

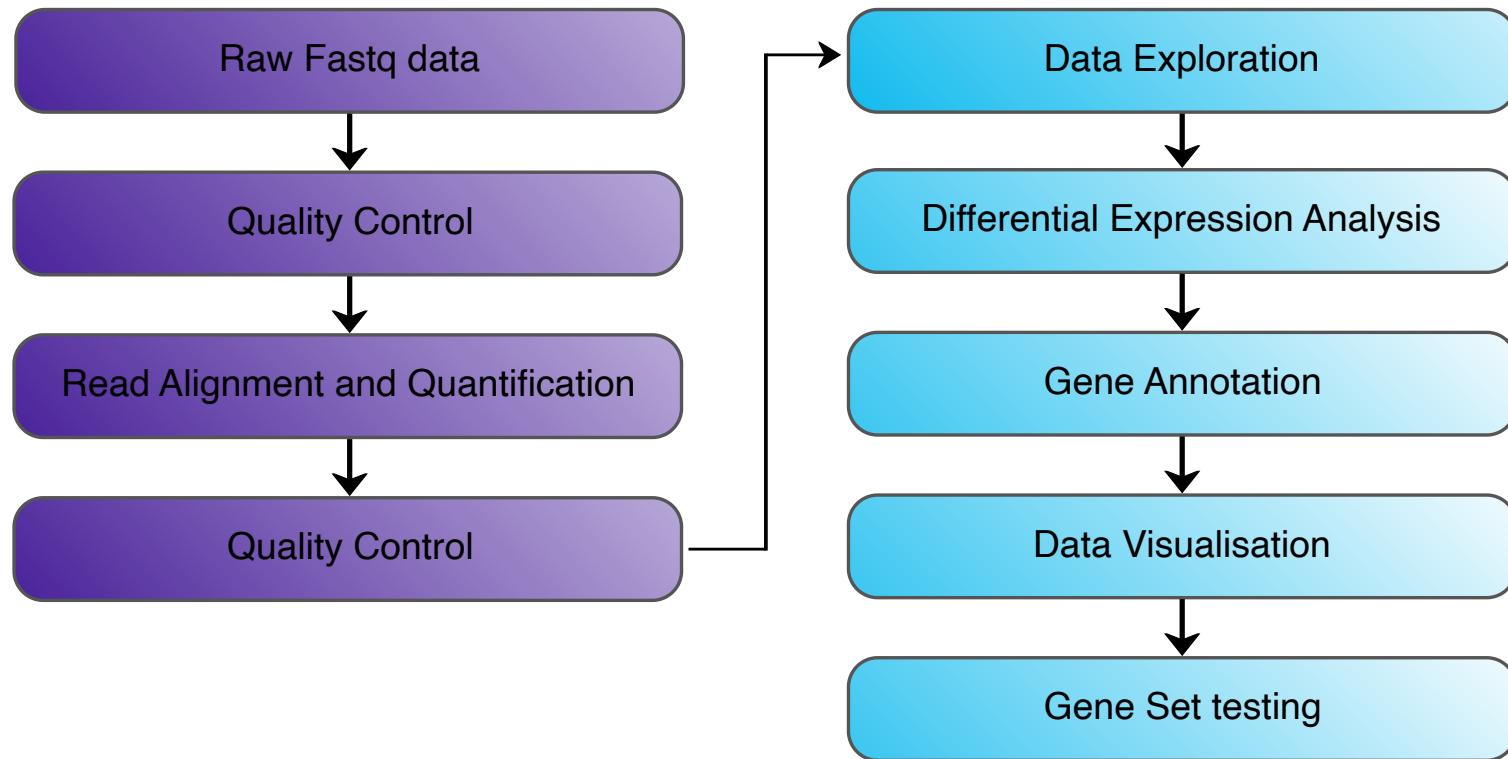
CANCER  
RESEARCH  
UK

CAMBRIDGE  
INSTITUTE

# Basic quality control with FastQC

March 2023

# Differential Gene Expression Analysis Workflow



# Fastq file format

## Fastq file format - Headers

```
@SRR7657883.sra.1 1 length=150
CTTGCGGTGTTCTCTGTTGCCTCAAGGATGGCCTTGAATTCCTCACAAAGGTTGTGTCGAATGTTTGACAGTTACTAATT
+SRR7657883.sra.1 1 length=150
AAAFJJJJJJJJJJJJAFJJFJFFFJJJJJJJJFJJFJAAAJJFJFFFJFJJ--AAJ--AA-AF-<-AF-7AAJ7J7-7--7-
@SRR7657883.sra.2 2 length=150
CCTGGTCATCGGGCACCTGGACACTTCACGGCCCAAGGCGGGCGAGACCCTGGGTTCTGTTGAGGCAGCTGGCTTTCCATA
+SRR7657883.sra.2 2 length=150
AAFFFJFJJJJ<JJJJJJFJJJJJJ<JJJJJJJJJJJJJJJJJJJJJJJJJJFAJJJJJJJJJJJJJJFJJAJJJJJJJAFJAFAAFA
@SRR7657883.sra.3 3 length=150
GTCCACTCTGCACTCCGAGCTTCTTCCCTGTCCAGGTAGCACCTTCGAGACGTGAAGATGTTAGAACGCCGCTTGGACT
+SRR7657883.sra.3 3 length=150
AAFFFJJJJJJJJJJJJJJJJJJJJJJJJFJJJJJJJJJJJJFJJJJJJJJJJJA7AFJFFJJFJJFJJFJJFJJFFFJ-AA
@SRR7657883.sra.4 4 length=150
CGGGGGATTGCAACACAATACCGTGCTGGTGGCAGCAGCATGAGCATCCCCCGTCGTGCACAATCCACCAACACGTGAATGCTGC
+SRR7657883.sra.4 4 length=150
AAAFFFFAJJ<AFJJJJJJFJAJJFAJFF<JAJJJJAFJJJJJF<FJFFJJJJJAJ7AJFJFFJJFJ7J7JFFAA<J-7-AFJ<F<
@SRR7657883.sra.5 5 length=150
CGGGAGTGCTCACTGTCACCCAAATCCACATCCCCCCCACCGTCTCCAGAGGTGTGGCGGAGAGCCGGACTCTTGGAAATCATCTC
