

# Understanding NGS raw data: FASTQ format, quality checking

Erika Souche

# Outline

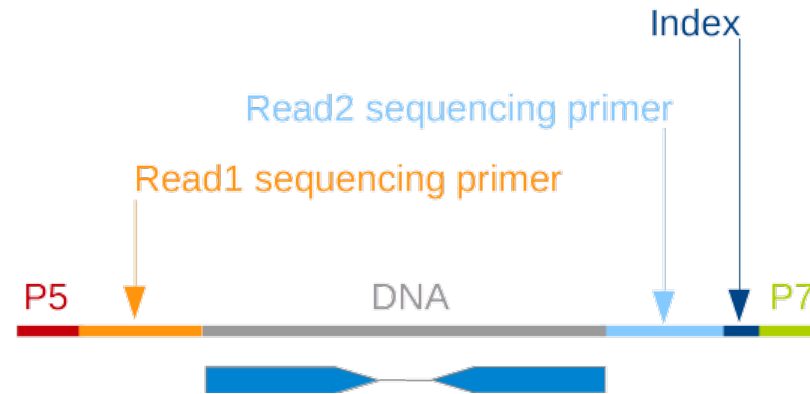
- Definitions
- FASTQ format
- Quality control
- Pre-processing

\*Focus on Illumina short read sequencing

# Outline

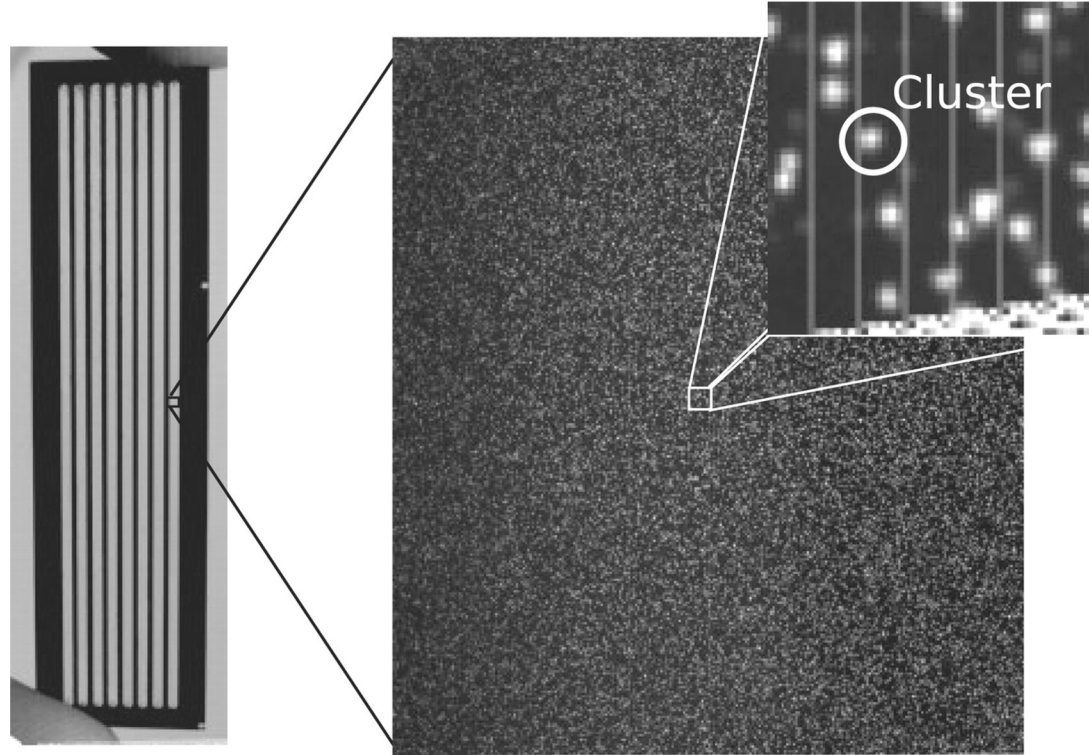
- Definitions
- FASTQ format
- Quality control
- Pre-processing

# Library



- DNA fragment
- Index/barcode
- Single End (SE) sequencing
- Paired End (PE) sequencing
- Insert size (=DNA fragment size)

# Illumina flowcell

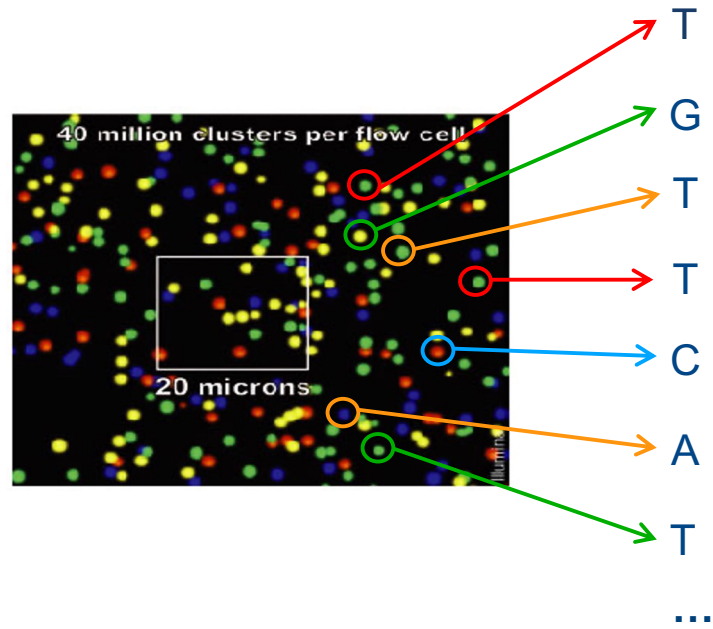


- Cluster

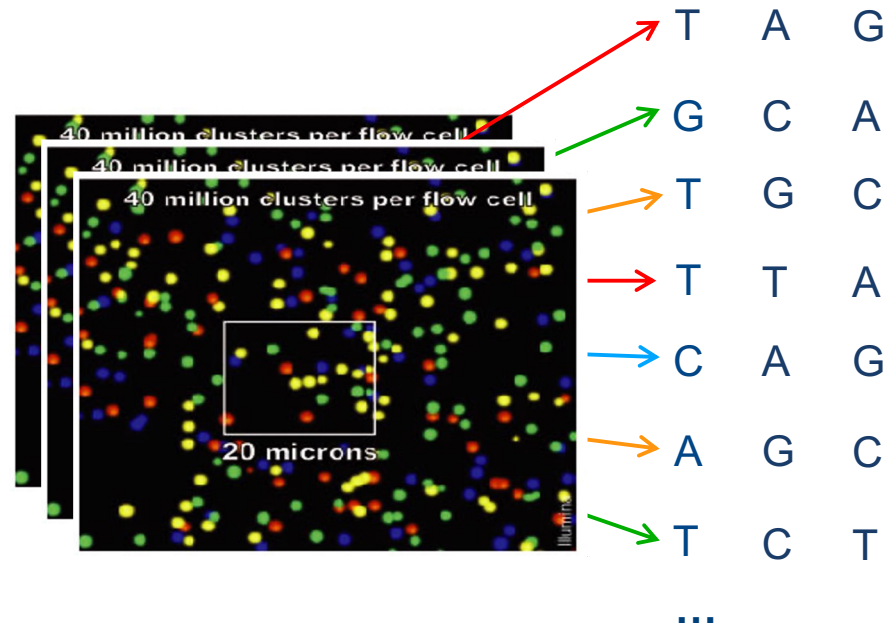
Bright spot on an image

Represents 1,000s of copies of the same DNA fragment

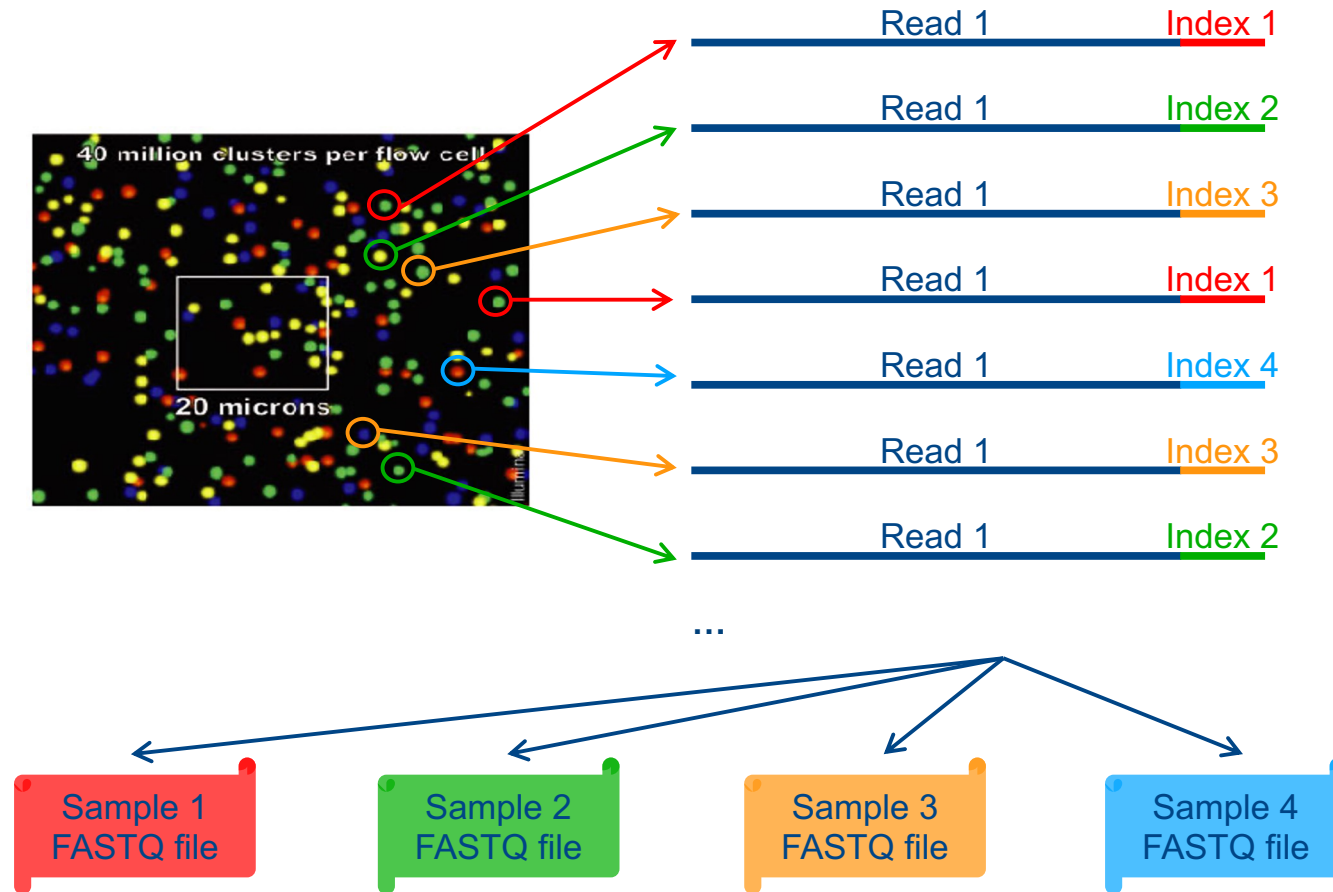
# Sequencing



# Sequencing



# Demultiplexing





# Outline

- Definitions
- **FASTQ format**
- Quality control
- Pre-processing

# FASTA file format

```
>1 dna:chromosome chromosome:GRCh38:1:1:248956422:1 REF
```

# Header 1

## Contig name

[illegible]

- Contig sequence

# FASTQ file format

The diagram illustrates the FASTQ file format structure. It shows three records, each consisting of four lines. The first record is highlighted with colored boxes: a green box for the header line, an orange box for the sequence line, a green box for the separator line, and a blue box for the quality line. Lines connect these boxes to labels on the right: 'Header 1' and 'Read name' for the first green box, 'Read sequence' for the orange box, 'Header 2' for the second green box, and 'Read quality' for the blue box. The second and third records are shown in plain text below the first.

