# Read mapping and alignment: SAM format, alignment, IGV

Luc Dehaspe
Erika Souche

GENOMICS
CORELEUVEN

# Overview



Sequencing

fastq — Sequence of nucleotides within a read?

bam — Genomic location of a read?

vcf — Genotype of a genomic location?

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping

look up read in reference sequence ….

… allowing mismatches

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence



DNA

Fragment

*TTACATGCATCT*

# Sequence Alignment & Mapping
## look up read in reference sequence



DNA

Fragment

*T T A C A T G C A T C T*

*A G A T G C A T C T A A*

Read1 **T T A C** *A T G C A T C T*
*A G A T G C A T C T A A*

TTAC
AATG

Read2 *T T A C A T G C A T C T*
*A G A T G C A T C T A A*

AGAT
ATCT

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence



DNA

Fragment

*TTACATGCATCT*

*AGATGCATGAA*

Read1 **TTAC** *ATGCATCT*

Read2 *TTACATGCATCT*

TTAC

AGAT

Reference

```
      1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16 17
chrN  C  G  T  T  A  C  A  T  G  C  A  T  C  T  A  A  G
```

TTAC

ATCT

GENOMICS
CORE LEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence



DNA

Fragment

*TTACATGCATCT*

Read1 **TTAC**ATGCATCT

Read2 TTACATGCATCT **AGAT**

Sam-like minimal description

```
Read1 chrN Forward 3
Read2 chrN Reverse 11
```

Reference

1 2 **3** 4 5 6 7 8 9 10 **11** 12 13 14 15 16 17

chrN CGTTACATGCATCTAAG

TTAC     ATCT

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

Fragment

*A T G C A T C T A A G*

Read3 ***A T G C*** *A T C T A A G*

ATGC

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|

Reference   chrN   C G T T A C A T G C A T C T A A G

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

Fragment

$A\ T\ G\ C\ A\ T\ C\ T\ A\ A\ G$

$T\ A\ C\ G\ T\ A\ G\ A\ T\ T\ C$

Read3  $A\ T\ G\ C$ $A\ T\ C\ T\ A\ A\ G$

$T\ A\ C\ G\ T\ A\ G\ A\ T\ T\ C$

| Read3 | chrN | Forward | 7 |

Or ?

| Read3 | chrN | Reverse | 9 |

Reference

 1  2  3  4  5  6  **7**  8  **9**  10 11 12 13 14 15 16 17

chrN  **C G T T A C A T G C A T C T A A G**

No unique mapping

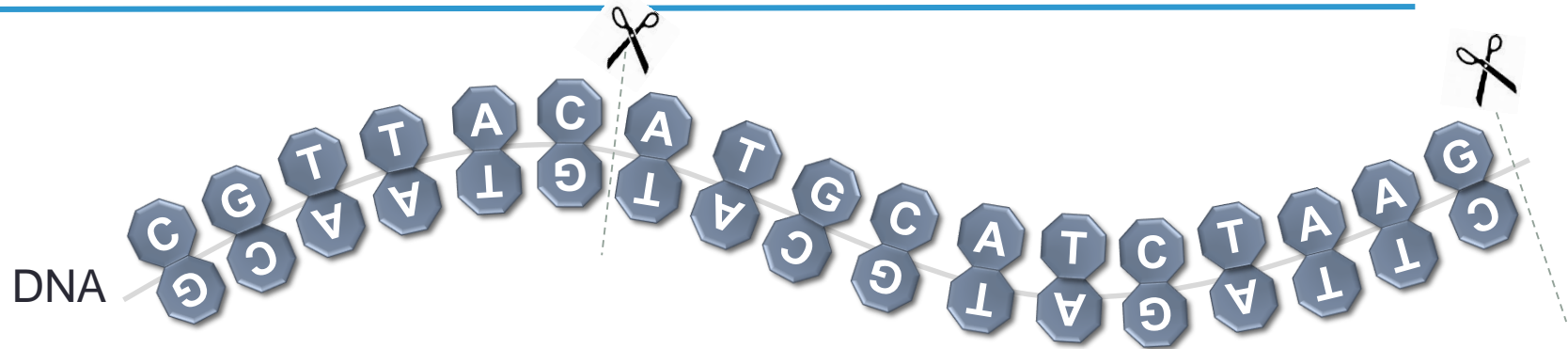| Read3 | chrN | Forward | 7 | **lowMQual** |

How can we improve mapping quality?

ATGC

GCAT

GCAT

ATGC

# Sequence Alignment & Mapping
## look up read in reference sequence



DNA

Fragment

*ATGCATCTAAG*

Read4 **ATGCATC** *TAAG*

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|

Reference    chrN  C G T T A C A T G C A T C T A A G

How can we improve mapping quality?

Strategy 1: READ ON

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

Fragment

$A\ T\ G\ C\ A\ T\ C\ T\ A\ A\ G$

Read4  **A T G C A T C** *T A A G*

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Reference chrN | C | G | T | T | A | C | A | T | G | C | A | T | C | T | A | A | G |

G A T G C A T      **No mapping**

How can we improve mapping quality?

Strategy 1: READ ON

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

Fragment

*A T G C A T C T A A G*
*ꓕ A Ɔ ꓨ ꓕ A ꓨ A ꓕ ꓕ Ɔ*

Read4 ***A T G C A T C*** *T A A G*
*Ɔ ꓕ ꓕ A ꓨ ꓕ A Ɔ ꓨ ꓕ A*

| 1 | 2 | 3 | 4 | 5 | 6 | **7** | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |

Reference   chrN  C G T T A C A T G C A T C T A A G

How can we improve mapping quality?

A T G C A T C
ꓕ A Ɔ ꓨ ꓕ A ꓨ          Unique mapping

Strategy 1: READ ON

Read3 chrN   Forward 7 **highMQual**

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

Fragment

$A\ T\ G\ C\ A\ T\ C\ T\ A\ A\ G$
$T\ A\ C\ G\ T\ A\ G\ A\ T\ T\ C$

Read5a **A T G C** A T C T A A G
Read5b C T A G A T G C A T **C T T A**

Reference

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|

chrN C G T T A C A T G C A T C T A A G

How can we improve mapping quality?