+SRR7657883.sra.5 5 length=150
AAAFFAFFFJJJJJJJJJJFJJJJJJJJJJ7FJJJJJJJJJJFJJFJJFJJAAAAJ<JF<JAJJJAJ<FJJJFFF-AF7AJ-<J
```

## Fastq file format - Sequences

```
@SRR7657883.sra.1 1 length=150
CTTGCGGTGTTCTCTGTTGCCTCAAGGATGGCCTTGAACCCCTACAAGGGTTGTGCGAATGTTTGACAGTTACTAATTCT
+SRR7657883.sra.1 1 length=150
AAAFJJJJJJJJJJJJJAFJJFJFFFJJJJJJJJFJJFJAAAJJFJFFFJJFJ--AAJ--AA-AF-<-AF-7AAJ7J7-7--7-
@SRR7657883.sra.2 2 length=150
CCTGGTCATGGGCACCTGGACACTTCACGGCCCAAGGCGGGCCGAGACCCCTGGGTTCTGTTGAGGCAGCTGGCTTTCCATA
+SRR7657883.sra.2 2 length=150
AAFFFJFJJJJ<JJJJJJFJJJJJJ<JJJJJJJJJJJJJJJJJJJJJJJJJJFAJJJJJJJJJJJJJJFJJAJJJJJJJAFJAFAAFA
@SRR7657883.sra.3 3 length=150
GTCCACTCTGCACTCCGCAGCTTCTTCCCTGTCCAGGTAGCACCTTCGAGACGTGAAGATGTTAGAACGCCGCTTGGACT
+SRR7657883.sra.3 3 length=150
AAFFFJJJJJJJJJJJJJJJJJJJJJJFJJJJJJJJJJJJJJFJJJJJJJJJJJA7AFJFFJJFJJFJJFJJFFFJ-AA
@SRR7657883.sra.4 4 length=150
CGGGGGATTGCAACACAATACCGTGTGGTGGGCAGCAGCATGAGCATGCCCGTGTGCACAATCCACCAACACGTGAATGCTGC
+SRR7657883.sra.4 4 length=150
AAAFFFAJJ<AFJJJJJJFJAJJFAJFF<JAJJJJAFJJJJJJF<FJFFJJJJJAJ7AJFJFFJJFJ7J7JFFAA<J-7-AFJ<F<
@SRR7657883.sra.5 5 length=150
CGGGAGTGTCACTGTCACCCCAAATCCACATCCCCCCCACCGTCTCCAGAGGTGTGGCCGGAGAGCCGGACTCTTGGAAATCATCTC
+SRR7657883.sra.5 5 length=150
AAAFAFFFJJJJJJJJJJFJJJJJJJJJJ7FJJJJJJJJFJJFJJFJJAAAAJ<JF<JAJJJAJ<FJJJFFF-AF7AJ-<J
```