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1069:110/1
GATCGCCATCCAGAAATCGGTGGCTTGCTCACTTTC
+
A>AAAAAAAAAAAAAAAAAAAA>>AA*:A>A:AAA3

@HWI-EAS337_4_FC_3038UAAXX:1:2:161:1891/1
GAACTGTTGCTGGAAGACTACAAAGCCTCCCTGAAA
+
A?A4?A>A>>A?<A>AA?AA=????;??6?????=<

@HWI-EAS337_4_FC_3038UAAXX:1:2:1041:124/1
GGTGATTTTCAGCCAACAGGCCAGATTGAGCGTAGCC
+
A:AAAAAAAAAAAA=AAA=AAAAAAAA<AA=AA9=

...
```

Header 1  
Read name  
Read sequence  
Header 2  
Read quality

# Read sequence & quality

@HWI-EAS337\_4\_FC\_3038UAAXX:1:2:1069:110/1

GATCGCCATCCAGAAATCGGTGGCTTGCTCACTTTC

Read sequence

+

A>AAAAAAAAAAAAAAAAAAAA>>AA\*:A>A:AAA3

Read quality

Probability of incorrect base call (P)

Phred Q =  $-10 \log_{10} P$

1 ASCII character per base

Phred Q =  $\text{dec}(\text{ASCII}) - 33$

## Regular ASCII Chart (character codes 0 - 127)

000	(nul)	016	(dle)	032	sp	048	0	064	@	080	P	096	`	112	p
001	@ (soh)	017	(dc1)	033	!	049	1	065	A	081	Q	097	a	113	q
002	(stx)	018	(dc2)	034	"	050	2	066	B	082	R	098	b	114	r
003	(etx)	019	(dc3)	035	#	051	3	067	C	083	S	099	c	115	s
004	(eot)	020	(dc4)	036	\$	052	4	068	D	084	T	100	d	116	t
005	(enq)	021	(nak)	037	%	053	5	069	E	085	U	101	e	117	u
006	(ack)	022	(syn)	038	&	054	6	070	F	086	V	102	f	118	v
007	(bel)	023	(etb)	039	'	055	7	071	G	087	W	103	g	119	w
008	(bs)	024	(can)	040	(	056	8	072	H	088	X	104	h	120	x
009	(tab)	025	(em)	041	)	057	9	073	I	089	Y	105	i	121	y
010	(lf)	026	(eof)	042	*	058	:	074	J	090	Z	106	j	122	z
011	(vt)	027	(esc)	043	+	059	;	075	K	091	[	107	k	123	{
012	(np)	028	(fs)	044	,	060	<	076	L	092	\	108	l	124	
013	(cr)	029	(gs)	045	-	061	=	077	M	093	]	109	m	125	}
014	(so)	030	(rs)	046	.	062	>	078	N	094	^	110	n	126	~
015	(si)	031	(us)	047	/	063	?	079	O	095	_	111	o	127	ó

# Read sequence & quality

@HWI-EAS337\_4\_FC\_3038UAAXX:1:2:1069:110/1

GATCGCCATCCAGAAATCGGTGGCTTGCTCACTTTC

Read sequence

+

A>AAAAAAAAAAAAAAAAAAAA>>AA\*:A>A:AAA3

Read quality

Probability of incorrect base call (P)

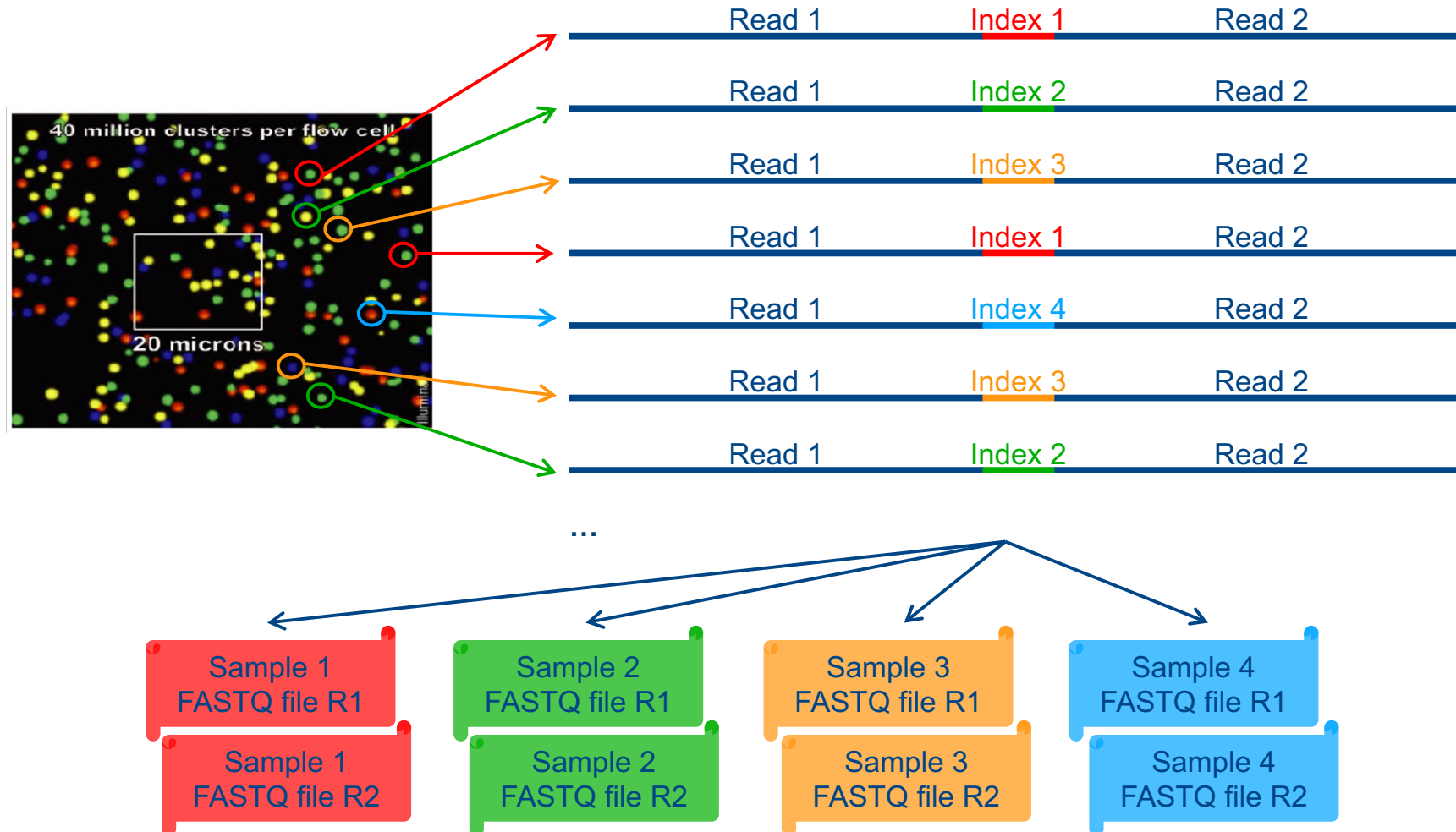
Phred Q =  $-10 \log_{10} P$

1 ASCII character per base

Phred Q =  $\text{dec}(\text{ASCII}) - 33$

Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

# Demultiplexing – PE sequencing



# FASTQ file format – PE sequencing

## Read 1

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1069:110/1
GATCGCCATCCAGAAATCGGTGGCTTGCTCACTTTC
+
A>AAAAAAAAAAAAAAAAAAAA>>AA*:A>A:AAA3
```

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:161:1891/1
GAACTGTTGCTGGAAGACTACAAAGCCTCCCTGAAA
+
A?A4?A>A>>A?<A>AA?AA=????;??6?????= <
```

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1041:124/1
GGTGATTTCAGCCAACAGGCCAGATTGAGCGTAGCC
+
A:AAAAAAAAAAAAAAAA=AAA=AAAAAAAA<AA=AA9=
```

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1114:113/1
GCTTTGTTTCGTGAAGCGTTATCGCTGGTGACATGGG
+
AAAAA?AAAAA?AA?A?AAAA>>AA>A:A>>7A>>A
```

...

## Read 2

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1069:110/2
AACGAAGACCGCGTCGTATTGTTCCAAAAGCGAATC
+
AAAAAAAAAAAAAAAAACAAAAAAAAAAAAAAAA>>>><
```

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:161:1891/2
CAATATTTTACGTTGCTAATGACAGTGAACAGACTT
+
AAAAA>;AA6AAAA?AAAA?A&AAA<?;AA><<5;<
```

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1041:124/2
CCGGAATTAAAGTCACCGTTGAGCATCCGGATAAAC
+
AAAA;AAAAAAAAAAAA?AAA=A?AAAAAAA6=>=><
```

```
@HWI-EAS337_4_FC_3038UAAXX:1:2:1114:113/2
CCGAAACAGACGCCAGCACCCGATCGGTGCCTGAC
+
AAAAA>AAAAAAAAAAAAAAAAAAAAAAAA(A==6=6<
```

...