C T T A — G C A T

G A A T — A T G C

**No mapping**

Strategy 2: READ BOTH SIDES

GENOMICS
CORELEUVEN

# Sequence Alignment & Mapping
## look up read in reference sequence

DNA

Fragment

$A\ T\ G\ C\ A\ T\ C\ T\ A\ A\ G$

Read5a **A T G C** A T C T A A G

Read5b G A T G C A T **C T T A**

Reference

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chrN | C | G | T | T | A | C | A | T | G | C | A | T | C | T | A | A | G |

How can we improve mapping quality?

Unique mapping

Strategy 2: READ BOTH SIDES

# Sequence Alignment & Mapping
## look up read in reference sequence



DNA

Fragment

$A\ T\ G\ C\ A\ T\ C\ T\ A\ A\ G$

Read5a  **A T G C** *A T C T A A G*

Read5b  *G A T G C A T* **C T T A**

Insert size = fragment size

|  1 | 2 | 3 | 4 | 5 | 6 | **7** | 8 | 9 | 10 | 11 | 12 | 13 | **14** | 15 | 16 | 17 |

Reference  chrN  `C G T T A C A T G C A T C T A A G`

Unique mapping

ATGC     TAAG

How can we improve mapping quality?

Strategy 2: READ BOTH SIDES

**Read5a** chrN F,**Paired,PairMapped,First**      7 hiMQ **chrN** 14  11
**Read5b** chrN R,**Paired,PairMapped,Second**    14 hiMQ **chrN**  7 -11

# Sequence Alignment & Mapping
# look up read in reference sequence



|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chrN | C | G | T | T | A | C | A | T | G | C | A | T | C | T | A | A | G |

Read5a  *A T G C* A T C T A A G
Read5b  T C G T A G C A T  *C T T A*

Read6a  *T T A G* A T G C A T G T A A C
Read6b  A A T C T A C G T A C  *G T T A*

Read7a  *C G T T* A C A T G C A T C
Read7b  G C A A T G T A C  *G A T G*

```
Read5a  chrN  F,Paired,PairMapped,First    7 hiMQ chrN 14   11
Read5b  chrN  R,Paired,PairMapped,Second  14 hiMQ chrN  7  -11

Read6a  chrN  R,Paired,PairMapped,Second  13 hiMQ chrN  2  -15
Read6b  chrN  F,Paired,PairMapped,First    2 hiMQ chrN 13   15

Read7a  chrN  F,Paired,PairMapped,First    1 hiMQ chrN 10   13
Read7b  chrN  R,Paired,PairMapped,Second  10 hiMQ chrN  1  -13
```

# Sequence Alignment & Mapping allowing mismatches



Variant

DNA

Fragment

*T T G C A T G C A T C T*

Read8 **T T G C** *A T G C A T C T*

| 1 | 2 | **3** | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

chrN C G T T A C A T G C A T C T A A G

T T G C

4 Matches or Mismatches

Read8 chrN  Forward 3 hiMQ **4M** **3A>G**
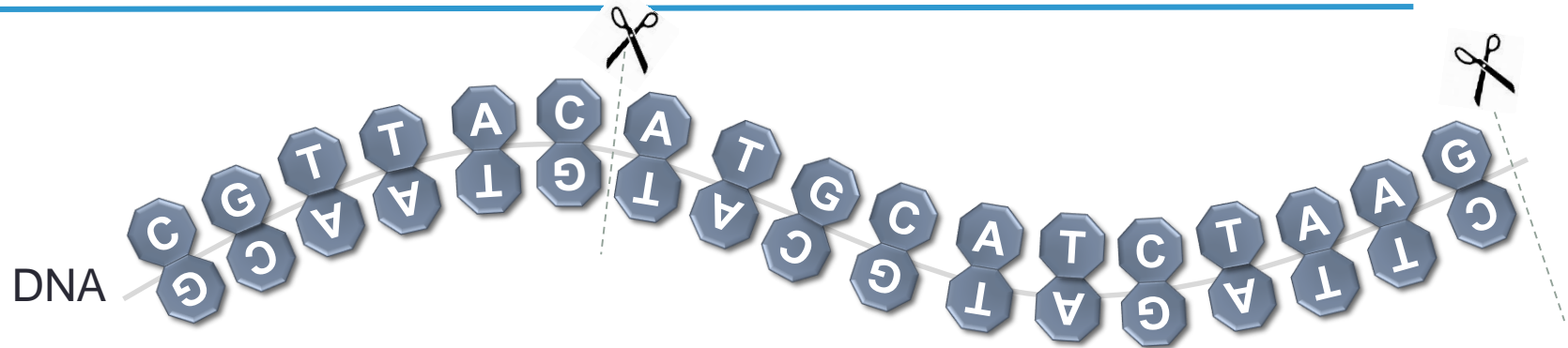
CIGAR string
Compact Ideosyncratic Gapped Alignment Report

GENOMICS CORELEUVEN

# Sequence Alignment & Mapping allowing mismatches

# Sequence Alignment & Mapping allowing mismatches

# Sequence Alignment & Mapping allowing mismatches

DNA

Fragment with sequencing artefact (e.g; adapter)

$$G\,G\,G\,G\,A\,T\,G\,C\,A\,T\,C\,T\,A\,A\,G$$
$$C\,C\,C\,C\,T\,A\,C\,G\,T\,A\,G\,A\,T\,T\,C$$

Read11

$$\textbf{\textit{G G G G A T G C}}\,ATCTAAG$$
$$CTTAGATGCATCGCC$$

Reference

```
        1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16 17
chrN    C  G  T  T  A  C  A  T  G  C  A  T  C  T  A  A  G
```

```
G G G G A T G C
C C C C T A C G
```

4 Soft clipped
4 Matches or Mismatches

```
Read11  chrN   Forward 7 hiMQ 4S4M
```

CIGAR

GENOMICS CORELEUVEN

# Sequence Alignment & Mapping
# SAM – BAM - CRAM

- BAM and CRAM files are compressed SAM files

  Not human readable, convert to SAM or use viewer (eg,IGV)

- CRAM smaller than BAM (40%-70%) but takes longer to read

- SAM format :

