

1 NGS Data formats and QC

1.1 Introduction

There are several file formats for storing Next Generation Sequencing (NGS) data. In this tutorial we will look at some of the most common formats for storing NGS reads and variant data. We will cover the following formats:

FASTQ - This format stores unaligned read sequences with base qualities

SAM/BAM - This format stores unaligned or aligned reads (text and binary formats)

CRAM - This format is similar to BAM but has better compression than BAM

VCF/BCF - Flexible variant call format for storing SNPs, indels, structural variations (text and binary formats)

Following this, we will work through some examples of converting between the different formats.

Further to understanding the different file formats, it is important to remember that all sequencing platforms have technical limitations that can introduce biases in your sequencing data. Because of this it is very important to check the quality of the data before starting any analysis, whether you are planning to use something you have sequenced yourself or publicly available data. In the latter part of this tutorial we will describe how to perform a QC assessment for your NGS data, and also suggest how to identify possible contamination.

1.2 Learning outcomes

On completion of the tutorial, you can expect to be able to:

- Describe the different NGS data formats available (FASTQ, SAM/BAM, CRAM, VCF/BCF)
- Perform conversions between the different data formats
- Perform a QC assessment of high throughput sequence data
- Identify possible contamination in high throughput sequence data

1.3 Tutorial sections

This tutorial comprises the following sections:

1. [Data formats](#)
2. [File conversion](#)
3. [QC assessment](#)
4. [Identifying contamination](#)

1.4 Authors

This tutorial was written by [Sara Sjonnebo](#) based on material from [Petr Danecek](#) and [Thomas Keane](#).

1.5 Running the commands from this tutorial

You can run the commands in this tutorial either directly from the Jupyter notebook (if using Jupyter), or by typing the commands in your terminal window.

1.5.1 Running commands on Jupyter

If you are using Jupyter, command cells (like the one below) can be run by selecting the cell and clicking *Cell -> Run* from the menu above or using *ctrl Enter* to run the command. Let's give this a try by printing our working directory using the *pwd* command and listing the files within it. Run the commands in the two cells below.



```
pwd
```



```
ls -l
```

1.5.2 Running commands in the terminal

You can also follow this tutorial by typing all the commands you see into a terminal window. This is similar to the "Command Prompt" window on MS Windows systems, which allows the user to type DOS commands to manage files.

To get started, select the cell below with the mouse and then either press control and enter or choose *Cell -> Run* in the menu at the top of the page.



```
echo cd $PWD
```

Now open a new terminal on your computer and type the command that was output by the previous cell followed by the enter key. The command will look similar to this:

```
cd /home/manager/pathogen-informatics-training/Notebooks/QC/
```

Now you can follow the instructions in the tutorial from here.

1.6 Let's get started!

This tutorial assumes that you have samtools, bcftools, bwa, Picard tools and Kraken installed on your computer. For download and installation instructions, please see:

- The [samtools website](#)
- The [bcftools website](#)
- The [bwa GitHub page](#)
- The [Picard tools GitHub page](#)
- The [Kraken website](#)

To check that you have installed these correctly, you can run the following commands:



```
samtools -help
```



```
bcftools -help
```



```
bwa
```



```
java -jar $PICARD -h
```



```
kraken -help
```

Where \$PICARD is an environmental variable set to point at picard.jar.

This should return the help message for samtools, bcftools bwa, picardtools and kraken respectively.

To get started with the tutorial, head to the first section: [Data formats](#)

2 Data formats for NGS data

Here we will take a closer look at some of the most common NGS data formats.

2.1 FASTA

The FASTA format is one of the most basic ways to store sequence data, and it can be used to store both nucleotide data and protein sequences. Each sequence in a FASTA file is represented by two parts, a header line and the actual sequence. The header always starts with the symbol “>” and is followed by information about the sequence, such as a unique identifier. The following lines show two sequences represented in FASTA format:

```
>Sequence_1
CTTGACGACTTGAAAAATGACGAAATCACTAAAAAACGTGAAAAATGAGAAATG
AAAATGACGAAATCACTAAAAAACGTGACGACTTGAAAAATGACCAC
>Sequence_2
CTTGAGACGAAATCACTAAAAAACGTGACGACTTGAAGTGAAAAATGAGAAATG
AAATCATGACGACTTGAAGTGAAAAAGTGAAAAATGAGAAATGAACGTGACGAC
AAAATGACGAAATCATGACGACTTGAAGTGAAAAATAAATGACC
```

2.1.1 Exercises

Q1: How many sequences are there in the fasta file data/example.fasta? (hint: is there a grep option you can use?)



2.2 FASTQ

FASTQ is a data format for raw unaligned sequencing reads. It is an extension to the FASTA file format, and includes a quality score for each base. For paired-end sequencing, two FASTQ files are produced. Have a look at the example below, containing two reads:

```
@ERR007731.739 IL16_2979:6:1:9:1684/1
CTTGACGACTTGAAAAATGACGAAATCACTAAAAAACGTGAAAAATGAGAAATG
+
BBCBCBBBBBBBABBBABBBBBBBBABBBBBBBBBBBBBBABAAAABBBBBB=>B
@ERR007731.740 IL16_2979:6:1:9:1419/1
AAAAAAAAAGATGTCATCAGCACATCAGAAAAGAAGGCAACTTTAAACTTTTC
+
BBABBBBABABAABABABBBBAAA>@B@BBAA@4AAA>.>BAA@779:AAA@A
```

We can see that for each read we get four lines:

1. The read metadata, such as the read ID. Starts with @ and, for paired-end Illumina reads, is terminated with /1 or /2 to show that the read is the member of a pair.
2. The read
3. Starts with + and optionally contains the ID again
4. The per base quality score

The quality scores are encoded as [ASCII characters](#). The first 32 codes are reserved for control characters which are not printable, and the 33rd is reserved for space. Neither of these can be used in the quality string, so we need to withdraw 33 from whatever the value of the quality character is. For example, the ASCII code of "A" is 65, so the corresponding quality is:

$$Q = 65 - 33 = 32$$

From this we can calculate the [Phred quality score](#) as:

$$P = 10^{-Q/10}$$

The Phred quality score is a measure of the quality of base calls. For example, a base assigned with a Phred quality score of 30 tells us that there is a 1 in 1000 chance that this base was called incorrectly.

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 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

The following simple perl command will print the quality score value for an ASCII character. Try changing the "A" to another character, for example one from the quality strings above (e.g. @, = or B).



```
perl -e 'printf "%d\n", ord("A")-33;'
```

Something to be aware of is that two different ways of calculating the quality scores have historically been in use. The standard Sanger variant uses the Phred score, while the Solexa pipeline earlier used a [different mapping](#).

2.2.1 Exercises

2: How many reads are there in the file example.fastq? (Hint: remember that @ is a possible quality score. Is there something else in the header that is unique?)



2.3 SAM/BAM

SAM (Sequence Alignment/Map) format was developed by the 1000 Genomes Project group in 2009 and is a unified format for storing read alignments to a reference genome. BAM is the compressed version of SAM. SAM/BAM format is the accepted standard format for storing NGS sequencing reads, base qualities, associated meta-data and alignments of the data to a reference genome. If no reference genome is available, the data can also be stored unaligned.

The files consist of a header section (optional) and an alignment section. The alignment section contains one record (a single DNA fragment alignment) per line describing the alignment between fragment and reference. Each record has 11 fixed columns and optional key:type:value tuples. Open the [SAM/BAM file specification document](#) as you may need to refer to it throughout this tutorial.

Now let us have a closer look at the different parts of the SAM/BAM format.

2.3.1 Header Section

Each line in the SAM header begins with an @, followed by a two-letter header record type code as defined in the [SAM/BAM format specification document](#). Each record type can contain meta-data captured as a series of key-value pairs in the format of 'TAG:VALUE'.

Read groups One useful record type is RG which can be used to describe each lane of sequencing. The RG code can be used to capture extra meta-data for the sequencing lane. Some common RG TAGs are:

- ID: SRR/ERR number
- PL: Sequencing platform
- PU: Run name
- LB: Library name
- PI: Insert fragment size
- SM: Individual/Sample
- CN: Sequencing centre

2.3.2 Exercises

From reading section 1.3 of the SAM specification, look at the following line from the header of the SAM/BAM file:

```
@RG ID:ERR003612 PL:ILLUMINA LB:g1k-sc-NA20538-TOS-  
1 PI:2000 DS:SRP000540 SM:NA20538 CN:SC
```

Q3: What does RG stand for?



Q4: What is the sequencing platform?



Q5: What is the sequencing centre?



Q6: What is the lane ID?



Q7: What is the expected fragment insert size?

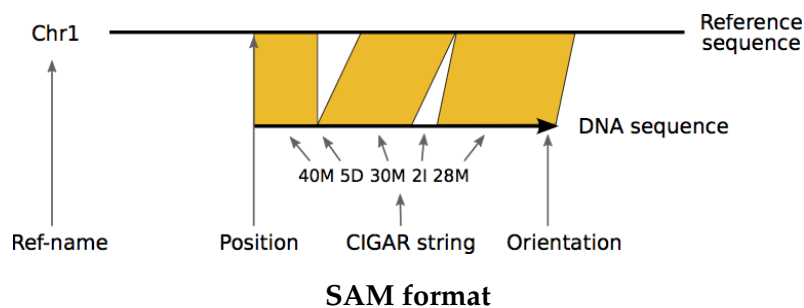


2.3.3 Alignment Section

The alignment section of SAM files contains one line per fragment alignment, which in turn contains the columns listed below. The first 11 columns are mandatory.

