

# Sequence format and pre-processing

Ramiro Logares, ICM, Barcelona





# Illumina short reads

- We receive files in fastq format from the sequencing center
- Two files per sample, 1 forward (R1) and 1 reverse (R2)
- Normally, reads overlap
- Depending on the library preparation, all reads are in the same direction (5'-3') or their directions are mixed in both R1 and R2
- We normally work with gzipped fastq files to save space



# How files look like when they are received

```
18S_BL060613_30r2-MSTAReuk_R2.fastq.gz 18S_BL091105_30-MSTAReuk_R2.fastq.gz 18S_BL130417_30-MSTAReuk_R2.fastq.gz
18S_BL060613_30r3_MSTAReuk_R1.fastq.gz 18S_BL091105_30r3_MSTAReuk_R1.fastq.gz 18S_BL130417_30r3_MSTAReuk_R1.fastq.gz
18S_BL060613_30r3_MSTAReuk_R2.fastq.gz 18S_BL091105_30r3_MSTAReuk_R2.fastq.gz 18S_BL130417_30r3_MSTAReuk_R2.fastq.gz
18S_BL060704_022-MSTAReuk_R1.fastq.gz 18S_BL091222_022-MSTAReuk_R1.fastq.gz 18S_BL130507_022-MSTAReuk_R1.fastq.gz
18S_BL060704_022-MSTAReuk_R2.fastq.gz 18S_BL091222_022-MSTAReuk_R2.fastq.gz 18S_BL130507_022-MSTAReuk_R2.fastq.gz
18S_BL060704_30-MSTAReuk_R1.fastq.gz 18S_BL091222_30-MSTAReuk_R1.fastq.gz 18S_BL130507_30-MSTAReuk_R1.fastq.gz
18S_BL060704_30-MSTAReuk_R2.fastq.gz 18S_BL091222_30-MSTAReuk_R2.fastq.gz 18S_BL130507_30-MSTAReuk_R2.fastq.gz
18S_BL060704_30r3_MSTAReuk_R1.fastq.gz 18S_BL091222_30r3_MSTAReuk_R1.fastq.gz 18S_BL130507_30r3_MSTAReuk_R1.fastq.gz
18S_BL060704_30r3_MSTAReuk_R2.fastq.gz 18S_BL091222_30r3_MSTAReuk_R2.fastq.gz 18S_BL130507_30r3_MSTAReuk_R2.fastq.gz
18S_BL060801_022-MSTAReuk_R1.fastq.gz 18S_BL100120_022-MSTAReuk_R1.fastq.gz 18S_BL130604_022-MSTAReuk_R1.fastq.gz
18S_BL060801_022-MSTAReuk_R2.fastq.gz 18S_BL100120_022-MSTAReuk_R2.fastq.gz 18S_BL130604_022-MSTAReuk_R2.fastq.gz
18S_BL060801_30r2-MSTAReuk_R1.fastq.gz 18S_BL100120_30-MSTAReuk_R1.fastq.gz 18S_BL130604_30-MSTAReuk_R1.fastq.gz
18S_BL060801_30r2-MSTAReuk_R2.fastq.gz 18S_BL100120_30-MSTAReuk_R2.fastq.gz 18S_BL130604_30-MSTAReuk_R2.fastq.gz
18S_BL060801_30r3_MSTAReuk_R1.fastq.gz 18S_BL100120_30r3_MSTAReuk_R1.fastq.gz 18S_BL130604_30r3_MSTAReuk_R1.fastq.gz
18S_BL060801_30r3_MSTAReuk_R2.fastq.gz 18S_BL100120_30r3_MSTAReuk_R2.fastq.gz 18S_BL130604_30r3_MSTAReuk_R2.fastq.gz
18S_BL060912_022-MSTAReuk_R1.fastq.gz 18S_BL100217_022-MSTAReuk_R1.fastq.gz 18S_BL130709_022-MSTAReuk_R1.fastq.gz
18S_BL060912_022-MSTAReuk_R2.fastq.gz 18S_BL100217_022-MSTAReuk_R2.fastq.gz 18S_BL130709_022-MSTAReuk_R2.fastq.gz
```



# fastq format

- Four sequences per line
  1. @sequence.ID
  2. ACTGACTGACTG # nucleotide sequence
  3. + (separator)
  4. Quality scores (Phred +33: normally 0-41)







@M02696:67:000000000-B44VG:1:1101:11781:1257 1:N:0:57

The first line, identifying the sequence, contains the following elements.

