

# NGS applications

Variant detection: SNPs, CNVs, structural variants, epigenetic variation

Elena Piñeiro - [epineiro@cnio.es](mailto:epineiro@cnio.es)

# Genomic Variants



# What are genomic variants?

- Genomic variants are **permanent** changes in the DNA sequence of an organism.
- They can emerge by different mechanisms:
  - Recombination during gametes formation.
  - Errors during the DNA replication.
  - External factors like radiation, viruses, transposons, tobacco, UV light.

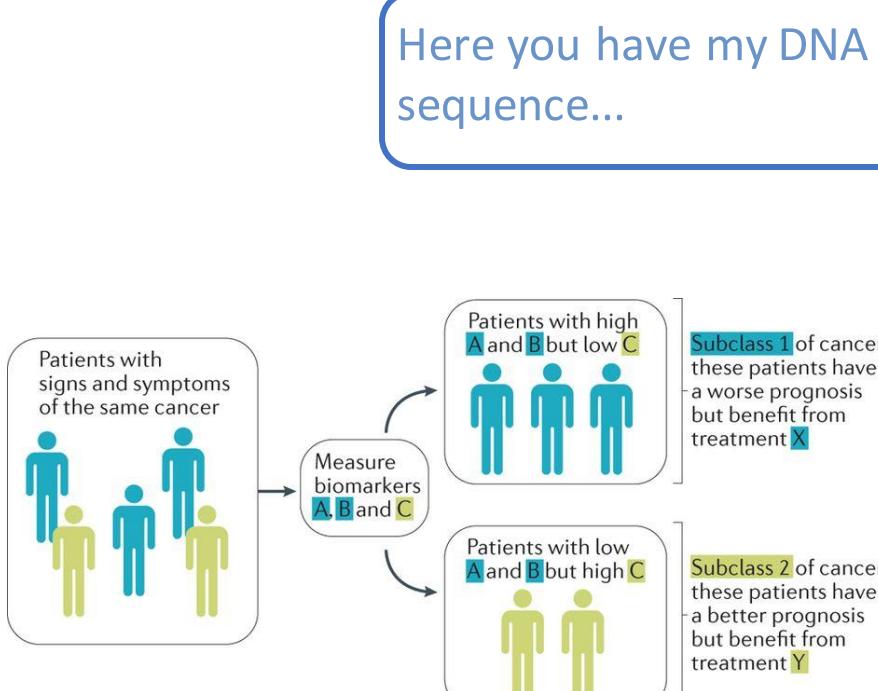
# Relevance of genomic variants

- They are a source of genomic variability, laying the foundations for the evolutionary mechanisms.
- Allow phenotypic differences between individuals (hair color, skin color, ...).
- Involved in diseases and drug response.

# Areas of application of variant

- Population and demographical studies.
- Phylogenetic studies.
- Epidemiological studies.
- Forensic applications.
- Clinical applications: disease susceptibility, drug response.

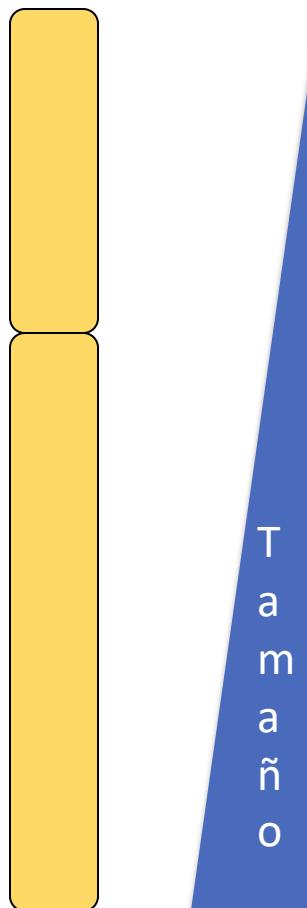
# Personalized Medicine



# Types of variants

- Different types of variants according to different criteria:
  - Variant size
  - Position of the variant in the DNA sequence
  - Consequence of the variant in transcription and translation
  - Clinical implication