1. QNAME: Query NAME of the read or the read pair
2. FLAG: Bitwise FLAG (pairing, strand, mate strand, etc.)
3. RNAME: Reference sequence NAME
4. POS: 1-Based leftmost POSition of clipped alignment
5. MAPQ: MAPping Quality (Phred-scaled)
6. CIGAR: Extended CIGAR string (operations: MIDNSHPX=)
7. MRNM: Mate Reference NaMe ('=' if same as RNAME)
8. MPOS: 1-Based leftmost Mate POSition
9. ISIZE: Inferred Insert SIZE
10. SEQ: Query SEQUENCE on the same strand as the reference
11. QUAL: Query QUALity (ASCII-33=Phred base quality)
12. OTHER: Optional fields

The image below provides a visual guide to some of the columns of the SAM format.



In a SAM file, this image representation could for instance translate to the following entry of 100 bases:

```
ERR005816.1408831 163 Chr1 19999970 23 40M5D30M2I28M = 20000147 213 GGTGGGTGGAT-
CACCTGAGATCGGGAGTTTGAGACTAGGTGG...
<=@A@??@=@A@A>@BAA@ABA:>@<>=BBB9@@2B3<=@A@...
```

2.3.4 Exercises

Let's have a look at `example.sam`. This file contains only a subset of the alignment section of a BAM-file that we will look closer at soon. Notice that we can use the standard UNIX operations like `cat` on this file.



```
cat data/example.sam
```

Q8: What is the mapping quality of `ERR003762.5016205`? (Hint: can you use `grep` and `awk` to find this?)



Q9: What is the CIGAR string for `ERR003814.6979522`? (Hint: we will go through the meaning of CIGAR strings in the next section)



Q10: What is the inferred insert size?



2.3.5 CIGAR string

Column 6 of the alignment is the CIGAR string for that alignment. The CIGAR string provides a compact representation of sequence alignment. Have a look at the table below. It contains the meaning of all different symbols of a CIGAR string:

Symbol	Meaning
M	alignment match or mismatch
=	sequence match
X	sequence mismatch
I	insertion to the reference
D	deletion from the reference
S	soft clipping (clipped sequences present in SEQ)

Symbol	Meaning
H	hard clipping (clipped sequences NOT present in SEQ)
N	skipped region from the reference
P	padding (silent deletion from padded reference)

Below are two examples describing the CIGAR string in more detail.

Example 1:

Ref: ACGTACGTACGTACGT

Read: ACGT- - - - ACGTACGA

Cigar: 4M 4D 8M

The first four bases in the read are the same as in the reference, so we can represent these as 4M in the CIGAR string. Next comes 4 deletions, represented by 4D, followed by 7 alignment matches and one alignment mismatch, represented by 8M. Note that the mismatch at position 16 is included in 8M. This is because it still aligns to the reference.

Example 2:

Ref: ACTCAGTG- - GT

Read: ACGCA- TGCAGTtagacgt

Cigar: 5M 1D 2M 2I 2M 7S

Here we start off with 5 alignment matches and mismatches, followed by one deletion. Then we have two more alignment matches, two insertions and two more matches. At the end, we have seven soft clippings, 7S. These are clipped sequences that are present in the SEQ (Query SEQUENCE on the same strand as the reference).

2.3.6 Exercises

Q11: What does the CIGAR from Q9 mean?



Q12: How would you represent the following alignment with a CIGAR string?

Ref: ACGT- - - - ACGTACGT

Read: ACGTACGTACGTACGT



2.3.7 Flags

Column 2 of the alignment contains a combination of bitwise FLAGS describing the alignment. The following table contains the information you can get from the bitwise FLAGS:

Hex	Dec	Flag	Description
0x1	1	PAIRED	paired-end (or multiple-segment) sequencing technology
0x2	2	PROPER_PAIR	each segment properly aligned according to the aligner
0x4	4	UNMAP	segment unmapped
0x8	8	MUNMAP	next segment in the template unmapped
0x10	16	REVERSE	SEQ is reverse complemented
0x20	32	MREVERSE	SEQ of the next segment in the template is reversed
0x40	64	READ1	the first segment in the template
0x80	128	READ2	the last segment in the template
0x100	256	SECONDARY	secondary alignment
0x200	512	QCFAIL	not passing quality controls
0x400	1024	DUP	PCR or optical duplicate
0x800	2048	SUPPLEMENTARY	supplementary alignment

For example, if you have an alignment with FLAG set to 113, this can only be represented by decimal codes $64 + 32 + 16 + 1$, so we know that these four flags apply to the alignment and the alignment is paired-end, reverse complemented, sequence of the next template/mate of the read is reversed and the read aligned is the first segment in the template.

Primary, secondary and supplementary alignments A read that aligns to a single reference sequence (including insertions, deletions, skips and clipping but not direction changes), is a **linear alignment**. If a read cannot be represented as a linear alignment, but instead is represented as a group of linear alignments without large overlaps, it is called a **chimeric alignment**. These can for instance be caused by structural variations. Usually, one of the linear alignments in a chimeric alignment is considered to be the **representative** alignment, and the others are called **supplementary**.

Sometimes a read maps equally well to more than one spot. In these cases, one of the possible alignments is marked as the **primary** alignment and the rest are marked as **secondary** alignments.

2.3.8 BAM

BAM (Binary Alignment/Map) format, is a binary version of SAM. This means that, while SAM is human readable, BAM is only readable for computers. BAM was developed for fast processing and random access. To achieve this, BGZF (Block GZIP) compression is used for indexing. BAM files can be viewed using samtools, and will then have the same format as a SAM file. The key features of BAM are:

- Can store alignments from most mappers
- Supports multiple sequencing technologies
- Supports indexing for quick retrieval/viewing
- Compact size (e.g. 112Gbp Illumina = 116GB disk space)
- Reads can be grouped into logical groups e.g. lanes, libraries, samples
- Widely supported by variant calling packages and viewers

2.3.9 Exercises

Since BAM is a binary format, we can't use the standard UNIX operations directly on this file format. **Samtools** is a set of programs for interacting with SAM and BAM files. Using the samtools view command, print the header of the BAM file:



```
samtools view -H data/NA20538.bam
```

Q13: What version of the human assembly was used to perform the alignments? (Hint: Can you spot this somewhere in the @SQ records?)



Q14: How many lanes are in this BAM file? (Hint: Do you recall what RG represents?)



Q15: What programs were used to create this BAM file? (Hint: have a look for the program record, @PG)



Q16: What version of bwa was used to align the reads? (Hint: is there anything in the @PG record that looks like it could be a version tag?)



The output from running samtools view on a BAM file without any options is a headerless SAM file. This gets printed to STDOUT in the terminal, so we will want to pipe it to something. Let's have a look at the first read of the BAM file:



```
samtools view data/NA20538.bam | head -n 1
```

Q17: What is the name of the first read? (Hint: have a look at the [alignment section](#) if you can't recall the different fields)

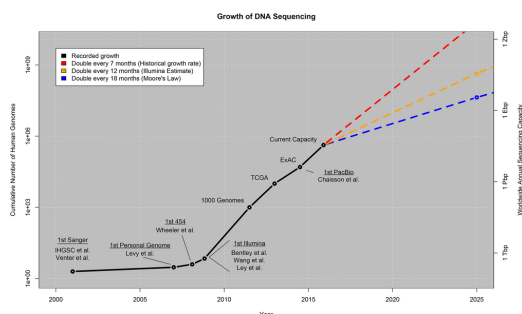


Q18: What position does the alignment of the read start at?



2.4 CRAM

Even though BAM files are compressed, they are still too large. Typically they use 1.5-2 bytes for each base pair of sequencing data that they contain, and while disk capacity is ever improving, increases in disk capacity are being far outstripped by sequencing technologies.



Growth of DNA sequencing

BAM stores all of the data, this includes every read base, every base quality, and it uses a single conventional compression technique for all types of data. Therefore, CRAM was designed for better compression of genomic data than SAM/BAM. CRAM uses three important concepts:

- Reference based compression
- Controlled loss of quality information
- Different compression methods to suit the type of data, e.g. base qualities vs. metadata vs. extra tags

The figure below displays how reference-based compression works. Instead of saving all the bases of all the reads, only the nucleotides that differ from the reference, and their positions, are kept.

Reference sequence: **ACGTACGTACGTACGTACGTACGTACGTAC**
 read 1: ACGTACGTACGTACGTACGT**G**C
 read 2: TACGTACG**C**ACGTACGTGCGTA
 read 3: CGTACG**C**ACGTACGTACGTACG
 read 4: TACGTACGTACGT**G**CGTACGTA
 read 5: CG**C**ACGTACGTACGTACGTACG
 read 6: TACGT**G**CGTACGTACGTAC

Reference sequence: **ACGTACGTACGTACGTACGTACGTACGTAC**
 read 1:**G**..
 read 2:**C**..
 read 3:**C**..
 read 4:**G**..
 read 5:**C**..
 read 6:**G**..