# Demultiplexing statistics

RUN	MACHINE	FLOWCELL	LANE	PROJECT	SAMPLE	BARCODE	CLUSTER.COUNT	BARCODE.COUNT	BARCODE.0.MISMATCH.COUNT	BARCODE.1.MISMATCH.COUNT	BARCODE.0.MISMATCH.PERC	BARCODE.1.MISMATCH.PERC	CLUSTER.COUNT	BASECALL.COUNT	YIELD.Q30.SUM	QSUM.SUM	LL.QUAL.AVG	BASECALL.Q30.PERC	CLUSTER.COUNT
191004	HiSeq2000	FCA	1	Project	GC085522	CGCATGAT+TCAGGCTT	5,552,039	5,552,039	5,552,039	NA	1	NA	151,682,423	1,121,511,878	1,076,542,969	42,813,220,820	38.17	0.96	0.04
191004	HiSeq2000	FCA	1	Project	GC085522_ID	CTGAAGCT+TATAGCCT	9,777	9,777	9,777	NA	1	NA	151,682,423	1,974,954	1,913,612	75,745,127	38.35	0.97	0
191004	HiSeq2000	FCA	1	Project	GC085523	CTTAGGAC+GTAGGAGT	5,814,208	5,814,208	5,814,208	NA	1	NA	151,682,423	1,174,470,016	1,129,817,062	44,881,213,892	38.21	0.96	0.04
191004	HiSeq2000	FCA	1	Project	GC085523_ID	CTGAAGCT+ATAGAGGC	17,079	17,079	17,079	NA	1	NA	151,682,423	3,449,958	3,316,987	131,549,660	38.13	0.96	0
191004	HiSeq2000	FCA	1	Project	GC085524	ATCCGGTA+TATCGGTC	6,290,933	6,290,933	6,290,933	NA	1	NA	151,682,423	1,270,768,466	1,221,984,994	48,547,800,901	38.20	0.96	0.04
191004	HiSeq2000	FCA	1	Project	GC085524_ID	CTGAAGCT+CCTATCCT	14,085	14,085	14,085	NA	1	NA	151,682,423	2,845,170	2,757,522	109,110,362	38.35	0.97	0
191004	HiSeq2000	FCA	1	default	Undetermined	unknown	11,130,696	712,364,544	712,364,544	NA	1	NA	151,682,423	2,248,400,592	2,030,960,115	82,104,595,067	36.52	0.9	0.07

SAMPLE	BARCODE	CLUSTER.COUNT	BARCODE.COUNT	T	BARCODE.0.MISMATCH.COUNT	BARCODE.1.MISMATCH.COUNT	BARCODE.0.MISMATCH.PERC	BARCODE.1.MISMATCH.PERC
GC085522	CGCATGAT+TCAGGCTT	5,552,039	5,552,039	5,552,039	NA	1	NA	NA
GC085522_ID	CTGAAGCT+TATAGCCT	9,777	9,777	9,777	NA	1	NA	NA
GC085523	CTTAGGAC+GTAGGAGT	5,814,208	5,814,208	5,814,208	NA	1	NA	NA
GC085523_ID	CTGAAGCT+ATAGAGGC	17,079	17,079	17,079	NA	1	NA	NA
GC085524	ATCCGGTA+TATCGGTC	6,290,933	6,290,933	6,290,933	NA	1	NA	NA
GC085524_ID	CTGAAGCT+CCTATCCT	14,085	14,085	14,085	NA	1	NA	NA
Undetermined	unknown	11,130,696	712,364,544	712,364,544	NA	1	NA	NA



# Demultiplexing statistics

RUN	MACHINE	FLOWC	LANE	PROJECT	SAMPLE	BARCODE	BARCODE.0.		BARCODE.1		BARCODE.		BARCODE.		BASECA		BASECALL.Q		CLUSTER.	
							CLUSTER.CO	BARCODE.CO	MISMATCH.C	MISMATCH.0	MISMAT	1.MISMAT	CLUSTER.COU	BASECALL.COU	LL.QUAL	UAL.ABOVE	COUNT.PE			
							UNT	UNT	OUNT	.COUNT	CH.PERC	CH.PERC	NT.SUM	LANE.NT	YIELD.Q30.SUM	QSUM.SUM	.AVG	30.PERC	RC	
191004	HiSeq2000	FCA	1	Project	GC085522	CGCATGAT+TCAGGCTT	5,552,039	5,552,039	5,552,039	NA	1	NA	151,682,423	1,121,511,878	1,076,542,969	42,813,220,820	38.17	0.96	0.04	
191004	HiSeq2000	FCA	1	Project	GC085522_ID	CTGAAGCT+TATAGCCT	9,777	9,777	9,777	NA	1	NA	151,682,423	1,974,954	1,913,612	75,745,127	38.35	0.97	0	
191004	HiSeq2000	FCA	1	Project	GC085523	CTTAGGAC+GTAGGAGT	5,814,208	5,814,208	5,814,208	NA	1	NA	151,682,423	1,174,470,016	1,129,817,062	44,881,213,892	38.21	0.96	0.04	
191004	HiSeq2000	FCA	1	Project	GC085523_ID	CTGAAGCT+ATAGAGGC	17,079	17,079	17,079	NA	1	NA	151,682,423	3,449,958	3,316,987	131,549,660	38.13	0.96	0	
191004	HiSeq2000	FCA	1	Project	GC085524	ATCCGGTA+TATCGGTC	6,290,933	6,290,933	6,290,933	NA	1	NA	151,682,423	1,270,768,466	1,221,984,994	48,547,800,901	38.20	0.96	0.04	
191004	HiSeq2000	FCA	1	Project	GC085524_ID	CTGAAGCT+CCTATCCT	14,085	14,085	14,085	NA	1	NA	151,682,423	2,845,170	2,757,522	109,110,362	38.35	0.97	0	
191004	HiSeq2000	FCA	1	default	Undetermined	unknown	11,130,696	712,364,544	712,364,544	unknown	1	NA	151,682,423	2,248,400,592	2,030,960,115	82,104,595,067	36.52	0.9	0.07	

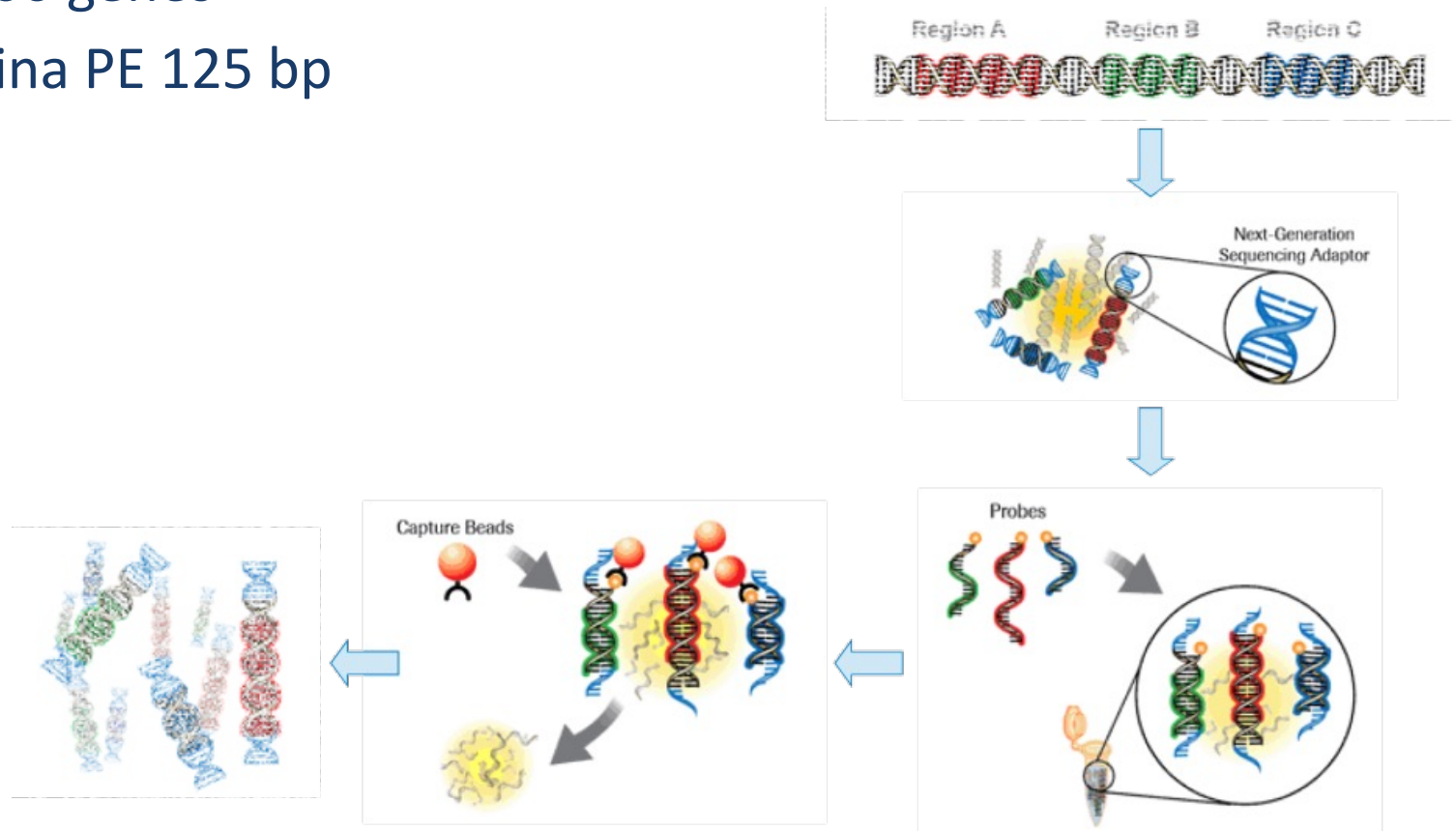
SAMPLE	CLUSTER.COUNT	BASECALL.COUNT	YIELD.Q30.SUM	QSUM.SUM	BASECALL.QUAL.AVG	BASECALL.QUAL.ABOVE.30.PERC	CLUSTER.COUNT.PERC
GC085522	151,682,423	1,121,511,878	1,076,542,969	42,813,220,820	38.17	0.96	0.04
GC085522_ID	151,682,423	1,974,954	1,913,612	75,745,127	38.35	0.97	0
GC085523	151,682,423	1,174,470,016	1,129,817,062	44,881,213,892	38.21	0.96	0.04
GC085523_ID	151,682,423	3,449,958	3,316,987	131,549,660	38.13	0.96	0
GC085524	151,682,423	1,270,768,466	1,221,984,994	48,547,800,901	38.20	0.96	0.04
GC085524_ID	151,682,423	2,845,170	2,757,522	109,110,362	38.35	0.97	0
Undetermined	151,682,423	2,248,400,592	2,030,960,115	82,104,595,067	36.52	0.9	0.07

