GBGGGGGGGGDGGDEDGGDGHHGDHHGHGBGG:GG

- Each quality score represents the probability that the corresponding nucleotide call is incorrect.
  - \_ This quality score is logarithmically based and is calculated as:

$$Q = -10 \times \log_{10} (P)$$

- P is the probability that a base call is erroneous

- $\text{Phred} = -10 * (\text{int})\log_{10}(p)$

- $p=0.1 \quad \text{Phred} = 10$

- $p=0.01 \quad \text{Phred} = 20$

- $p=0.001 \quad \text{Phred} = 30$

**Higher Phred Score  
Higher Confidence**

# Phred Score Encoding

- Translation of Phred score to single ASCII letter
- Based on standard ASCII table

Different quality encoding scales exist, but the most commonly one used is fastq-sanger, which is the scale output by Illumina since mid-2011.

0	NUL	17	C1	33	!	50	2	67	C
1	SOH	18	DC2	34	"	51	3	68	D
2	STX	19	DC3	35	#	52	4	69	E
3	ETX	20	DC4	36	\$	53	5	70	F
4	EOT	21	NAK	37	%	54	6	71	G
5	ENQ	22	SYN	38	&	55	7	72	H
6	ACK	23	ETB	39	'	56	8	73	I
7	BEL	24	CAN	40	(	57	9	74	J
8	BS	25	EM	41	)	58	:	75	K
9	HT	26	SUB	42	*	59	;	76	L
10	LF	27	ESC	43	+	60	<	77	M
11	VT	28	FS	44	,	61	=	78	N
12	FF	29	GS	45	-	62	>	79	O
13	CR	30	RS	46	.	63	?	80	P
14	SO	31	US	47	/	64	@	81	Q
15	SI	32	{SPACE}	48	0	65	A	82	R
16	DLE			49	1	66	B	83	S

# QC metrics we can work with:

## Phred Scores



### How the Sequencer Performed

- At different cycles
- Across different locations
- For different reads

## Library Composition

Adapter      Insert      Adapter

### The Nature of our Sequenced Reads

- Biases
- Contaminants
- Duplication

## Mapping



### Where our Sequencing Reads Come From

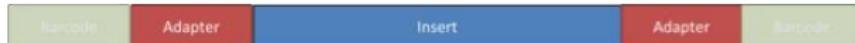
- Species (plural or singular!)
- Region (repetitive or unique)

# Universal QC Metrics

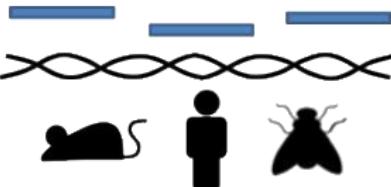
- Demultiplexing



- Base Call Quality



- Adapter Content

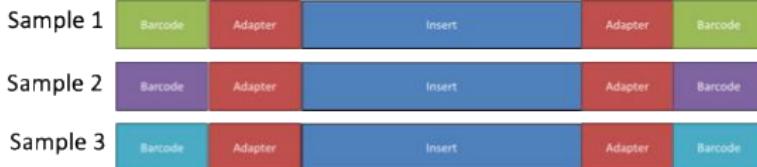


- Mapping Quality

# Demultiplexing: Expectation



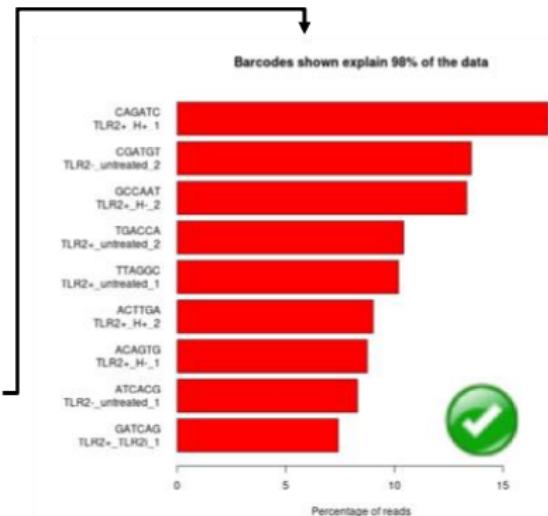
Only the barcodes we assigned to samples should be present—**no others**



etc...

What barcode sequences have we found?

█ Expected Barcode    █ Unknown Sequence



What could unknown barcode sequences mean?

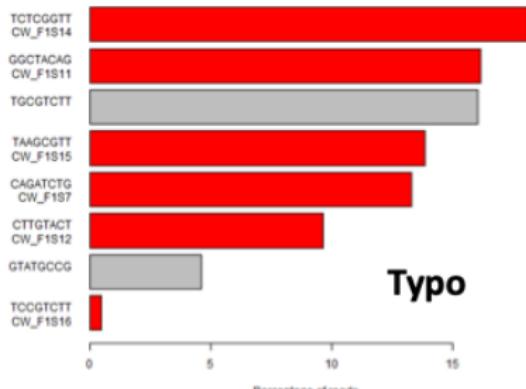


# Demultiplexing: Unknown

■ Expected Barcode

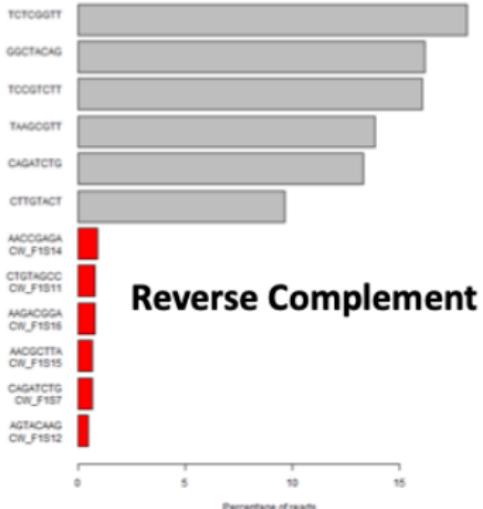
■ Unknown Sequence

Barcodes shown explain 92% of the data



Typo

Barcodes shown explain 91% of the data



Reverse Complement



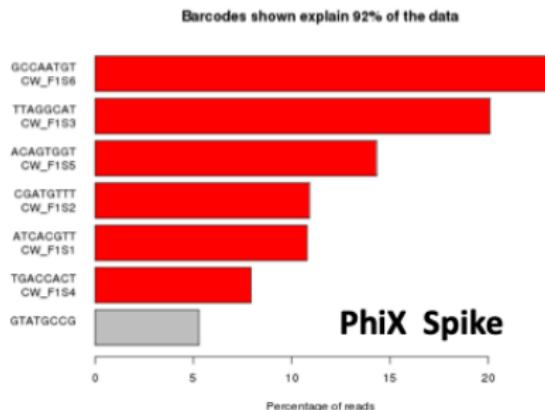
Human Error is a really common source of barcode issues



# Demultiplexing: Unknown

■ Expected Barcode

■ Unknown Sequence



PhiX is a well characterized bacteriophage that is used as a control in Illumina sequencing runs.

PhiX is typically added for low-complexity libraries. To add sequence diversity.

Improves flow cell clustering and reduces sequencing errors.



Sometimes more technical in nature...