@<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<x-pos>:<y-pos>:<UMI> <read>:<is filtered>:<control number>:<index>

Table 1 FASTQ File Elements

Element	Requirements	Description
@	@	Each sequence identifier line starts with @.
<instrument>	Characters allowed: a–z, A–Z, 0–9 and underscore	Instrument ID.
<run number>	Numerical	Run number on instrument.
<flowcell ID>	Characters allowed: a–z, A–Z, 0–9	
<lane>	Numerical	Lane number.
<tile>	Numerical	Tile number.
<x_pos>	Numerical	X coordinate of cluster.
<y_pos>	Numerical	Y coordinate of cluster.
<UMI>	Restricted characters: A/T/G/C/N	Optional, appears when UMI is specified in sample sheet. UMI sequences for Read 1 and Read 2, seperated by a plus [+].
<read>	Numerical	Read number. 1 can be single read or Read 2 of paired-end.
<is filtered>	Y or N	Y if the read is filtered (did not pass), N otherwise.
<control number>	Numerical	0 when none of the control bits are on, otherwise it is an even number. On HiSeq X and NextSeq systems, control specification is not performed and this number is always 0.
<index>	Restricted characters: A/T/G/C/N	Index of the read.



# Sanger Phred quality scores

Phred quality scores are logarithmically linked to error probabilities		
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%

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

Q = Phred quality scores

P = base calling error probability



# Calculating Phred scores

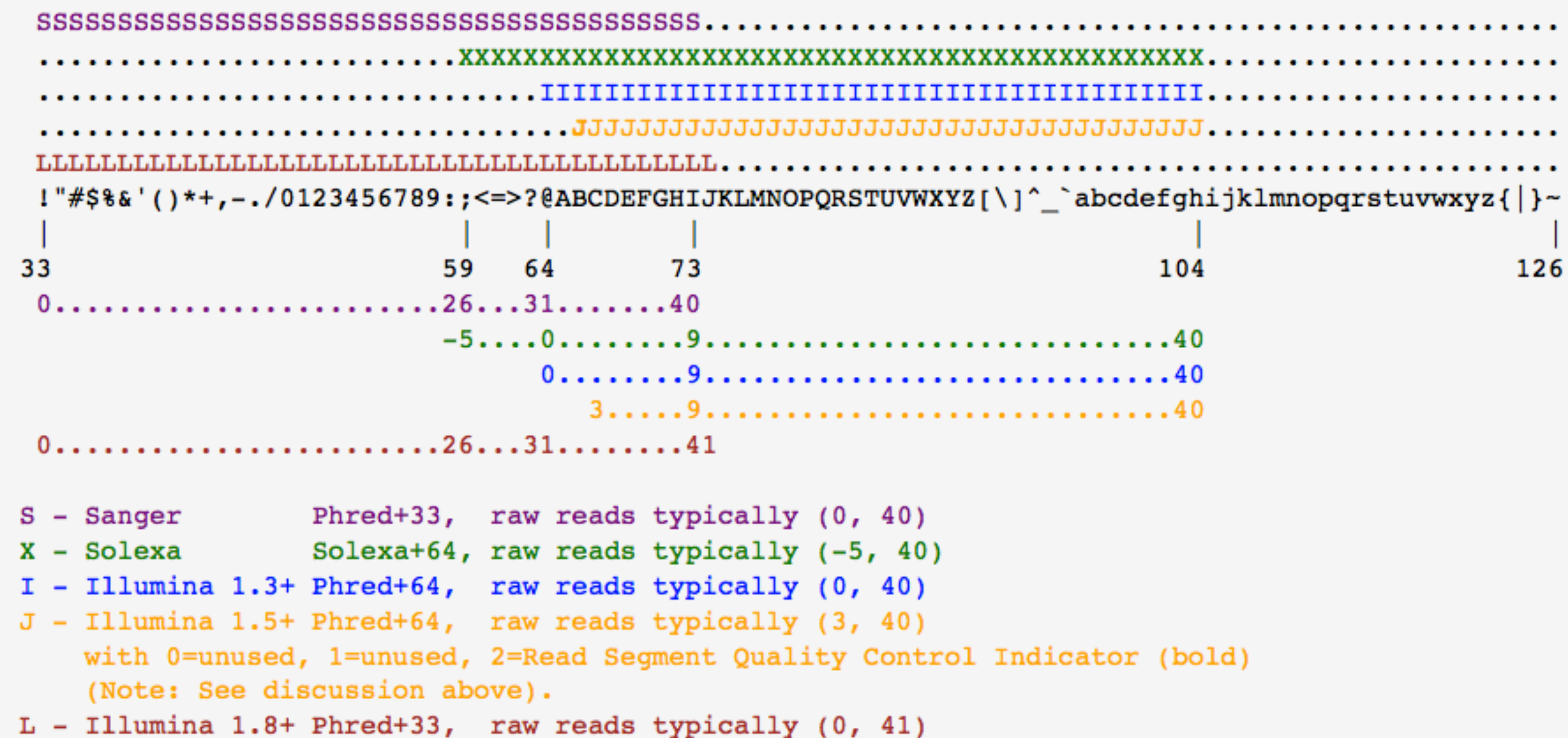
- To determine quality scores, Phred first calculates several parameters related to peak shape and peak resolution at each base.
- Phred then uses these parameters to look up a corresponding quality score in huge lookup tables.
- These lookup tables were generated from sequence traces where the correct sequence was known, and are hard coded in Phred; different lookup tables are used for different sequencing chemistries and machines.