# Outline

- Definitions
- FASTQ format
- **Quality control**
- Pre-processing

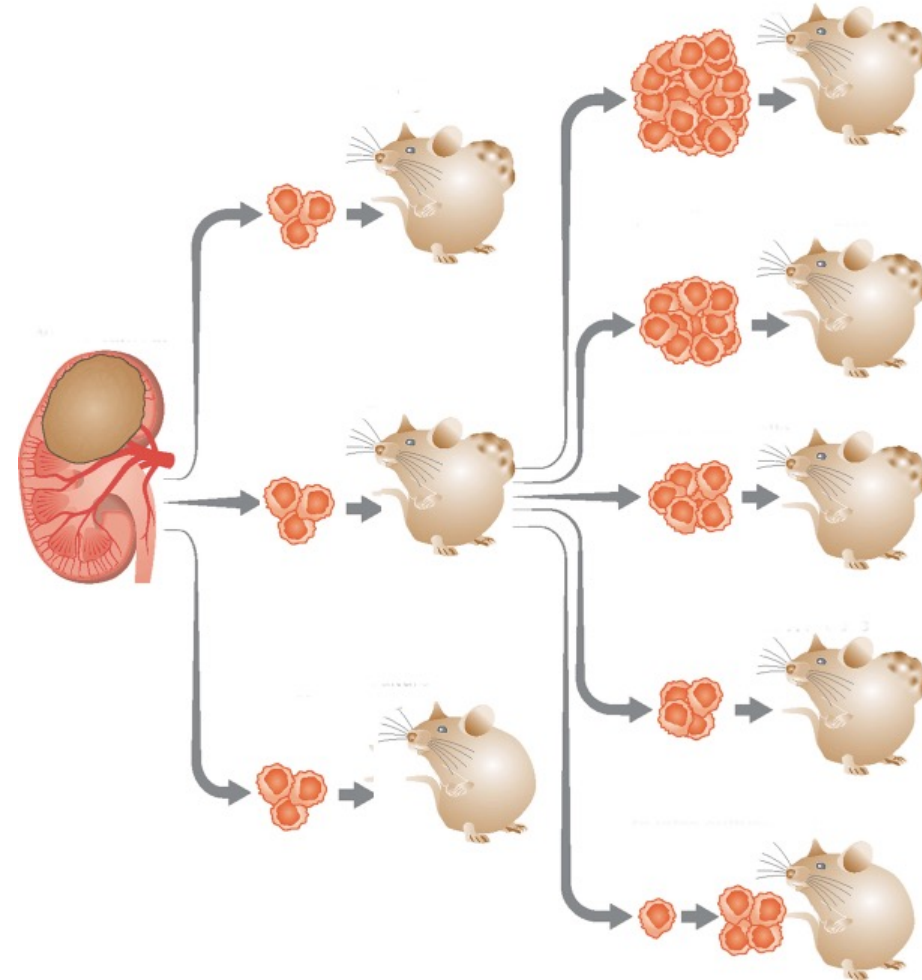
# Example 1 – targeted capture sequencing

- Finding the genetic cause of a disease
  - ~ 6,000 genes
  - Illumina PE 125 bp



# Example 2 – targeted amplicon sequencing

- Fingerprinting of xenografts
  - 31 SNPs
  - Illumina PE 150 bp

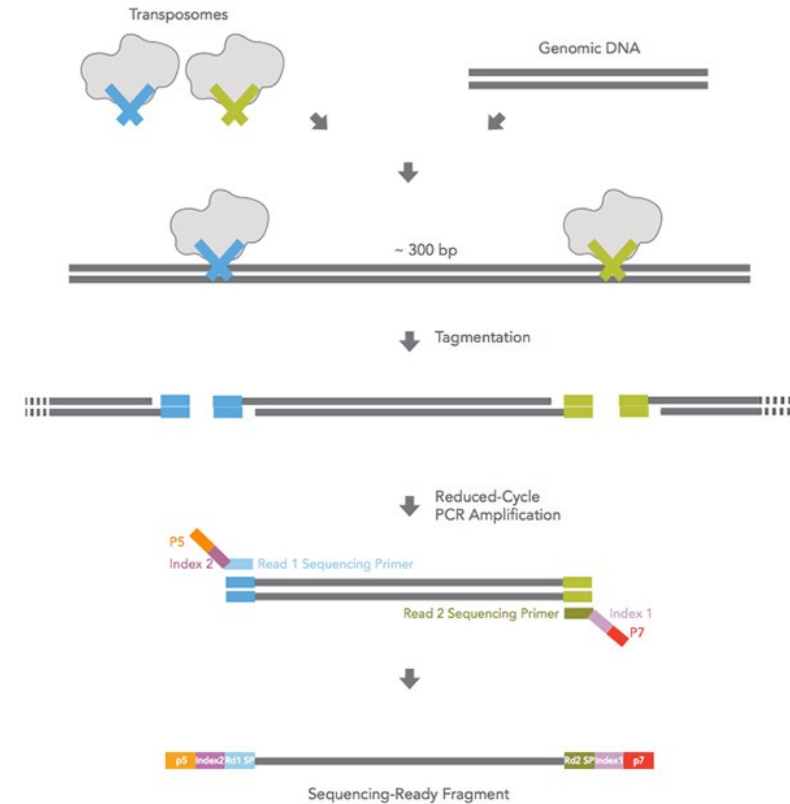


Trace platform, <http://www.uzleuven-kuleuven.be/lki/trace/>

Figure adapted from Peter Hohenstein, EMBO Molecular Medicine: 5 (1), 2013

# Example 3 – whole genome sequencing

- Predicting bacterial resistance
  - Whole Genome Sequencing (WGS)
    - *Mycobacterium tuberculosis* (BWGS)
  - Illumina PE 300 bp



# Need for quality control

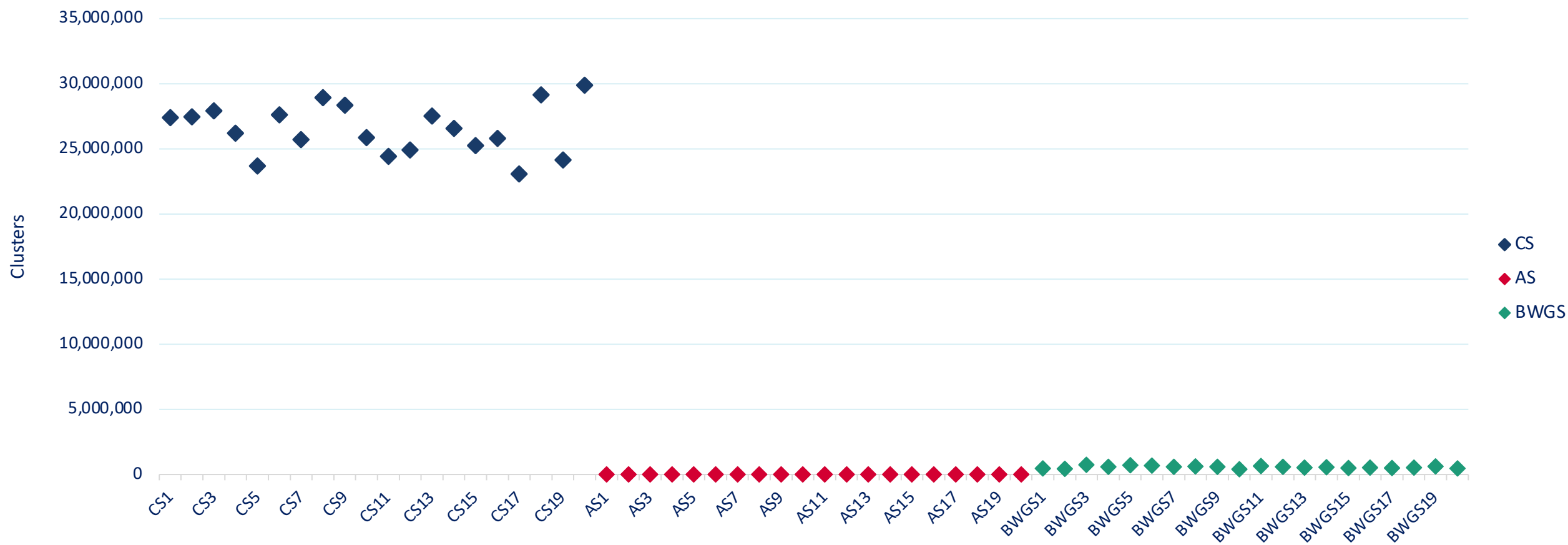
- Complex workflow
- Cost efficient assay
- Ensure sequenced data is OK
  - Enough data ?
  - Correct content ?
  - Good quality data ?

# Number of reads

Capture Sequencing (CS)  
18,013,929 bp

Amplicon Sequencing (AS)  
1,835 bp

Bacterial WGS (BWGS)  
2,221,315 bp

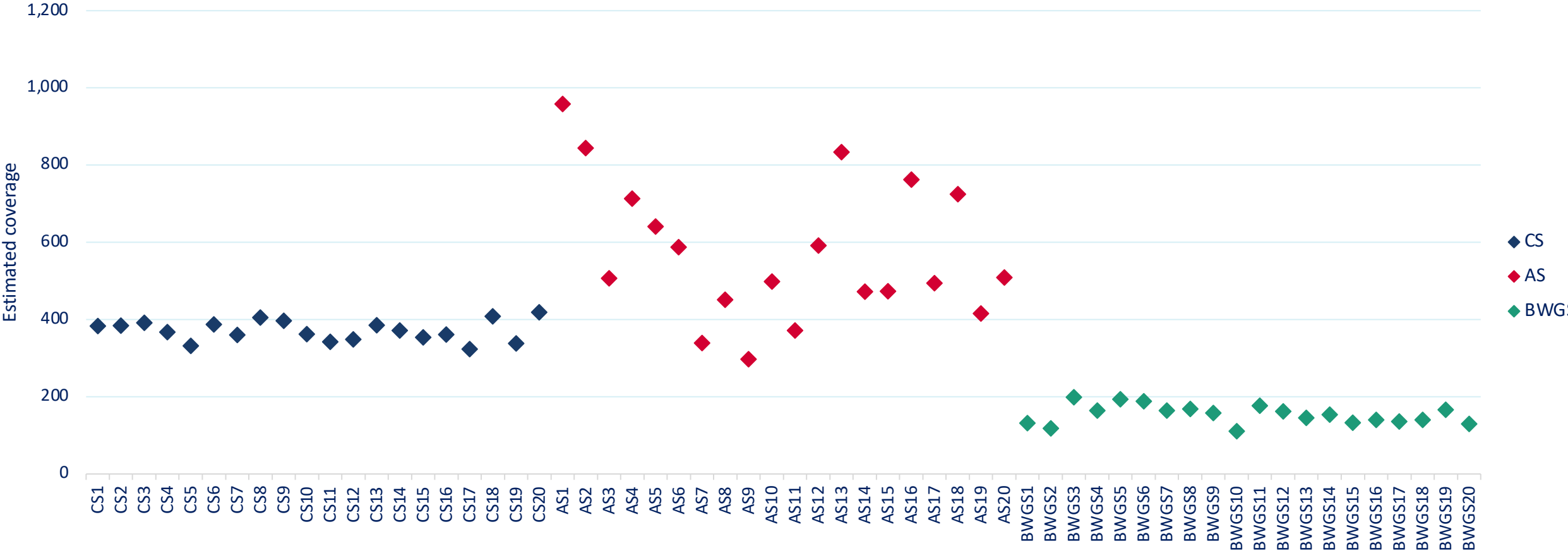


# Estimated mean coverage

Capture Sequencing (CS)  
18,013,929 bp

Amplicon Sequencing (AS)  
1,835 bp

Bacterial WGS (BWGS)  
2,221,315 bp





# FASTQ files – QC

- FastQC
  - Check Phred quality scores
  - Check GC content
  - Check read content
  - ...