Chr1    Mapping Quality    Sequence, Base Quality    Alignment scores
AS:best, XS:2nd best

```
R9184:51295  163  1  693897  20  126M     = 694037   266  TAA... 3BB...  NM:i:0 MD:Z:126     AS:i:126 XS:i:121
R9184:51295  83   1  694037  20  126M     = 693897  -266  ACA... ;EE...  NM:i:0 MD:Z:126     AS:i:126 XS:i:126
R5802:69397  99   1  948309  60  126M     = 948364   178  TGA... :A3...  NM:i:1 MD:Z:124A1  AS:i:124 XS:i:24
R5802:69397  147  1  948364  60  58M3I65M = 948309  -178  GTG... GGC...  NM:i:3 MD:Z:123    AS:i:114 XS:i:47
```

Read id    position    CIGAR    Read pair info    Mismatch info

163: Forward,Paired,PairMapped,Second
83:  Reverse,Paired,PairMapped,First
99:  Forward,Paired,PairMapped,First
147: Reverse,Paired,PairMapped,Second
*https://broadinstitute.github.io/picard/explain-flags.html*

bio-bwa.sourceforge.net
samtools.sourceforge.net

# Duplicate removal

DNA

Fragment wih PCR error

$$T\ T\ \textcolor{red}{G}\ C\ A\ T\ G\ C\ A\ T\ C\ T$$

Duplicated erroneous fragment

$$T\ T\ \textcolor{red}{G}\ C\ A\ T\ G\ C\ A\ T\ C\ T$$

Reference

chrN C G T T A C A T G C A T C T A A G

```
1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16 17
```

```
TTGC      ATCT

TTGC      ATCT
```

# Base quality score recalibration (BQSR)

- Corrects systematic errors made by the sequencer when it estimates the quality score of each base call

- Performed on BAM, not on FASTQ
  - Requires genomic location of base
  - Ignores genomic locations where variants known to occur frequently
  - Considers each remaining variant an error

- Uses machine learning to characterize regions where more/less errors found than predicted by sequencer

- Example: *any base call that comes after AA in a read should have its quality score reduced by 1%*

*Source: https://gatkforums.broadinstitute.org/gatk/discussion/44/base-quality-score-recalibration-bqsr*

# Viewing a BAM file

- Without graphical user interface
  - o Samtools – suite of tools for handling SAM, BAM, CRAM

- With graphical user interface
  - o IGV – Integrative Genome Viewer

https://github.com/samtools/samtools
http://software.broadinstitute.org/software/igv/

# Samtools

- ## View

  - ### Read mapping information

```
samtools view sample.bam 1:11131116-11133317 | less

D00210:1282:CD2J0ANXX:5:1303:1159:51350 163     1       11131905
60      126M    =       11132084        305
GACTGCCTTCTCCAACCACCAACGAGACAGCTACAGCACCTCCAGCACTCCCCACCAATCTCTCTGCACAGCACCTGC
TGCCATCTGCCAGGATAGATACTGATTGCCCACCATCCCTCAGCAGAA
@=>BBECDCFEADBCFDDFDD@F@CGCFDEFEAFBHFBFDFFDDGFDGECDADDFADCBFEFEFEHFDGBGCAGEEIG
FHFDDBFEIGEDGFBAAGCBAFEIBCBIFDDDFDDBGEDFFAGDAC?B  MC:Z:126M
BD:Z:MMNNNNNMLMMMNNMMNNLNNLLLMMNMKNONMNLONNMNMMNNOONMNMNJJNMNNMMNMMMMMNNNMLOON
MNMNNNNNNNNNMNNNNOOMNLNNNLNNNNNMMNNJNMNNNNNJMMNOONONM    MD:Z:126
BI:Z:QQRRRQORQQQQQRRQQRQQRRQQQQRRQRRRQRQRRQQQRQQRRRQQRQQPPRQQRRRRQQQQQRQQQQRRQ
QQRRQRRQORRRQRQORRQRRQQRRQRRRQRRQQOPRQQRRRQPRQRRRQRRQ    NM:i:0  MQ:i:60
AS:i:126        XS:i:19 RG:Z:GC065340.run.181130.HiSeq2500.FCA.lane5-389E55F4-
49A02F07  PG:Z:MarkDuplicates-746E271E
```

# Samtools

- Pileup

  o Bases observed at each position

```
samtools mpileup sample.bam –f genome.fa –r 1:11131116-11133317

1     11132152    T     56
.$...............,...,,..,.,...........,..,.,...,,...,..^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@

1     11132153    C     57
...............^..,...,,..,.,...........,..,.,...,,...,...^].^].
GFoGGHmGDFoo^mHnFHHHGFHGHEGFDF/FFGGAGGHGHGFHGHHGGDHCCAAAA

1     11132154    A     57
...............,...,,..,.,...........,..,.,...,,...,......
DDgDDDhD2DgiUhEiDEEEECEDEDDEDE<CADE?2DFDEDDDDEEDDDED>??AA
```

# Samtools

- Pileup
  - Bases observed at each position

```
samtools mpileup sample.bam –f genome.fa –r 1:11131116-11133317

1       11132152    T       56
.$...............^].,...,,..,.,.............,..,.,...,,...,..^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@
```

Chromosome

GENOMICS
CORELEUVEN

# Samtools

- Pileup

  o Bases observed at each position

```
samtools mpileup sample.bam -f genome.fa -r 1:11131116-11133317

1     11132152     T     56
.$............,...,,..,.,...........,..,.,...,,...,..^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@
```

Position

Chromosome

# Samtools

- Pileup
  - Bases observed at each position

```
samtools mpileup sample.bam -f genome.fa -r 1:11131116-11133317

1     11132152     T      56
.$...............,..,,..,.,..........,..,.,...,,...,..^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@
```

Reference

Position

Chromosome

GENOMICS CORELEUVEN

# Samtools

- Pileup
  - Bases observed at each position

```
samtools mpileup sample.bam -f genome.fa -r 1:11131116-11133317

1      11132152     T       56
.$..............,...,,...,.,..........,..,.,...,,...,..^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@
```

Depth

Reference

Position

Chromosome

# Samtools

- Pileup
  - Bases observed at each position

```
samtools mpileup sample.bam -f genome.fa -r 1:11131116-11133317

1       11132152     T       56
.$...............^].,..,,..,.,............,..,.,...,,...,..^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@
```