# Basic Statistics



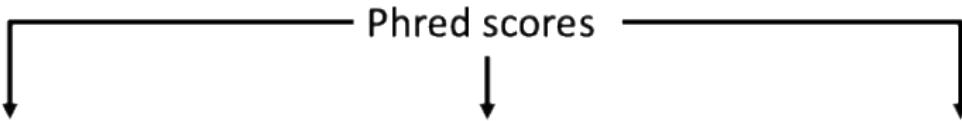
## Basic Statistics

Measure	Value
Filename	Mov10_oe_1_subset.fq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	305900
Sequences flagged as poor quality	0
Sequence length	100
%GC	47

# Base Call Quality: Expectations



Illumina Sequencers are technically reliable, so we expect confident calls



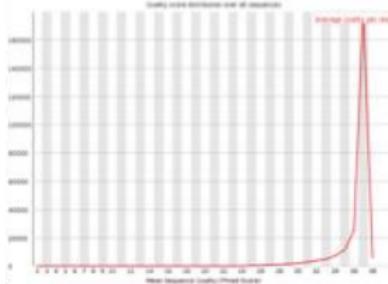
Per base/cycle



FastQC Report

Per base sequence quality

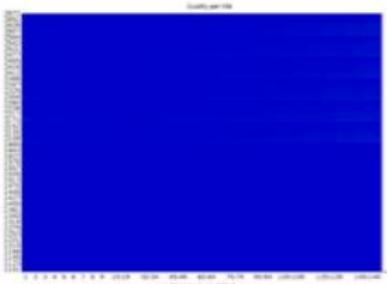
Per Read



FastQC Report

Per sequence quality scores

Location



FastQC Report

Per tile sequence quality



# Per base sequence quality

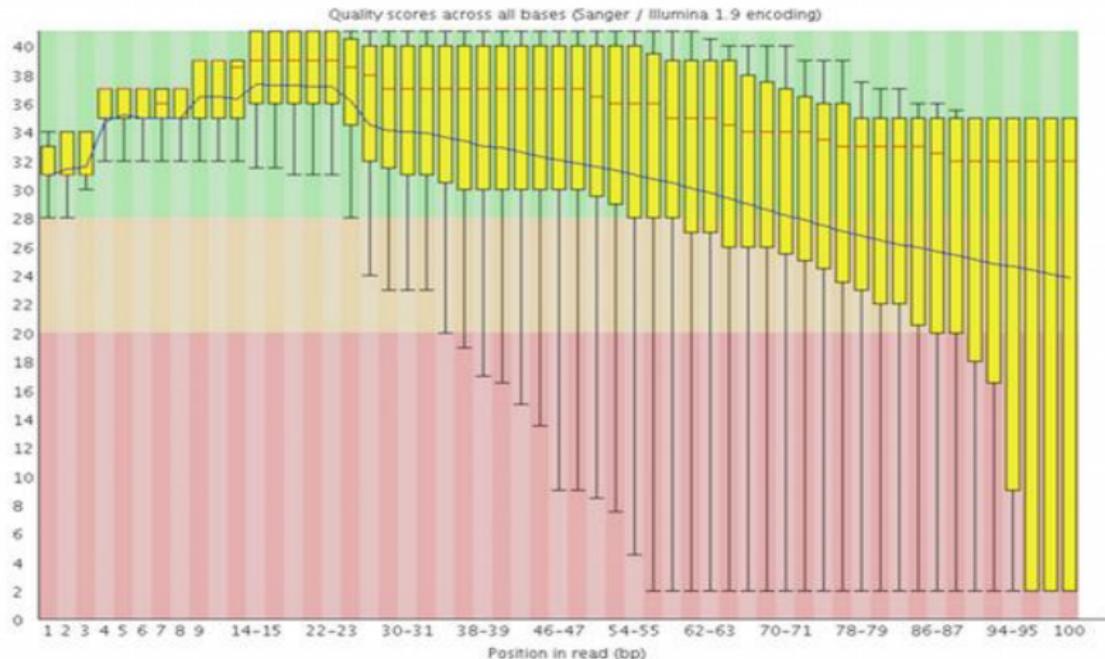


*Distribution of Phred score at each position across all reads*



# Per base sequence quality

## ✖ Per base sequence quality

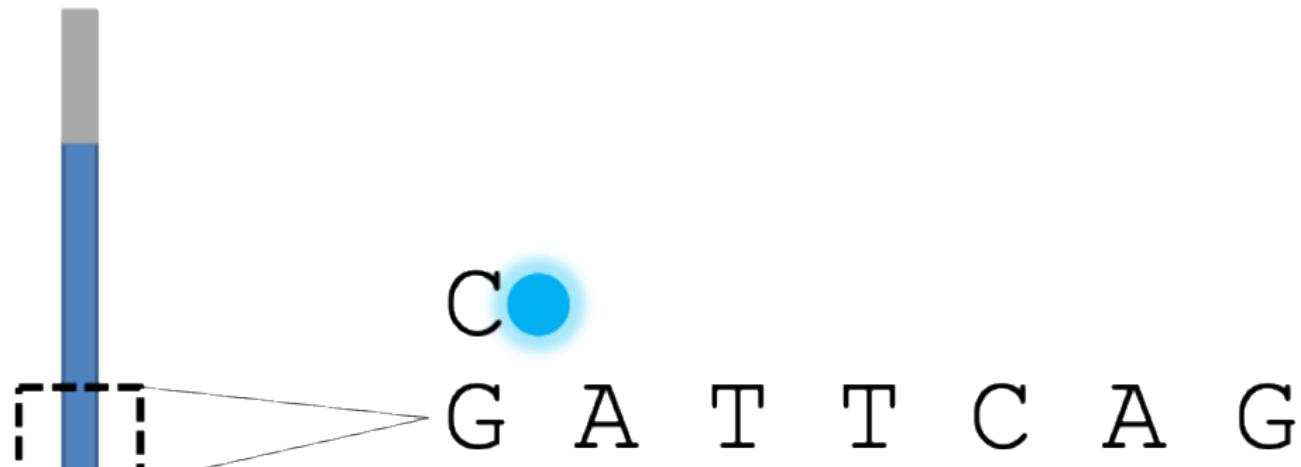