In lossless (no information is lost) mode a CRAM file is 60% of the size of a BAM file, so archives and sequencing centres are now moving from BAM to CRAM.

Since samtools 1.3, CRAM files can be read in the same way that BAM files can. We will look closer at how you can convert between SAM, BAM and CRAM formats in the next section.

2.5 Indexing

To allow for fast random access of regions in BAM and CRAM files, they can be indexed. The files must first be coordinate-sorted. This can be done using **samtools sort**. If no options are supplied, it will by default sort by the left-most position.



```
samtools sort -o data/NA20538_sorted.bam data/NA20538.bam
```

Now we can use **samtools index** to create an index file (.bai) for our sorted BAM file:



```
samtools index data/NA20538_sorted.bam
```

To look for reads mapped to a specific region, we can use **samtools view** and specify the region we are interested in as: `RNAME[:STARTPOS[-ENDPOS]]`. For example, if we wanted to look at all the reads mapped to a region called chr4, we could use:

```
samtools view alignment.bam chr4
```

To look at the region on chr4 beginning at position 1,000,000 and ending at the end of the chromosome, we can do:

```
samtools view alignment.bam chr4:1000000
```

And to explore the 1001bp long region on chr4 beginning at position 1,000 and ending at position 2,000, we can use:

```
samtools view alignment.bam chr4:1000-2000
```

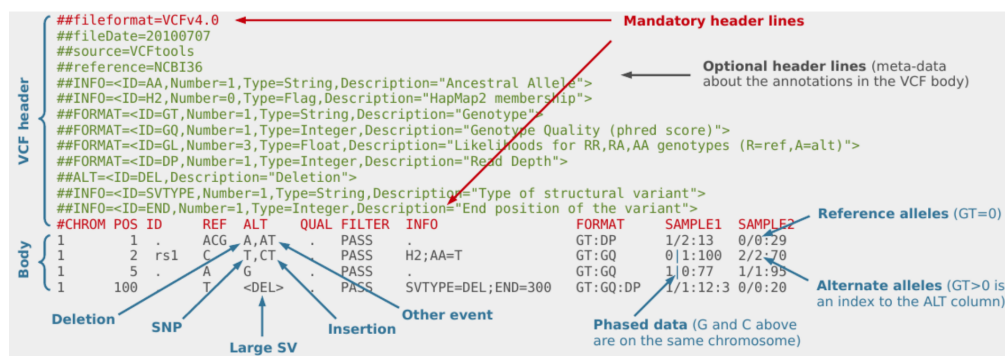
2.5.1 Exercises

Q19: How many reads are mapped to region 20025000-20030000 on chromosome 1?



2.6 VCF/BCF

The VCF file format and its binary version BCF were introduced to store variation data. VCF consists of tab-delimited text and is parsable by standard UNIX commands which makes it flexible and user-extensible. The figure below provides an overview of the different components of a VCF file:



VCF format

2.6.1 VCF header

The VCF header consists of meta-information lines (starting with `##`) and a header line (starting with `#`). All meta-information lines are optional and can be put in any order, except for *fileformat*. This holds the information about which version of VCF is used and must come first.

The meta-information lines consist of key=value pairs. Examples of meta-information lines that can be included are `##INFO`, `##FORMAT` and `##reference`. The values can consist of multiple fields enclosed by `<>`. More information about these fields is available in the [VCF specification](#).

2.6.2 Header line

The header line starts with `#` and consists of 8 required fields:

1. **CHROM:** an identifier from the reference genome
2. **POS:** the reference position
3. **ID:** a list of unique identifiers (where available)
4. **REF:** the reference base(s)
5. **ALT:** the alternate base(s)
6. **QUAL:** a phred-scaled quality score
7. **FILTER:** filter status
8. **INFO:** additional information

If the file contains genotype data, the required fields are also followed by a **FORMAT** column header, and then a number of sample IDs. The **FORMAT** field specifies the data types and order. Some examples of these data types are:

- **GT:** Genotype, encoded as allele values separated by either `/` or `|`
- **DP:** Read depth at this position for this sample
- **GQ:** Conditional genotype quality, encoded as a phred quality

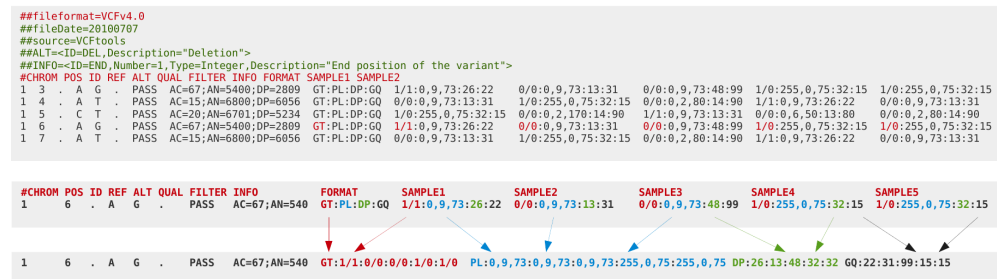
2.6.3 Body

In the body of the VCF, each row contains information about a position in the genome along with genotype information on samples for each position, all according to the fields in the header line.

2.6.4 BCF

VCF can be compressed with BGZF (bgzip) and indexed with TBI or CSI (tabix), but even compressed it can still be very big. For example, a compressed VCF with 3781 samples of human data will be 54 GB for chromosome 1, and 680 GB for the whole genome.

VCFs can also be slow to parse, as text conversion is slow. The main bottleneck is the “FORMAT” fields. For this reason the BCF format, a binary representation of VCF, was developed. In BCF files the fields are rearranged for fast access. The following images show the process of converting a VCF file into a BCF file.



Bcftools comprises a set of programs for interacting with VCF and BCF files. It can be used to convert between VCF and BCF and to view or extract records from a region.

bcftools view Let’s have a look at the header of the file 1kg.bcf in the data directory. Note that bcftools uses **-h** to print only the header, while samtools uses **-H** for this.



```
bcftools view -h data/1kg.bcf
```

Similarly to BAM, BCF supports random access, that is, fast retrieval from a given region. For this, the file must be indexed:



```
bcftools index data/1kg.bcf
```

Now we can extract all records from the region 20:24042765-24043073, using the **-r** option. The **-H** option will make sure we don’t include the header in the output:



```
bcftools view -H -r 20:24042765-24043073 data/1kg.bcf
```

bcftools query The versatile **bcftools query** command can be used to extract any VCF field. Combined with standard UNIX commands, this gives a powerful tool for quick querying of VCFs. Have a look at the usage options:



```
bcftools query -h
```

Let’s try out some useful options. As you can see from the usage, **-l** will print a list of all the samples in the file. Give this a go:



```
bcftools query -l data/1kg.bcf
```

Another very useful option is **-s** which allows you to extract all the data relating to a particular sample. This is a [common option](#) meaning it can be used for many bcftools commands, like `bcftools view`. Try this for sample HG00131:



```
bcftools view -s HG00131 data/1kg.bcf | head -n 50
```

The format option, **-f** can be used to select what gets printed from your query command. For example, the following will print the position, reference base and alternate base for sample HG00131, separated by tabs:



```
bcftools query -f'%POS\t%REF\t%ALT\n' -s HG00131 data/1kg.bcf | head
```

Finally, let's look at the **-i** option. With this option we can select only sites for which a particular expression is true. For instance, if we only want to look at sites that have at least 2 alternate alleles, we can use the following expression (piped to `head` to only show a subset of the output):



```
bcftools query -f'%CHROM\t%POS\n' -i 'AC[0]>2' data/1kg.bcf | head
```

We use **-i** with the expression `AC[0]>2`. `AC` is an info field that holds the `__a__llele __c__ount`. Some fields can hold multiple values, so we use `[0]>2` to indicate that we are looking for the first value (this is zero indexed, and hence starts at 0 instead of 1), and that this value should be greater than 2. To format our output, we use **-f** to specify that we only want to print the chromosome name and the position.

There is more information about expressions on the [bcftools manual page](#)

2.6.5 Exercises

Now, try and answer the following questions about the file `1kg.bcf` in the data directory. For more information about the different usage options you can open the [bcftools query manual page](#) in a new tab.

Q20: What version of the human assembly do the coordinates refer to?



Q21: How many samples are there in the BCF?