## Fastq file format - Third line

```
@SRR7657883.sra.1 1 length=150
CTTGCGGTGTTCTCTGTTGCCTCAAGGATGGCCTTGACTTCCTCACAAAGGTTGTGCGAATGTTTGACAGTTACTAATT
+SRR7657883.sra.1 1 length=150
AAAFJJJJJJJJJJJJJAFJJFJFFFJJJJJJJJFJJFJAAAJJFJFFFJFJJ--AAJ--AA-AF-<-AF-7AAJ7J7-7--7-
@SRR7657883.sra.2 2 length=150
CCTGGTCATCGGGCACCTGGACACTTCACGGCCCAAGGCCGGCCGAGACCCCTGGGTTCTGTTGAGGCAGCTGGCTTTCCATA
+SRR7657883.sra.2 2 length=150
AAFFFJFJJJJ<JJJJJJFJJJJJJ<JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJFAJJJJJJJJJJJJJJFJJAJJJJJJJAFJAFAAFA
@SRR7657883.sra.3 3 length=150
GTCCACTCTGCACTCCGAGCTTCTTCCCTGTCCAGGTAGCACCTTCGAGACGTGAAGATGTTAGAACGCCGCTTGGACT
+SRR7657883.sra.3 3 length=150
AAFFFJJJJJJJJJJJJJJJJJJJJJJJJFJJJJJJJJJJJJJJFJJJJJJJJJA7AFJFFJJFJJFJJFJJFJJFFFJ-AA
@SRR7657883.sra.4 4 length=150
CGGGGGATTGCAACACAATACCGTGTGGTGGCAGCAGCATGAGCATCCCCCGTCGTGCACAATCCACCAACACGTGAATGCTGC
+SRR7657883.sra.4 4 length=150
AAAFFF AJJ<AFJJJJJJFJAJJFAJFF<JAJJJJAFJJJJJF<FJFFJJJJJAJ7AJFJFFJJFJ7J7JFFAA<J-7-AFJ<F<
@SRR7657883.sra.5 5 length=150
CGGGAGTGTCACTGTCACCCAAATCCACATCCCCCCCACCGTCTCCAGAGGTGTGGCGGAGAGCCGGACTCTGGAATCATCTC
+SRR7657883.sra.5 5 length=150
AAAFFFJJJJJJJJJJFJJJJJJJJJJ7FJJJJJJJJFJJFJJFJJAAAAJ<JF<JAJJJAJ<FJJFFF-AF7AJ-<J
```

