



# IX International Course of Massive Data Analysis FOR GENOMICS

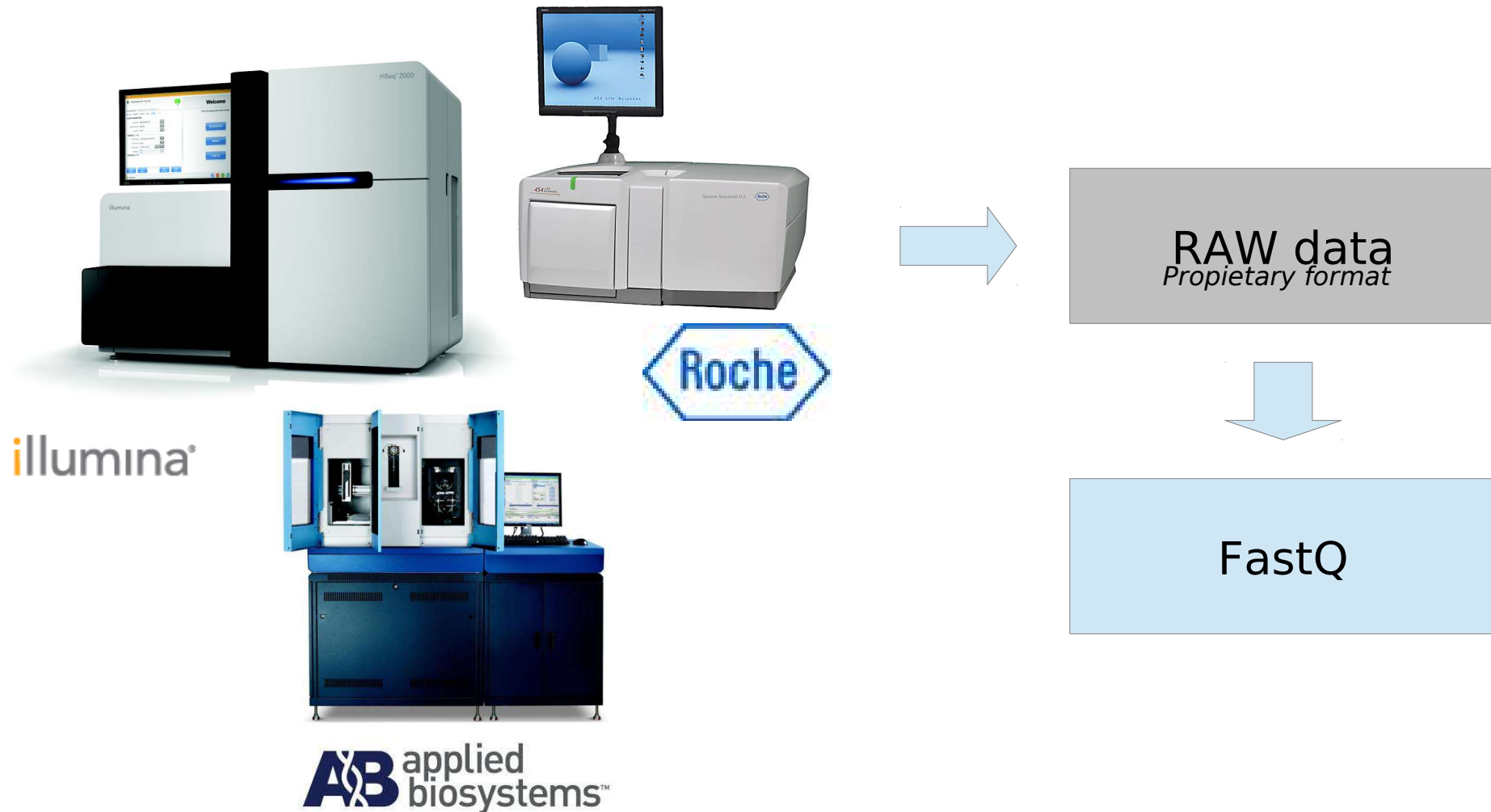


# Contents

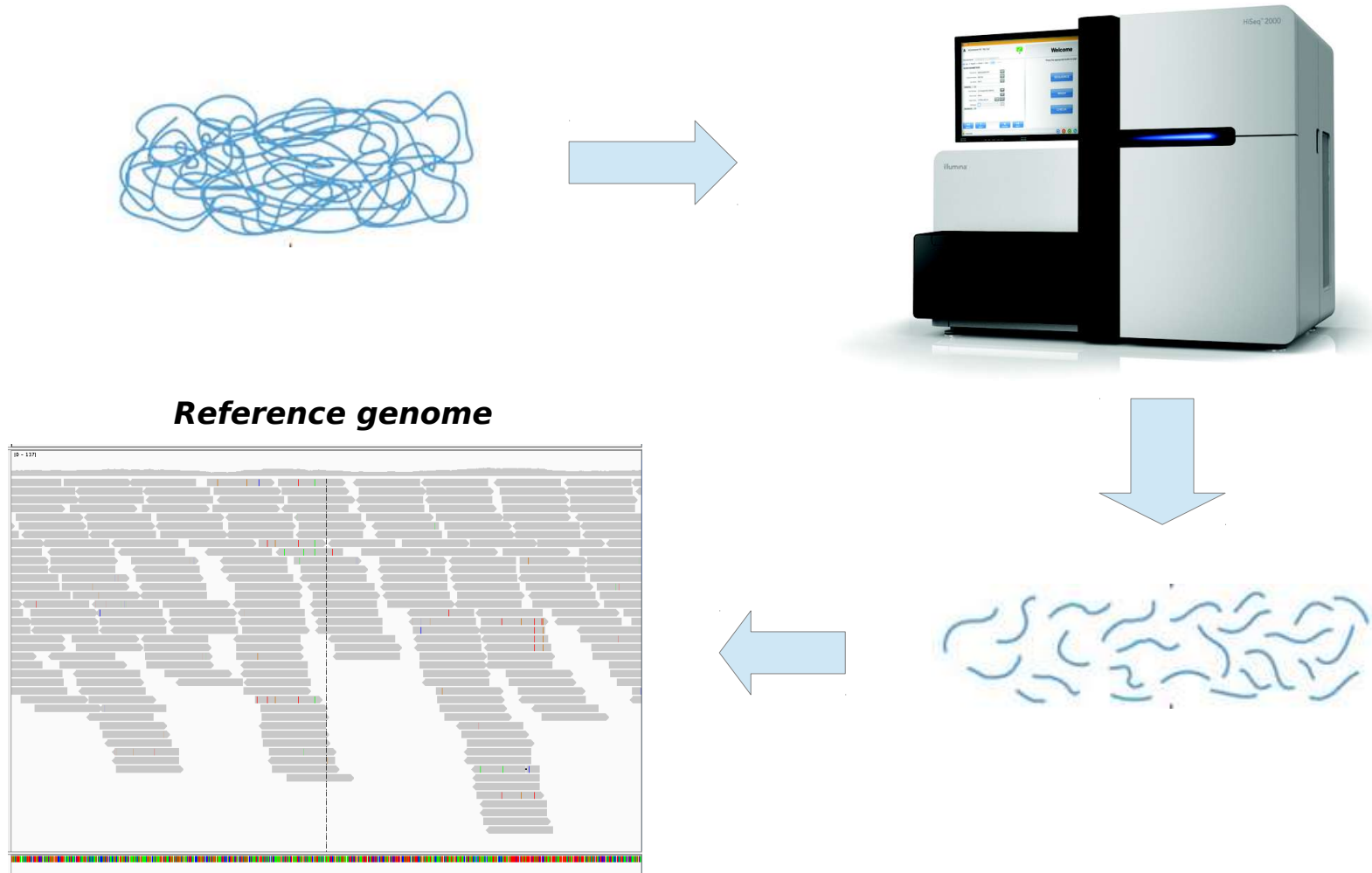
---

- Data formats
  - ▣ Sequence capture
  - ▣ Fasta and fastq formats
  - ▣ Sequence quality encoding
- Quality Control
  - ▣ Evaluation of sequence quality
  - ▣ Quality control tools
  - ▣ Identification of artifacts & filtering
- Practical session

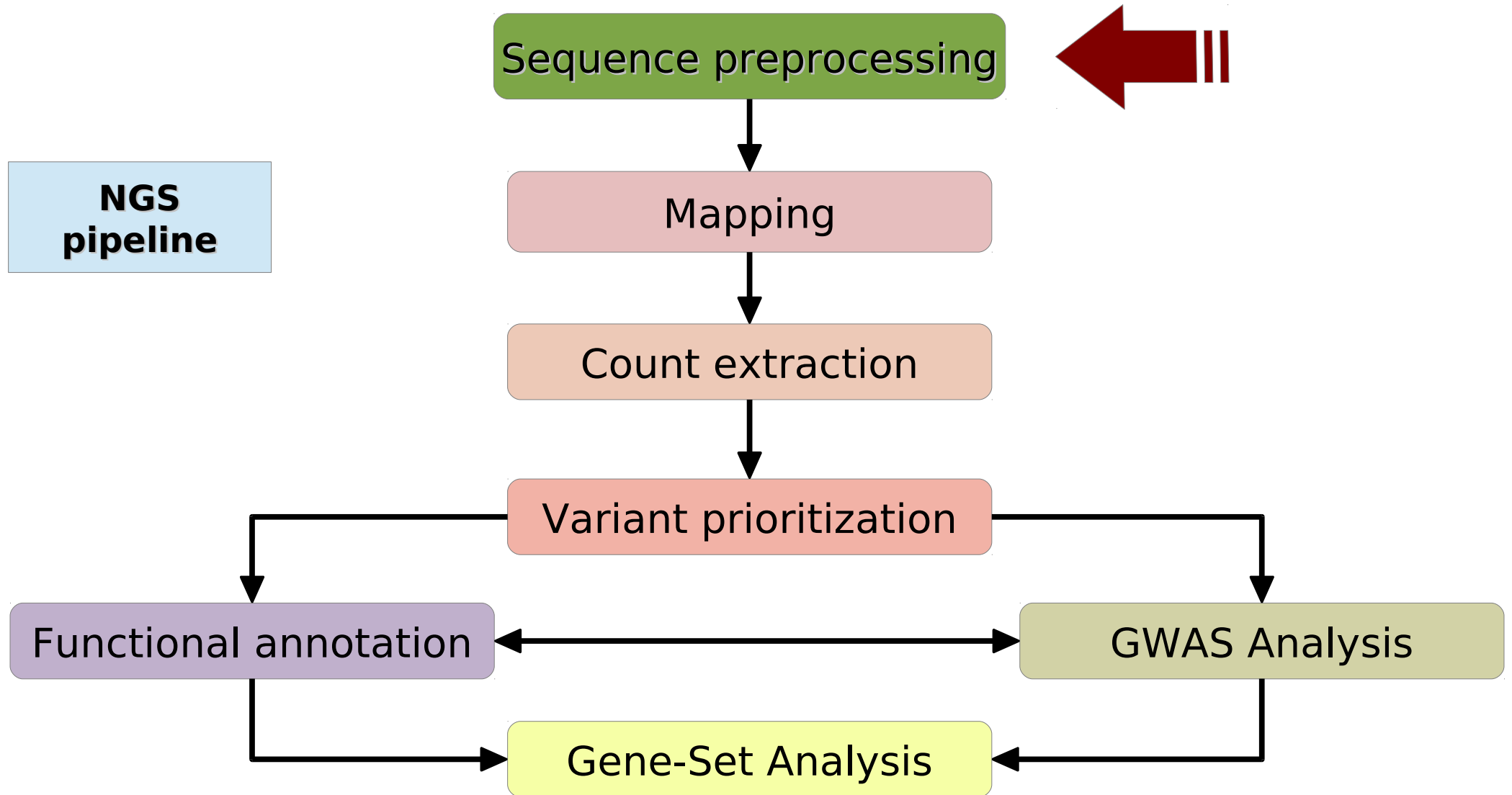
# Sequence capture



# Genome sequencing



# Where are we?



# From sequencers to digital data

- **What structure does the data have?**
  - Text-based formats (easy to use!)
  - If not compressed, it can be huge
- **Different data formats:**
  - Different sequencers output different files (sff, fasta, csfasta, qual file, fastq...)
  - There are some data formats widely accepted (e.g. FastQ format)