Q22: What is the genotype of the sample HG00107 at the position 20:24019472? (Hint: use the combination of **-r, **-s**, and **-f** options)**



Q23: How many positions are there with more than 10 alternate alleles? (Hint: use the `-i` filtering option)



Q24: In how many positions does HG00107 have a non-reference genotype and a read depth bigger than 10? (Hint: you can use pipes to combine bcftools queries)



2.6.6 gVCF

Often it is not enough to know variant sites only. For instance, we don't know if a site was dropped because it matches the reference or because the data is missing. We sometimes need evidence for both variant and non-variant positions in the genome. In gVCF format, blocks of reference-only sites can be represented in a single record using the "INFO/END" tag. Symbolic alleles (`<*>`) are used for incremental calling:

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	Sample
19	9902	.	G	<*>	.	.	END=9915;MinDP=0	PL:DP	0,0,0:0
19	9916	.	C	<*>	.	.	END=9922;MinDP=5	PL:DP	0,15,137:5
19	9923	.	G	<*>	.	.	END=9948;MinDP=10	PL:DP	0,30,214:10
19	9949	.	G	A,<*>	.	.	DP=28	PL:DP	0,60,255,78,255,255:27
19	9950	.	C	<*>	.	.	END=9958;MinDP=28	PL:DP	0,84,255:28
19	9959	.	G	T,<*>	.	.	DP=34	PL:DP	0,82,255,99,255,255:34
19	9960	.	C	<*>	.	.	END=9969;MinDP=34	PL:DP	0,102,255:34

Symbolic "unobserved" allele
Represents any other possible alternate allele

A block of 10 sites with
at least 34 reference reads

Genotype likelihoods
for CC, C*, **

gVCF

2.6.7 Exercises

Q25: In the above example, what is the size of the reference-only block starting at position 9923?



Q26: For the same block, what is the first base?



Q27: How many reference reads does the block have?



The answers to the questions on this page can be found [here](#).

Now continue to the next section of the tutorial: [File conversion](#).

You can also return to the [index page](#).

3 File conversion

In this section we are going to look at how to convert from one file format to another. There are many tools available for converting between file formats, and we will use some of the most common ones: samtools, bcftools and Picard.

3.1 SAM to BAM

To convert from SAM to BAM format we are going to use the **samtools view** command. In this instance, we would like to include the SAM header, so we use the **-h** option:



```
samtools view -h data/NA20538.bam > data/NA20538.sam
```

Now, have a look at the first ten lines of the SAM file. They should look like they did in the previous section when you viewed the BAM file header.



```
head data/NA20538.sam
```

Well that was easy! And converting SAM to BAM is just as straightforward. This time there is no need for the **-h** option, however we have to tell samtools that we want the output in BAM format. We do so by adding the **-b** option:



```
samtools view -b data/NA20538.sam > data/NA20538_2.bam
```

Samtools is very well documented, so for more usage options and functions, have a look at the [samtools manual](#).

3.2 BAM to CRAM

From samtools version 1.3, support for CRAM format was introduced. This means that the samtools view command can also be used to convert a BAM file to CRAM format. In the data directory there is a BAM file called yeast.bam that was created from *S. cerevisiae* Illumina sequencing data. There is also a reference genome in the directory, called *Saccharomyces_cerevisiae*.EF4.68.dna.toplevel.fa. For the conversion, an index file (.fai) must first be created. This can be done using **samtools faidx**. However, as we will see, samtools will generate this file on the fly when we specify a reference file using the **-F** option.

To convert to CRAM, we use the **-C** option to tell samtools we want the output as CRAM, and the **-T** option to specify what reference file to use for the conversion. We also use the **-o** option to specify the name of the output file. Give this a go:



```
samtools view -C \  
    -T data/Saccharomyces_cerevisiae.EF4.68.dna.toplevel.fa \  
    -o data/yeast.cram data/yeast.bam
```

Have a look at what files were created:



```
ls -l data
```

As you can see, this has created an index file for the reference genome called `Saccharomyces_cerevisiae.EF4.68.dna.toplevel.fa.fai` and the CRAM file `yeast.cram`.

Q1: Since CRAM files use reference-based compression, we expect the CRAM file to be smaller than the BAM file. What is the size of the CRAM file?



Q2: Is your CRAM file smaller than the original BAM file?



```
ls -l data
```

To convert CRAM back to BAM, simply change `-C` to `-b` and change places for the input and output CRAM/BAM:

```
samtools view -b -T data/Saccharomyces_cerevisiae.EF4.68.dna.toplevel.fa \
-o data/yeast.bam data/yeast.cram
```

3.3 FASTQ to SAM

As mentioned in the previous section of this tutorial, SAM format is mainly used to store alignment data. However, in some cases we may want to store unaligned data in SAM format and for this we can use the picard tools **FastqToSam** application. Picard tools is a Java application that comes with a number of useful options for manipulating high-throughput sequencing data. Apart from FASTQ to SAM, we won't go into any detail about Picard tools in this tutorial, but feel free to explore it on the [Picard tools website](#). To convert the FASTQ files of lane 13681_1#18 to unaligned SAM format, run:



```
java -jar $PICARD FastqToSam F1=data/13681_1#18_1.fastq.gz \
F2=data/13681_1#18_2.fastq.gz \
O=data/13681_1#18.sam SM=13681_1#18
```

Where `$PICARD` should contain the path to the `picard.jar` file, as described on the [index page](#).

From here you can go on and convert the SAM file to BAM and CRAM, as described previously. There are also multiple options for specifying what metadata to include in the SAM header. To see all available options, run:



```
java -jar $PICARD FastqToSam -h
```

3.4 CRAM to FASTQ

Although it is possible to convert CRAM to FASTQ directly using the `samtools fastq` command, for many applications we need the fastq files to be ordered. For this reason, we will first use `samtools collate`, which will produce a collated BAM file. The reference file and its index file that was created when we converted BAM to CRAM is required for this as well.



```
samtools collate data/yeast.cram data/yeast.collated
```

The newly produced BAM file will be called `yeast.collated.bam`. Let's use this to create two FASTQ files, one for the forward reads and one for the reverse reads:



```
samtools fastq -1 data/yeast.collated_1.fastq \  
-2 data/yeast.collated_2.fastq data/yeast.collated.bam
```

For further information and usage options, have a look at the [samtools manual page](#).

3.5 VCF to BCF

As we saw in the previous section, `bcftools` comprises a set of programs for interacting with VCF/BCF files. In a similar way that `samtools view` can be used to convert between SAM, BAM and CRAM, **`bcftools view`** can be used to convert between VCF and BCF. To convert the file called `1kg.bcf` to a compressed VCF file called `1kg.vcf.gz`, run:



```
bcftools view -O z -o data/1kg.vcf.gz data/1kg.bcf
```

The `-O` option allows us to specify in what format we want the output, compressed BCF (b), uncompressed BCF (u), compressed VCF (z) or uncompressed VCF (v). With the `-o` option we can select the name of the output file.

Have a look at what files were generated (the options `-lrt` will list the files in reverse chronological order):



```
ls -lrt data
```

As you can see, this also generated an index file, `1kg.bcf.csi`.

To convert a VCF file to BCF, we can run a similar command. If we want to keep the original BCF, we need to give the new one a different name so that the old one is not overwritten:



```
bcftools view -O b -o data/1kg_2.bcf data/1kg.vcf.gz
```

The answers to the questions on this page can be found [here](#).

Now continue to the next section of the tutorial: [QC assessment of NGS data](#).

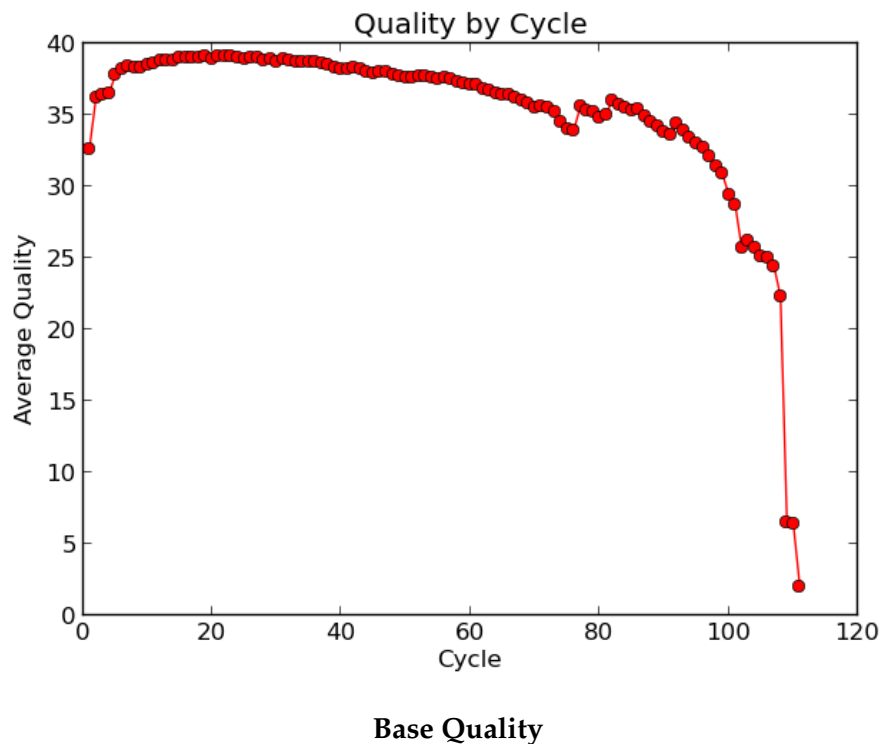
You can also return to the [index page](#).

4 QC assessment of NGS data

As mentioned previously, QC is an important part of any analysis. In this section we are going to look at some of the metrics and graphs that can be used to assess the QC of NGS data.