Depth

Aligned bases

Reference                    . = reference base forward strand
                             , = reference base reverse strand

Position

Chromosome

GENOMICS
CORELEUVEN

# Samtools

- Pileup
  - Bases observed at each position

```
samtools mpileup sample.bam -f genome.fa -r 1:11131116-11133317

1     11132152    T      56
.$...............^].^].
CEAgFEFgE;EggKgChEFCCEEDECEEBEE=EEEECEECECEEEECCEE6CAAA@
```

Depth

Aligned bases
. = reference base forward strand
, = reference base reverse strand

Reference

Position

Chromosome

Base qualities

# IGV



Load Genome from File

# IGV



Load from File
(Select BAM file)

# Example 1

# Example 1 – targeted sequencing

# Example 1 – targeted sequencing

# Example 2

# Example 2 – amplicon sequencing

# Example 3

# Example 3 – WGS

# Example 3 – WGS

# Example 3 – WGS

# BAM file quality control

- Did I select/sequence what I wanted to ?

- What is the mean coverage ?

- How much of my region of interest is
  - Covered at 30X ?
  - Not covered at all ?

- Is the coverage even ?

- …

# Picard tools

- Set of (command line) tools to
  - Manipulate NGS data
    - SAM/BAM/CRAM
    - VCF/BCF
  - Compute metrics

https://broadinstitute.github.io/picard/

# Alignment metrics

- Picard CollectAlignmentSummaryMetrics



```
java -jar ~/bin/picard.jar CollectAlignmentSummaryMetrics
I=sample.bam O=sample_sum_metrics.txt R=ref.fasta
```

GENOMICS
CORELEUVEN

# Insert size

- Picard CollectInsertSizeMetrics



```
java -jar ~/bin/picard.jar CollectInsertSizeMetrics I=sample.bam
O=sample_insert_size_metrics.txt H=sample_insert_size_histogram.pdf
```

GENOMICS
CORE LEUVEN

# WGS metrics

- CollectWgsMetrics
  - Mean coverage



```
java -jar picard.jar CollectWgsMetrics I=sample.bam
O=sample_wgs_metrics.txt R=reference_sequence.fasta
```