# Fasta format

- Two lines per sequence:
  - 1. Header lines starts with “>” followed by a sequence ID
  - 2. Sequence (string of nt or peptides)

```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]  
LCLYTHIGRNIYYGSYLYSETWNTGIMLLITMATAFMGYVLPWGQMSFWGATVITNLFSAIPYIGTNLV  
EWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIPFHPYYTIKDFLG  
LLILILLLLLLLALLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRSVPNKLGGVLALFLSIVIL  
GLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYPYTIIGQMASILYFSIILAFPLIAGX  
IENY
```

```
>BBTBSCRYR  
tgcaccaaaccatgtctaaagctggaaccaaattactttctttgaagacaaaaactttca  
aggccgccactatgacagcgattgcgactgtgcagatttccacatgtacctgagccgctg  
caactccatcagagtggaaggaggcacctgggctgtgtatgaaaggccaattttgctgg  
gtacatgtacatcctaccccgggcgagtatcctgagtaccagcactggatgggcctcaa
```

- Typical file extensions (.fasta, .fa, .fna, .fnn, .faa, ...)

# Fastq format

- We could say “it is a fasta with **qualities**”:
  - 1. Header (like the fasta but starting with “@”)
  - 2. Sequence (string of nt)
  - 3. “+” and sequence ID (optional)
  - 4. Encoded quality of the sequence

```
@SEQ_ID  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT  
+  
!'!*((( (***) )%%%++) (%%%) .1***-+*') **55CCF>>>>>CCCCCCC65
```



# Quality codification

## □ Phred quality score

□ Error probability

□ ASCII encoded

□ Phred +33

- Sanger [0,40]
- Illumina 1.8 [0,41]
- Illumina 1.9 [0,41]

□ Phred +64

- Illumina 1.3 [0,40]
- Illumina 1.5 [3,40]

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32;	Space	64	40	100	&#64;	@	96	60	140	&#96;	`
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33;	!	65	41	101	&#65;	A	97	61	141	&#97;	a
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34;	"	66	42	102	&#66;	B	98	62	142	&#98;	b
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35;	#	67	43	103	&#67;	C	99	63	143	&#99;	c
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	\$	68	44	104	&#68;	D	100	64	144	&#100;	d
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	%	69	45	105	&#69;	E	101	65	145	&#101;	e
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	&	70	46	106	&#70;	F	102	66	146	&#102;	f
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39;	'	71	47	107	&#71;	G	103	67	147	&#103;	g
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40;	(	72	48	110	&#72;	H	104	68	150	&#104;	h
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	)	73	49	111	&#73;	I	105	69	151	&#105;	i
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	*	74	4A	112	&#74;	J	106	6A	152	&#106;	j
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	+	75	4B	113	&#75;	K	107	6B	153	&#107;	k
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	,	76	4C	114	&#76;	L	108	6C	154	&#108;	l
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45;	-	77	4D	115	&#77;	M	109	6D	155	&#109;	m
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46;	.	78	4E	116	&#78;	N	110	6E	156	&#110;	n
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47;	/	79	4F	117	&#79;	O	111	6F	157	&#111;	o
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48;	0	80	50	120	&#80;	P	112	70	160	&#112;	p
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49;	1	81	51	121	&#81;	Q	113	71	161	&#113;	q
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50;	2	82	52	122	&#82;	R	114	72	162	&#114;	r
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51;	3	83	53	123	&#83;	S	115	73	163	&#115;	s
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52;	4	84	54	124	&#84;	T	116	74	164	&#116;	t
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	5	85	55	125	&#85;	U	117	75	165	&#117;	u
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	6	86	56	126	&#86;	V	118	76	166	&#118;	v
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	7	87	57	127	&#87;	W	119	77	167	&#119;	w
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56;	8	88	58	130	&#88;	X	120	78	170	&#120;	x
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57;	9	89	59	131	&#89;	Y	121	79	171	&#121;	y
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58;	:	90	5A	132	&#90;	Z	122	7A	172	&#122;	z
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59;	;	91	5B	133	&#91;	[	123	7B	173	&#123;	{
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60;	<	92	5C	134	&#92;	\	124	7C	174	&#124;	
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61;	=	93	5D	135	&#93;	]	125	7D	175	&#125;	}
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62;	>	94	5E	136	&#94;	^	126	7E	176	&#126;	~
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63;	?	95	5F	137	&#95;	_	127	7F	177	&#127;	DEL

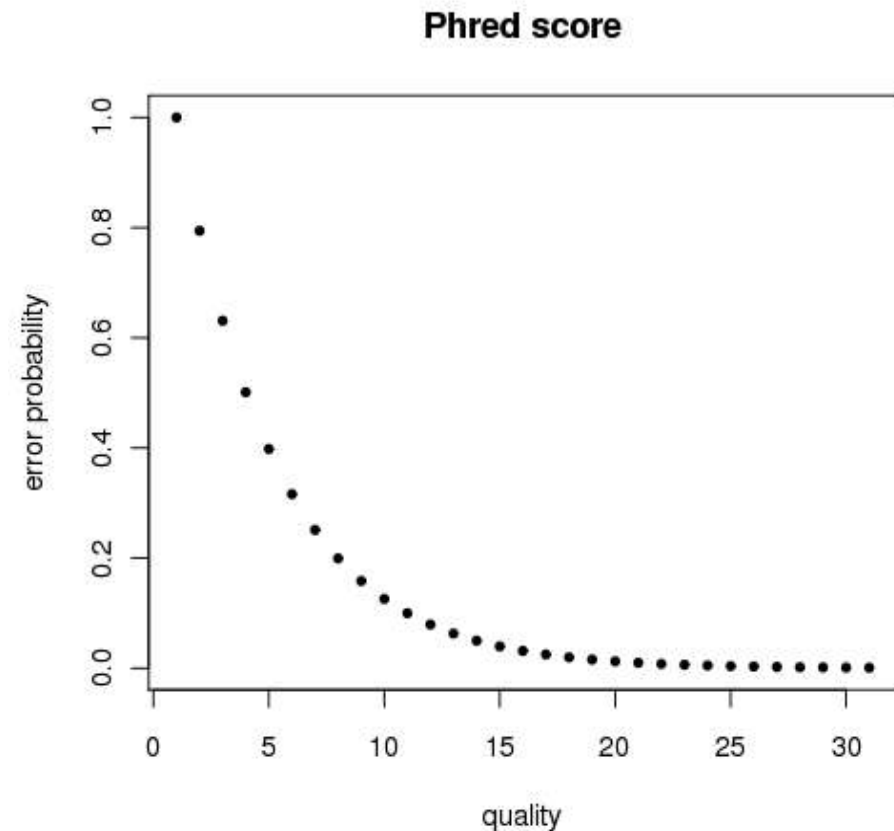
Source: [www.LookupTables.com](http://www.LookupTables.com)

# Quality codification

- Phred scores

$$Q = -10 \log_{10} P \quad \longleftrightarrow \quad P = 10^{-\frac{Q}{10}}$$

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%
50	1 in 100000	99.999%



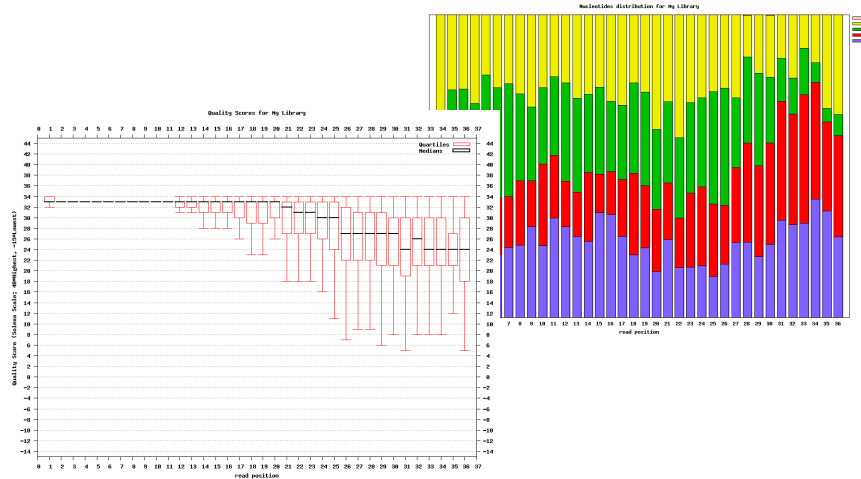
# Sequence quality evaluation

---

- If we evaluate our sequence in depth ...  
... we will know how reliable our results are
- **Problem:**
  - ▣ **Huge files** → Need of a tool to do it

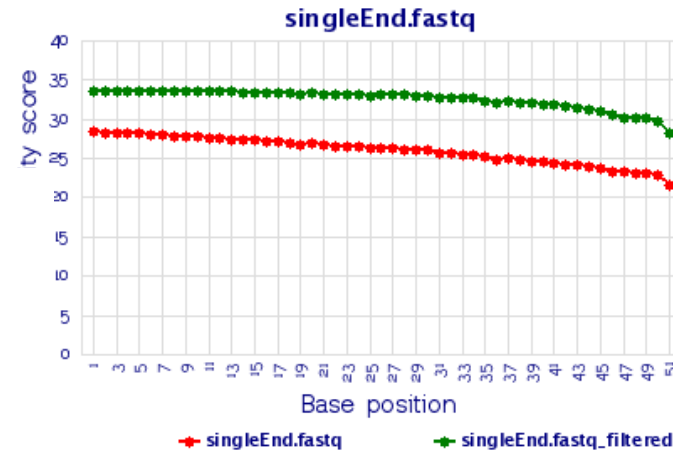
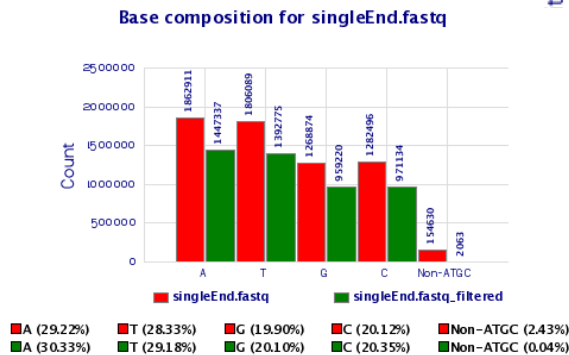
# Sequence quality evaluation