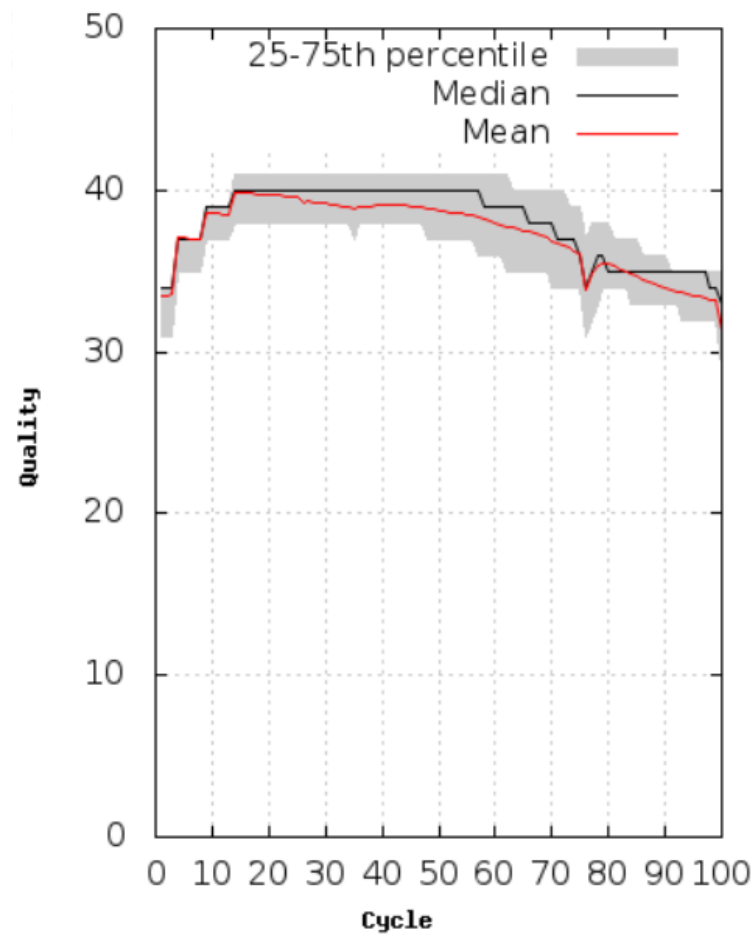
4.1 Base quality

[Illumina sequencing](#) technology relies on sequencing by synthesis. One of the most common problems with this is **dephasing**. For each sequencing cycle, there is a possibility that the replication machinery slips and either incorporates more than one nucleotide or perhaps misses to incorporate one at all. The more cycles that are run (i.e. the longer the read length gets), the greater the accumulation of these types of errors gets. This leads to a heterogeneous population in the cluster, and a decreased signal purity, which in turn reduces the precision of the base calling. The figure below shows an example of this.

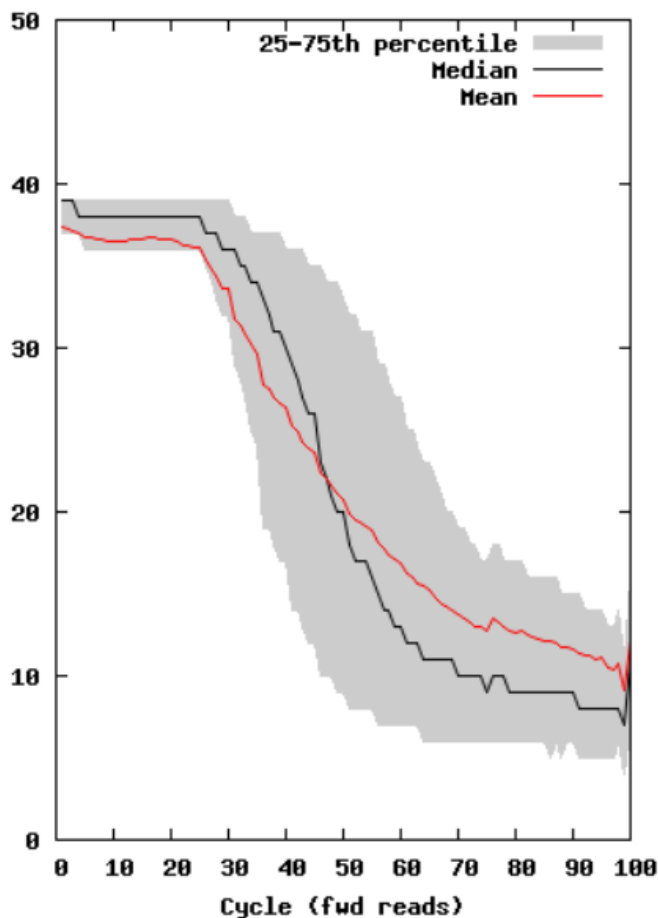


Because of dephasing, it is possible to have high-quality data at the beginning of the read but really low-quality data towards the end of the read. In those cases you can decide to trim off the low-quality reads, for example using a tool called [Trimmomatic](#).

The figures below shows an example of a high-quality read data (top) and a poor quality read data (bottom).



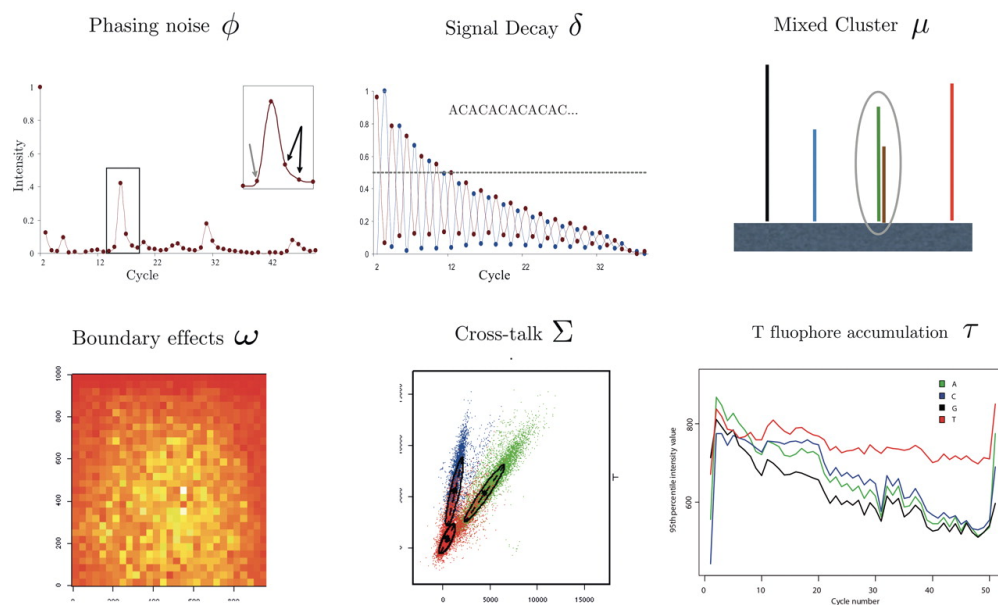
High-quality read data



Poor quality read data

4.2 Other base calling errors

There are several different reasons for a base to be called incorrectly, as shown in the figure below. **Phasing noise** and **signal decay** is a result of the dephasing issue described above. During library preparation, **mixed clusters** can occur if multiple templates get co-located. These clusters should be removed from the downstream analysis. **Boundary effects** occur due to optical effects when the intensity is uneven across each tile, resulting in lower intensity toward the edges. **Cross-talk** occurs because the emission frequency spectra for each of the four dyes partly overlap. Finally, for early chemistries **T fluorophore accumulation** was an issue, where incomplete cleavage of the dye coupled to thymine lead to an accumulation the nucleotide.