GENOMICS
CORELEUVEN
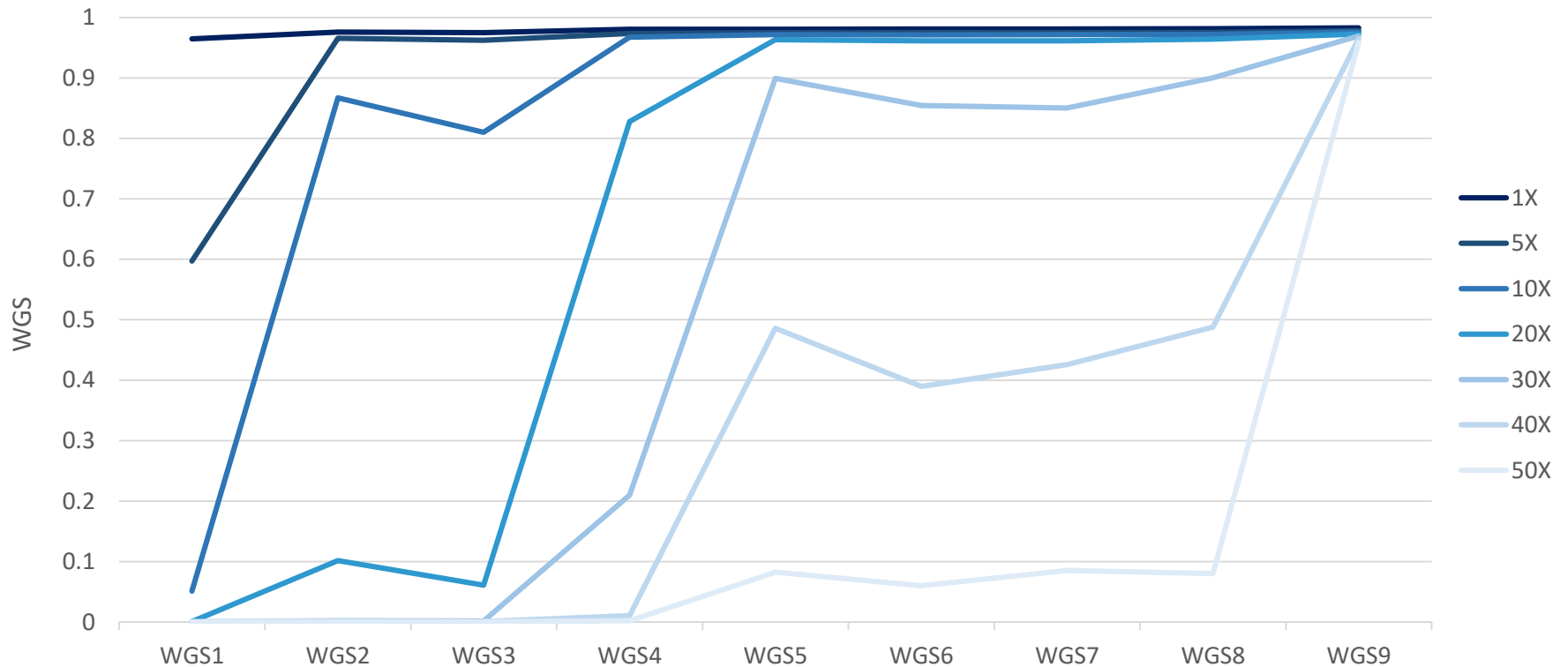
# WGS metrics

- CollectWgsMetrics
  - Raw vs informative coverage

# WGS metrics
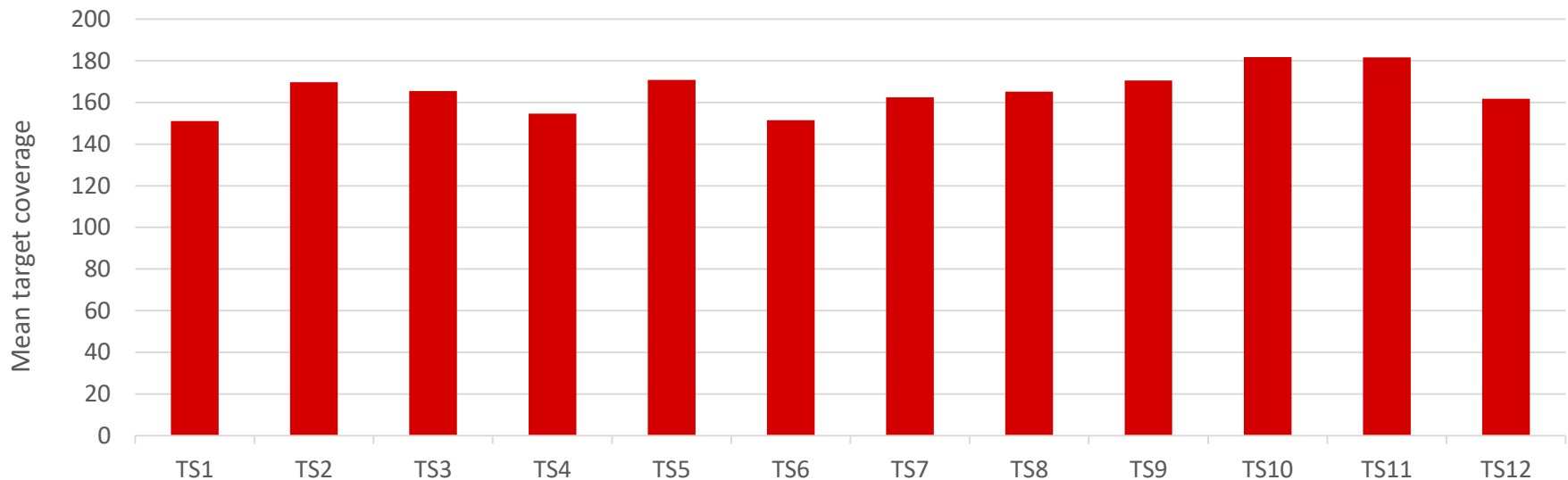
- CollectWgsMetrics
  - Excluded bases

# WGS metrics

- ## CollectWgsMetrics
  - Proportion of WGS covered at 1-50X

GENOMICS
CORELEUVEN

# Targeted metrics

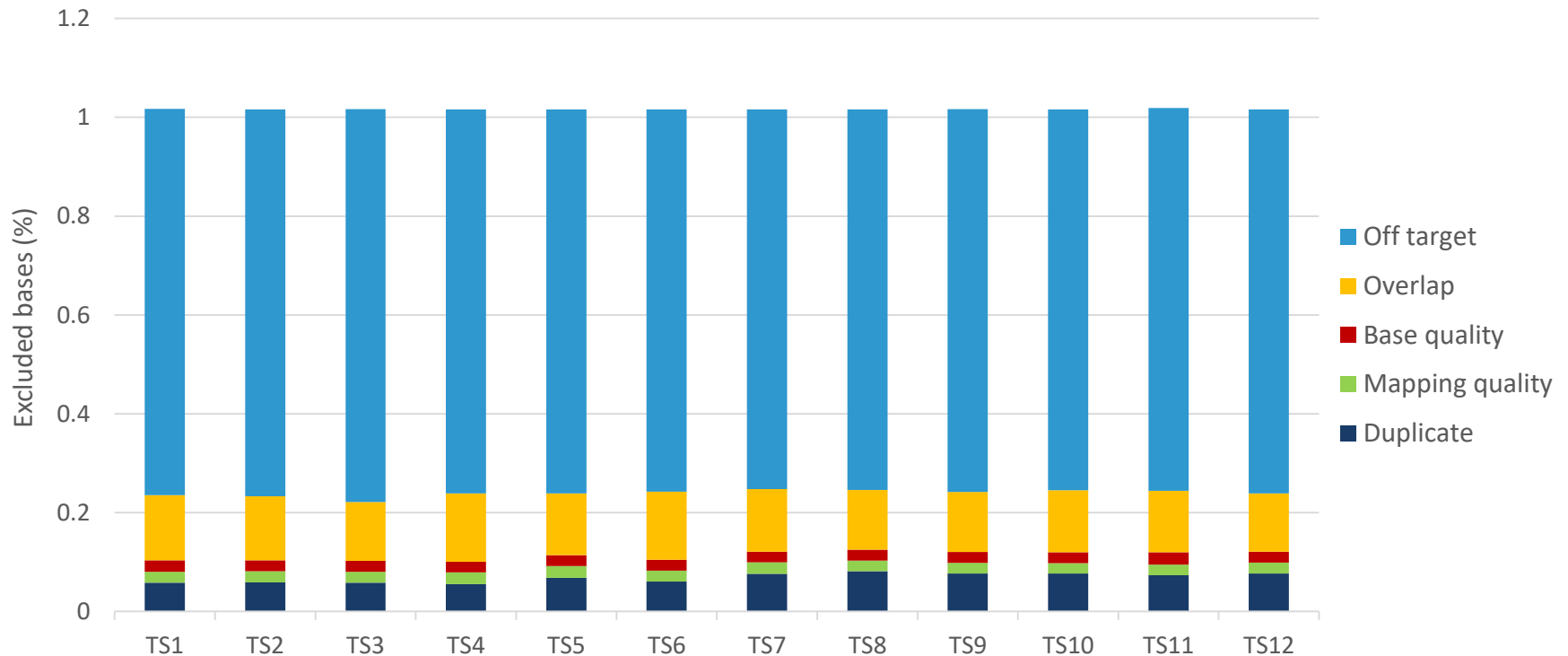- ## CollectHsMetrics