# Fastq file format - Quality Scores

```
@SRR7657883.sra.1 1 length=150
CTTGCGGTGTTCTCTGTTGCCTCAAGGATGGCCTTGAATTCCTACAAGGGTTGTCGAATGTTTGACAGTTACTAATT
+SRR7657883.sra.1 1 length=150
AAAFAFJJJJJJJJJJJJJJFJFJFFFJJJJJJJJFJJFJAAAJJFJFFFJJFJ--AAJ--AA-AF-<-AF-7AAJ7J7-7--7-
@SRR7657883.sra.2 2 length=150
CCTGGTCATCGGGCACCTGGACACTTCACGGCCCAAGGCGGGCCGAGACCCTGGGTTCTGTTGAGGCAGCTGGCTTCCATA
+SRR7657883.sra.2 2 length=150
AAFFFJFJJJJ<JJJJJJFJJJJ<JJJJJJJJJJJJJJJJJJJJJJJJJJJJFAJJJJJJJJJJJJJJFJJAJJJJJJJAFJAFAAFA
@SRR7657883.sra.3 3 length=150
GTCCACTCTGCACTCCGCAGCTTCTTCCCTGTCCAGGTAGCACCTTCGAGACGTGAAGATGTTAGAACGCCGCTTGGACT
+SRR7657883.sra.3 3 length=150
AAFFFJJJJJJJJJJJJJJJJJJJJJJJJJJFJJJJJJJJJJJJJJFJJJJJJJJJJJA7AFJFFJJFJJFJJFJJFFFJ-AA
@SRR7657883.sra.4 4 length=150
CGGGGGATTGCAACACAATACCGTGCTGGTGGCAGCAGCATGAGCATCCCCCGTCGTGCACAATCCACCAACACGTGAATGCTGC
+SRR7657883.sra.4 4 length=150
AAAFFFAJJ<AFJJJJJJFJAJJFAJFF<JAJJJJAFJJJJF<FJFFJJJJJAJ7AJFJFFJJFJ7J7JFFAA<J-7-AFJ<F<
@SRR7657883.sra.5 5 length=150
CGGGAGTGCTCACTGTCACCCAAATCCACATCCCCCCCACCGTCTCCAGAGGTGTGGCGGAGAGCCGGACTCTTGAATCATCTC
+SRR7657883.sra.5 5 length=150
AAAFAFFFJJJJJJJJJJFJJJJJJJJJJ7FJJJJJJJJFJJFJJFJJAAAAJ<JF<JAJJJAJ<FJJFFF-AF7AJ-<J
```

# (Phred) Quality Scores

Sequence quality scores are transformed and translated p-values

- Sequence bases are called after image processing (base calling)
  - Each base in a sequence has a *p-value* associated with it
  - p-values range from 0-1 (e.g.: 0.05, 0.01, 1e-30)
  - p-value of 0.01 inferred as 1 in 100 chance that called base is wrong

# (Phred) Quality Scores ...

How do we assign p-values to bases in the fastq file?

- P-vales can be many characters long (e.g.: 0.000005)
- Transform to Phred quality scores  $Q$
- $Q = -10(\log_{10}P)$  (e.g.: 0.01 = Q value of 20, 0.001 = Q value of 30)
- Translate  $Q$  values to ASCII characters (adding 33) (Q value of 30 = ?, Q value of 40 = !)

Dec	Hex	Chr	Dec	Hex	Chr	Dec	Hex	Chr	Dec	Hex	Chr
0	00	NUL	32	20	Space	64	40	@	96	60	'
1	01	SOH	33	21	!	65	41	A	97	61	a
2	02	STX	34	22	"	66	42	B	98	62	b
3	03	ETX	35	23	#	67	43	C	99	63	c
4	04	EOT	36	24	\$	68	44	D	100	64	d
5	05	ENQ	37	25	%	69	45	E	101	65	e
6	06	ACK	38	26	&	70	46	F	102	66	f
7	07	BEL	39	27	'	71	47	G	103	67	g
8	08	BS	40	28	(	72	48	H	104	68	h
9	09	HT	41	29	)	73	49	I	105	69	i
10	0A	LF	42	2A	*	74	4A	J	106	6A	j
11	0B	VT	43	2B	+	75	4B	K	107	6B	k
12	0C	FF	44	2C	,	76	4C	L	108	6C	l
13	0D	CR	45	2D	-	77	4D	M	109	6D	m
14	0E	SO	46	2E	.	78	4E	N	110	6E	n
15	0F	SI	47	2F	/	79	4F	O	111	6F	o
16	10	DLE	48	30	0	80	50	P	112	70	p
17	11	DC1	49	31	1	81	51	Q	113	71	q
18	12	DC2	50	32	2	82	52	R	114	72	r
19	13	DC3	51	33	3	83	53	S	115	73	s
20	14	DC4	52	34	4	84	54	T	116	74	t
21	15	NAK	53	35	5	85	55	U	117	75	u
22	16	SYN	54	36	6	86	56	V	118	76	v
23	17	ETB	55	37	7	87	57	W	119	77	w
24	18	CAN	56	38	8	88	58	X	120	78	x
25	19	EM	57	39	9	89	59	Y	121	79	y
26	1A	SUB	58	3A	:	90	5A	Z	122	7A	z
27	1B	ESC	59	3B	;	91	5B	[	123	7B	{
28	1C	FS	60	3C	<	92	5C	\	124	7C	
29	1D	GS	61	3D	=	93	5D	]	125	7D	}
30	1E	RS	62	3E	>	94	5E	^	126	7E	~
31	1F	US	63	3F	?	95	5F	_	127	7F	DEL