# Phred encoding in different sequencers





# Removing primers

- The sequences received from the sequencing center may contain primers used to amplify them
- Primers need to be removed as they normally contain ambiguous positions that can interfere with DADA2
- DADA2 assumes primers have been removed



# Checking if sequences have primers in unix

- We know the primer sequence expected in the reads
- Fw :CCAGCA[ACGT]C[ACGT]GCGGTAATTCC
- Rv: ACTTTCGTTCTTGAT[AGCT][AGCT][AGCT]
  - Ambiguities are included to match sequences
- We use zgrep to match the primer against the gzipped sequences

```
[rlogares@marbits raw]$ ls
BL100525E-MSTAREuk_R1.fastq.gz  BL100706E-MSTAREuk_R1.fastq.gz  BL100914E-MSTAREuk_R1.fastq.gz  primers_R1_in_reads
BL100525E-MSTAREuk_R2.fastq.gz  BL100706E-MSTAREuk_R2.fastq.gz  BL100914E-MSTAREuk_R2.fastq.gz  primers_R2_in_reads
BL100622E-MSTAREuk_R1.fastq.gz  BL100803E-MSTAREuk_R1.fastq.gz  clipping_primers.sh
BL100622E-MSTAREuk_R2.fastq.gz  BL100803E-MSTAREuk_R2.fastq.gz  cutadapt.o40252
[rlogares@marbits raw]$ zgrep -c --color CCAGCA[ACGT]C[ACGT]GCGGTAATTCC BL100525E-MSTAREuk_R1.fastq.gz
23843
[rlogares@marbits raw]$

[rlogares@marbits raw]$ zgrep -c --color ACTTTCGTTCTTGAT[AGCT][AGCT][AGCT] BL100525E-MSTAREuk_R2.fastq.gz
23856
[rlogares@marbits raw]$
```

Forward primer

Reverse primer



- As several counts are given, we inspect the sequences visually

```
[rlogares@marbits raw]$ zgrep --color CCAGCA[ACGT]C[ACGT]GCGGTAATTCC BL100525E-MSTAreuk_R1.fastq.gz
```

```
CCAGCACCTGCGGTAATTCCGGCTCCTTCAGCCTGAGGTAGAATTGTTGTAGTTAAAACGCTCGTAGTTGGATTTTGTAAAGAGTTTTGTGTGTGTTGGTTGCGTATATATTCGTATATTCGTGATT
CTTCATGCCACTTTTATACTGATTGTGGATAATTTTCGGATTATTTGCAATATTACTGTGAGAAAAAGAGTGCGCTTAAGGGCGGCTTTATGCTAAGATCATTTAGCATGGAATAAACATAACGG
CCAGCACCTGCGGTAATTCCGGCTCCTTCAGCCTGAGGTAGAATTGTTGTAGTTAAAACGCTCGTAGTTGGATTTTGTAAAGAGTTTTGTGTGTGTTGGTTGCGTATATATTCGTATATTCCTGATT
CTTCATGCCACTTTTATACTGATTGTGGATAATTTTCGGATTATTTGCAATATTACTGTGAGAAAAAGAGTGCGCTTAAGGGCGGCTTTATGCTAAGATCATTTAGCATGGAATAAACATAACGG
CCAGCACCTGCGGTAATTCCGGCTCCTTCAGCCTGAGGTAGAATTGTTGTAGTTAAAACGCTCGTAGTTGGATTTTGTAAAGAGTTTTGTGTGTGTTGGTTGCGTATATATTCGTATATTCGTGATT
CTTCATGCCACTTTTATACTGATTGTGGATAATTTTCGGATTATTTGCAATATTACTGTGAGAAAAAGAGTGCGCTTAAGGGCGGCTTTATGCTAAGATCATTTAGCATGGAATAAACATAACGG
CCAGCACCTGCGGTAATTCCGGCTCCTTCAGCCTGATGTAGAATTGTTGTAGTTAAAACGCTCGTAGTTGGATTTTGTAAAGAGTTTTGTGTGTGTTGGTTGCGTATATATTCGTATATTCGTGATT
CTTCATGCCACTTTTATACTGATTGTGGATAATTTTCGGATTATTTGCAATATTACTGTGAGAAAAAGAGTGCGCTTAAGGGCGGCTTTATGCTAAGATCATTTAGCATGGAATAAACATAACGG
CCAGCAGCCGCGGTAATTCCGGCTCCTTCAGCCTGAGGTAGAATTGTTGTAGTTAAAACGCTCGTAGTTGGATTTTGTAAAGAGTTTTGTGTGTGTTGGTTGCGTATATATTCGTATATTCGTGATT
CTTCATGCCACTTTTATACTGATTGTGGATAATTTTCGGATTATTTGCAATATTACTGTGAGAAAAAGAGTGCGCTTAAGGGCGGCTTTATGCTAAGATCATTTAGCATGGAATAAACATAACGG
CCAGCAGCCGCGGTAATTCCGGCTCCTTCAGCCTGAGGTAGAATTGTTGTAGTTAAAACGCTCGTAGTTGGATTTTGTAAAGAGTTTTGTGTGTGTTGGTTGCGTATATATTCGTATATTCGTGATT
CTTCATGCCACTTTTATACTGATTGTGGATAATTTTCGGATTATTTGCAATATTACTGTGAGAAAAAGAGTGCGCTTAAGGGCGGCTTTATGCTAAGATCATTTAGCATGGAATAAACATAACGG
[rlogares@marbits raw]$
```