Clusters get out of sync over length of the read

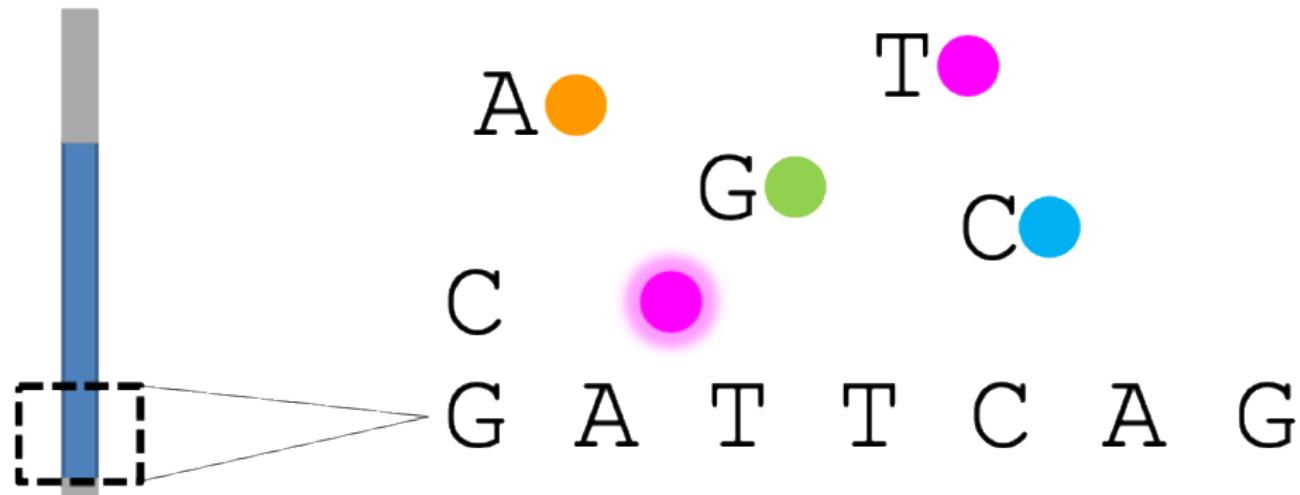
# Sequence by Synthesis



# Sequence by Synthesis



# Sequence by Synthesis



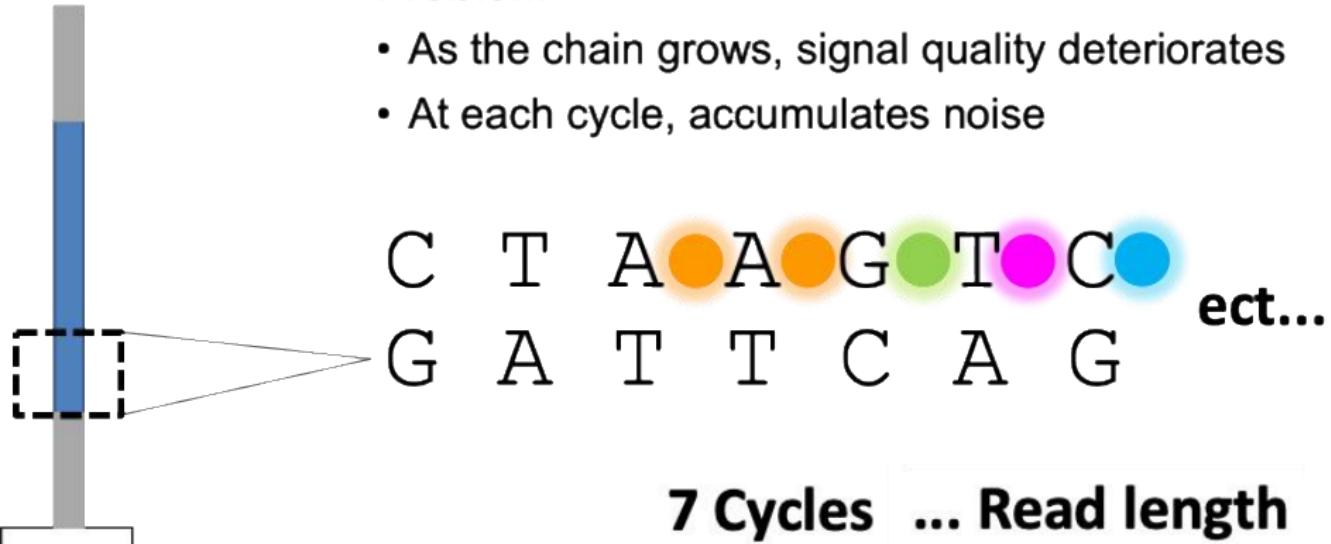
**1 Cycle**

A photo is taken....

# Sequence by Synthesis

## Problem

- As the chain grows, signal quality deteriorates
- At each cycle, accumulates noise



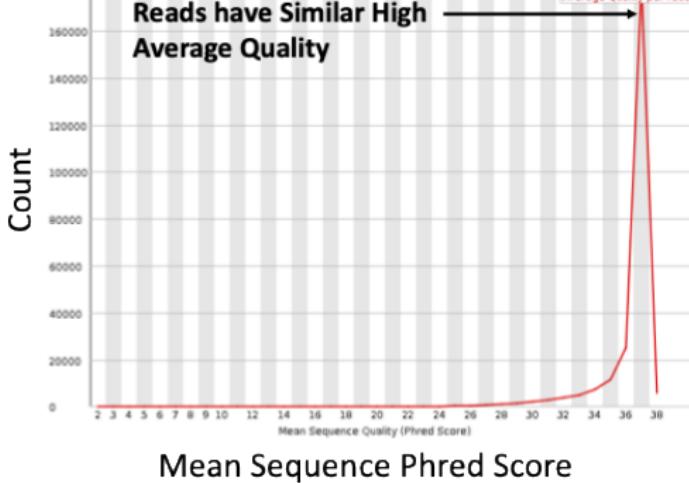
# Per sequence quality scores



Quality score distribution over all sequences

Reads have Similar High  
Average Quality

Average Quality per read



What to expect?

- A single peak centered at high quality (Q30+)

# Per sequence quality scores



Quality score distribution over all sequences

Reads have Similar High Average Quality

Average Quality per read

Count



Mean Sequence Phred Score



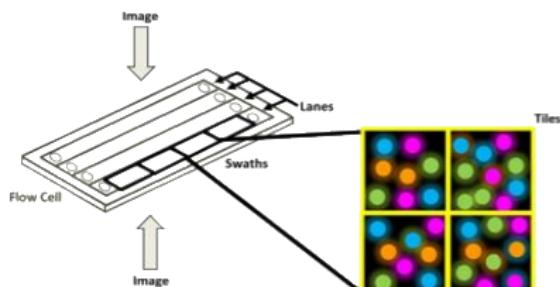
Quality score distribution over all sequences