# Classification of variants according to size



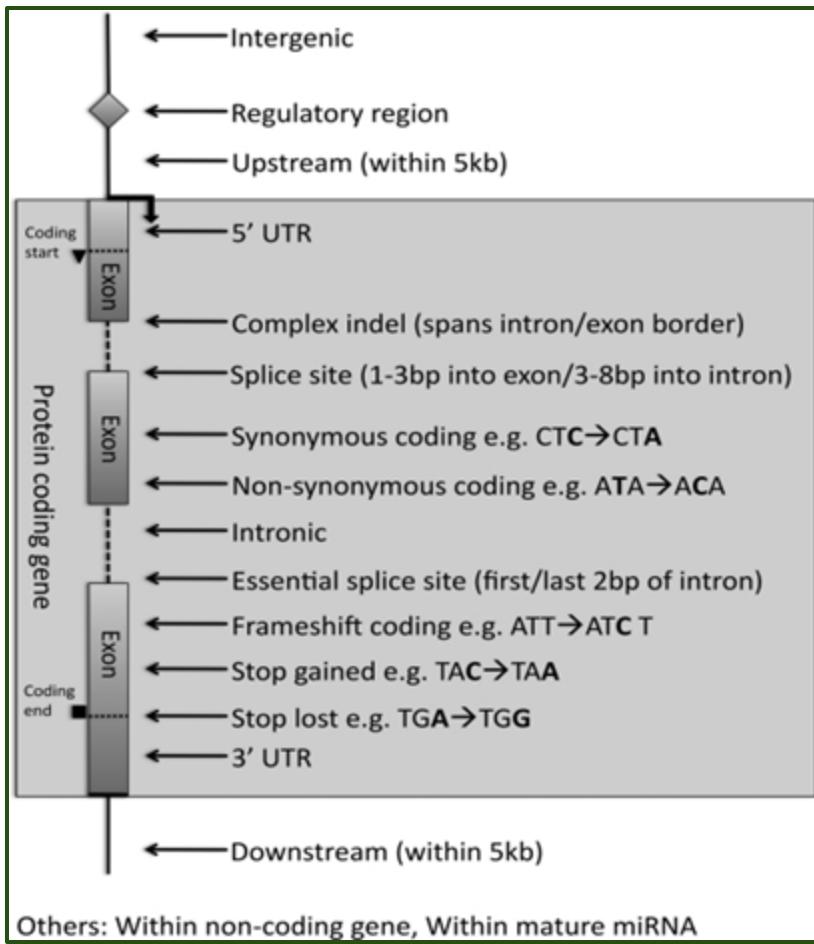
SNVs Single Nucleotide Variation  
Indels Small insertions and deletions

VNTRs (Micro, Minisatellites)  
Variable number of tandem repetitions

CNVs Copy number variation  
Translocations, Inversions

Aneuploidies

# Classification of variants according to the sequence position



- Intergenic
- In regulatory regions
- Upstream
- Downstream
- In genes
  - Untranslated regions: 5'UTR y 3'UTR
  - Exons
  - Introns
  - Splicing sites

# Consequences of SNV in protein coding regions

**SNV: Single-Nucleotide Variant**

**Wild-type sequence (reference)**  
ATCTTCAGCCT**AAA**GATGAAGTT

**Transition (green)**  
ATCTTCAGCCT**GAA**GATGAAGTT

**Transversion (blue)**  
ATCTTCAGCCT**CAA**GATGAAGTT

The diagram illustrates the consequences of point mutations at the DNA, mRNA, and protein levels. It shows four categories: No mutation, Silent, Nonsense, and Missense. The DNA level shows mutations from TTC to TTT, ATC, TCC, or TGC. The mRNA level shows AAG, AAA, UAG, AGG, or ACG. The protein level shows Lysine (Lys), Lysine (Lys), Stop, Arginine (Arg), or Threonine (Thr). The legend indicates that basic amino acids are purple and polar ones are green.