# QC is important

Check for any problems before we put time and effort into analysing potentially bad data

- Start with FastQC
  - Quick
  - Outputs an easy to read html report



We run fastQC from the terminal with the command

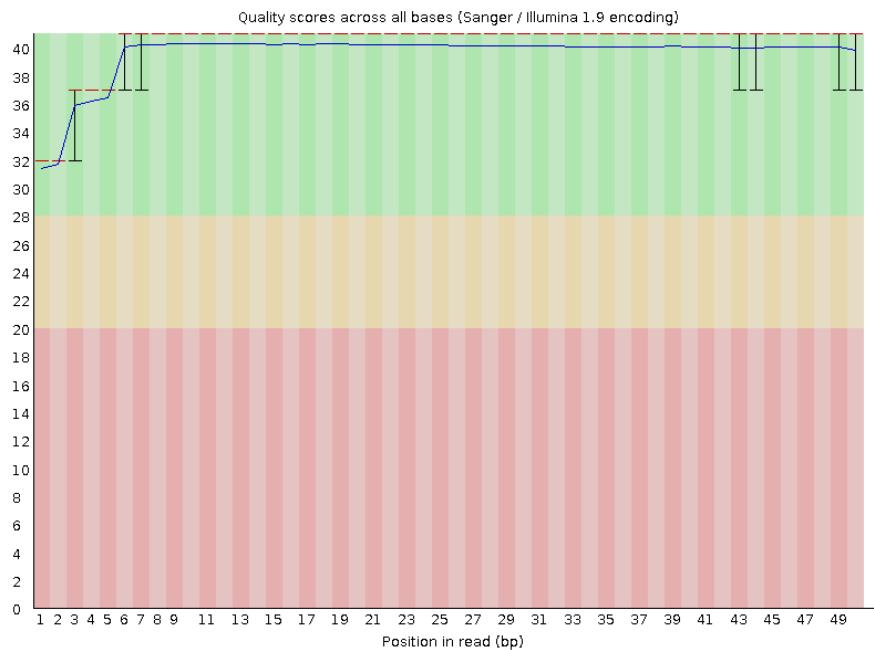
```
fastqc <fastq>
```

but there are lots of other parameters which you can find to tailor your QC by typing