  - ### Mean target coverage



```
java -jar picard.jar CollectHsMetrics I=sample.bam
O=sample_hs_metrics.txt R=reference_sequence.fasta
BAIT_INTERVALS=bait.interval_list
TARGET_INTERVALS=target.interval_list
```
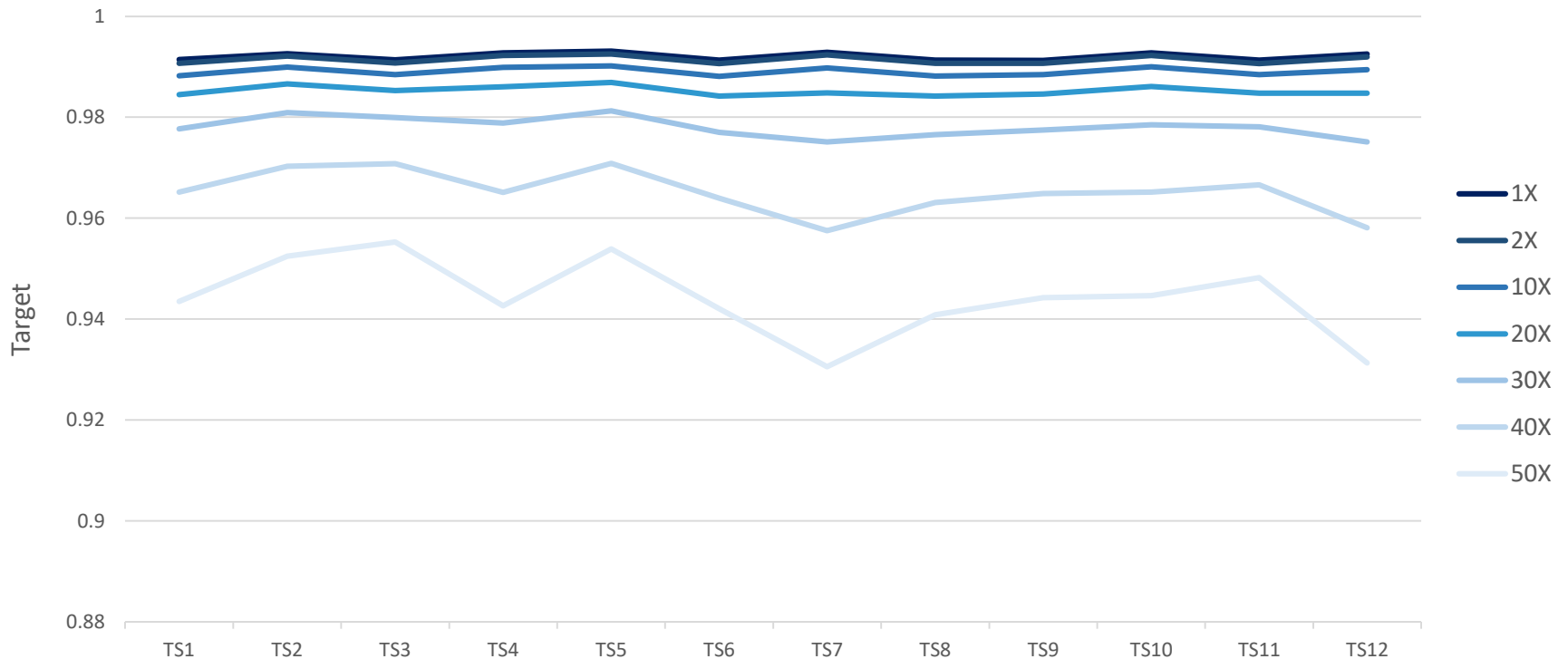
# Targeted metrics

- CollectHsMetrics
  - Excluded bases

# Targeted metrics

- CollectHsMetrics
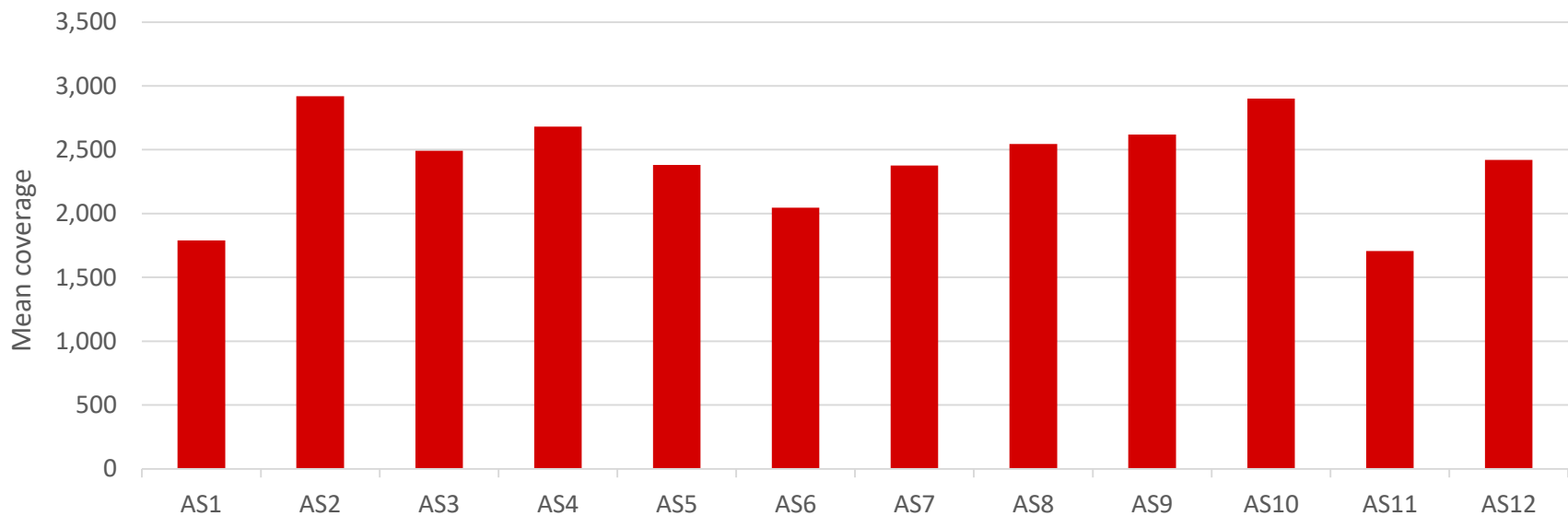  - Proportion of target covered at 1-50X