**Chemical structures of the four bases:**

- adenine: purine, nitrogenous base, two-ring system
- guanine: purine, nitrogenous base, two-ring system
- cytosine: pyrimidine, nitrogenous base, one-ring system
- thymine: pyrimidine, nitrogenous base, one-ring system

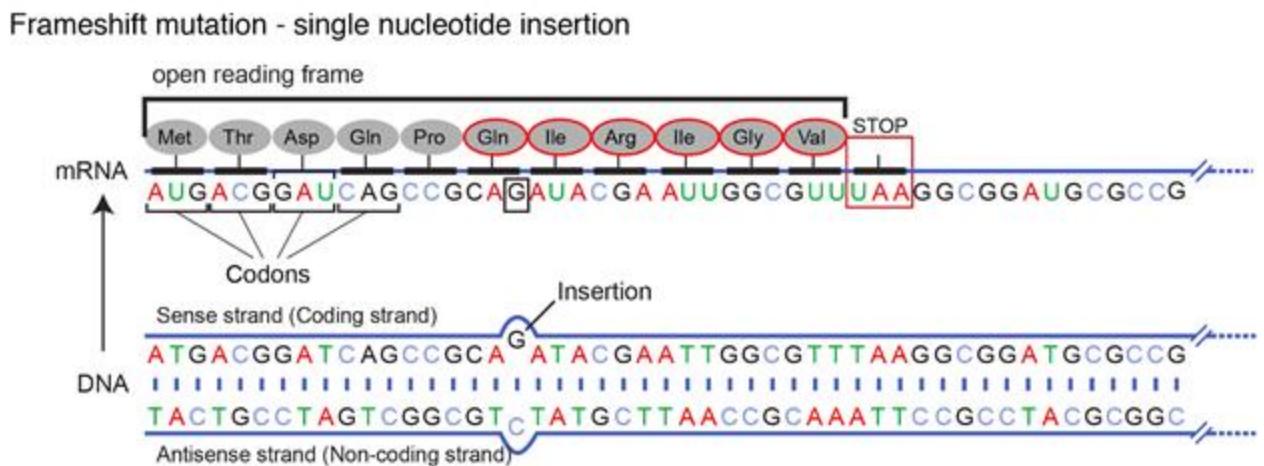
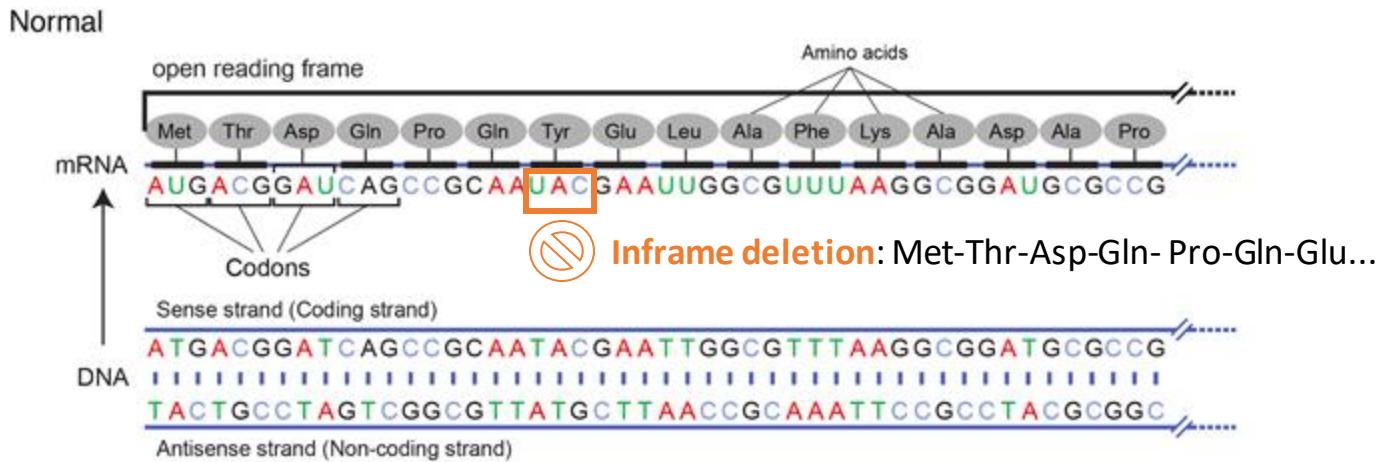
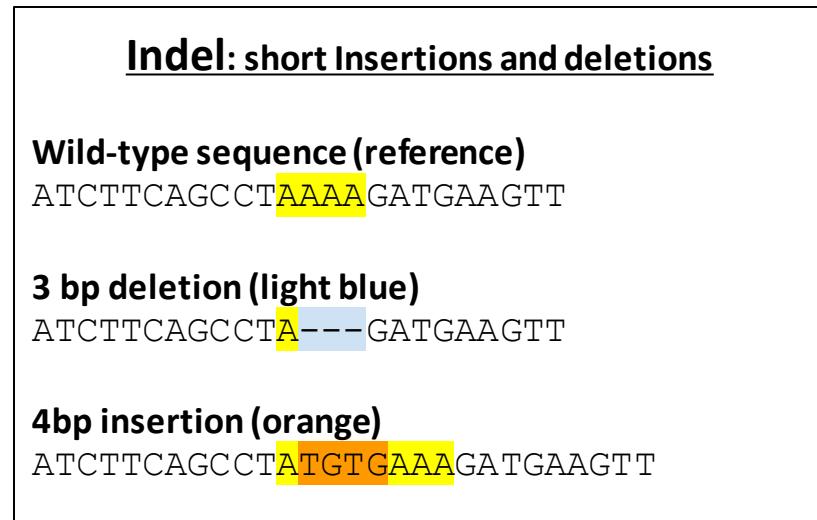
**Mutation types:**