- Quality control tools:
  - **Fastx-toolkit**



[http://hannonlab.cshl.edu/fastx\\_toolkit/download.html](http://hannonlab.cshl.edu/fastx_toolkit/download.html)

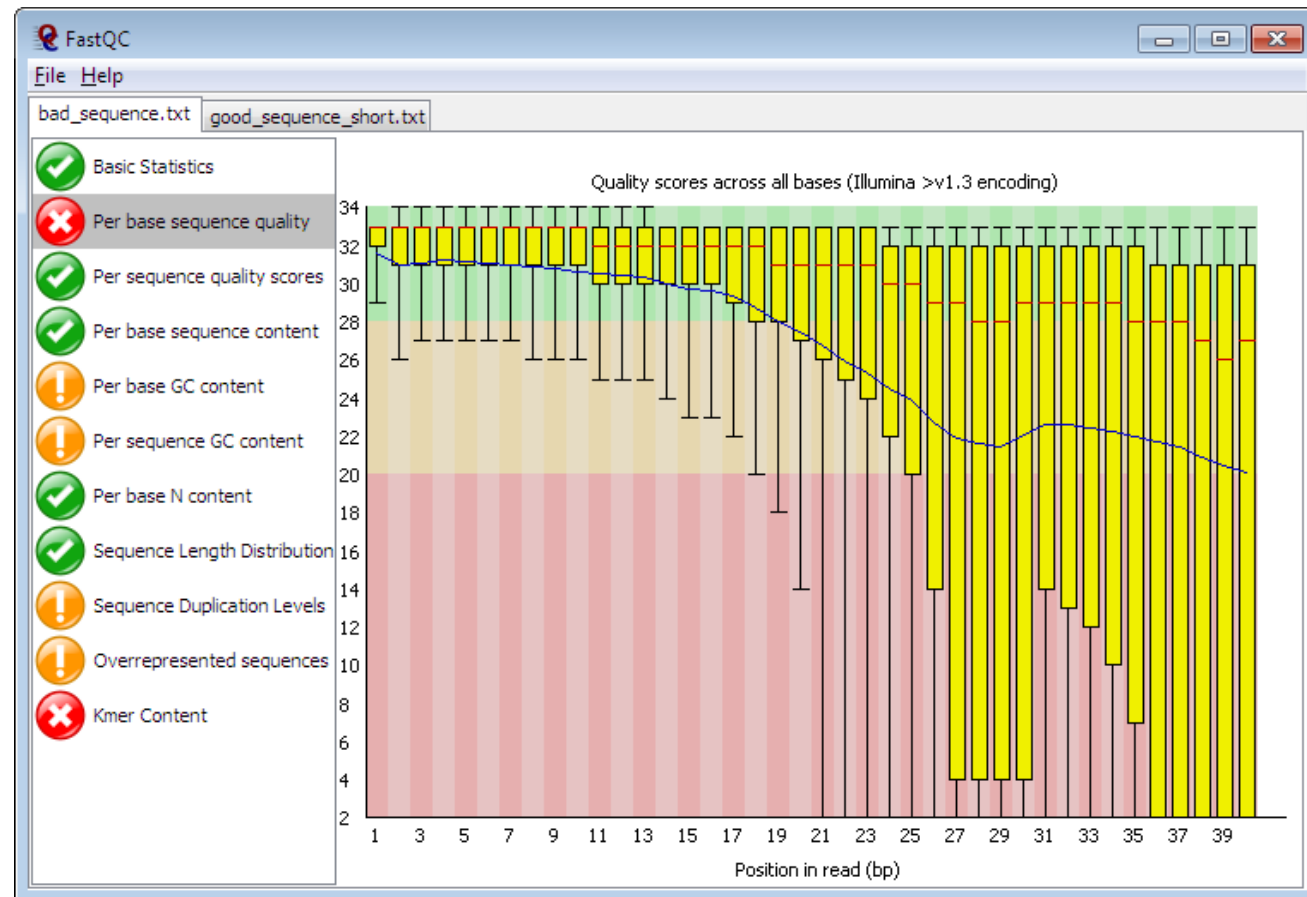
- **NGS QC Toolkit**



<http://www.nipgr.res.in/ngsqctoolkit.html>

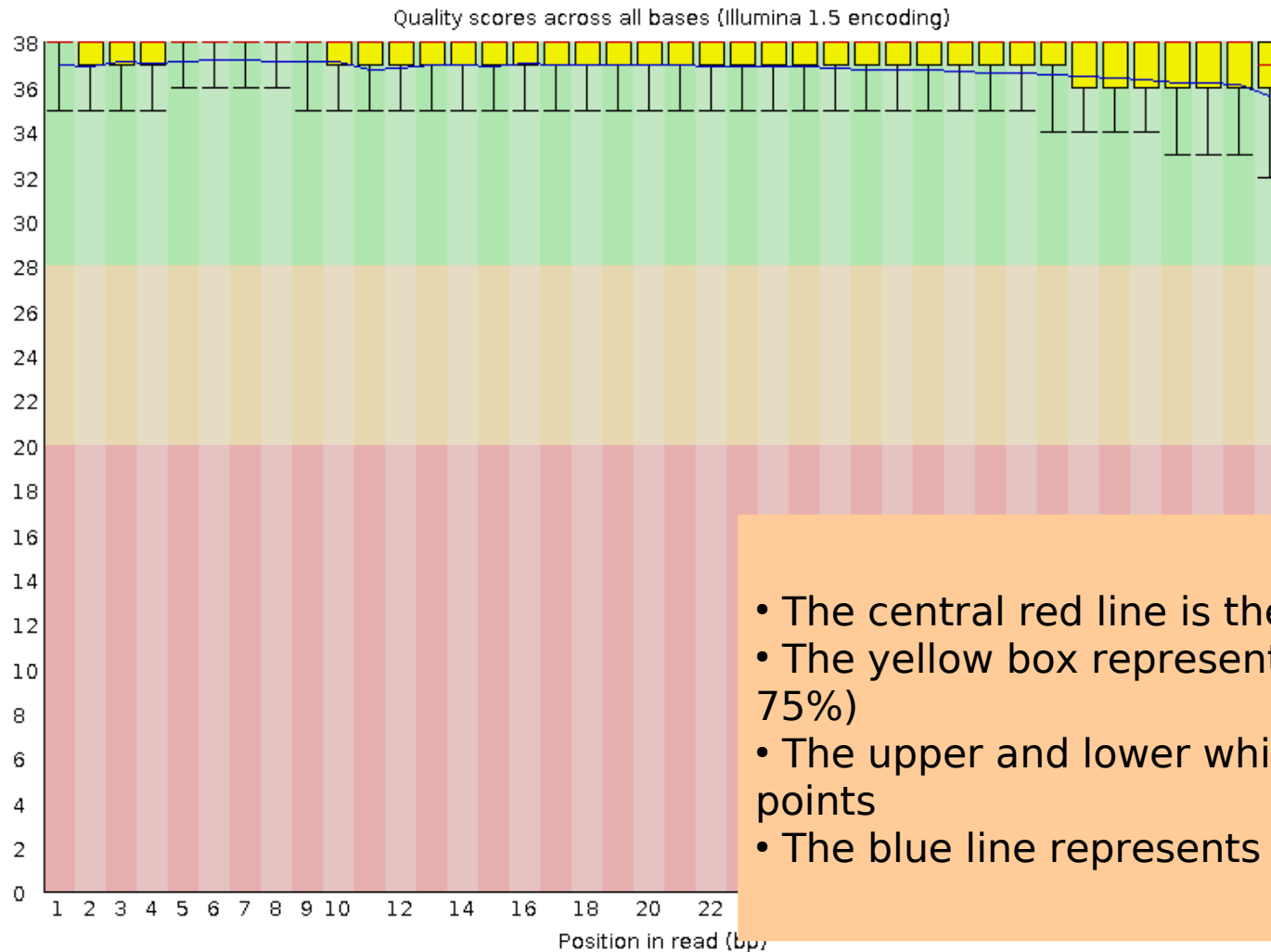
# Sequence quality evaluation

- Other quality control tool: **FastQC**



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

# Sequence quality per base position

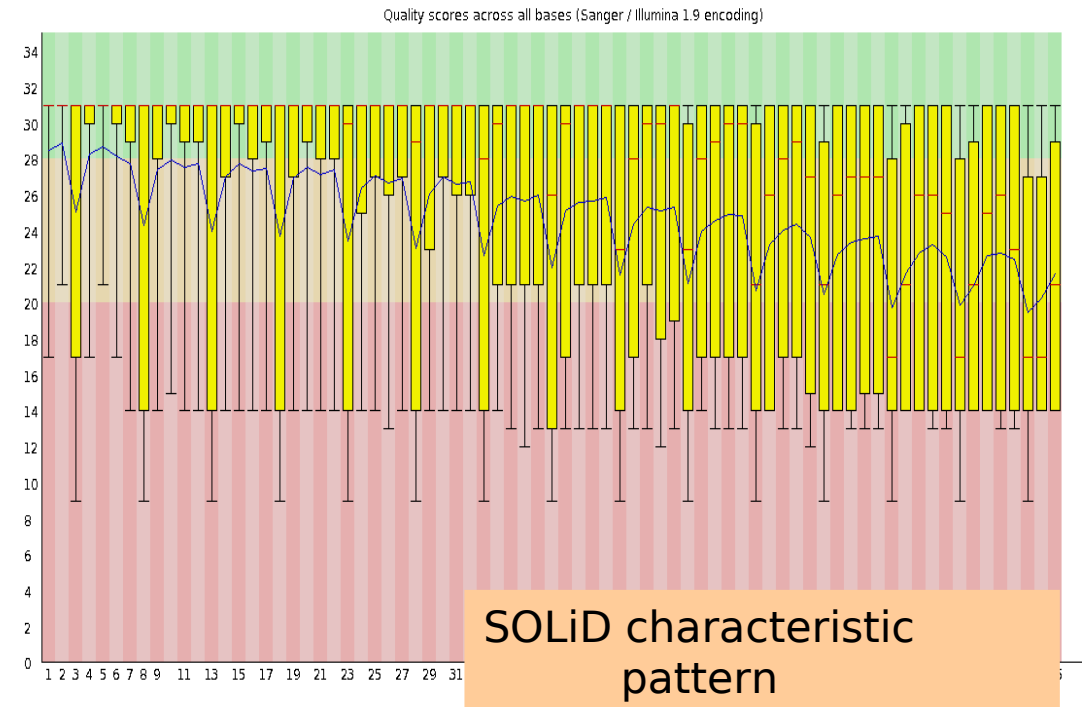
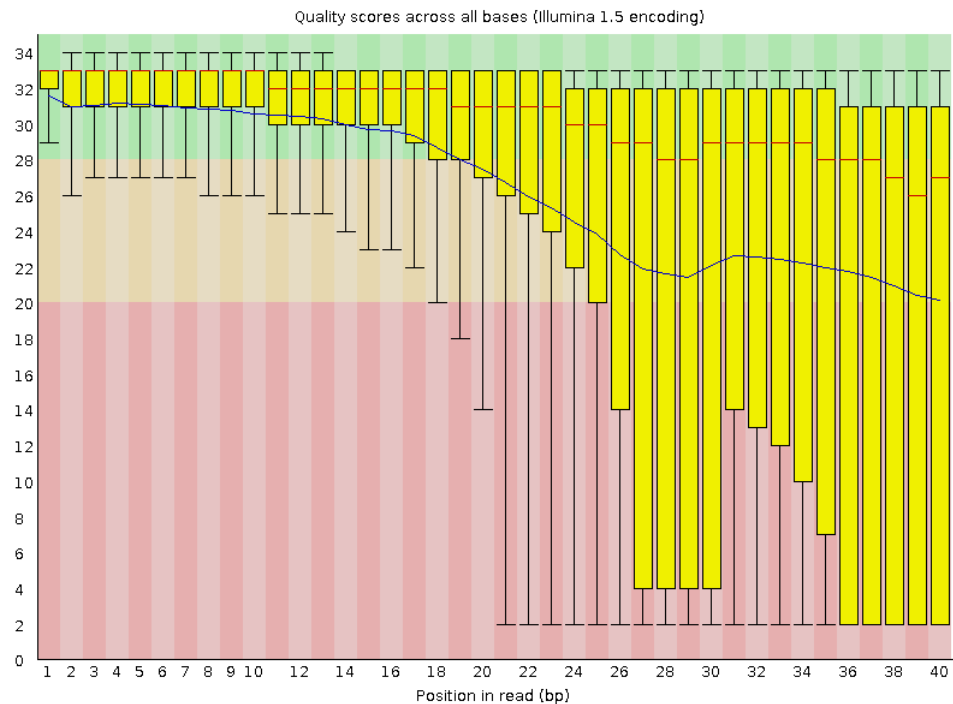