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

# FastQC

- Summary report

```
fastqc -o result sample.R1.fastq.gz  
fastqc -o result sample.R2.fastq.gz
```

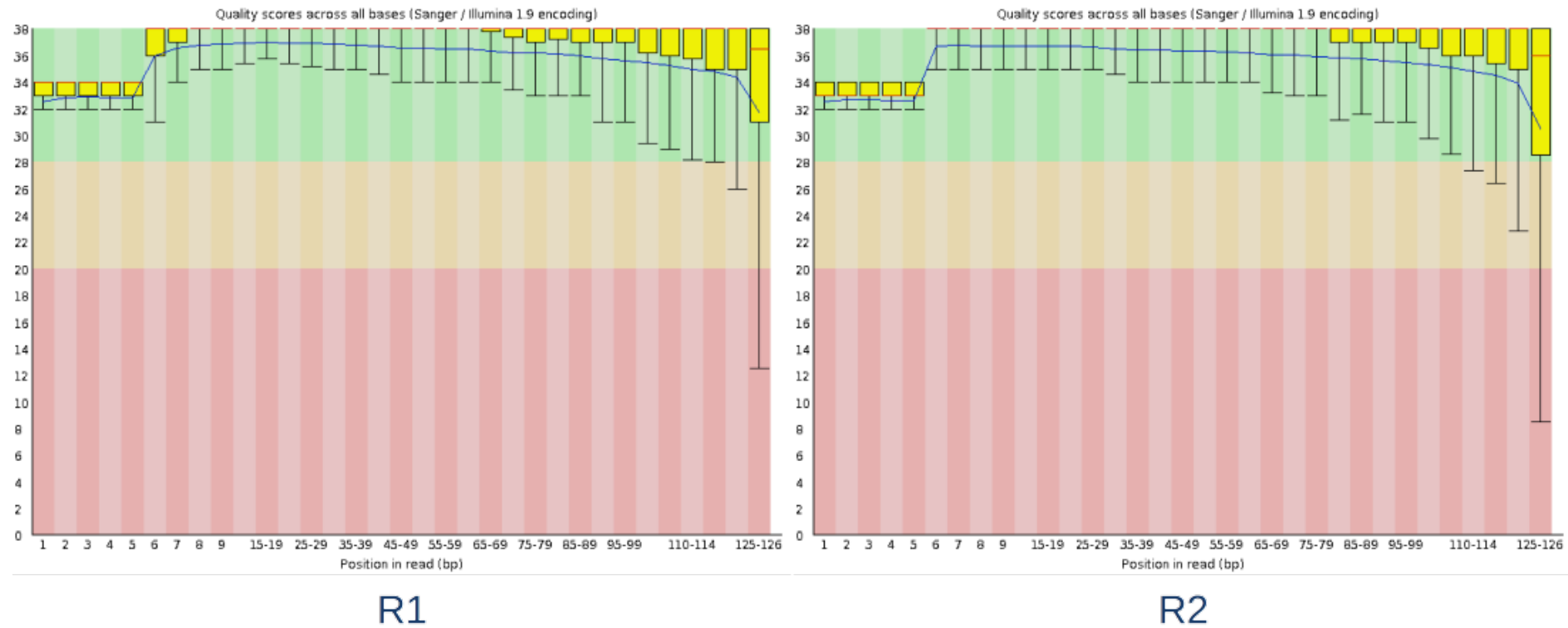
## FastQC Report

### Summary

- ✓ [Basic Statistics](#)
- ✓ [Per base sequence quality](#)
- ✓ [Per sequence quality scores](#)
- ✓ [Per base sequence content](#)
- ✓ [Per base GC content](#)
- ✗ [Per sequence GC content](#)
- ✓ [Per base N content](#)
- ✓ [Sequence Length Distribution](#)
- ✓ [Sequence Duplication Levels](#)
- ✓ [Overrepresented sequences](#)
- ! [Kmer Content](#)

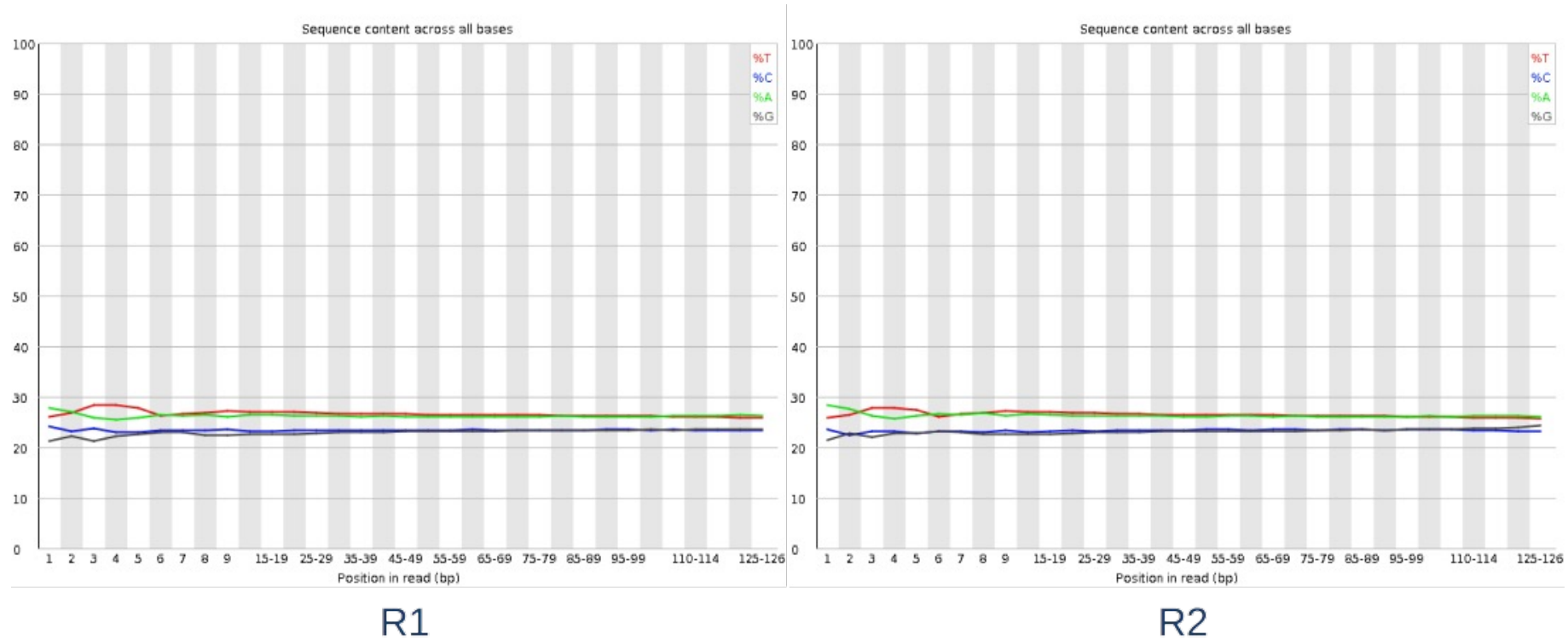
# FastQC – Phred quality score by position