# Targeted metrics

- ## CollectHsMetrics
  - ### Capture efficiency

| Sample | Selected bases | Fold enrichment | Fold_80_base penalty |
|--------|----------------|-----------------|----------------------|
| TS1    | 0.95           | 101.21          | 1.86                 |
| TS2    | 0.94           | 100.95          | 1.91                 |
| TS3    | 0.94           | 100.03          | 1.84                 |
| TS4    | 0.94           | 101.63          | 1.89                 |
| TS5    | 0.94           | 100.50          | 1.88                 |
| TS6    | 0.94           | 101.53          | 1.87                 |
| TS7    | 0.95           | 102.72          | 2.08                 |
| TS8    | 0.95           | 101.88          | 2.01                 |
| TS9    | 0.95           | 102.28          | 2.03                 |
| TS10   | 0.95           | 102.85          | 2.09                 |
| TS11   | 0.95           | 102.67          | 2.06                 |
| TS12   | 0.95           | 101.78          | 2.07                 |

# Amplicon metrics

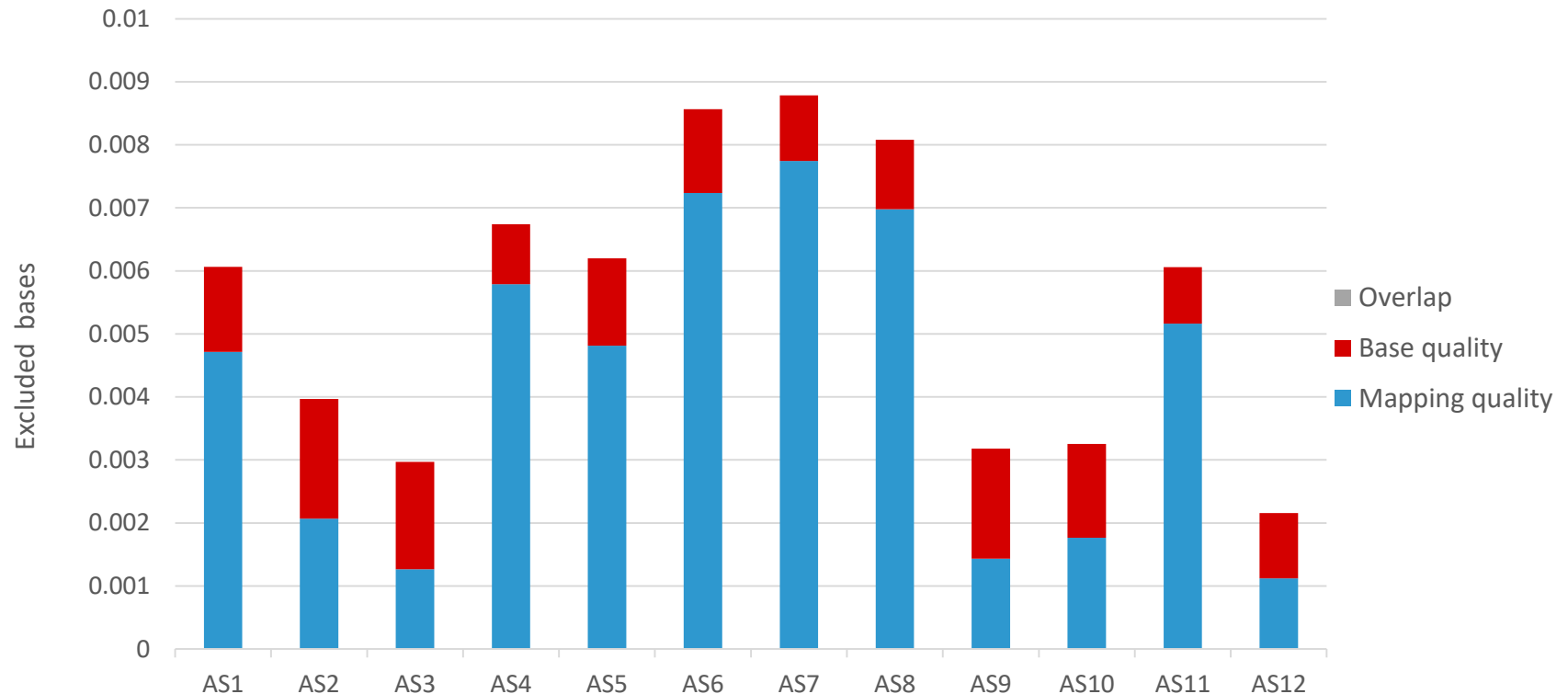- ## CollectTargetedPcrMetrics

  - ### Mean target coverage



```
java -jar picard.jar CollectTargetedPcrMetrics I=input.bam
O=sample_pcr_metrics.txt R=reference_sequence.fasta
AMPLICON_INTERVALS=amplicon.interval_list
TARGET_INTERVALS=targets.interval_list
```