- Good data
- Consistent
- High quality along the read

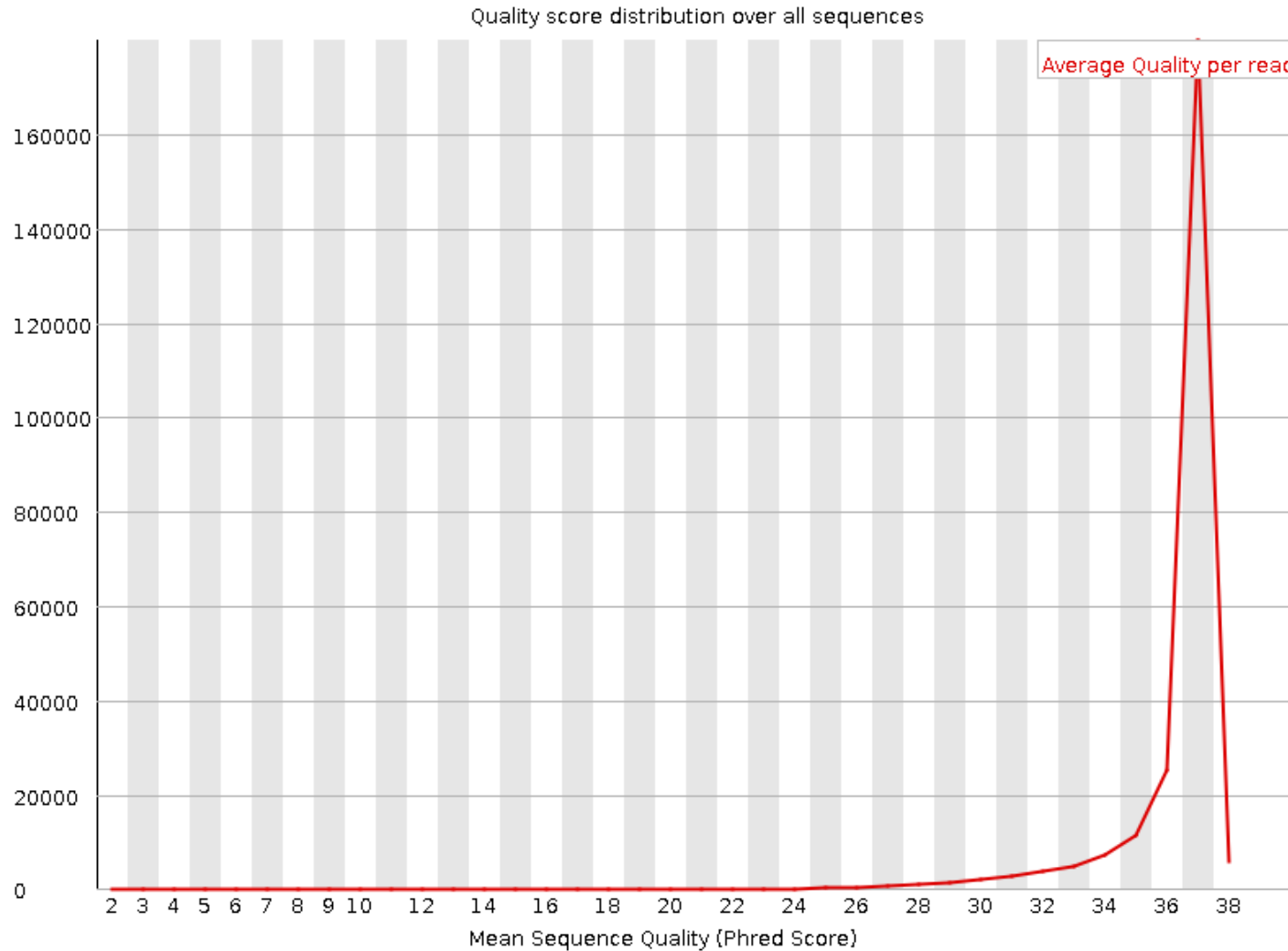
- The central red line is the median value
- The yellow box represents the inter-quartile range (25-75%)
- The upper and lower whiskers represent the 10% and 90% points
- The blue line represents the mean quality