- Example 1 – targeted capture sequencing



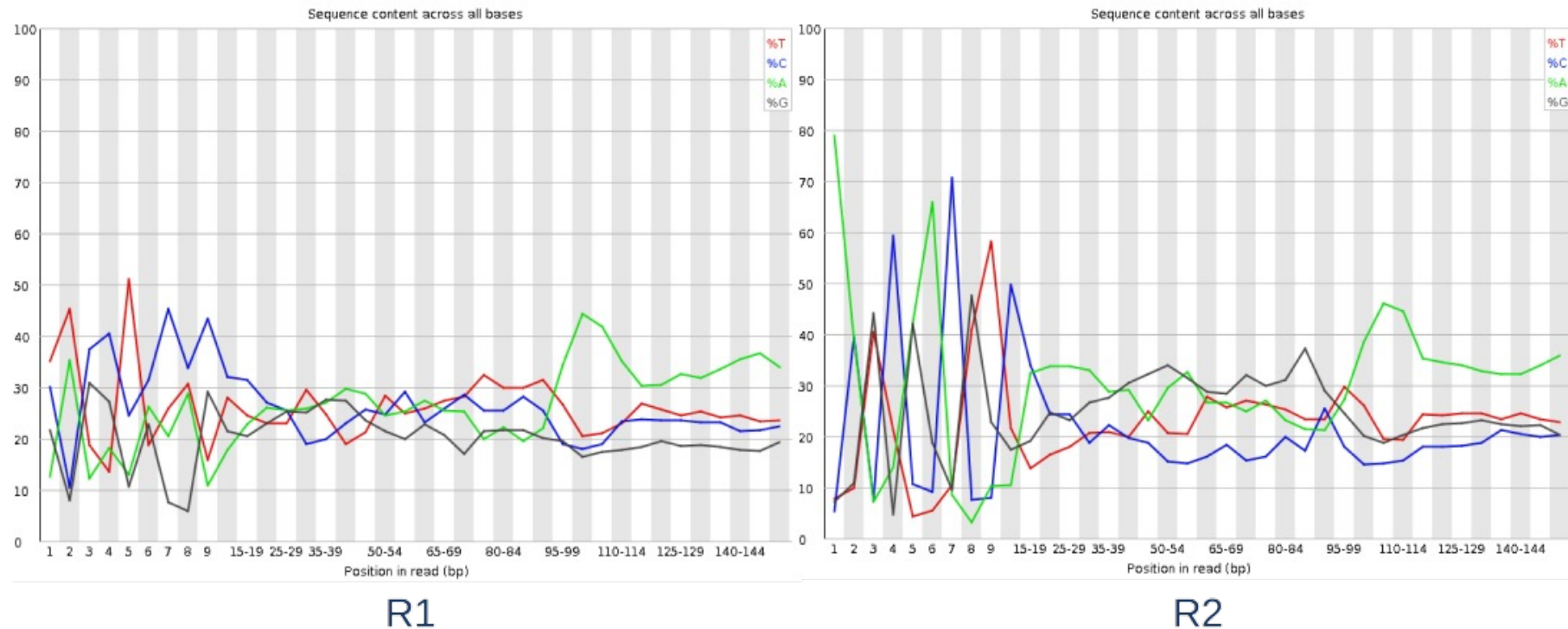
# FastQC – Base content by position

- Example 1 – targeted capture sequencing
  - G-C 25-26%
  - A-T 24-25%



# FastQC – Base content by position

- Example 2 – targeted amplicon sequencing
  - G-C 25-27%
  - A-T 24-24%



# FastQC – Over-represented sequences

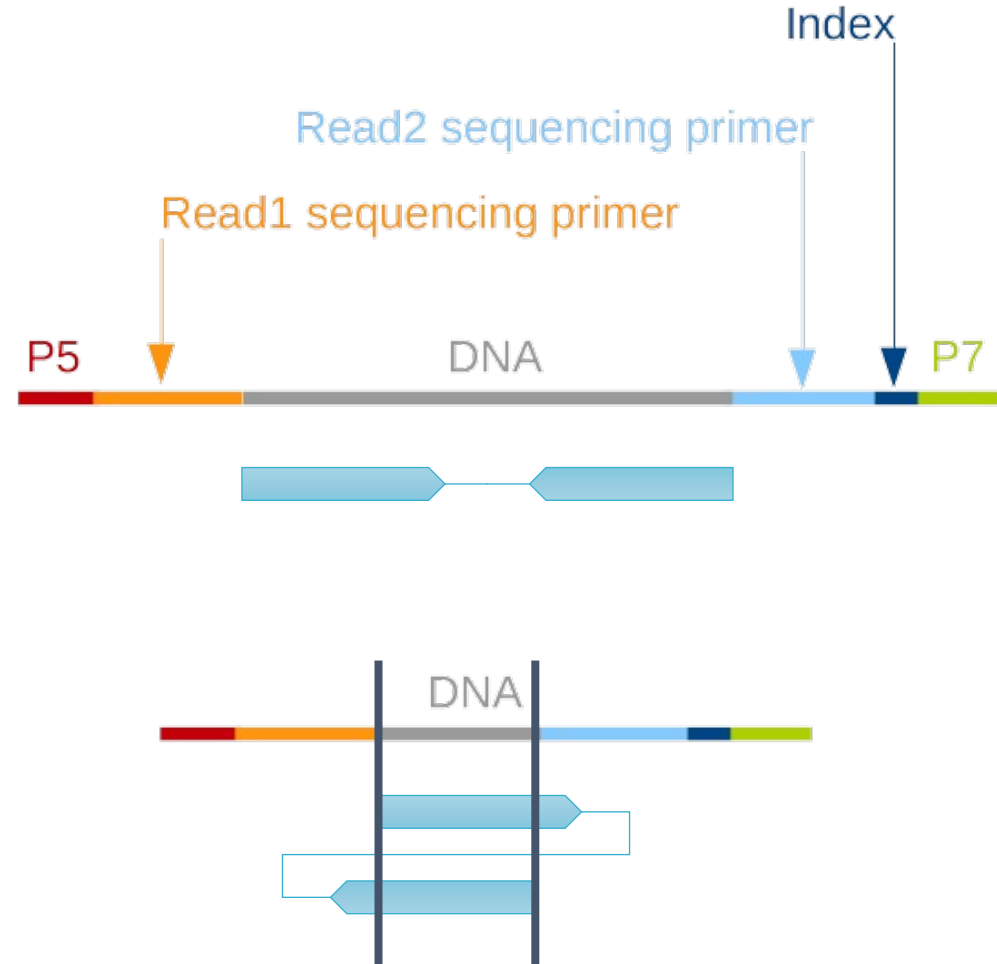
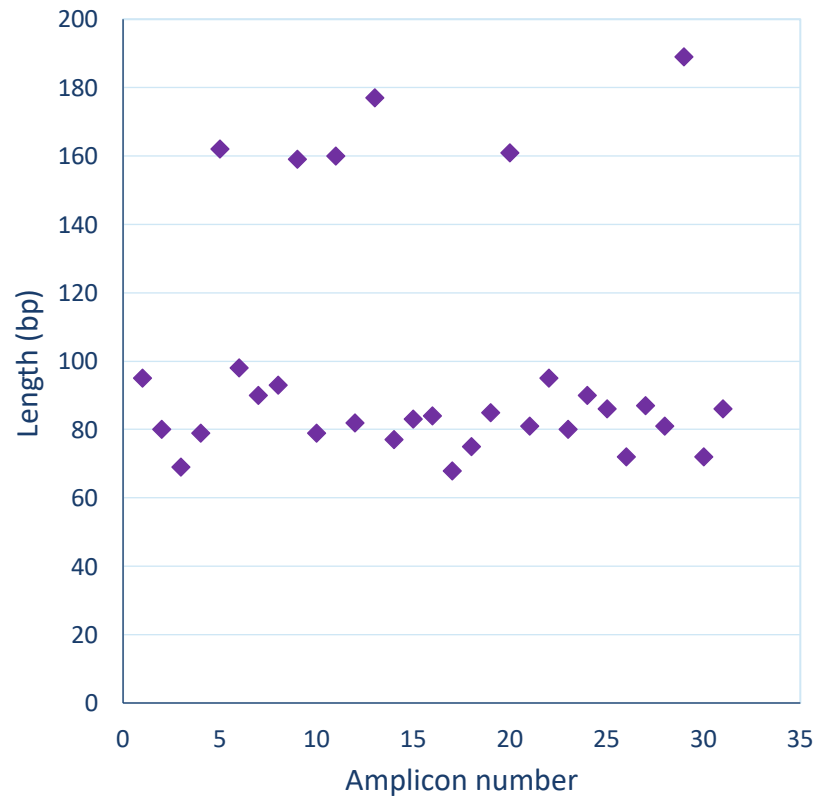
- Example 2 – targeted amplicon sequencing

## ✖ Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GTGATCTCCAACCTTTGACCTGACCGTCGCTTAGATCGGAAGAGCACACGT	92	5.2421652421652425	No Hit
CAGTGACACTAGTCTGCAACAAACGCCACTTAGATCGGAAGAGCACACGT	88	5.014245014245014	No Hit
TTCTCACCCTGCCCTGACCGTCGCTTAGATCGGAAGAGCACACGTGAA	52	2.9629629629629632	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
TTCTCACCCTGCCCTGCAACAAACGCCACTTAGATCGGAAGAGCACACGT	49	2.792022792022792	No Hit
CTCACATCAGCCTGACCTGACCGTCGCTTAGATCGGAAGAGCACACGTCT	31	1.7663817663817662	Illumina Multiplexing PCR Primer 2.01 (100% over 21bp)
CAGCCTCTGCTCTACCTGACCTGACCGTCGCTTAGATCGGAAGAGCACA	26	1.4814814814814816	No Hit
TTCTCACCCTGCCCTGCACTCAATCATCGTCTCTAGATCGGAAGAGCAC	25	1.4245014245014245	No Hit
TTCTCACCCTGCCCTGCACTCTCTCACCTCCACCTGCACTCTCTCTCAC	22	1.2535612535612535	No Hit
CTCACATCAGCCTGACACAACCTTAGGACCACTTGAATAGAGAGCCTCAGT	22	1.2535612535612535	No Hit
TTCTCACCCTGCCCTGCACTCTCTCACCTCTCTCTGCACTCAATCATC	21	1.1965811965811968	No Hit
GACCAGAAGAACCTGACCTGACCGTCGCTTAGATCGGAAGAGCACACGTC	21	1.1965811965811968	No Hit
CTCACATCAGCCTGACCGTCGCTTAGATCGGAAGAGCACACGTCTGAACT	19	1.0826210826210827	Illumina Multiplexing PCR Primer 2.01 (100% over 26bp)
TCCTCCCTCTTGATGTGACCGTCGCTTAGATCGGAAGAGCACACGTCTGA	14	0.7977207977207977	Illumina Multiplexing PCR Primer 2.01 (100% over 23bp)
GTACAGCTGCACTGTGAAGATCGGAAGAGCACACGTCTGAACTCCAGTCA	14	0.7977207977207977	Illumina Multiplexing PCR Primer 2.01 (100% over 33bp)
TTCTCACCCTGCCCTGCACCATGAATGTTTTTATAAAAAGGCTGTTGGC	12	0.6837606837606838	No Hit
AGGTAAGTGACAGTTTGCTCATGGGAAAGGAGATAGATCGGAAGAGCACA	12	0.6837606837606838	No Hit
GTGATCTCCAACCTTTGACCGTCGCTTAGATCGGAAGAGCACACGTCTGAA	11	0.6267806267806267	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
CAAGAGCTCAGAGGAGGAAGCTGTCAGAGATCGGAAGAGCACACGTCTGA	11	0.6267806267806267	Illumina Multiplexing PCR Primer 2.01 (100% over 23bp)
TTGTACTTGACCTGACCGTCGCTTAGATCGGAAGAGCACACGTCTGAA	11	0.6267806267806267	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
CTCACATCAGCCTGACACTTTAAGTCGGGAGTCAGAAAGTACCAAGGAG	9	0.5128205128205128	No Hit
TTGGTGATCATGTGTTGTGTGTGTGGGGGAAGTTGAGTAGATCG	9	0.5128205128205128	No Hit
TTCTCACCCTGCCCTGCACTCGATAATCAATACATAATATTCAATAATT	9	0.5128205128205128	No Hit
GTACAGCTGGTACAAGAACCAGATCGGAAGAGCACACGTCTGAACTCCAG	9	0.5128205128205128	Illumina Multiplexing PCR Primer 2.01 (100% over 30bp)
TTCTCACCCTGACCTGACCGTCGCTTAGATCGGAAGAGCACACGTCTGAA	8	0.4558404558404558	Illumina Multiplexing PCR Primer 2.01 (100% over 24bp)
TTGGTGATCATGTGTTGTGTGTGGGGGAAGTTGAGTAGATCGGAAG	8	0.4558404558404558	No Hit
CATCTGCATGGTGATCTGGGCTCTGTAGTGGTGGCTGCAAGAGGTGCT	8	0.4558404558404558	No Hit
CATTTCCATTGCCAACCGAGTCCATTGTGCACAGTATGAAGACAGCACAT	8	0.4558404558404558	No Hit
GATGTTCAAGGATTCCAGTTAGGTGAGTAAACCTTGATCAGTCACTAT	7	0.39886039886039887	No Hit
AGGTAAGTGACAGTTTGCTCAGGGAAAGTGTGAGATTGGATTCTTTAAAC	7	0.39886039886039887	No Hit
TGGCCTTGACAACAGATCGGAAGAGCACACGTCTGAACTCCAGTCACAG	7	0.39886039886039887	TruSeq Adapter, Index 2 (97% over 35bp)
TTCTCACCCTGCCCTGCACTCTCTCACCTCACTCTCATCACTCCACC	7	0.39886039886039887	No Hit
ACACTGGGCTAGACACTCGTATGGTTGTATGGGTTTCTCTTCTTAGAGA	7	0.39886039886039887	No Hit
TTGTACTTGATCTGGGCGCATCGTTTCTTTTCAAGTTGTGGATAGCAC	7	0.39886039886039887	No Hit
AAGAGCCTGCTGACCGTCGCTTAGATCGGAAGAGCACACGTCTGAACTC	6	0.3418803418803419	Illumina Multiplexing PCR Primer 2.01 (100% over 27bp)

# FastQC – Over-represented sequences

- Example 2 – targeted amplicon sequencing



# Contamination check

- Compare reads to various libraries
- Any library can be searched against
- Output proportion of reads with hit(s) to libraries