- Transitions:** Changes between purines (A to G) or pyrimidines (C to T).
- Transversions:** Changes between a purine and a pyrimidine (A to C/T, G to C/T).

|               | No mutation | Point mutations |          |          |
|---------------|-------------|-----------------|----------|----------|
|               |             | Silent          | Nonsense | Missense |
| DNA level     | TTC         | TTT             | ATC      | TCC      |
| mRNA level    | AAG         | AAA             | UAG      | AGG      |
| protein level | Lys         | Lys             | STOP     | Arg      |
|               |             |                 |          |          |
|               |             |                 |          |          |
|               |             |                 |          |          |

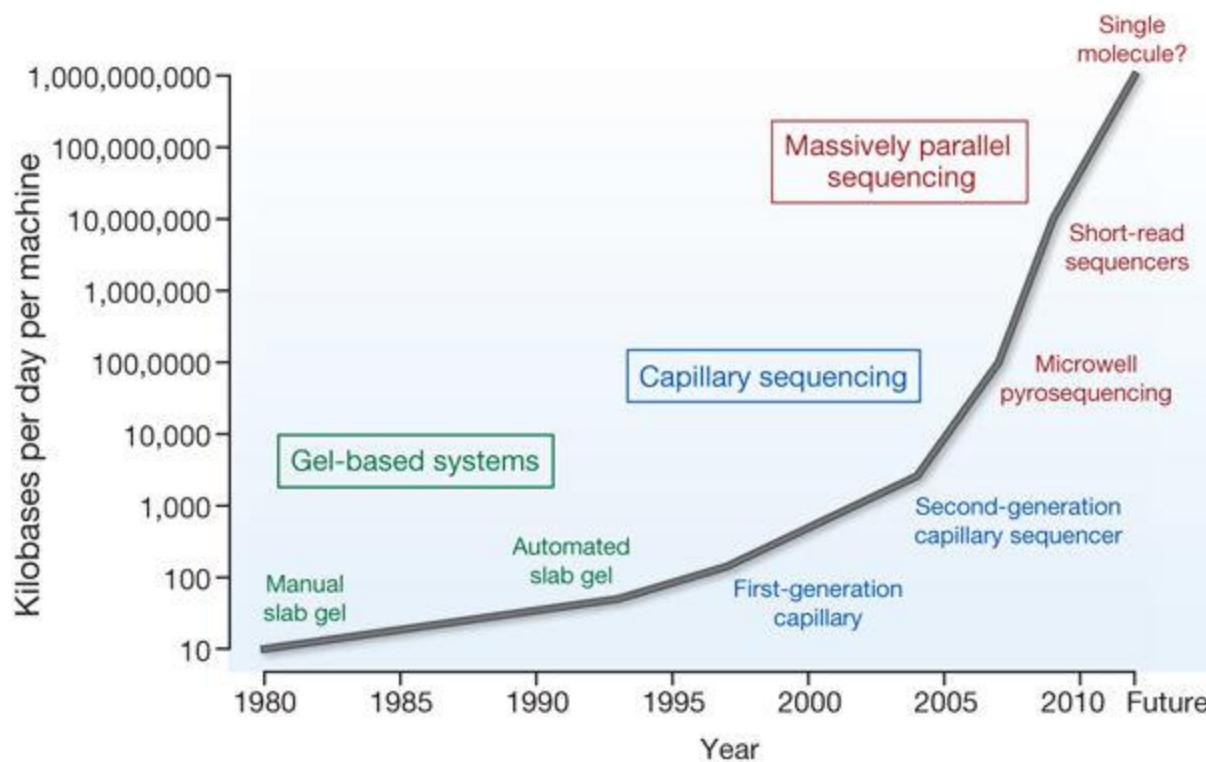
Stop loss, start gain, start loss

# Consequence of INDEL in protein coding regions

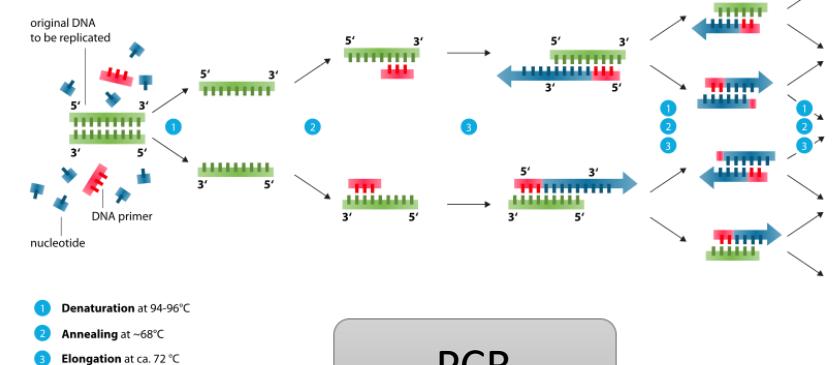