# We use cutadapt to remove primers

- Program: <https://cutadapt.readthedocs.io/en/stable/>
- Runs in unix
- Cutadapt will search for primers in R1 and R2 sequences and remove them
- It can also remove all sequences where primers have not been found



```
# Running cutadapt in a loop (NB: use arrays if you have a cluster)
```

```
for i in $(ls *fastq.gz | cut -f 1 -d - | uniq); \
do cutadapt -g CCAGCASCYGC GGTAATTCC -G ACTTTCGTTCTTGATYRR \
-m 100 -M 350 --match-read-wildcards --pair-filter=both -q 10 \
-o $i-MSTAREuk_R1.clipped.fastq.gz -p $i-MSTAREuk_R2.clipped.fastq.gz \
$i-MSTAREuk_R1.fastq.gz $i-MSTAREuk_R2.fastq.gz; done

# -g ADAPTER, --front=ADAPTER
#           Sequence of an adapter ligated to the 5' end (paired
#           data: of the first read). The adapter and any
#           preceding bases are trimmed. Partial matches at the 5'
#           end are allowed. If a '^' character is prepended
#           ('anchoring'), the adapter is only found if it is a
#           prefix of the read.
# Paired-end options:
# The -A/-G/-B/-U options work like their -a/-b/-g/-u counterparts, but
# are applied to the second read in each pair.
#
# -G ADAPTER          5' adapter to be removed from second read in a pair.
# -m LENGTH, --minimum-length=LENGTH
#           Discard reads shorter than LENGTH. Default: 0
# -M LENGTH, --maximum-length=LENGTH
#           Discard reads longer than LENGTH. Default: no limit
# --match-read-wildcards
#           Interpret IUPAC wildcards in reads. Default: False
# --pair-filter=(any|both)
#           Which of the reads in a paired-end read have to match
#           the filtering criterion in order for the pair to be
#           filtered. Default: any
# -q [5'CUTOFF,]3'CUTOFF, --quality-cutoff=[5'CUTOFF,]3'CUTOFF
#           Trim low-quality bases from 5' and/or 3' ends of each
#           read before adapter removal. Applied to both reads if
#           data is paired. If one value is given, only the 3' end
#           is trimmed. If two comma-separated cutoffs are given,
#           the 5' end is trimmed with the first cutoff, the 3'
#           end with the second.
# -o output file R1
# -p FILE, --paired-output=FILE
#           Write second read in a pair to FILE.
```

ambiguities are interpreted

There are several additional options



- After cutadapt, sequences are ready to be used in dada2
- It is good to double check primers are gone using the same zgrep command used before
- We don't analyse the overall quality of the sequences, as this will be done later with dada2
- We only remove entire sequences that look very wrong with cutadapt
- It is important to consider whether sequences come from sequencers with 4- or 2- color chemistries, as this will change cutadapt parameters
- There are alternative tools, such as Trimmomatic



THE END