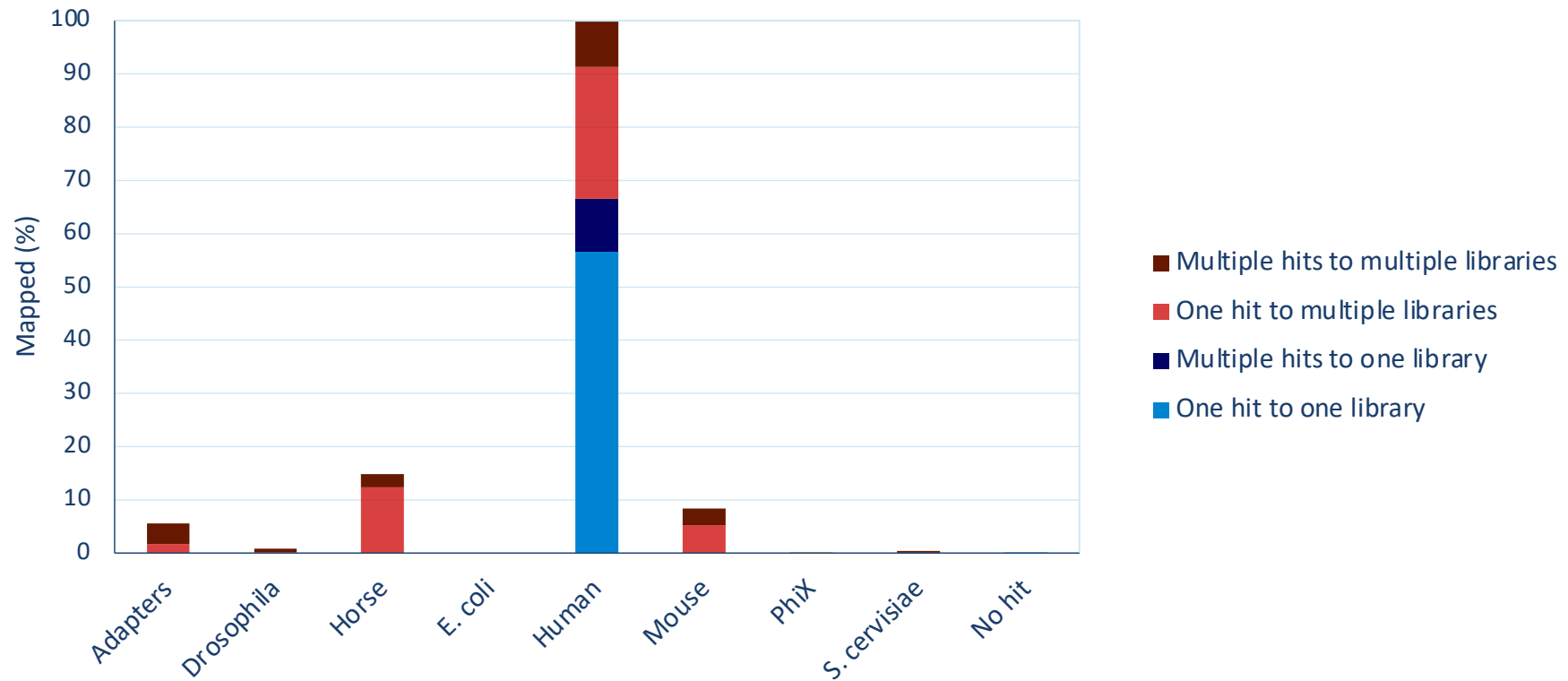
```
fastq_screen --subset 100000 --conf fastq_screen.conf --aligner bowtie2 --  
outdir result --nohits sample.R1.fastq.gz
```

[https://www.bioinformatics.babraham.ac.uk/projects/fastq\\_screen/](https://www.bioinformatics.babraham.ac.uk/projects/fastq_screen/)



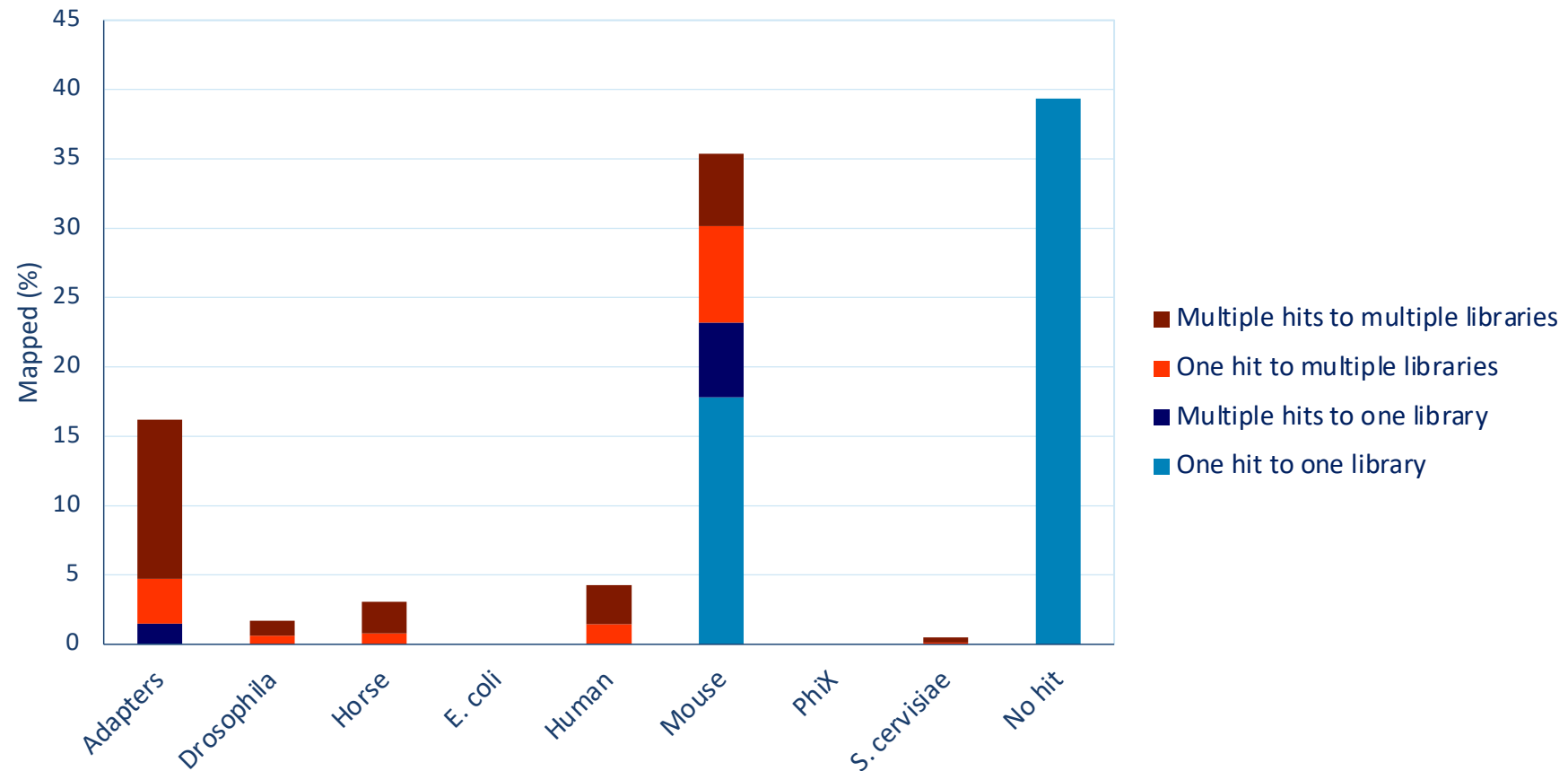
# Contamination check

- Example 1 – targeted capture sequencing



# Contamination check

- Example 2 – targeted amplicon sequencing



# Contamination check

- Taxonomic sequence classification system
  - Build (custom) database
  - Compare k-mer from reads to database

```
kraken2-build --download-taxonomy --db database
```

```
kraken2-build --download-library library --db database
```

```
kraken2-build --build --db database --minimizer-spaces 0
```

```
kraken2 --db database --paired sample.R1.fastq.gz sample.R2.fastq.gz -  
report kraken2Report.txt --use-names > kraken2.output.txt
```

<https://ccb.jhu.edu/software/kraken2/>

# Contamination check

- Example 3 – bacterial WGS

Fragments covered by clade (%)	Fragments covered by clade	Fragments assigned to taxon	Rank code	NCBI taxonomic ID	scientific name
0.04	243	243 U		0 unclassified	
99.96	632527	0 R		1 root	
99.96	632527	0 R1		131567 cellular organisms	
99.96	632491	121 D		2 Bacteria	
99.55	629954	1 D1		1783272 Terrabacteria group	
99.53	629777	0 P		201174 Actinobacteria	
99.53	629777	17 C		1760 Actinobacteria	
99.52	629754	20 O		85007 Corynebacteriales	
99.52	629733	74 F		1762 Mycobacteriaceae	
99.51	629657	5911 G		1763 Mycobacterium	
98.56	623663	5295 G1		77643 Mycobacterium tuberculosis complex	
97.69	618141	613920 S		1773 Mycobacterium tuberculosis	
0.04	227	78 S		78331 Mycobacterium canettii	
0.01	37	35 S		1768 Mycobacterium kansasii	
0.03	176	0 P		1239 Firmicutes	
0.03	176	0 C		91061 Bacilli	
0.03	172	0 O		186826 Lactobacillales	
0.03	172	0 F		1300 Streptococcaceae	
0.03	172	7 G		1301 Streptococcus	
0.03	161	155 S		1313 Streptococcus pneumoniae	
0	2	1 S		28037 Streptococcus mitis	
0	1	0 S		257758 Streptococcus pseudopneumoniae	

# Outline

- Definitions
- FASTQ format
- Quality control
- Pre-processing

# Pre-processing ?

- Process FASTQ files prior to further analysis
  - Remove reads from other species
  - Trim adapters
  - Clip low quality bases
  - Merge overlapping reads from same DNA fragment
  - ...

# Clipping & trimming

- Sequencing adapters and primers
- Poor quality bases at the ends of reads
- Ns from ends of reads
- Remove low complexity reads

```
fastq-mcf -H -X -o sample_filtered.R1.fastq.gz -o  
/sample_filtered.R2.fastq.gz adapters.fa sample.R1.fastq sample.R2.fastq
```

<https://github.com/ExpressionAnalysis/ea-utils/blob/wiki/FastqMcf.md>

# Adapter clipping & trimming

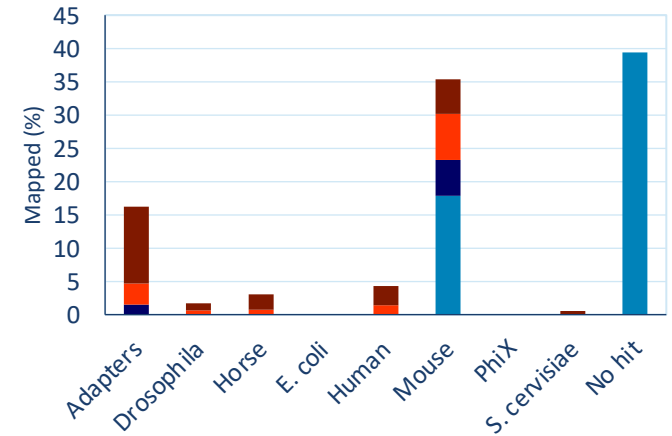
- Example 2 – targeted amplicon sequencing

- Input

- 2 fastq files of 1,834 reads each

- Outputs

- 2 fastq files of 1,801 reads each
    - List of adapter found



Adapter TruSeq\_Adapter\_Index\_1 : counted **1038** at the 'end' of 'sample.R1.fastq'

...

Adapter Illumina\_Single\_End\_Sequencing\_Primer\_3p : counted **1046** at the 'end' of 'sample.R2.fastq'

...