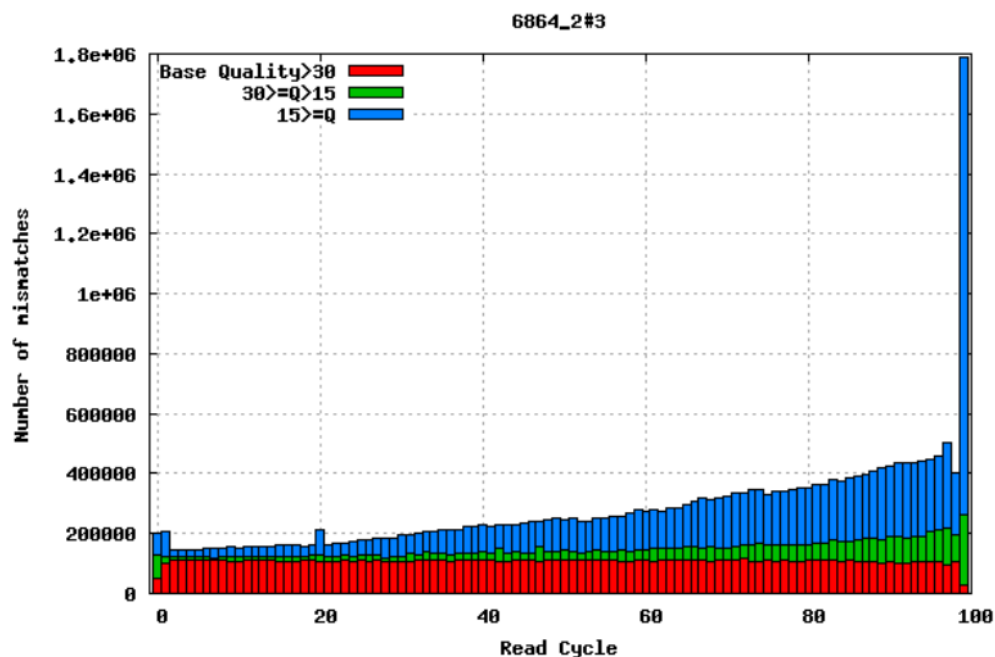


Base Calling Errors

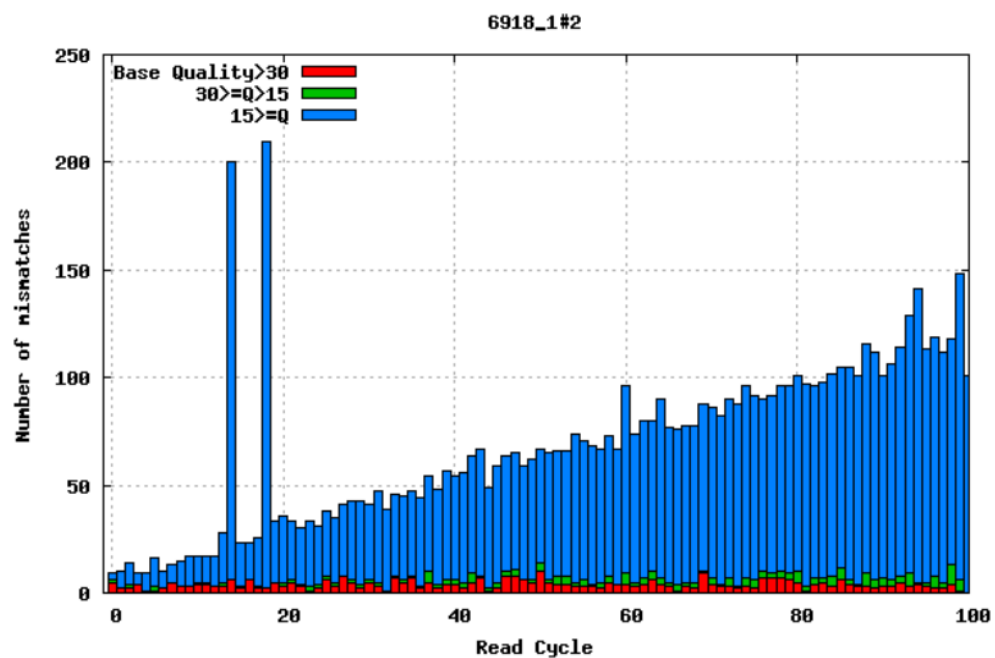
Base-calling for next-generation sequencing platforms, doi: [10.1093/bib/bbq077](https://doi.org/10.1093/bib/bbq077)

4.3 Mismatches per cycle

Aligning reads to a high-quality reference genome can provide insight to the quality of a sequencing run by showing you the mismatches to the reference sequence. This can help you detect cycle-specific errors. Mismatches can occur due to two main causes, sequencing errors and differences between your sample and the reference genome, which is important to bear in mind when interpreting mismatch graphs. The figures below show an example of a good run (top) and a bad one (bottom). In the first figure, the distribution of the number of mismatches is even between the cycles, which is what we would expect from a good run. However, in the second figure, two cycles stand out with a lot of mismatches compared to the other cycles.



Good run

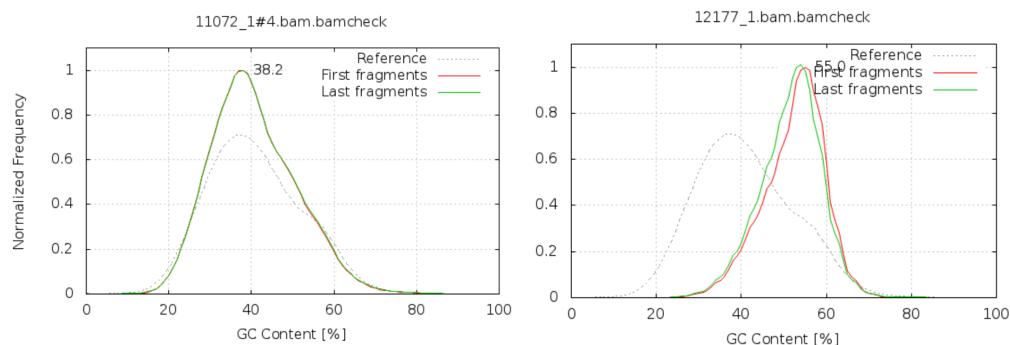


Poor run

4.4 GC bias

It is a good idea to compare the GC content of the reads against the expected distribution in a reference sequence. The GC content varies between species, so a shift in GC content like the one

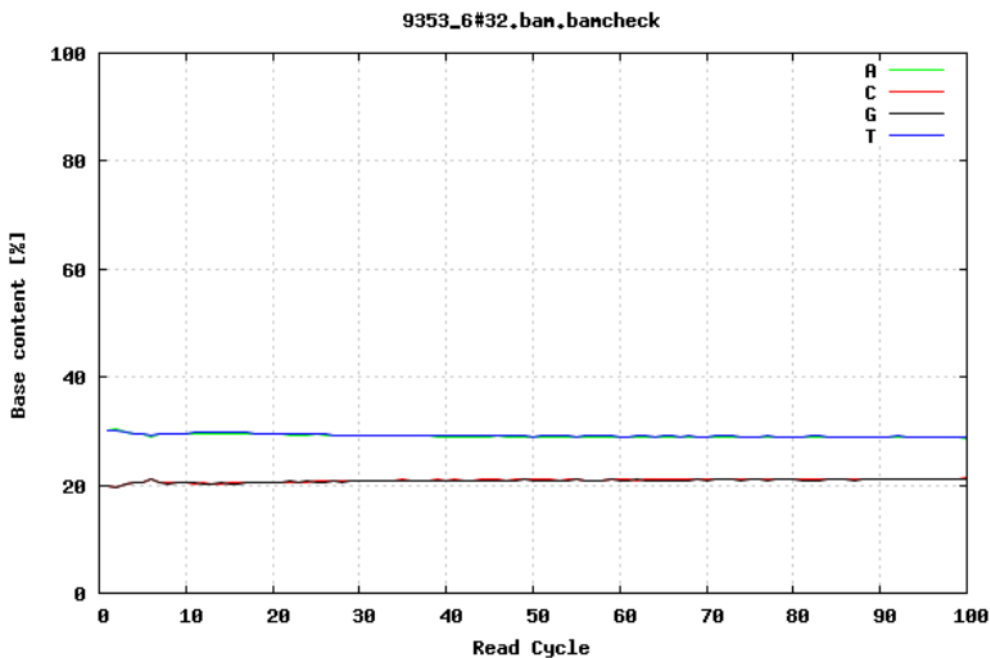
seen below could be an indication of sample contamination. In the left image below, we can see that the GC content of the sample is about the same as for the reference, at ~38%. However, in the right figure, the GC content of the sample is closer to 55%, indicating that there is an issue with this sample.



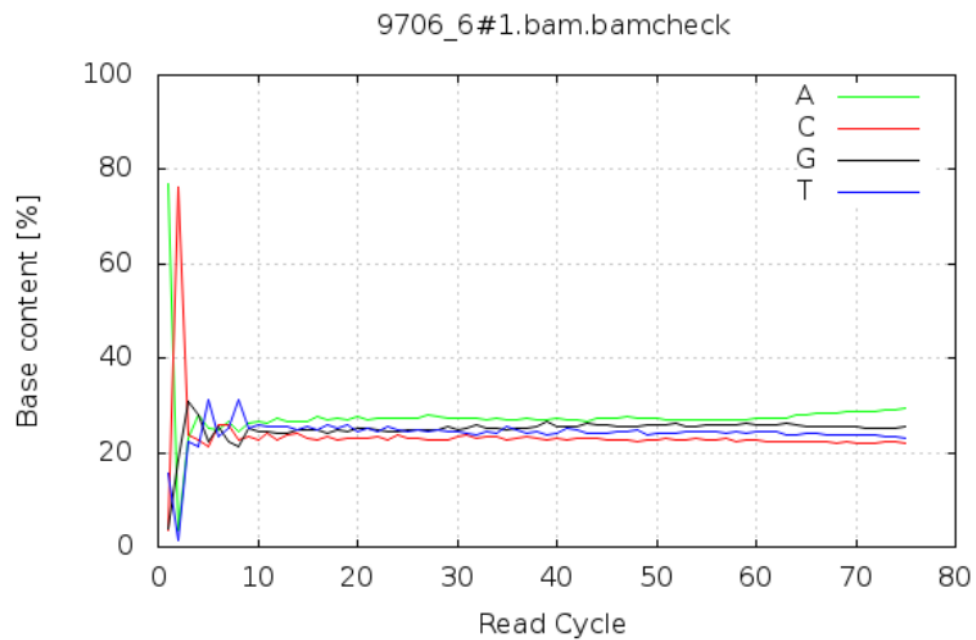
GC Bias

4.5 GC content by cycle

Looking at the GC content per cycle can help detect if the adapter sequence was trimmed. For a random library, it is expected to be little to no difference between the different bases of a sequence run, so the lines in this plot should be parallel with each other like in the first of the two figures below. In the second of the figures, the initial spikes are likely due to adapter sequences that have not been removed.