Poorer Quality Read Subsets

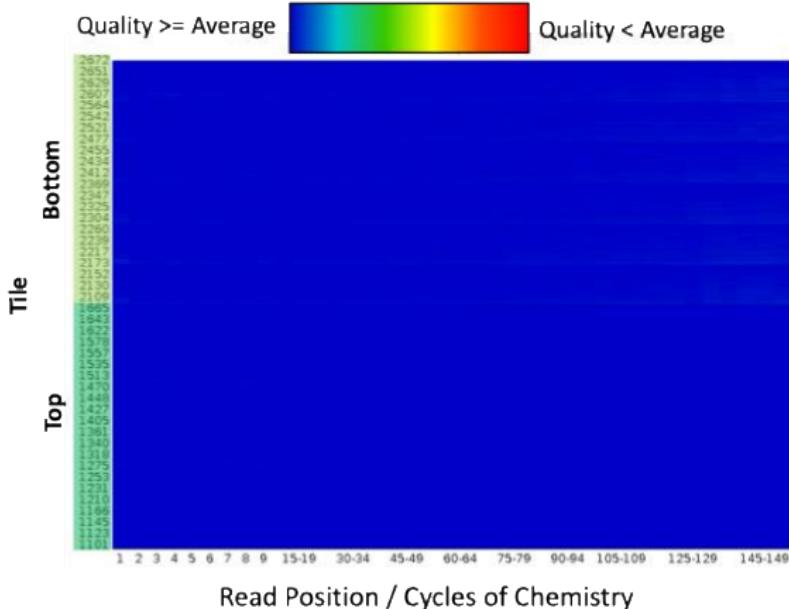
Average Quality per read



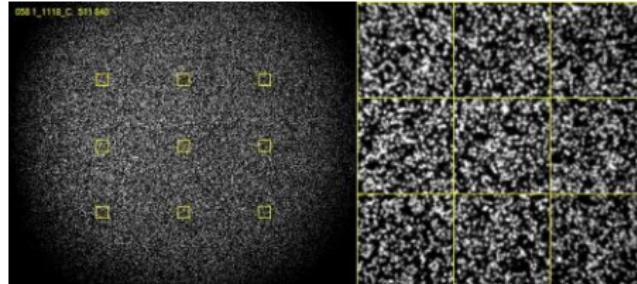
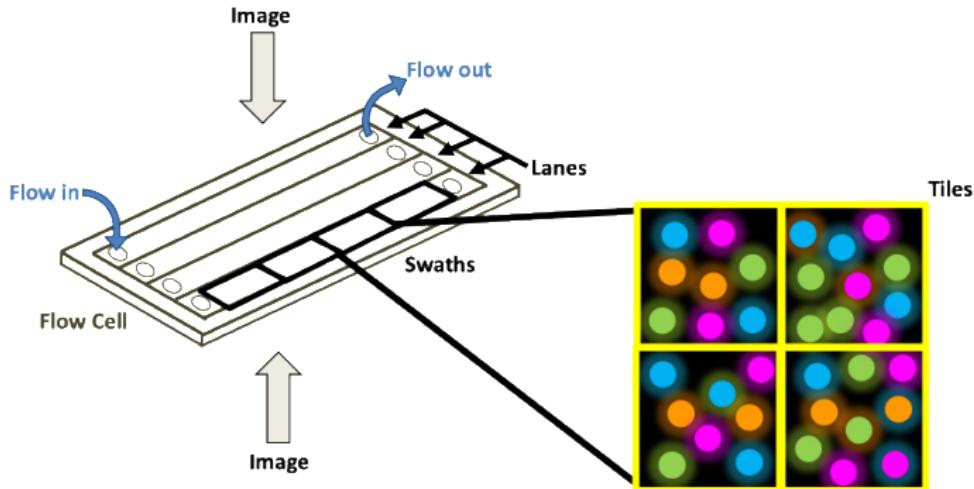
# Positional Quality



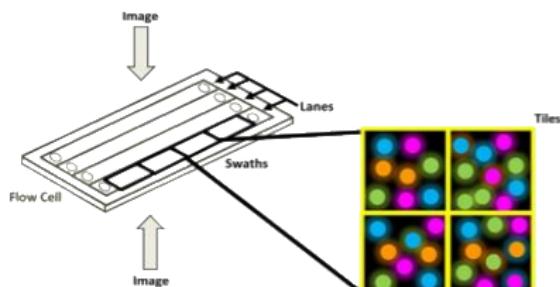
A good tile plot is a blue square



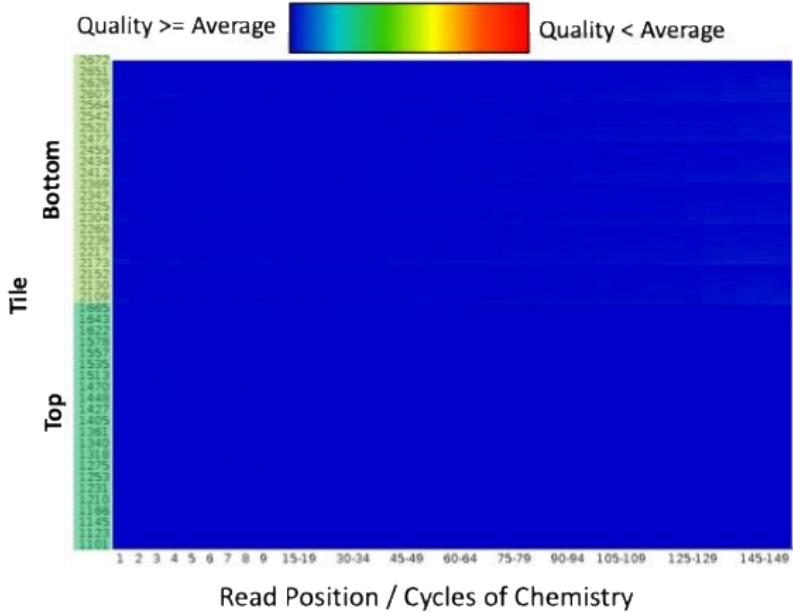
# Flow Cell Imaging



# Positional Quality

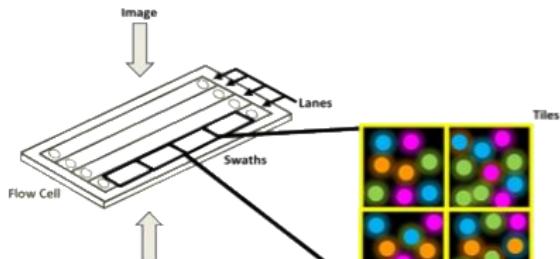


A good tile plot is a blue square

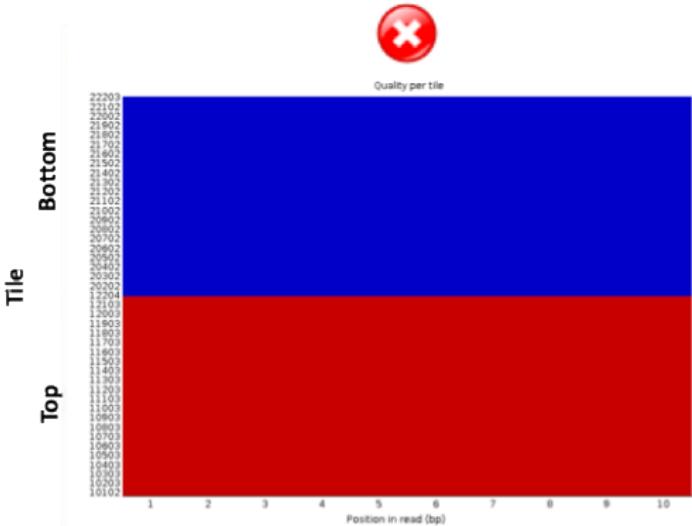


Read Position / Cycles of Chemistry

# Positional Quality



A good tile plot is a blue square



Focusing Fail!