# Sequence quality per base position

- Bad data
- High variance
- Quality decrease with length



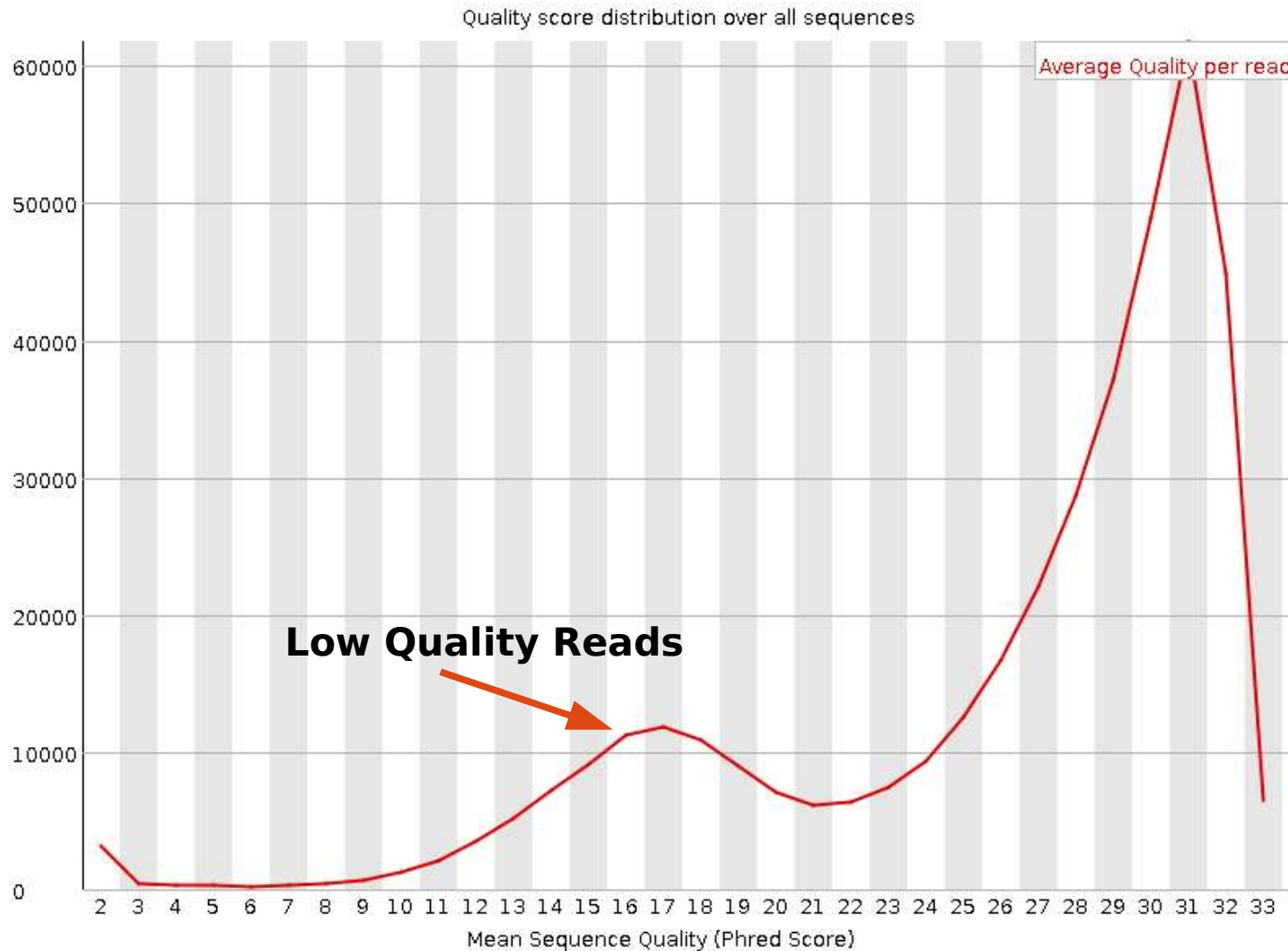
# Per sequence quality distribution



- Good data
- Most are high-quality sequences



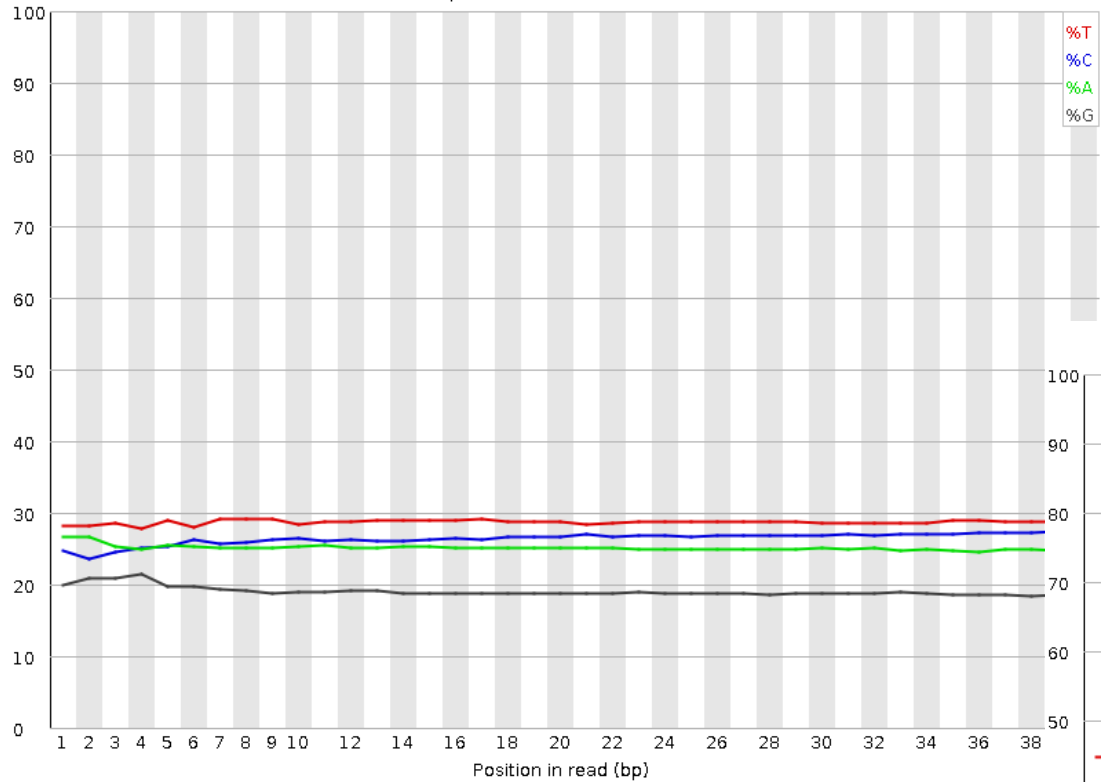
# Per sequence quality distribution



- Bad data
- Non-uniform distribution

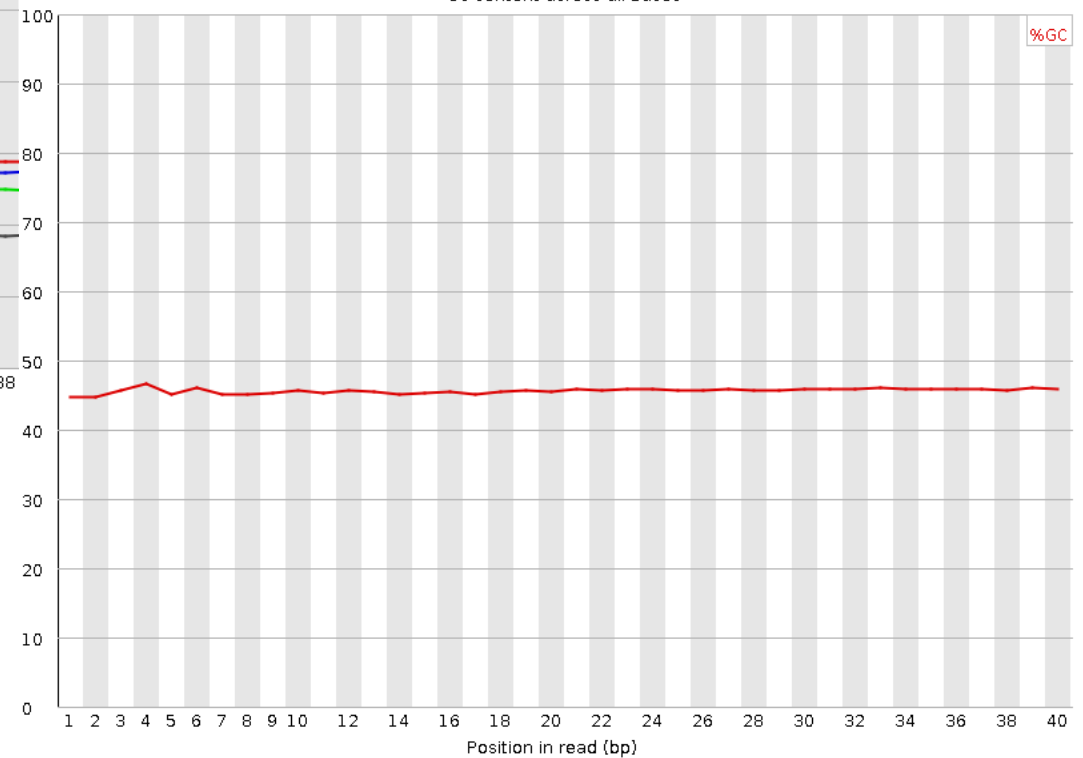
# Per base sequence content

Sequence content across all bases

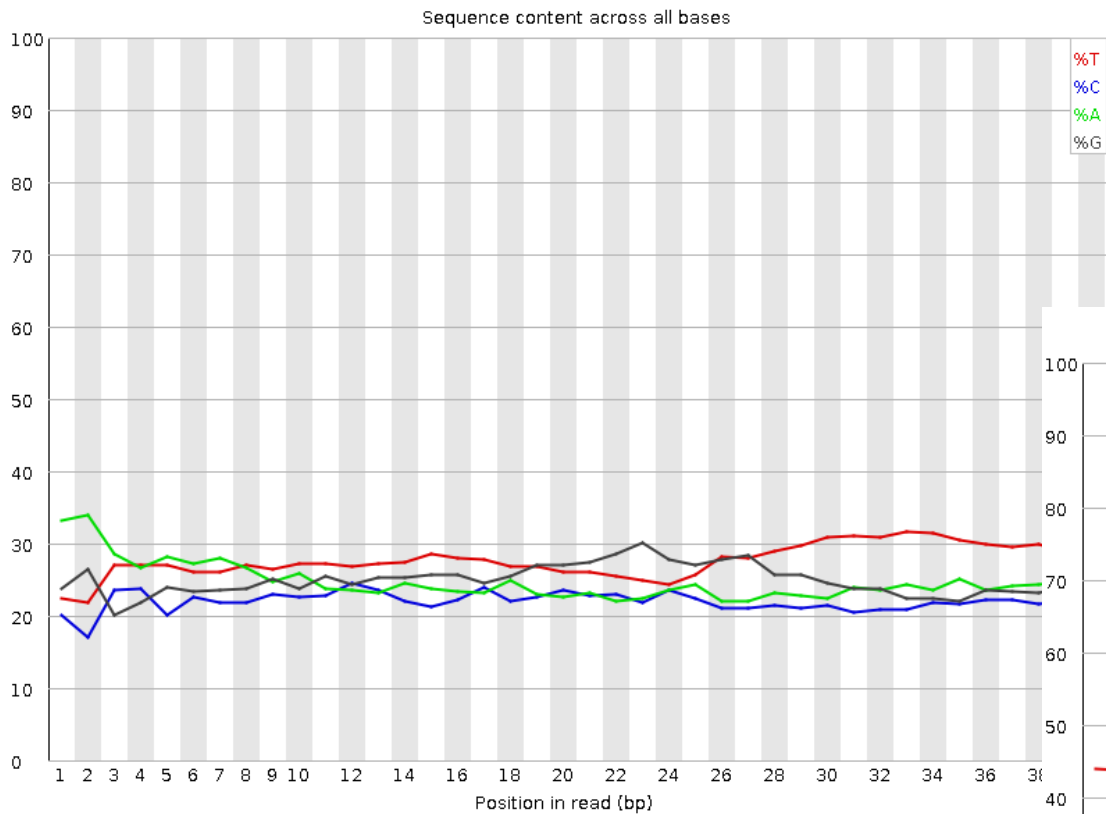


- Good data
  - Smooth over length
  - Organism dependent (GC)