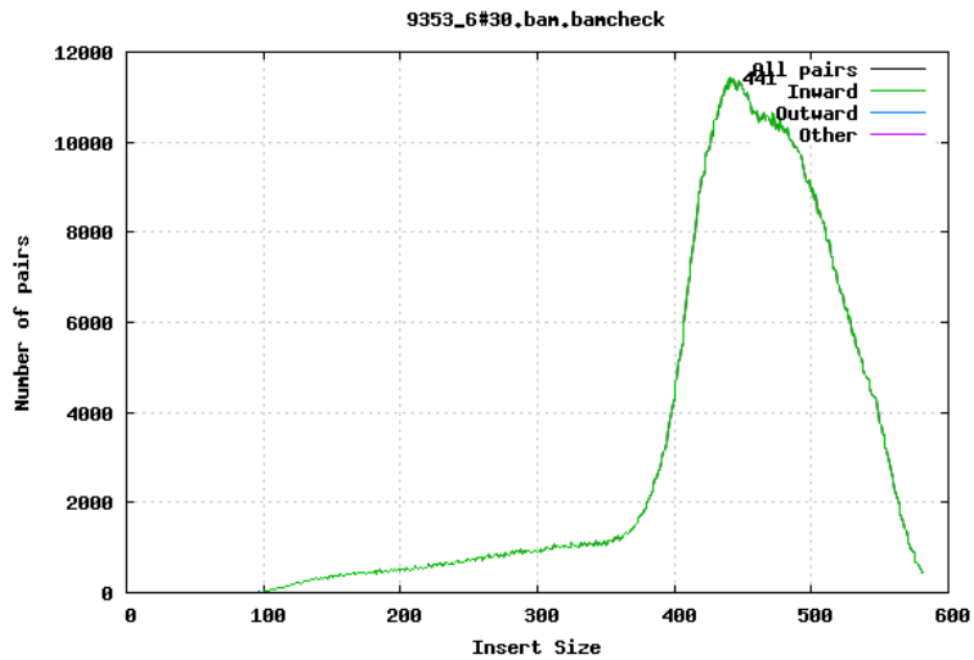
Good run



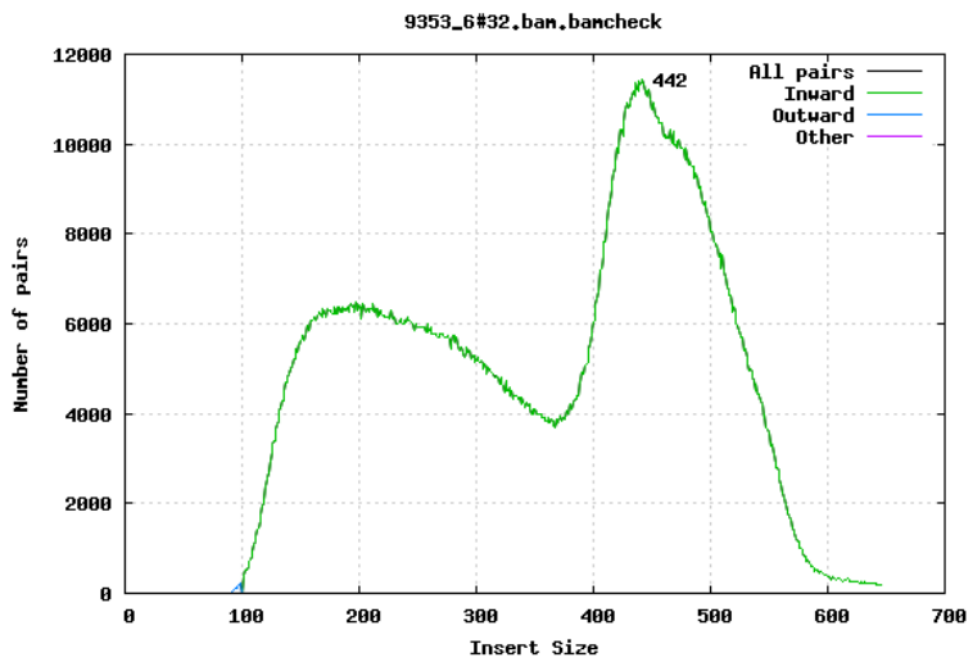
Poor run

4.6 Insert size

For paired-end sequencing the size of DNA fragments also matters. In the first of the examples below, the insert size peaks around 440 bp. In the second however, there is also a peak at around 200 bp. This indicates that there was an issue with the fragment size selection during library prep.



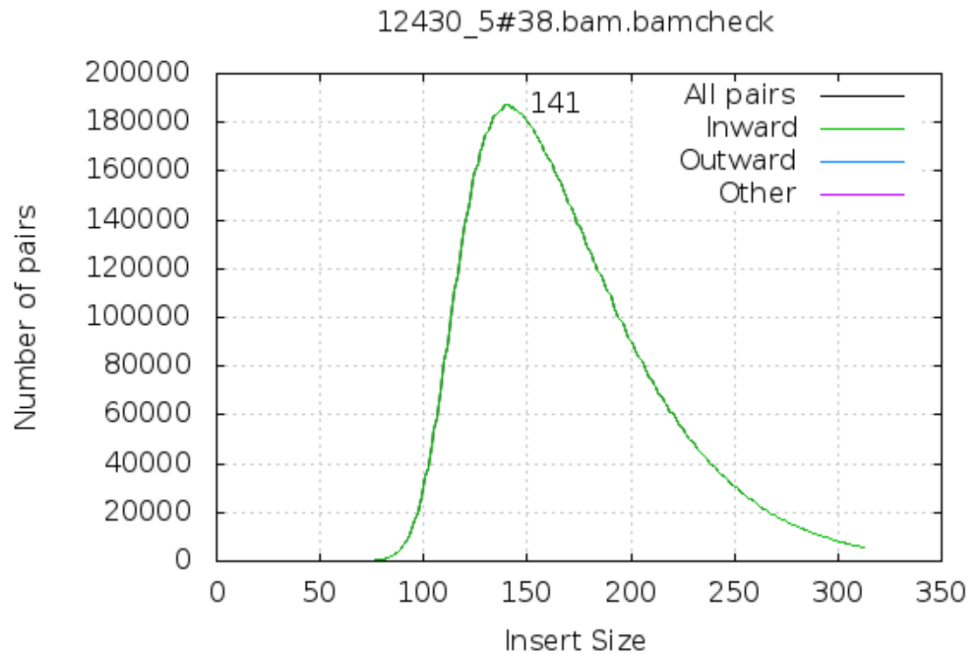
Good run



Poor run

4.6.1 Exercises

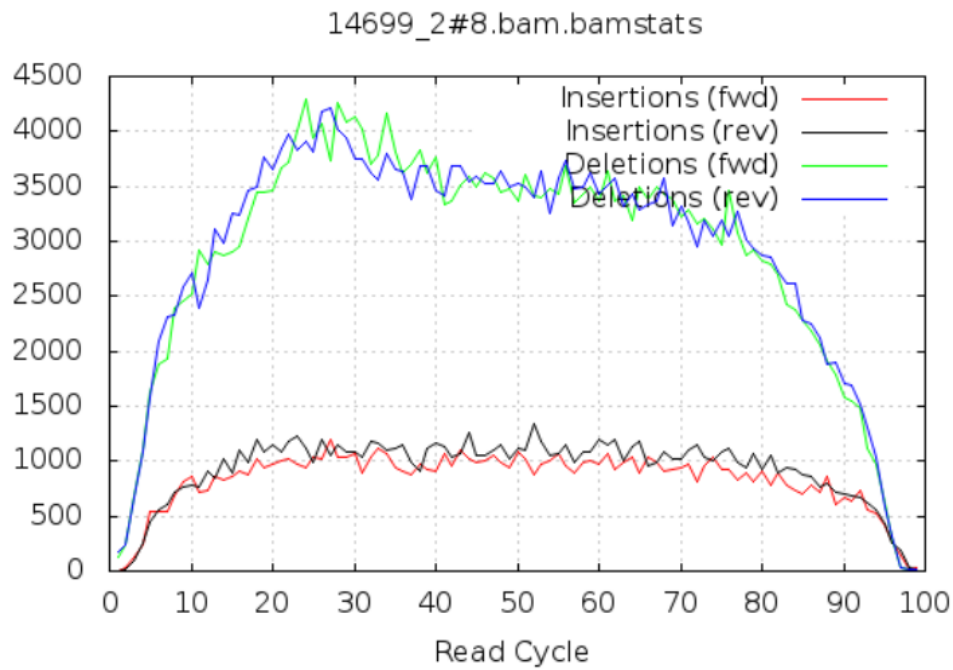
Q1: The figure below is from a 100bp paired-end sequencing. Can you spot any problems?