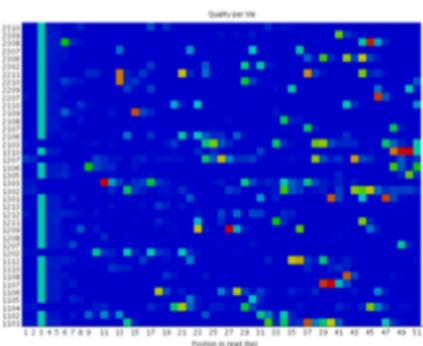
# Positional Quality



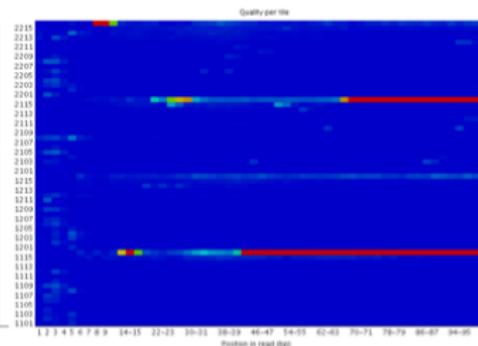
## More Examples of Positional Fails

Position Specific Patterning

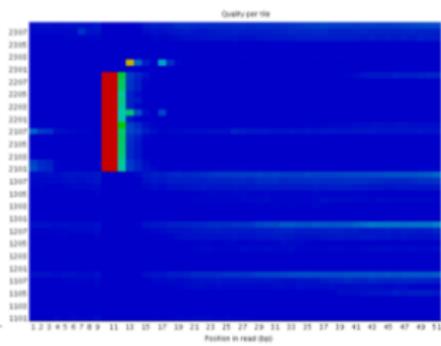
Random Pattern



Permanent



Transient



Overloading of Flow Cell

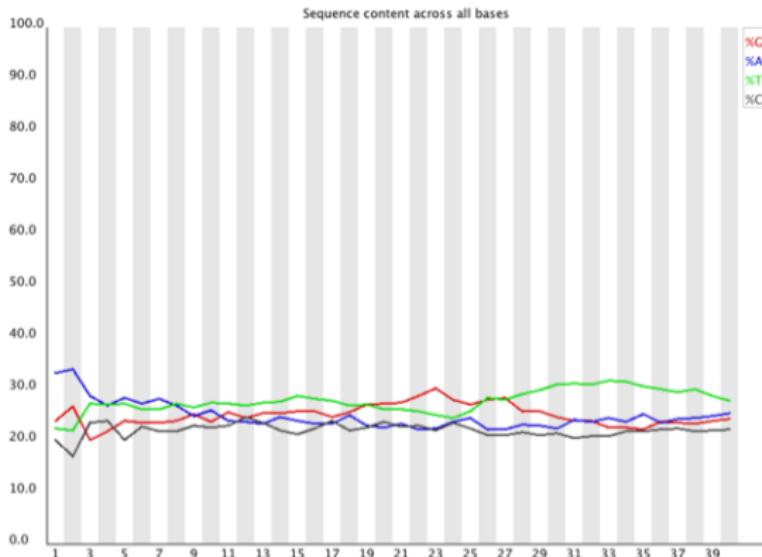


Obstruction

Bubble

# Library Base Composition

Shows proportion of A, T, G, C's at each position



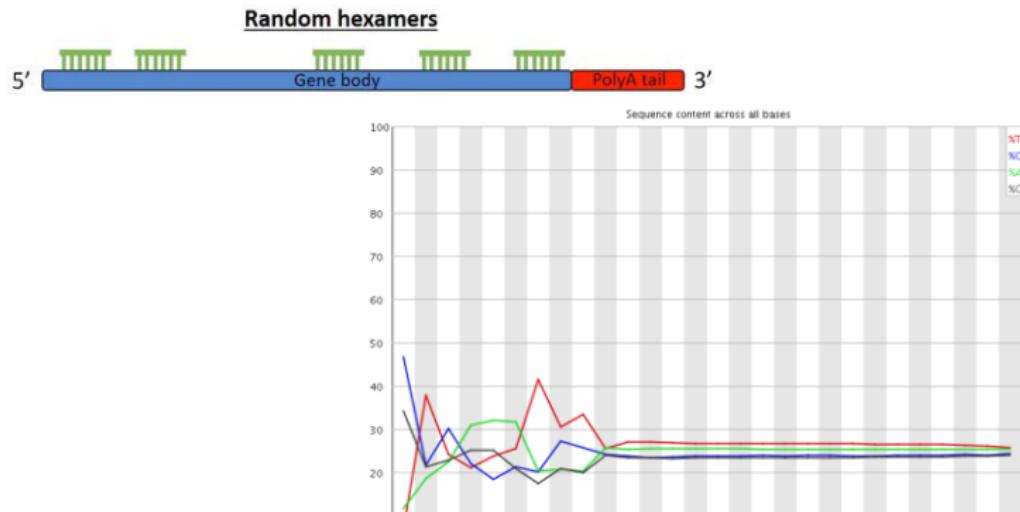
- For every chemistry cycle we can look at the number of ATGC
- The composition should be the same for all cycles



# Library Base Composition

Always gives a FAIL for RNA-seq data. This is because the first 10-12 bases result from the ‘random’ hexamer priming that occurs during RNA-seq library preparation.

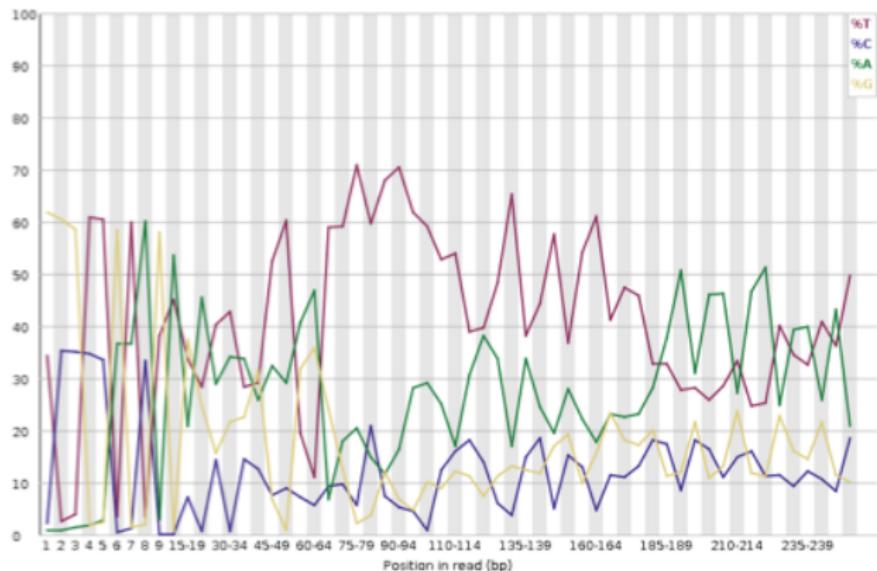
This priming is not as random as we might hope giving an enrichment in particular bases for these initial nucleotides.