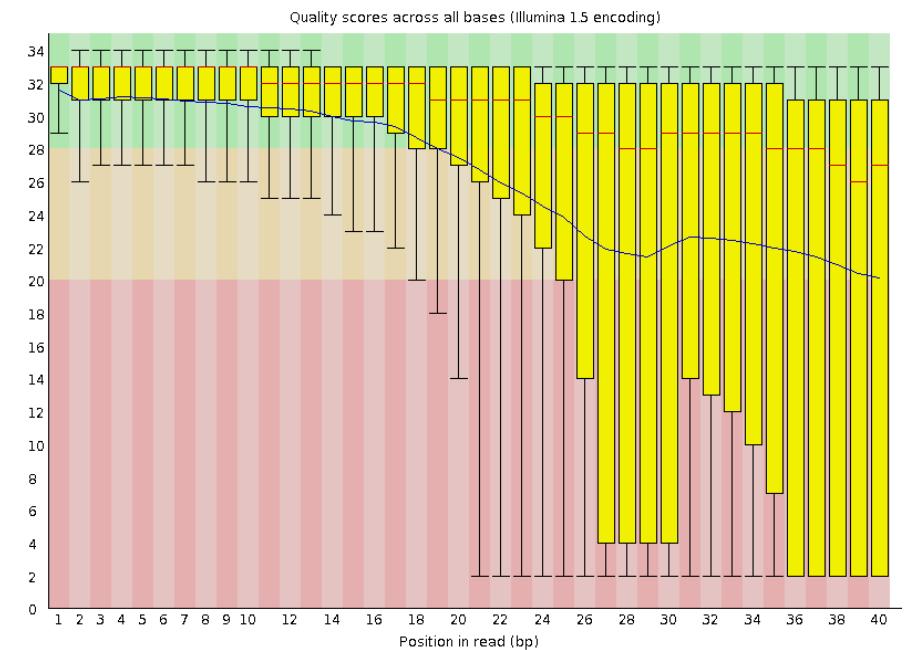
```
fastqc -h
```

# Per base sequence quality

Good Data

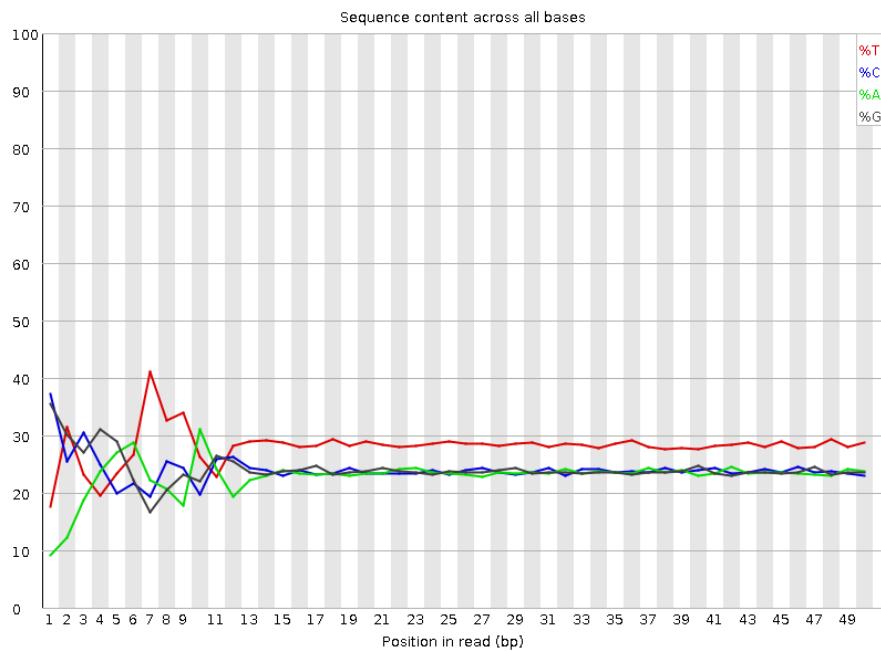


Bad Data

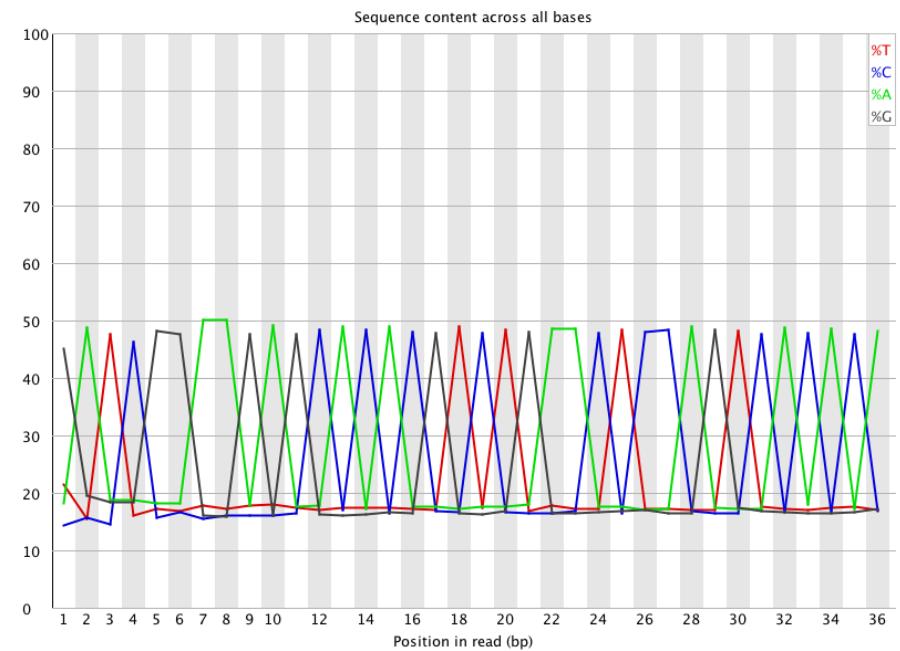


# Per base sequence content

Good Data

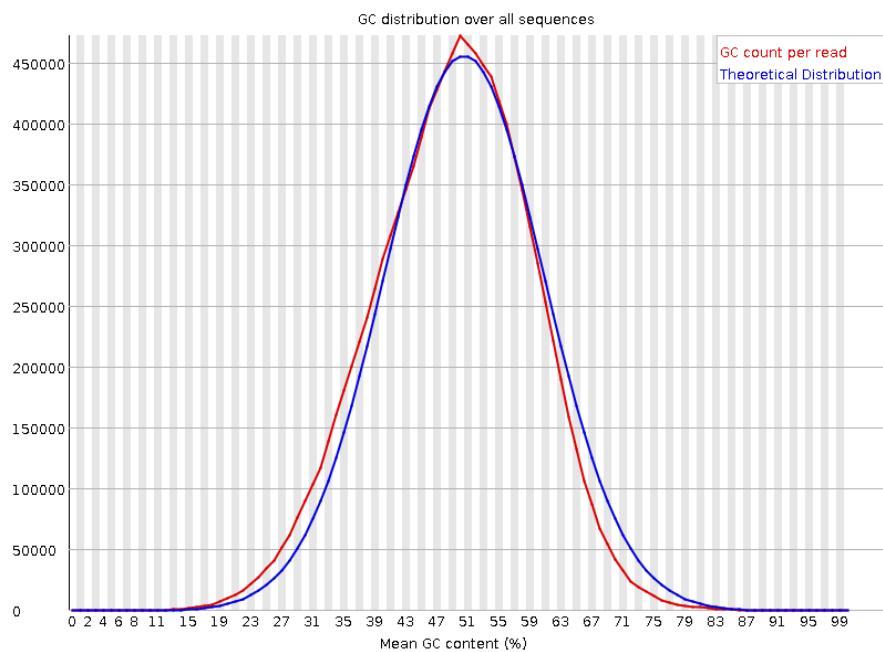


Bad Data

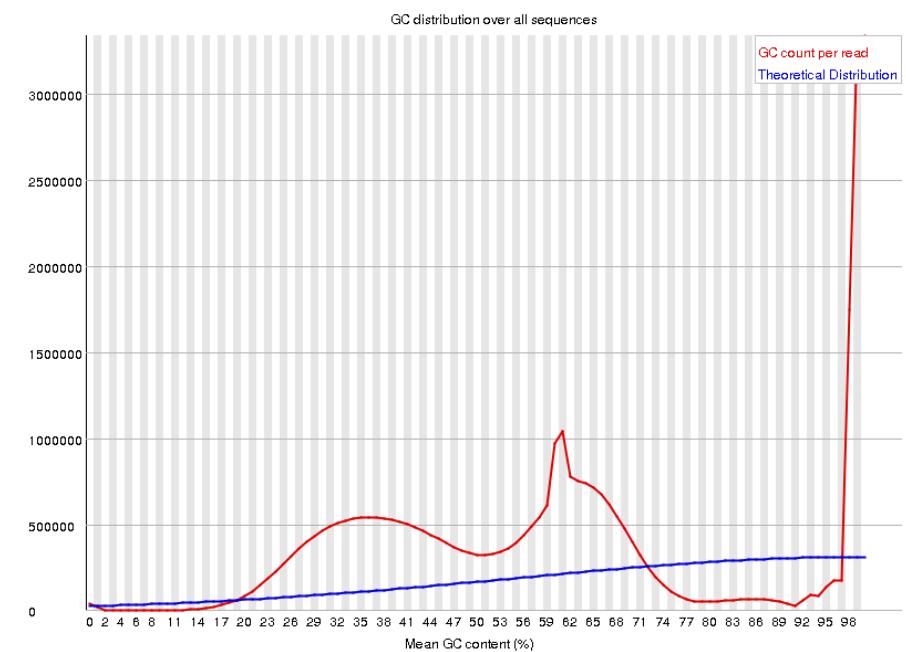


# Per sequence GC content

Good Data

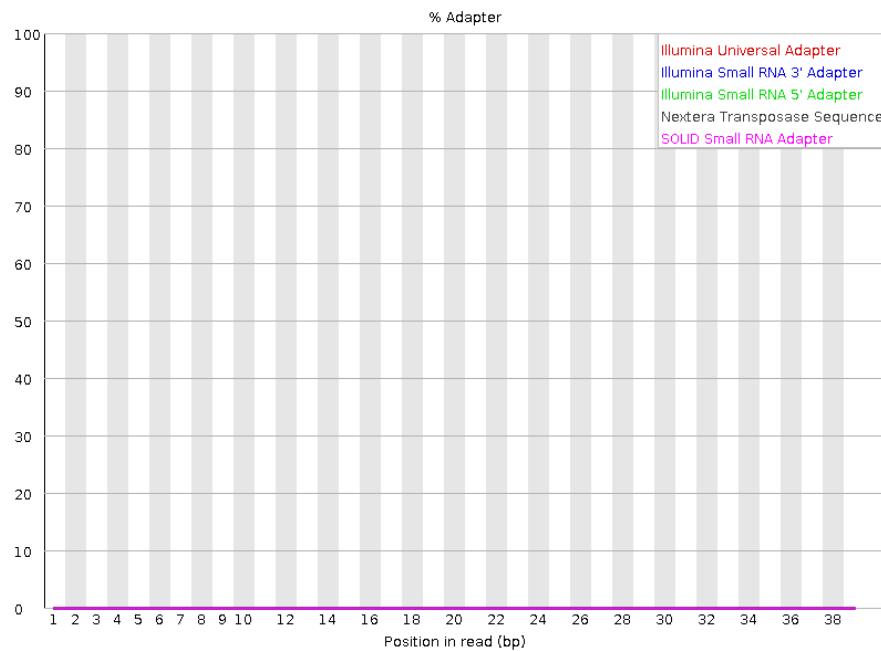