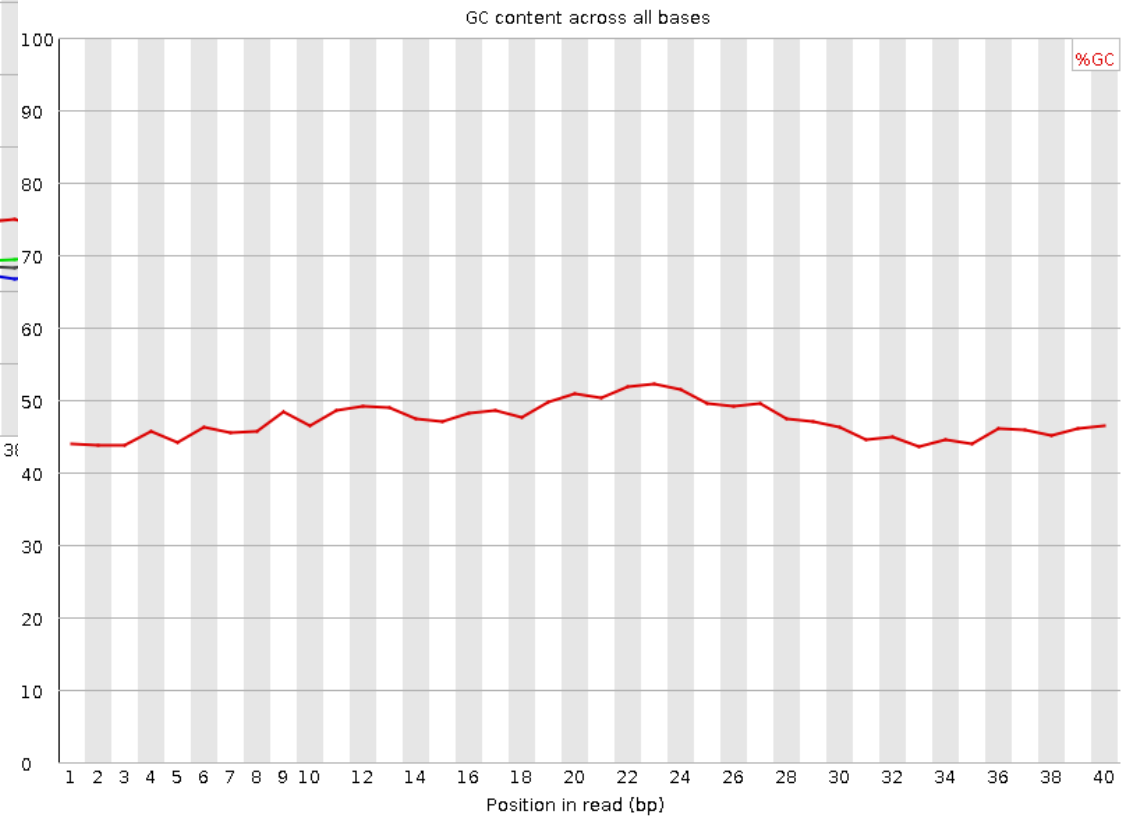
GC content across all bases



# Per base sequence content



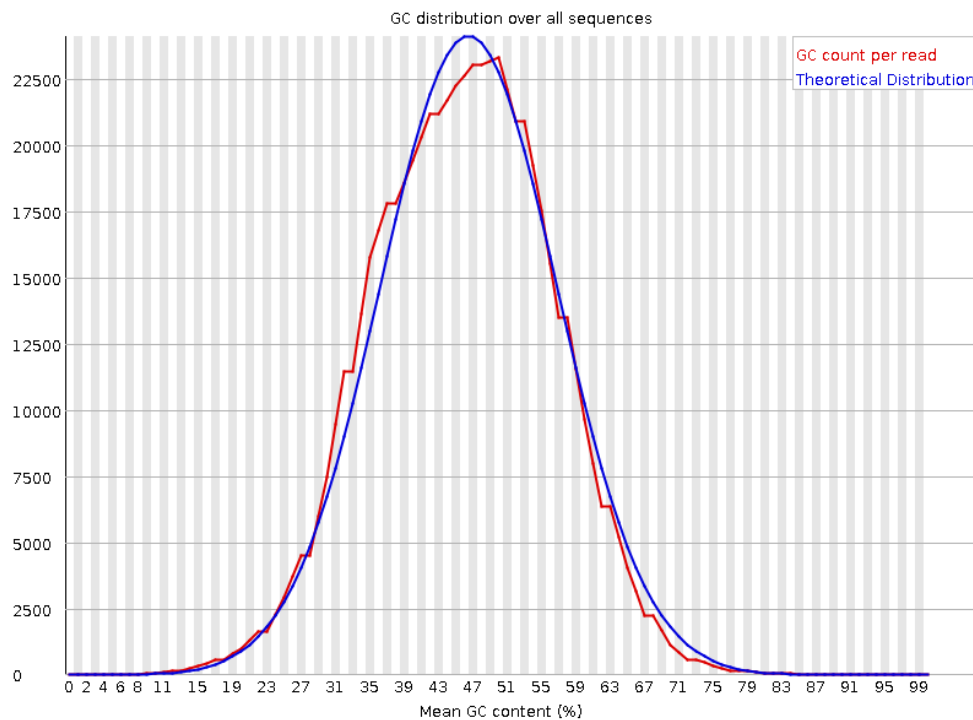
- Bad data
- Sequence position bias



# Per sequence GC content

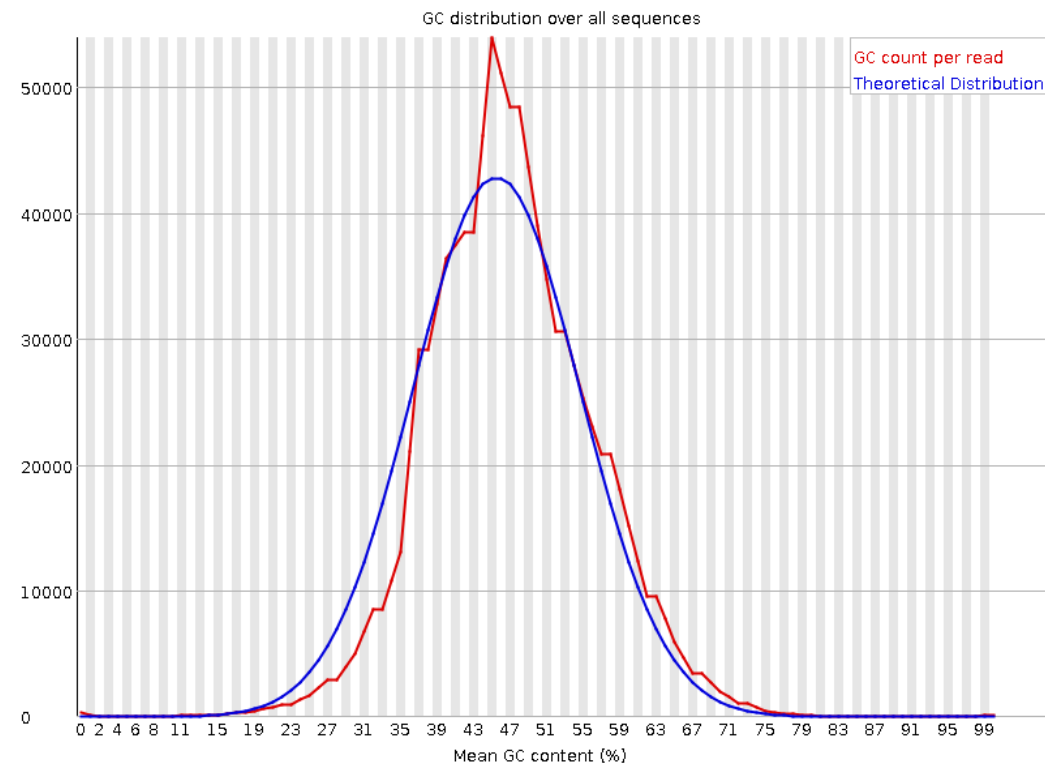
## Good data

- Fits with expected
- Organism dependent



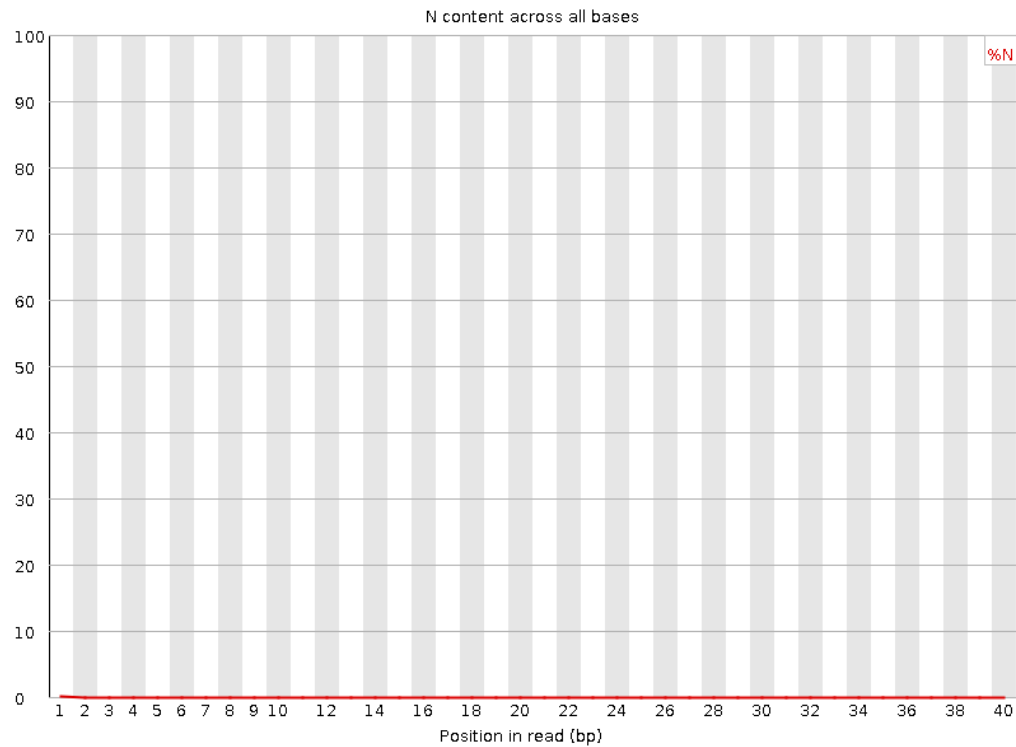
## Bad data

- Does not fit with expected
- Library contamination?