# Library Base Composition: Bias

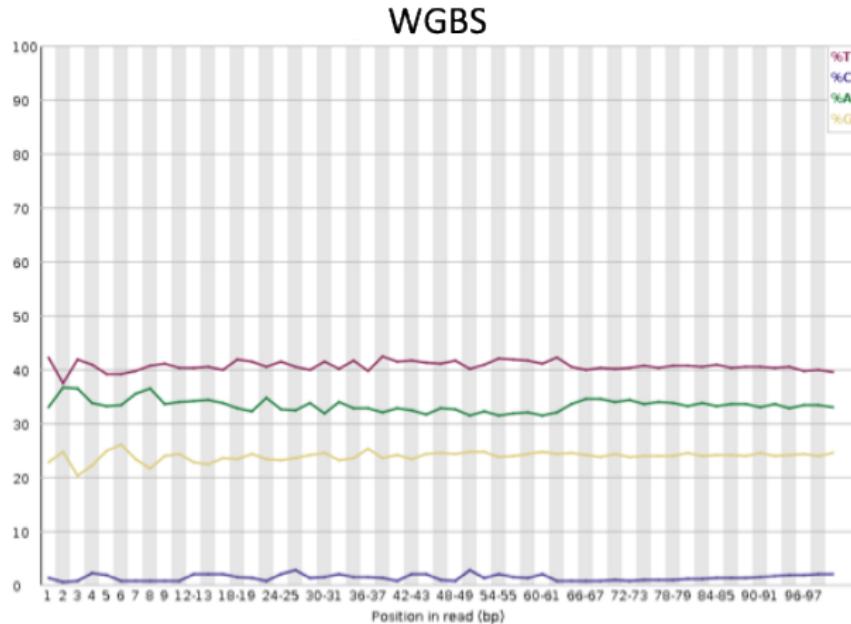
Amplicon



Very low diversity



# Library Base Composition: Bias



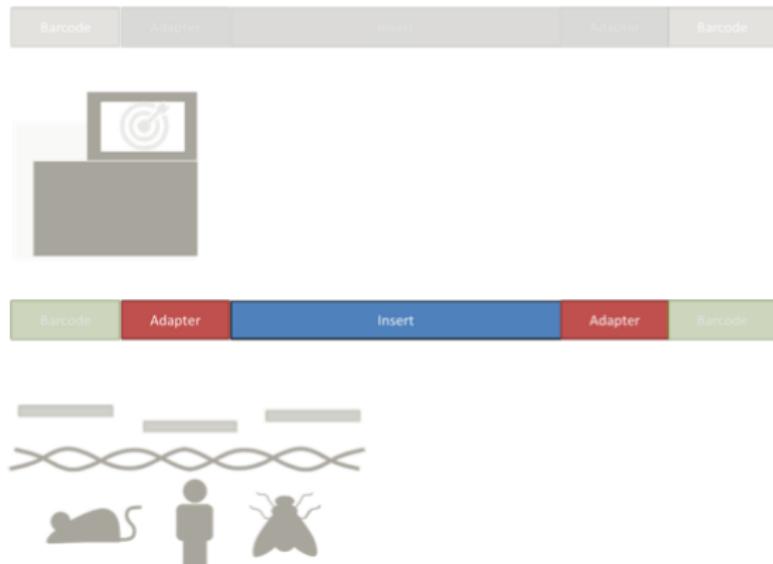
Bisulphite  
treated – C  
is converted  
to T

Consistent disproportional expression of bases



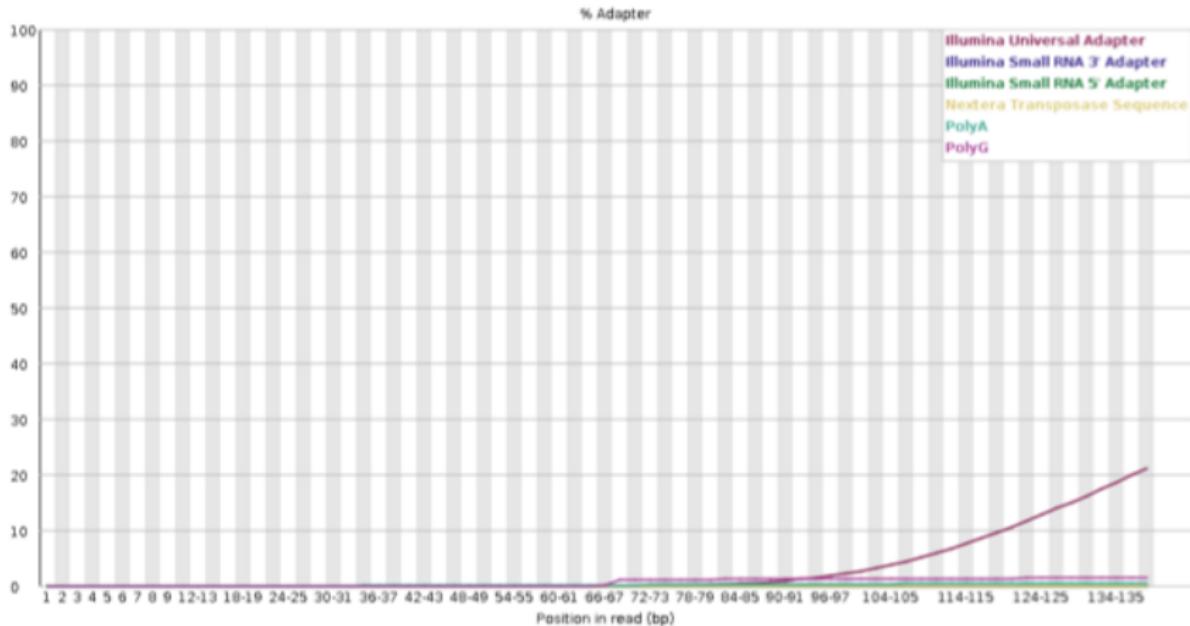
# Universal QC Metrics

- Demultiplexing
- Base Call Quality
- Adapter Content
- Mapping Quality





# Measuring Read-through Adapters



# Adapter Content

*Read length = 200bp*



*DNA insert = 150 bp*

# Library Dependent QC Metrics

From the Base Sequence:

- GC Content



- Base Composition

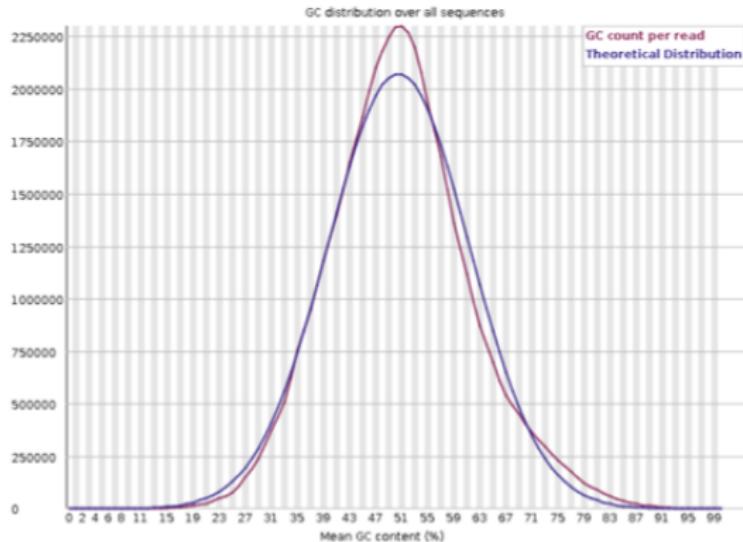
GATC

- Duplication

GATCTACGAGTTACGATCAGT  
GATCTACGAGTTACGATCAGT  
GATCTACGAGTTACGATCAGT  
GATCTACGAGTTACGATCAGT



# Library GC Content

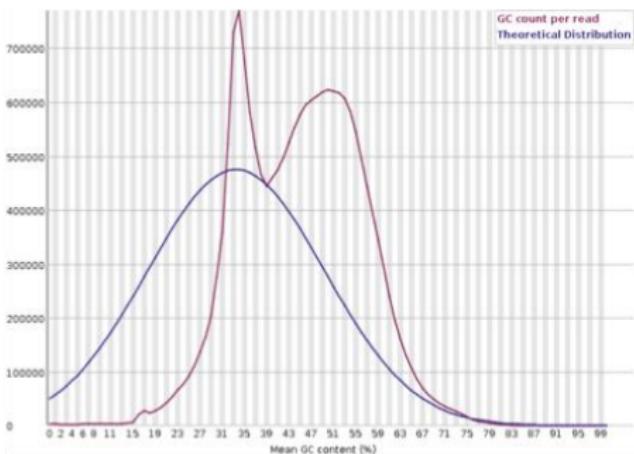


- Generic summary of library composition at a read level
- Expect a normally distributed set of values centred on the overall GC content