Bad Data

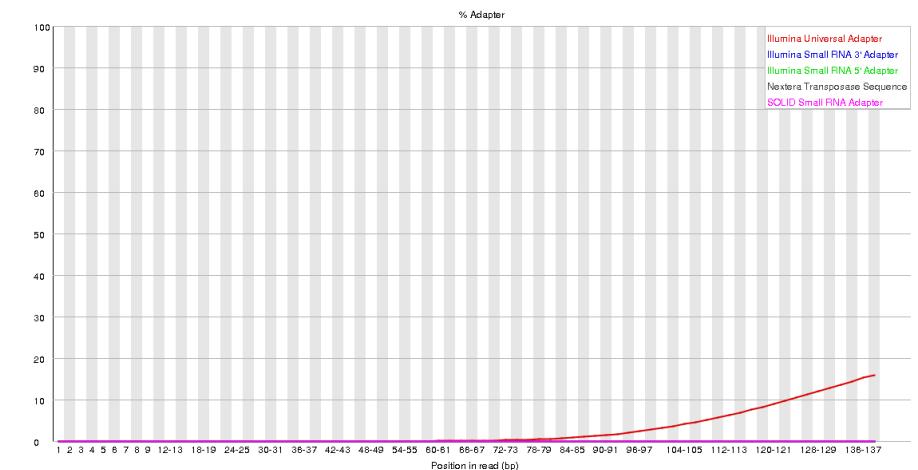


# Adaptor content

Good Data



Bad Data



## And now onto the exercise...

- Go to: <https://ushers.bio.cam.ac.uk/guacamole2>
- Log on with YOUR credentials that were emailed to you

# A quick intro to the environment

- The terminal is just a text based version of the operating system
- We will look at an example with side by side GUI and text file system...
- You use commands instead of mouse clicks - commands are case-sensitive and can be followed by arguments with spaces
  - cd
  - pwd
  - ls
  - flags - e.g. ls -a
  - the directory structure is like a tree, you can go back with cd ..
  - Up arrows to get through history
  - tab complete to avoid errors
  - More to look at the files and q to exit
  - ctrl-c