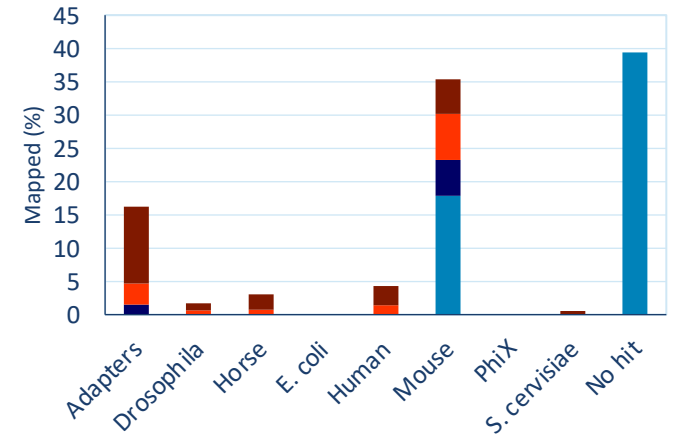
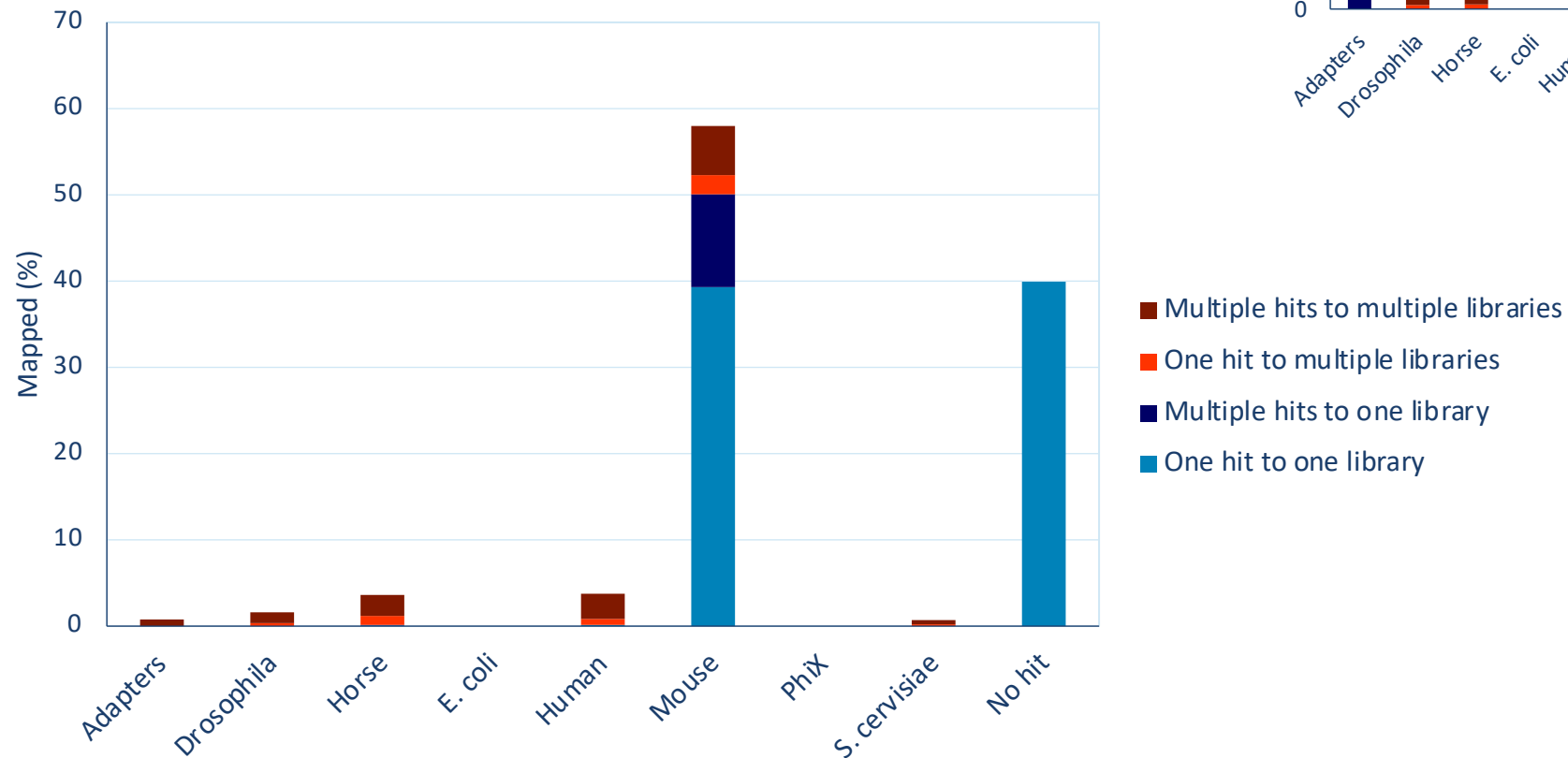
Total reads: 1801

Too short after clip: 33



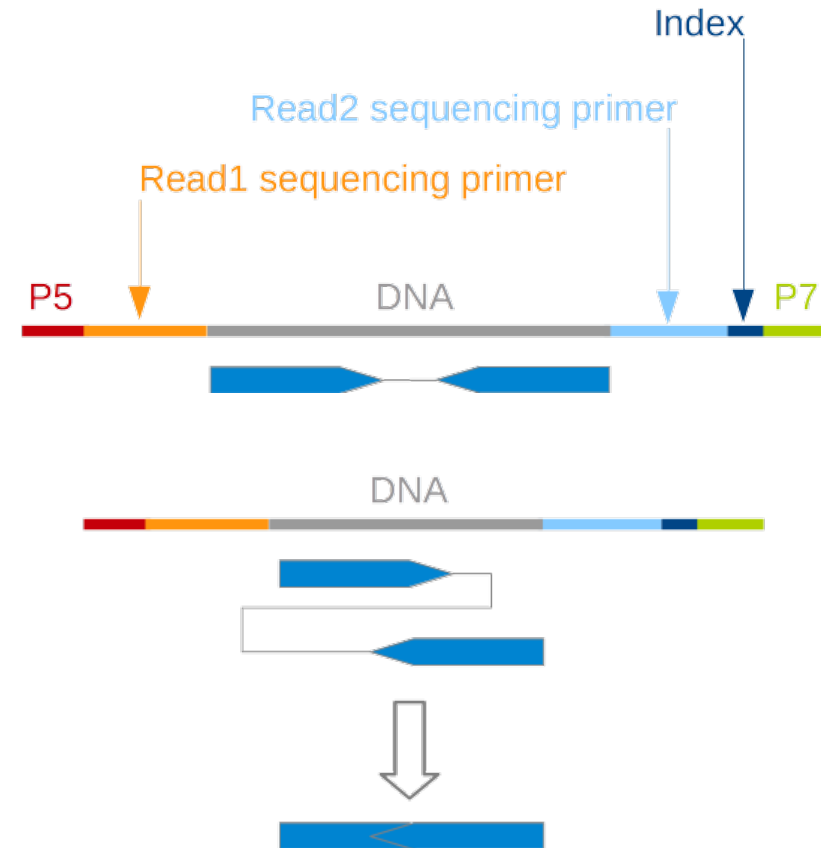
# Adapter clipping & trimming

- Example 2 – targeted amplicon sequencing



# Merge overlapping reads

- Merge paired-end reads
- Keep DNA insert only

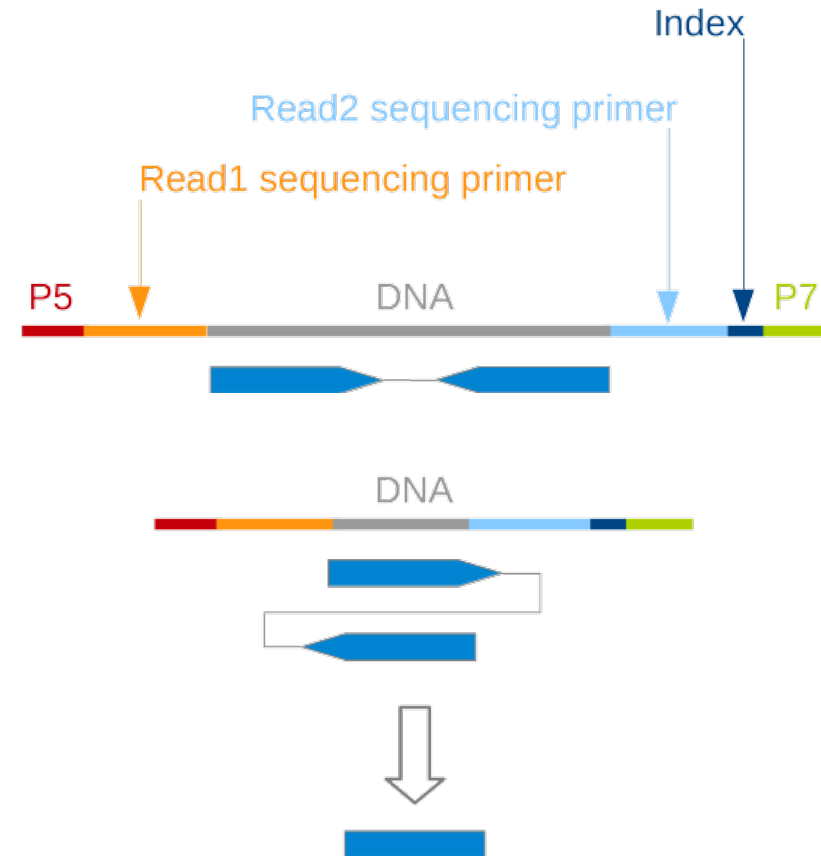


```
flash2 --max-overlap=250 --min-overlap=20 --allow-outies -d result -o  
sample.flashed sample.R1.fastq.gz sample.R2.fastq.gz > flash.log
```

<https://github.com/dstreett/FLASH2>

# Merge overlapping reads

- Merge paired-end reads
- Keep DNA insert only



```
flash2 --max-overlap=250 --min-overlap=20 --allow-outies -d result -o  
sample.flashed sample.R1.fastq.gz sample.R2.fastq.gz > flash.log
```

<https://github.com/dstreett/FLASH2>

# Read selection

- Select reads from sequenced organism

Fragments covered by clade (%)	Fragments covered by clade	Fragments assigned to taxon	Rank code	NCBI taxonomic ID	scientific name
99.52	629733	74F		1762	Mycobacteriaceae
99.51	629657	5911G		1763	Mycobacterium
98.56	623663	5295G1		77643	Mycobacterium tuberculosis complex
97.69	618141	613920S		1773	Mycobacterium tuberculosis
0.04	227	78S		78331	Mycobacterium canettii
0.01	37	35S		1768	Mycobacterium kansasii

```
grep "organism" kraken2.output.txt | cut -f2 > reads.list
```

```
seqtk subseq sample.R1.fastq.gz reads.list | gzip - >  
sample.selected.R1.fastq.gz
```

```
seqtk subseq sample.R2.fastq.gz reads.list | gzip - >  
sample.selected.R2.fastq.gz
```

<https://ccb.jhu.edu/software/kraken2/>

<https://github.com/lh3/seqtk>

# Summary

- Quality control can be done at the raw data level
- Pipeline has to be tailored for analysis
- Pre-processing might be required depending on the analysis

# Questions?

Erika Souche