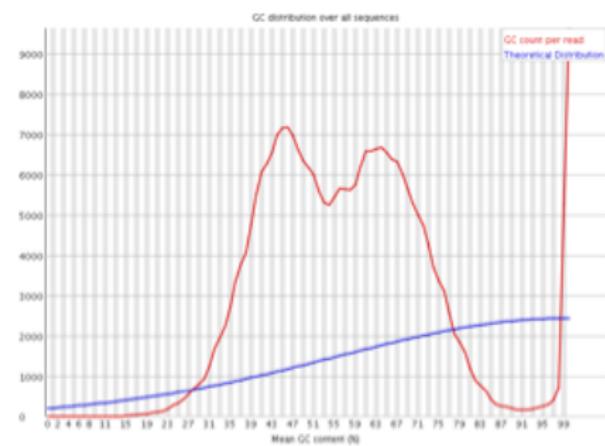


# Per sequence GC content

This plot would indicate some type of over-represented sequence with the sharp peaks, indicating either contamination or a highly over-expressed gene.



*Single contamination*



*Broad contamination*

# Duplication

If the exact same sequence appears more than once it could be...

Technical:

ATCCGAGCTATTGGCGAGCTGCCAGTTACG

ATCCGAGCTATTGGCGAGCTGCCAGTTACG

ATCCGAGCTATTGGCGAGCTGCCAGTTACG

Coincidental:

ATCCGAGCTATTGGCGAGCTGCCAGTTACG

ATCCGAGCTATTGGCGAGCTGCCAGTTACG

ATCCGAGCTATTGGCGAGCTGCCAGTTACG

- PCR duplicates

- Deep sequencing
- Highly present sequences
- Restricted diversity libraries



# Overrepresented Sequences

- Extreme duplication
- Displays the sequences (at least 20 bp) that occur in more than 0.1% of the total number of sequences.
- The exact same sequence is a significant proportion of the whole library (which might not be duplicated overall)

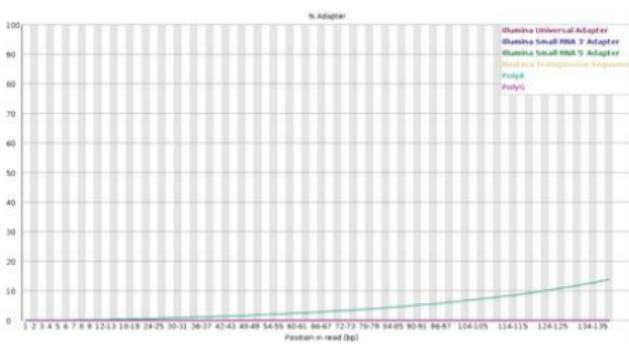
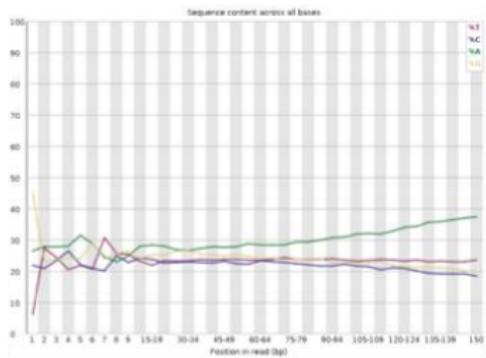
# Low complexity sequences

- Poly A
- Poly G
- Poly N

## Overrepresented Sequences: Poly A



## PolyA (or PolyT) – Common in RNA-Seq



# Overrepresented Sequences: Poly N



## PolyN – Quality too poor to make any calls

Sequence	Count	Percentage
TCCCCCCCCCCCCCCCCCCCCCCCCCCCC	462344	1.070097045533307
GNNNNNNNNNNNNNNNNNNNNNNNNNN	232540	0.5382147642627897
ANNNNNNNNNNNNNNNNNNNNNNNNN	127291	0.29461553090984244
CNNNNNNNNNNNNNNNNNNNNNNNN	87792	0.20319493671694688
TTNNNNNNNNNNNNNNNNNNNNNNNN	85181	0.19715176672688003
GANNNNNNNNNNNNNNNNNNNNNNNN	48918	0.11322090753507845

## Low complexity sequences

- Poly A
- Poly G
- Poly N

Why are these a problem?

- Poor alignment
- N's will break alignment step

*Indication that read should be trimmed or discarded*

# Overrepresented Sequences: Adapter Dimers

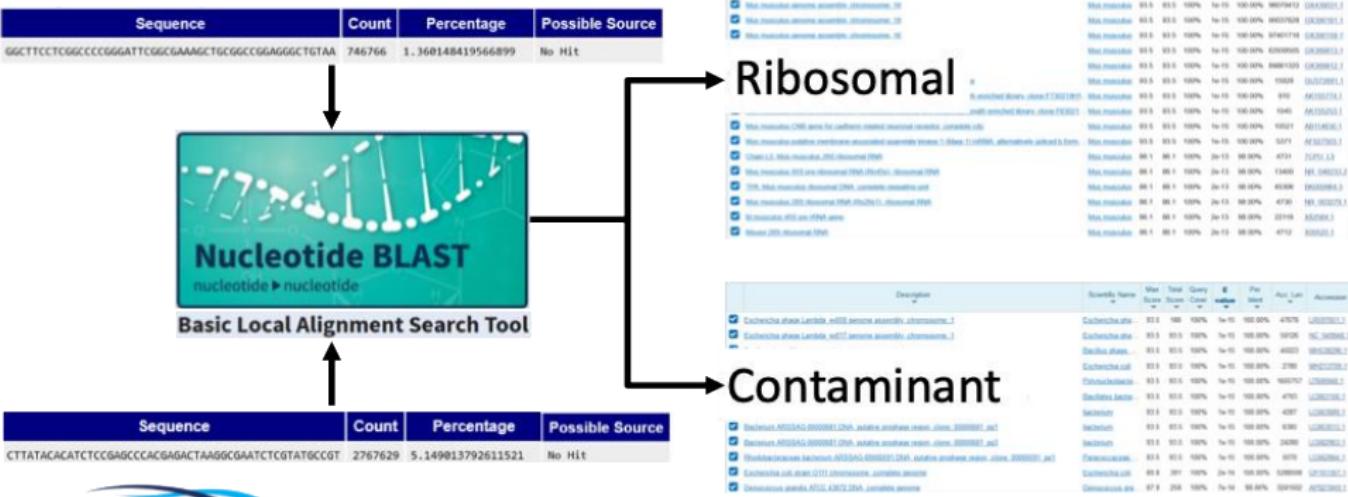


## ! Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GATCGGAAGAGCACACGTCTGAACCCAGTCACCTTGAAATCTCGTATGC	17957	0.14359551756800035	TruSeq Adapter, Index 12 (100% over 50bp)

# Overrepresented Sequences: Specific Sequences

- Other potential sources...