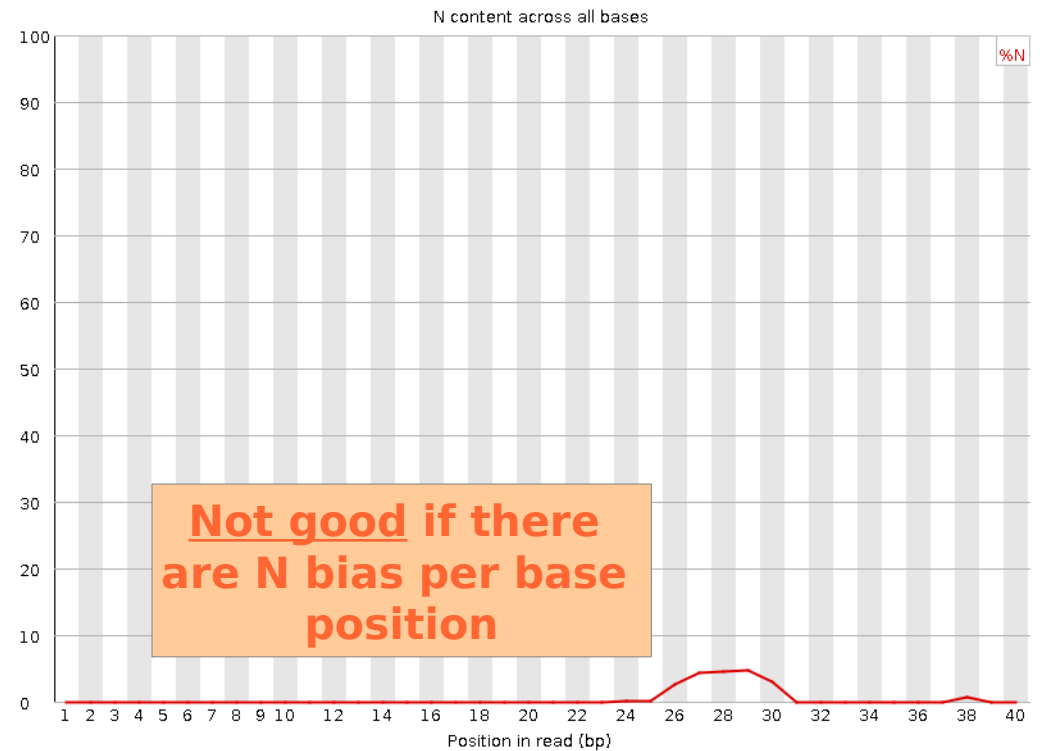


# Per base N content

□ Good data

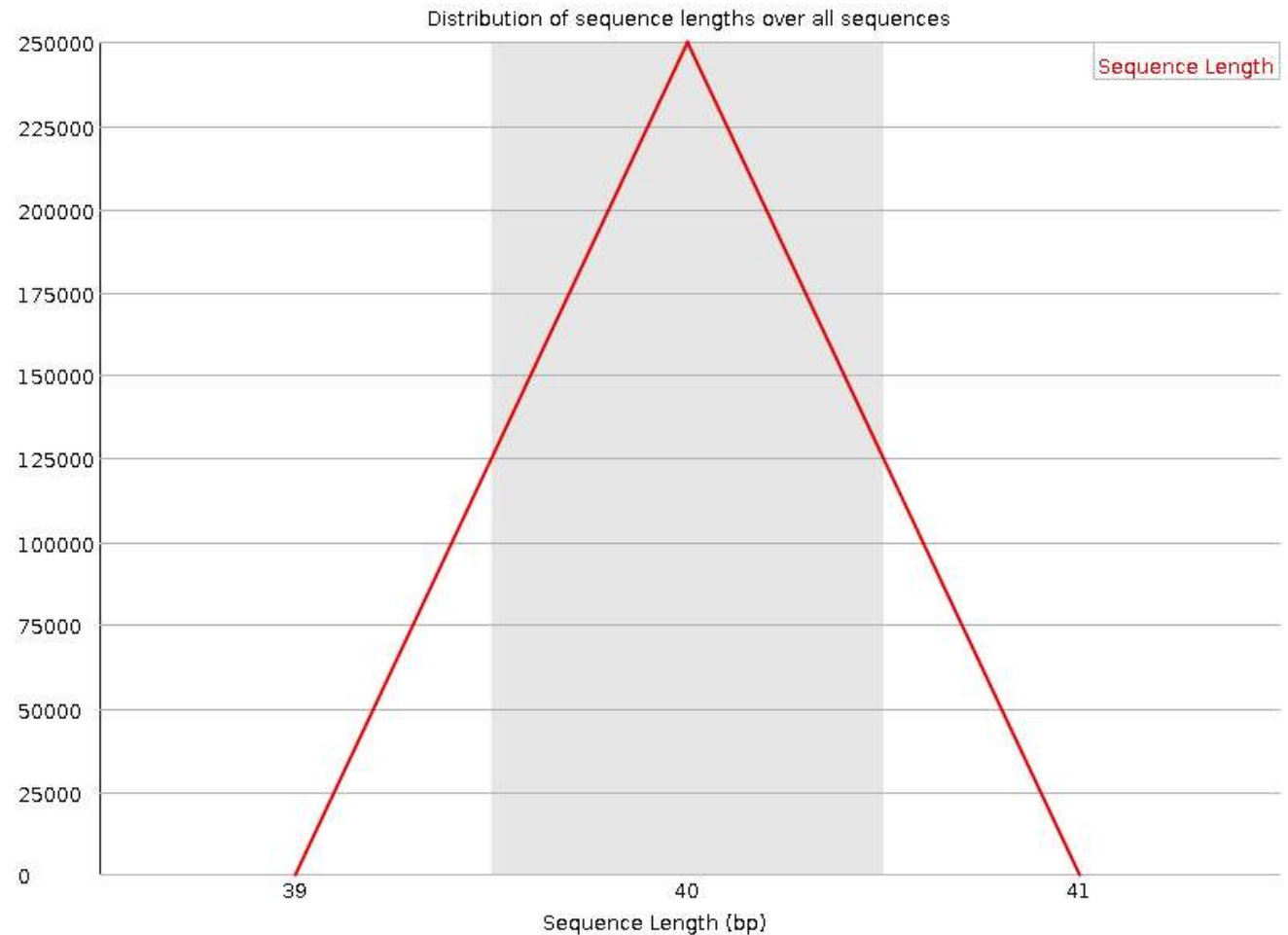


□ Bad data



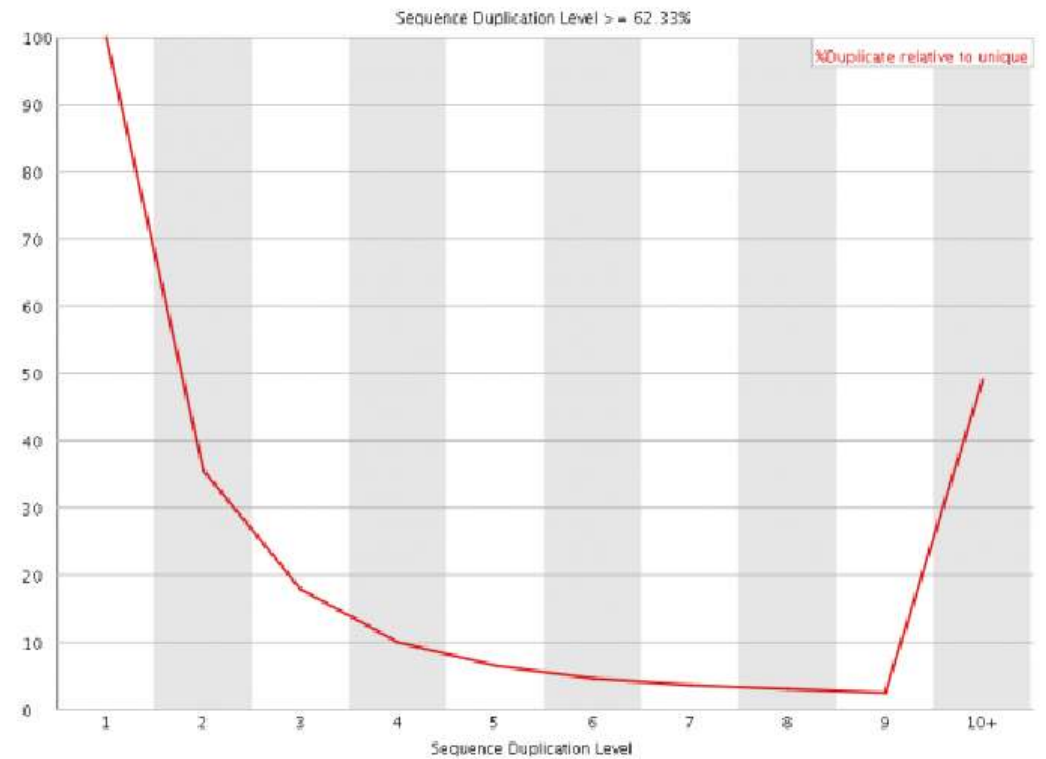
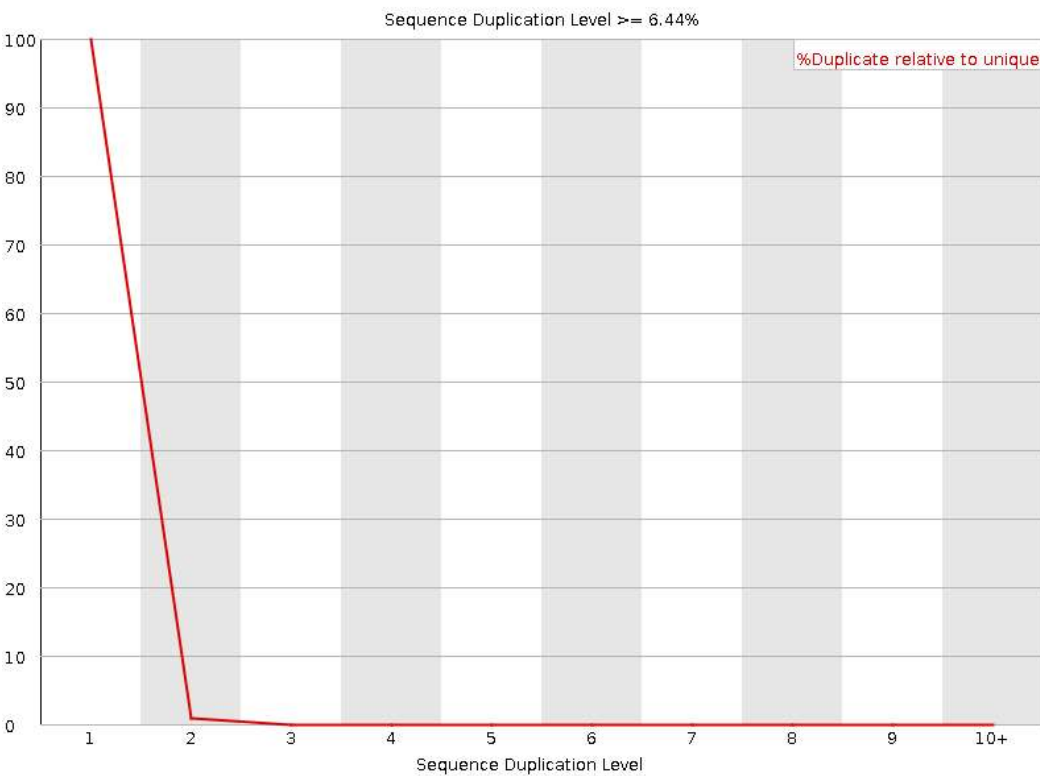
# Sequence length distribution

- Just descriptive:
  - Some sequencers output sequences of different length (e.g. 454)



# Sequence duplication levels

- In **transcriptomics**, you expect higher number of duplicated sequences.
- In **genomics** you should be worried if this happens → PCR artifact?



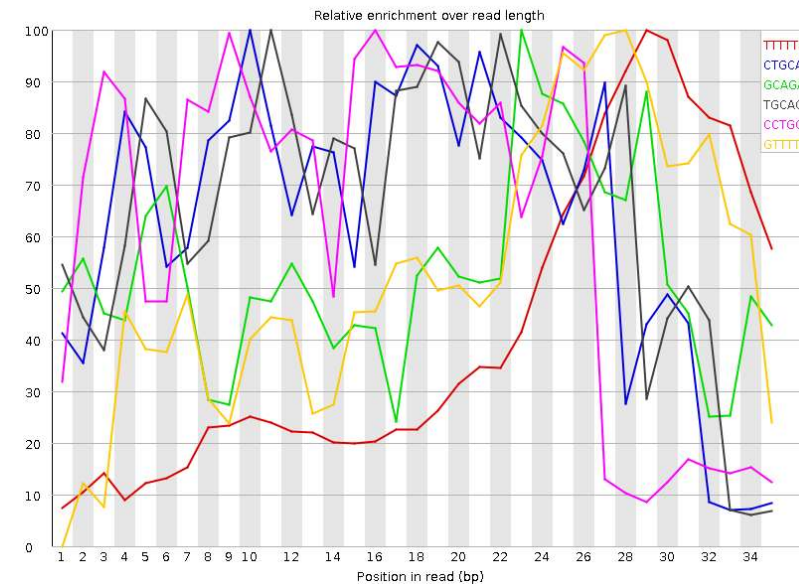
# Overrepresented sequences & Kmer content

- Question:
  - If we obtain the exact same sequences too many times  
→ **Do we have a problem?**

- Answer:
  - **Sometimes !**

Sequence	Count	Percentage
AGAGTTTTATCGCTTCCATGAC GCAGAAGTTAACACTTTC	2065	0.522403918155876 3
GATTGGCGTATCCAACCTGCA GAGTTTTATCGCTTCCATG	2047	0.517850276254275 4

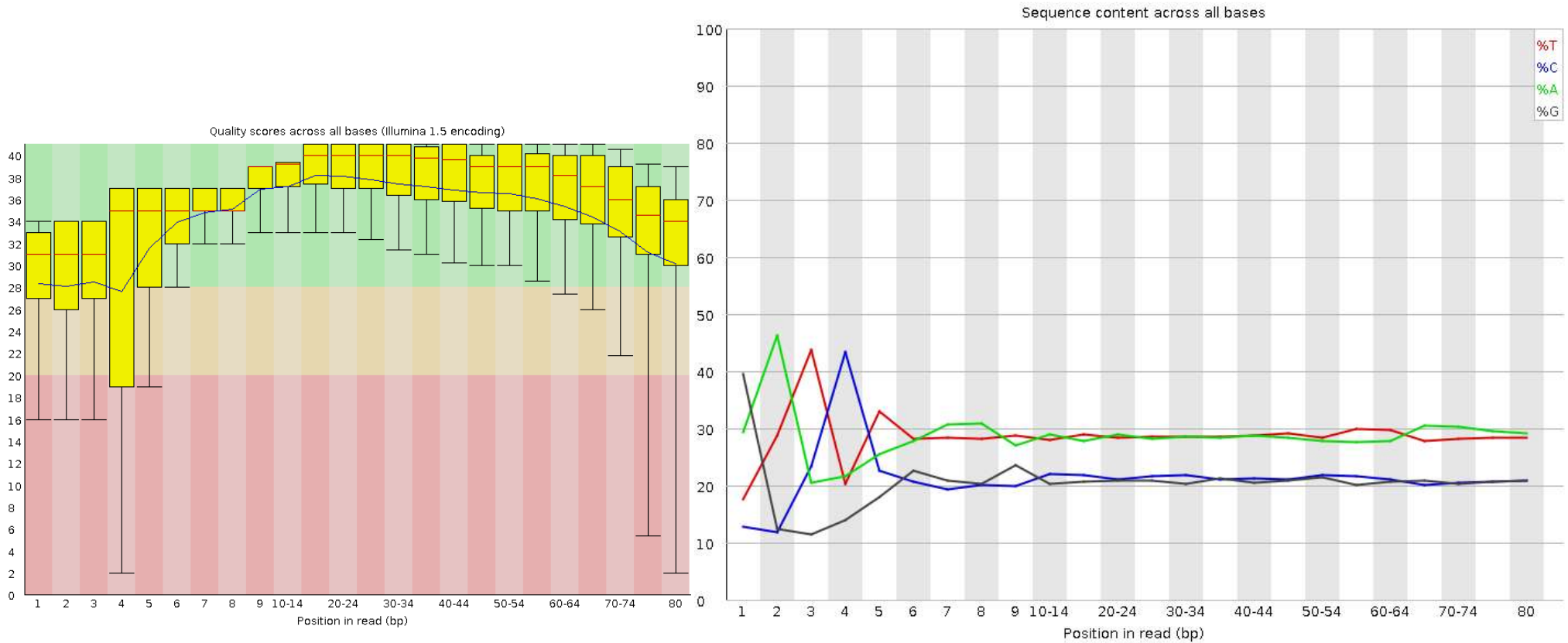
- Examples:
  - PCR primers, adapters ...