# Amplicon metrics

- CollectTargetedPcrMetrics
  - Excluded bases

# Overview



Sequencing

fastq — Sequence of nucleotides within a read?
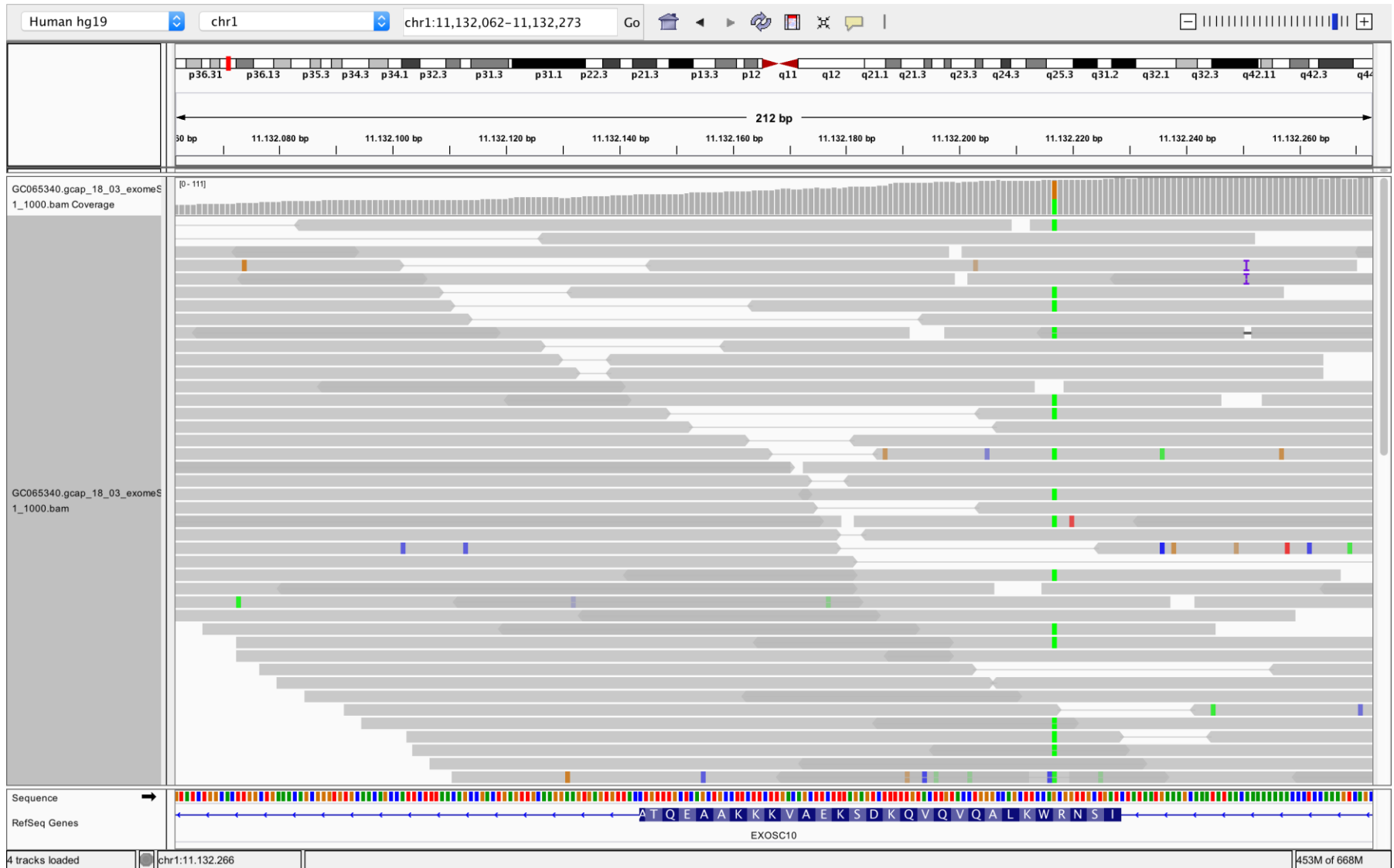
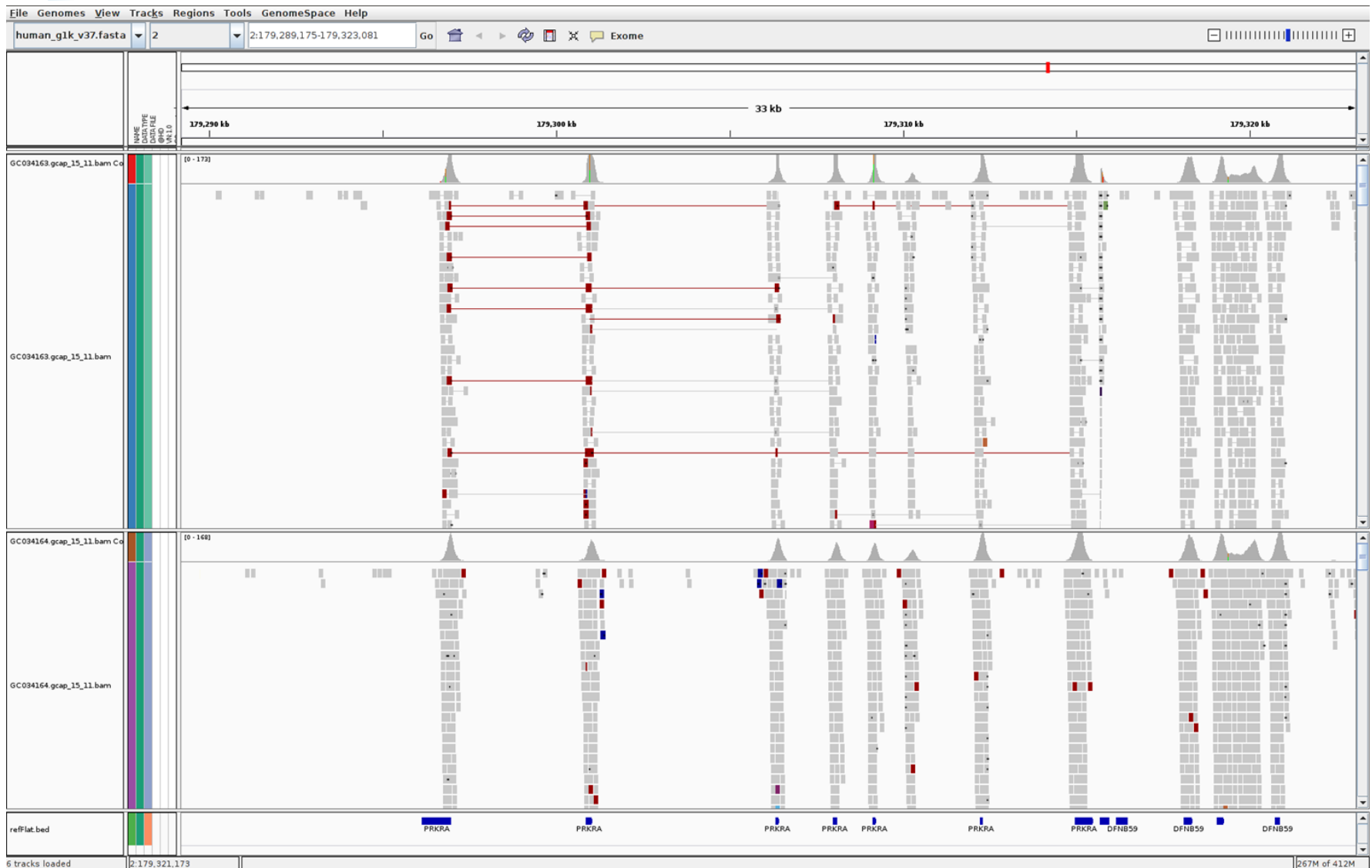bam — Genomic location of a read?

vcf — Genotype of a genomic location?

# Questions?

Luc Dehaspe
Erika Souche

GENOMICS CORELEUVEN

# Example 1

# Special cases

# Special cases