# **FASTQC is a decision tool**

**FASTQC helps us decide:**

- Do we trim? (quality drop, adapters present)
- How aggressively do we trim?
- Do we filter reads entirely
- Is the dataset salvageable

**Real sequencing data is almost never perfect!**

# Run fastqc with the Environment Module System:

```
module load gcc/13.3.0-xp3epy7  
module load fastqc/0.12.1-qxseug5
```

# Class Exercise #1

- Run FASTQC on all FASTQ files in `raw_fastq`. FASTQC allows you to redirect your output into a specified location with the `-o` parameter. Be sure to use this parameter in your final code.

If successful, you will see the following outputs inside of the `fastqc` folder:

```
Irrel_kd_1.subset_fastqc.html Irrel_kd_3.subset_fastqc.h  
Irrel_kd_1.subset_fastqc.zip Irrel_kd_3.subset_fastqc.z  
Irrel_kd_2.subset_fastqc.html Mov10_oe_1.subset_fastqc.h  
Irrel_kd_2.subset_fastqc.zip Mov10_oe_1.subset_fastqc.z
```

## How to use fastqc?

```
fastqc --help
```

### SYNOPSIS

```
fastqc seqfile1 seqfile2 .. seqfileN
```

```
fastqc [-o output dir]  
[--(no)extract]  
[-f fastq|bam|sam]  
[-c contaminant file]  
seqfile1 .. seqfileN
```

# FASTQC Outputs

For each individual FASTQ file that is input to FastQC, there are **two output files that are generated.**

- ① The first is **an HTML file** which is a self-contained document with various graphs embedded into it. Each of the graphs evaluate different quality aspects of our data, we will discuss in more detail in this lesson.
- ② Alongside the HTML file is **a zip file**. This file contains the different plots from the report as separate image files but also contains data files which are designed to be easily parsed to allow for a more detailed and automated evaluation of the raw data on which the QC report is built.

# Special Note

We are running FASTQC interactively. This is running on the login node relatively quickly. This is because this alignment for these FASTQ files was only performed for a small portion of the chromosome 1. Later on, this will take a lot longer. Therefore, you will need to generate a script.

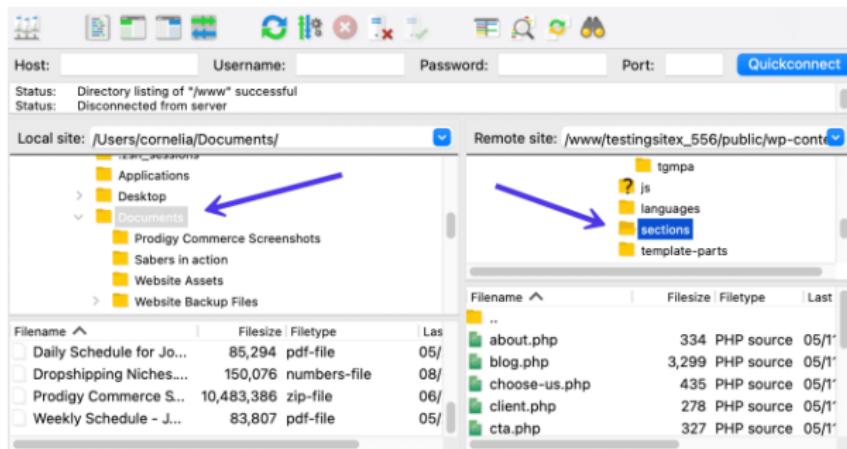
Running Parameters for FASTQC:

- 10G of memory is required
- 1 node, 2 tasks

# Viewing the HTML report from FASTQC

All of the following are solutions that allow students to transfer files between remote (i.e. VACC) and local (i.e. your laptop) servers.

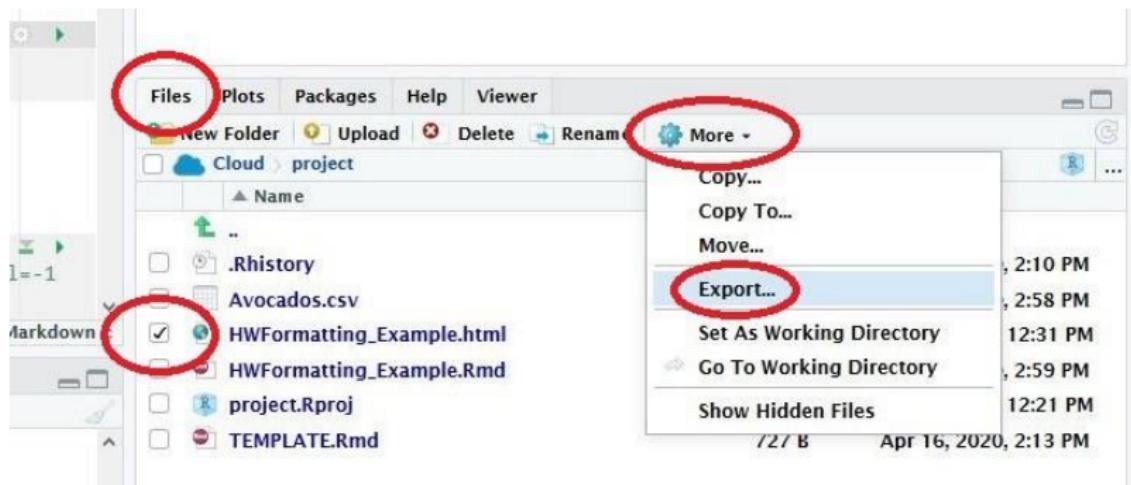
An FTP application such as Filezilla



# Viewing the HTML report from FASTQC

## RStudio (via VACC-OOD)

You can export it or simply view it using RStudio



# Viewing the HTML report from FASTQC

In File Explorer on OpenOnDemand, use the “Download” button

The screenshot shows the OpenOnDemand File Explorer interface. At the top, there is a toolbar with buttons for 'Go To...', 'Open in Terminal', 'New File', 'New Dir', 'Upload', 'Show Dotfiles', and 'Show Owner/Mode'. Below the toolbar, the path '/users/appl/shussain/' is displayed. On the left, a sidebar titled 'Home Directory' lists various folders: Desktop, Documents, Downloads, Music, Pictures, Public, Templates, Videos, WorkflowOptimization, crimson\_files, links, ondemand, ood\_data, test, and tmp. The main area displays a list of files and directories with columns for 'name', 'size', and 'modified date'. A red box highlights the 'Download' button in the toolbar. The table data is as follows:

name	size	modified date
..	<dir>	10/07/2015
Desktop	<dir>	10/07/2015
Documents	<dir>	10/07/2015
Downloads	<dir>	10/07/2015
Music	<dir>	10/07/2015
Pictures	<dir>	10/07/2015
Public	<dir>	10/07/2015
Templates	<dir>	10/07/2015
Videos	<dir>	10/07/2015
WorkflowOptimization	<dir>	09/02/2015
crimson_files	<dir>	01/11/2016

# Class Exercise #2: Assessing Universal Metrics



Grab the following folder from the location below.

/gpfs1/cl/mmg3320/course\_materials/FASTQC\_example

Please answer questions in **ClassExercise2-FASTQC-020926.docx**

# Citation

*This lesson has been developed by members of the teaching team at the [Harvard Chan Bioinformatics Core \(HBC\)](#). These are open access materials distributed under the terms of the [Creative Commons Attribution license](#) (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.*