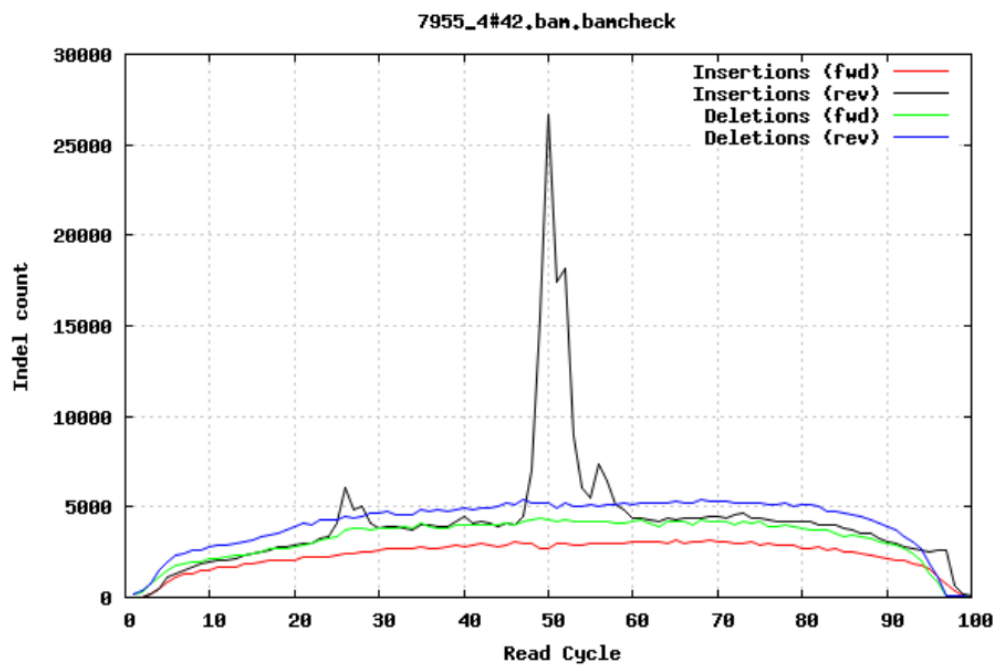
Q1

4.7 Insertions/Deletions per cycle

Sometimes, air bubbles occur in the flow cell, which can manifest as false indels. The spike in the second image provides an example of how this can look.



Good run



Poor run

4.8 Generating QC stats

Now let's try this out! We will generate QC stats for two lanes of Illumina paired-end sequencing data from yeast. We will use the bwa mapper to align the data to the [Saccromyces cerevisiae genome](#), followed by samtools stats to generate the stats.

Read pairs are usually stored in two separate FASTQ files so that n-th read in the first file and the n-th read in the second file constitute a read pair. Can you devise a quick sanity check that reads in these two files indeed form pairs? The files must have the same number of lines and the naming of the reads usually suggests if they form a pair. The location of the files is:

```
data/lane1/s_7_1.fastq
data/lane1/s_7_2.fastq
```



Let's have a look at the script we are going to run to create the mappings:



```
cat create_mapping.sh
```

The script contains several commands, some are combined together using pipes. (UNIX pipes is a very powerful and elegant concept which allows us to feed the output of one command into the next command and avoid writing intermediate files. If you are not comfortable with UNIX, consider having a go at the UNIX tutorial).

Now run the script to create the mappings and stats:



```
./create_mapping.sh
```

The script will produce the BAM file lane1.sorted.bam and a matching index file:



```
ls -alrt data
```

Now we will use **samtools stats** to generate the stats for the primary alignments. The option **-f** can be used to filter reads with specific tags, while **-F** can be used to *filter out* reads with specific tags. The following command will include only primary alignments:



```
samtools stats -F SECONDARY data/lane1.sorted.bam \
> data/lane1.sorted.bam.bchk
```

Have a look at the first 41 lines of the statistics file that was generated:



```
head -n 41 data/lane1.sorted.bam.bchk
```

This file contains a number of useful stats that we can use to get a better picture of our data, and it can even be plotted with **plot-bamstats**, as you will see soon. First let's have a closer look at some of the different stats. Each part of the file starts with a # followed by a description of the section and how to extract it from the file. Let's have a look at all the sections in the file:



```
grep ^'#' data/lane1.sorted.bam.bchk | grep 'Use'
```


4.8.1 Summary Numbers (SN)

This initial section contains a summary of the alignment and includes some general statistics. In particular, you can see how many bases mapped, and how much of the genome that was covered.

Now look at the output and try to answer the questions below.

Q2: What is the total number of reads?



Q3: What proportion of the reads were mapped?



Q4: How many pairs were mapped to a different chromosome?



Q5: What is the insert size mean and standard deviation?



Q6: How many reads were paired properly?



Finally, we will create some QC plots from the output of the stats command using the command **plot-bamstats** which is included in the samtools package:



```
plot-bamstats -p data/lane1-plots/ data/lane1.sorted.bam.bchk
```

Now in your web browser open the file lane1-plots/index.html to view the QC information.

Q7: How many reads have zero mapping quality?



Q8: Which of the first fragments or second fragments are higher base quality on average?



The answers to the questions on this page can be found [here](#).

Now continue to the next section of the tutorial: [Identifying contamination](#)
Alternatively, you can return to the [previous section](#) or the [index page](#).

5 Identifying contamination

It is always a good idea to check that your data is from the species you expect it to be. A very useful tool for this is [Kraken](#). In this tutorial we will go through how you can use Kraken to check your samples for contamination.

Note if using the Sanger cluster: Kraken is run as part of the automatic qc pipeline and you can retrieve the results using the `pf qc` script. For more information, run `pf man qc`.

5.1 Setting up a database

To run Kraken you need to either build a database or download an existing one. The standard database is very large (33 GB), but thankfully there are some smaller, pre-built databases available. To download the smallest of them, the 4 GB MiniKraken. If you don't already have a kraken database set up, run:



```
wget https://ccb.jhu.edu/software/kraken/dl/minikraken_20171019_4GB.tgz
```

Then all you need to do is un-tar it:



```
tar -zxvf minikraken_20171019_4GB.tgz
```

This version of the database is constructed from complete bacterial, archaeal, and viral genomes in RefSeq, however it contains only around 3 percent of the kmers from the original kraken database (more information [here](#)). If the pre-packaged databases are not quite what you are looking for, you can create your own customized database instead. Details about this can be found [here](#).

Note if using the Sanger cluster: There are several pre-built databases available centrally on the Sanger cluster. For more information, please contact the Pathogen Informatics team.

5.2 Running Kraken

To run Kraken, you need to provide the path to the database you just created. By default, the input files are assumed to be in FASTA format, so in this case we also need to tell Kraken that our input files are in FASTQ format, gzipped, and that they are paired end reads:



```
kraken -db ./minikraken_20171013_4GB -output kraken_results \
-fastq-input -gzip-compressed -paired \
data/13681_1#18_1.fastq.gz data/13681_1#18_2.fastq.gz
```

The five columns in the file that's generated are:

1. "C"/"U": one letter code indicating that the sequence was either classified or unclassified.
2. The sequence ID, obtained from the FASTA/FASTQ header.
3. The taxonomy ID Kraken used to label the sequence; this is 0 if the sequence is unclassified.

4. The length of the sequence in bp.
5. A space-delimited list indicating the LCA mapping of each k-mer in the sequence.

To get a better overview you can create a kraken report:



```
kraken-report -db ./minikraken_20171013_4GB -print_header \  
kraken_results > kraken-report
```

5.3 Looking at the results

Let's have a closer look at the kraken_report for the sample. If for some reason your kraken-run failed there is a prebaked kraken-report at data/kraken-report



```
head -n 20 kraken-report
```

The six columns in this file are:

1. Percentage of reads covered by the clade rooted at this taxon
2. Number of reads covered by the clade rooted at this taxon
3. Number of reads assigned directly to this taxon
4. A rank code, indicating (U)nclassified, (D)omain, (K)ingdom, (P)hylum, (C)lass, (O)rder, (F)amily, (G)enus, or (S)pecies. All other ranks are simply '-'.
5. NCBI taxonomy ID
6. Scientific name

5.4 Exercises

Q1: What is the most prevalent species in this sample?



Q2: Are there clear signs of contamination?



Q3: What percentage of reads could not be classified?



5.5 Heterozygous SNPs

For bacteria, another thing that you can look at to detect contamination is if there are heterozygous SNPs in your samples. Simply put, if you align your reads to a reference, you would expect any snps to be homozygous, i.e. if one read differs from the reference genome, then the rest of the reads that map to that same location will also do so:

Homozygous SNP

Ref: CTTGAGACGAAATCACTAAAAAACGTGACGACTTG
Read1: CTTGAGtCG
Read2: CTTGAGtCGAAA
Read3: GAGtCGAAATCACTAAAA
Read4: GtCGAAATCA

But if there is contamination, this may not be the case. In the example below, half of the mapped reads have the T allele and half have the A.

Heterozygous SNP

Ref: CTTGAGACGAAATCACTAAAAAACGTGACGACTTG
Read1: CTTGAGtCG
Read2: CTTGAGaCGAAA
Read3: GAGaCGAAATCACTAAAA
Read4: GtCGAAATCA

Note if using the Sanger cluster: Heterozygous SNPs are calculated as part of the automated QC pipeline. The result for each lane is available in the file heterozygous_snps_report.txt.

Congratulations! You have reached the end of this tutorial. You can find the answers to all the questions of the tutorial [here](#).

To revisit the previous section [click here](#). Alternatively you can head back to the [index page](#)