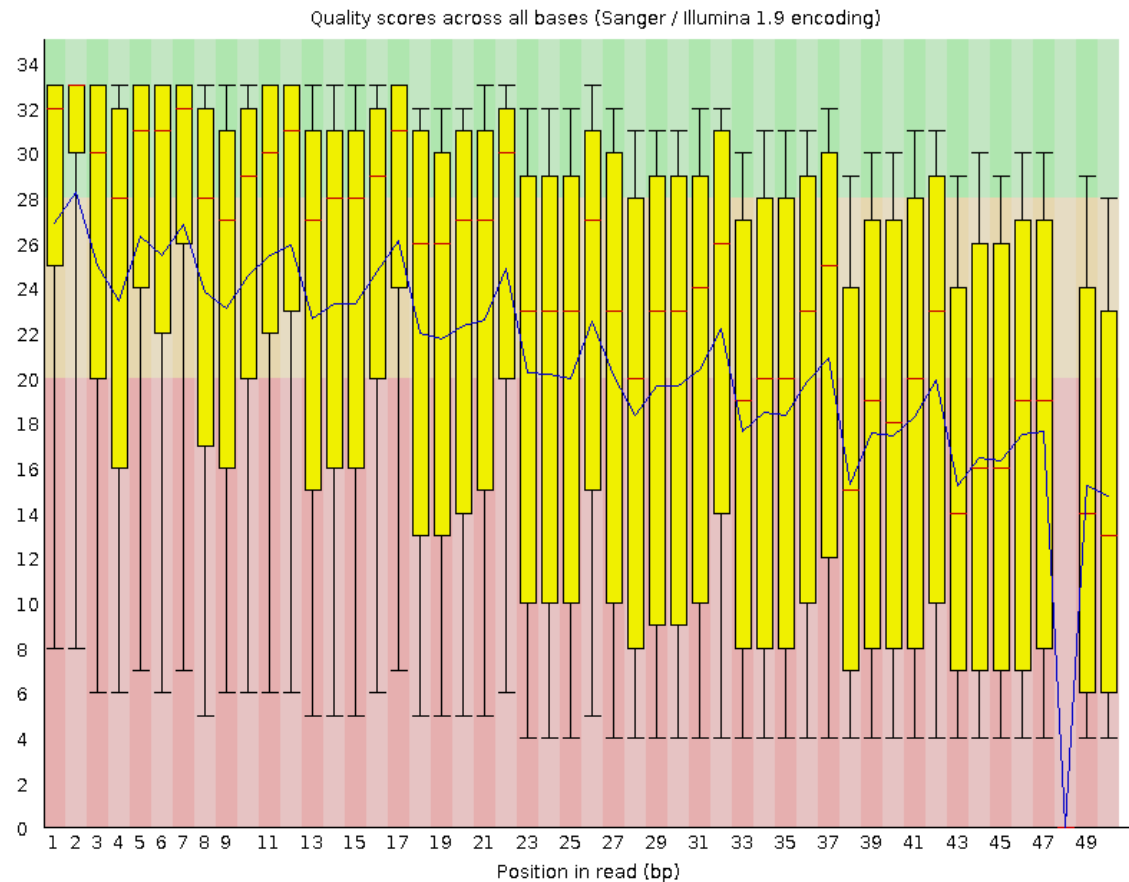
# Typical artifacts

## □ Sequence adapters

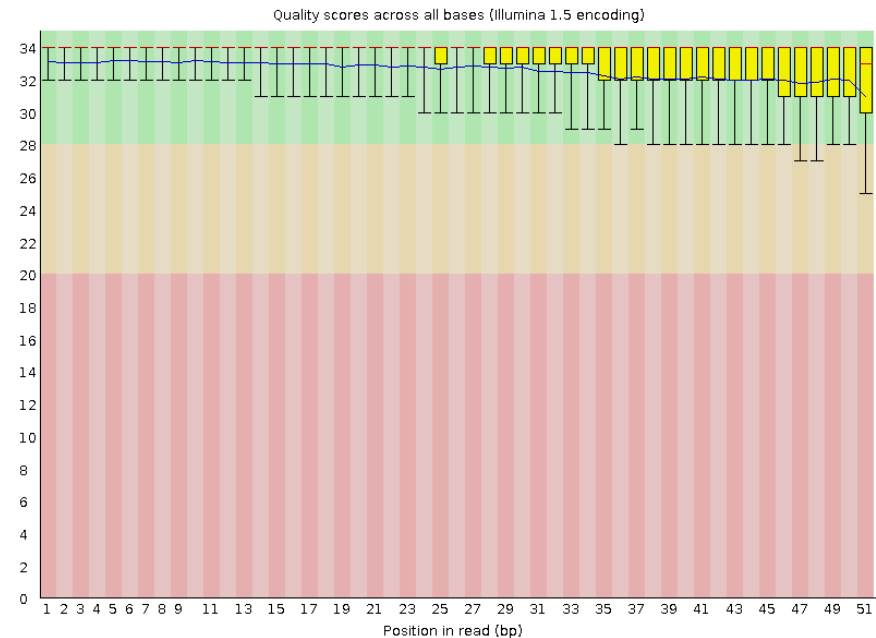
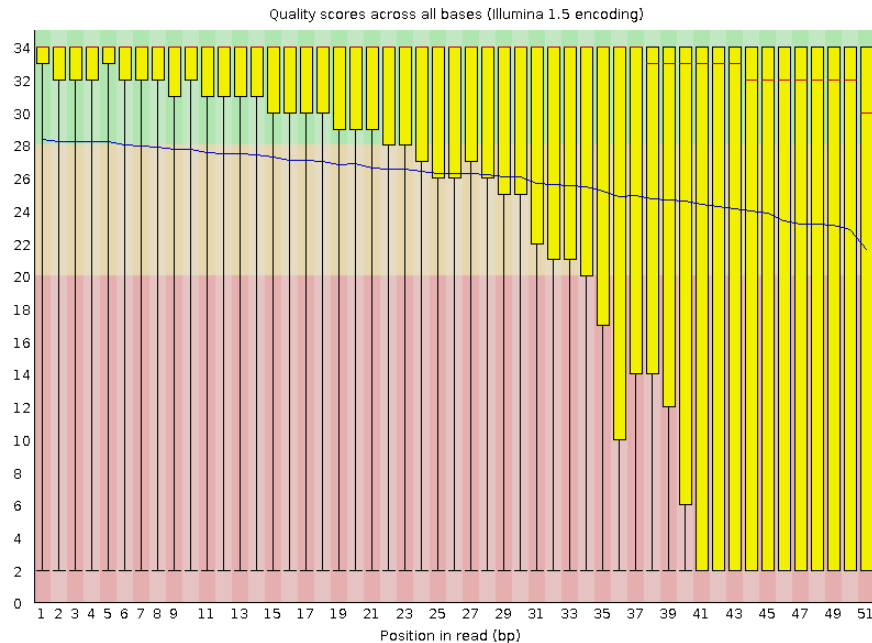


# Typical artifacts

- Platform dependent



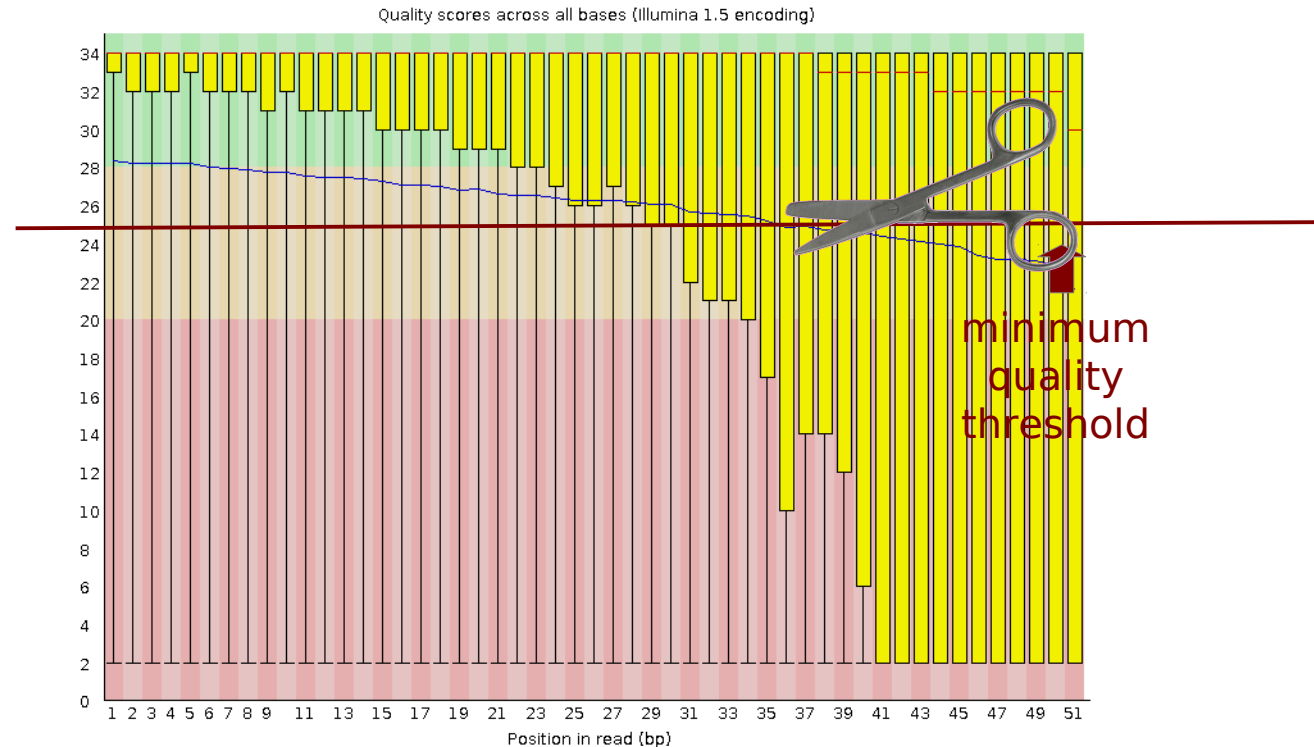
# Filtering & trimming



- Removing bad quality data will improve our confidence on downstream analysis

# Filtering & trimming

- Sequence filtering
  - Mean quality
  - Read length
  - Read length after trimming
  - Percentage of bases above Q
  - Adapter trimming
  - Adapter reads



# Filtering & trimming

- Sequence filtering tools
  - Fastx-toolkit
  - Galaxy (<https://main.g2.bx.psu.edu/>)
  - SeqTK (<https://github.com/lh3/seqtk>)
  - Cutadapt (<http://code.google.com/p/cutadapt/>)
  - And more....

# Practical: FastQC & Fastx-toolkit

---

- Use **FastQC** to see your starting state.
- Use **Fastx-toolkit** to optimize different datasets and then visualize the result with FastQC to prove your success!

**Hints:** Try trimming, clipping and quality filtering.

*Go to the tutorial and try the exercises...*