# Variant Detection with NGS

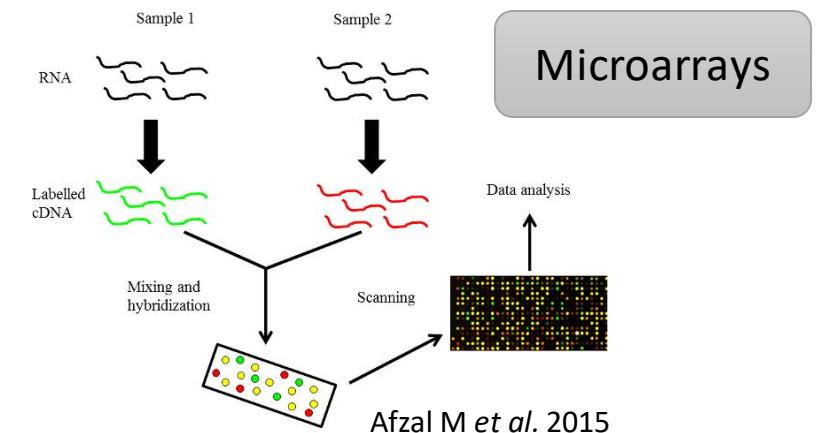
# Variant detection



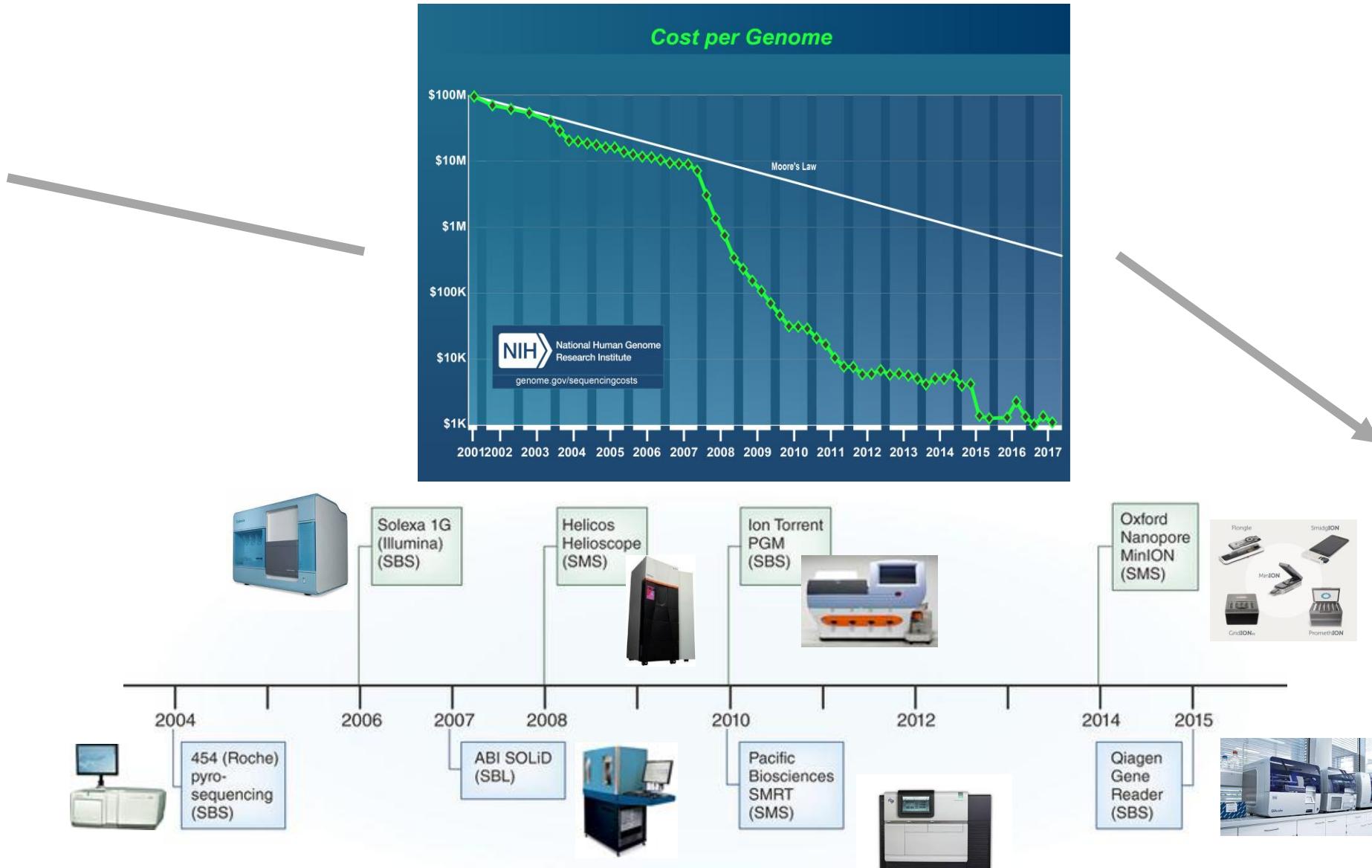
Polymerase chain reaction - PCR



Microarrays

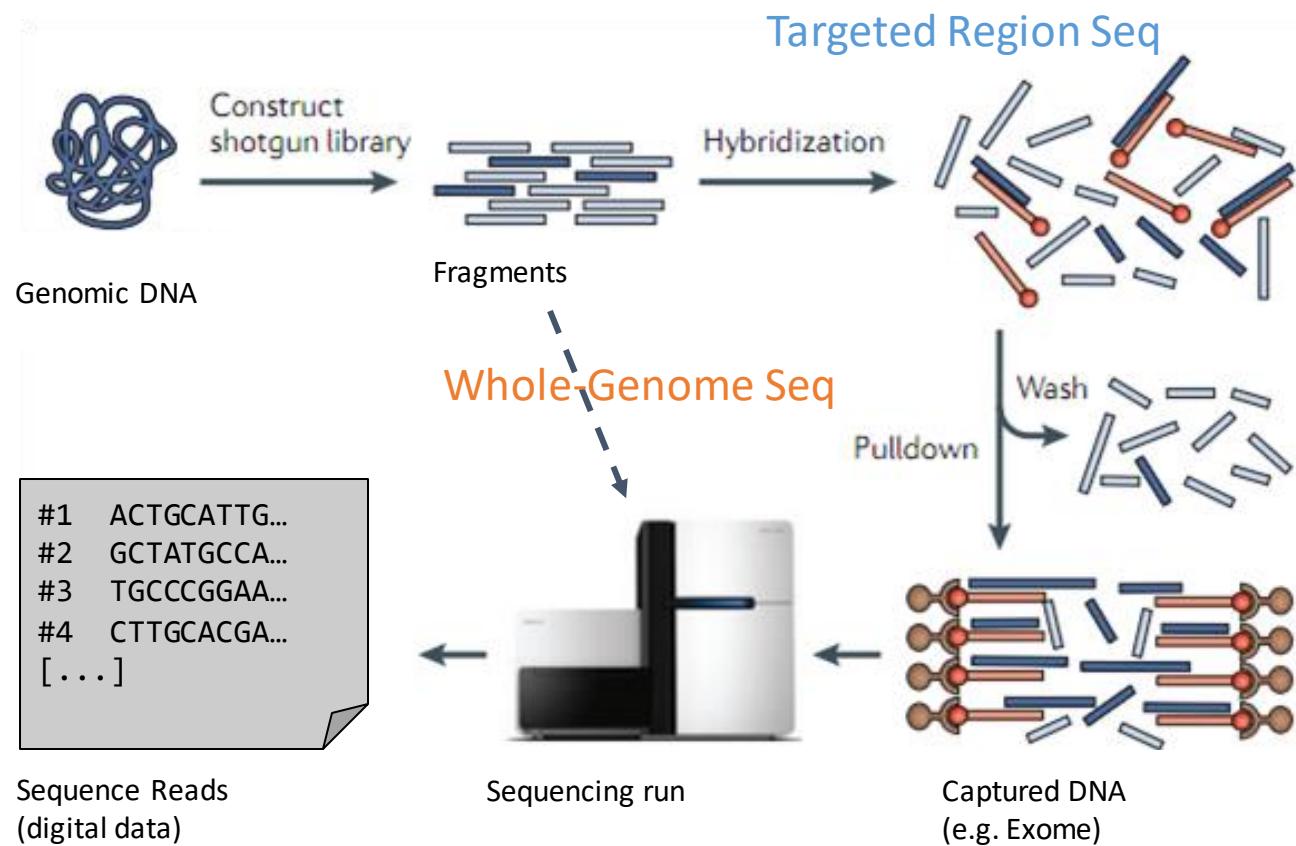


# Sequencing cost has been coming down



Adapted from Mardis ER. Nat Protoc 2017

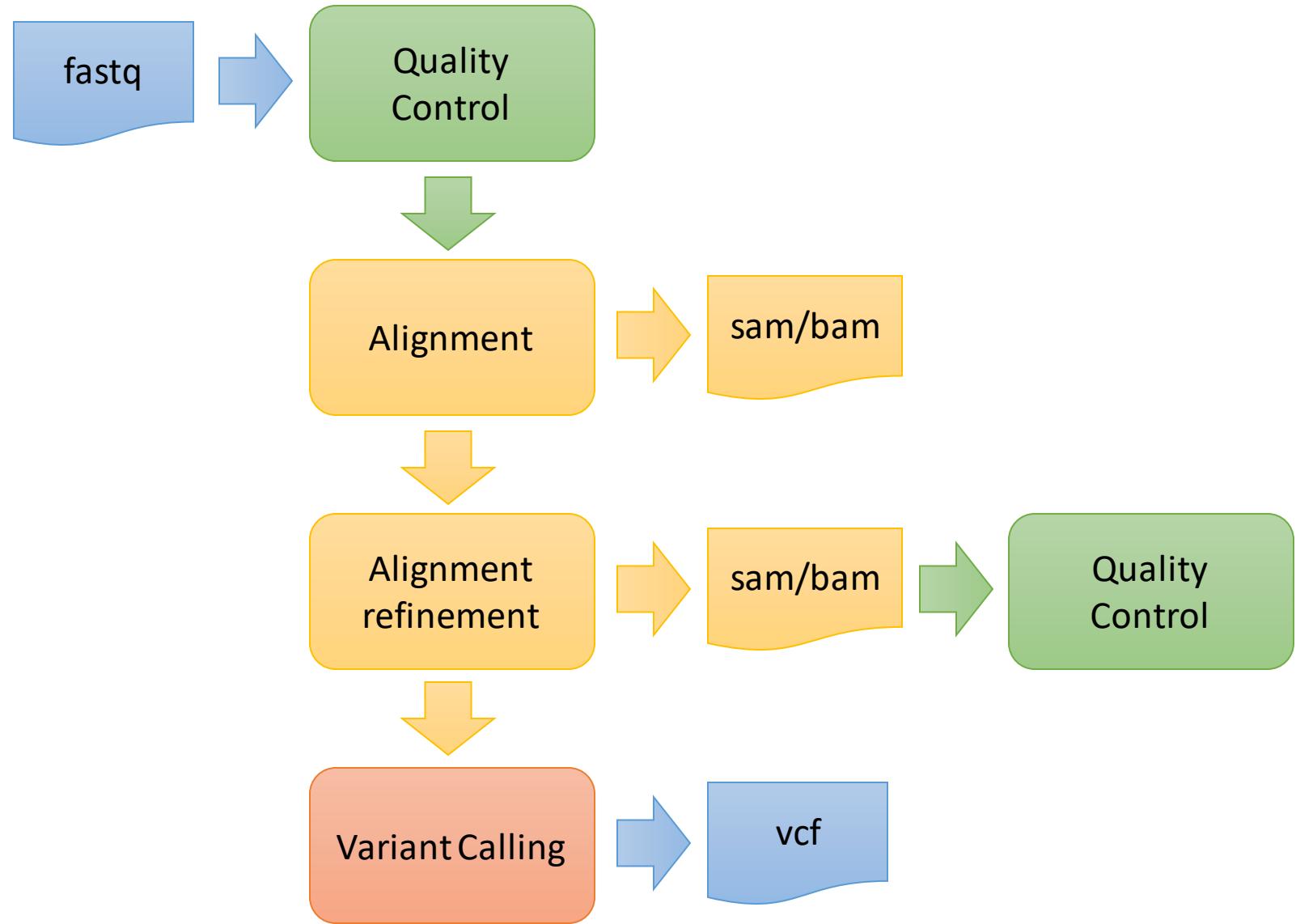
# DNA Sequencing data generation



The sequence reads belong from the ends of the original fragment.

# Steps in bioinformatic analysis

# General steps



fastq  
file

Fasta & Quality scores

# FASTA

- Typical extensions: .fasta, .fas, .fa, .fna, .fsa
- Each sequence is composed by at least two consecutive lines:
  - ">" Sequence name and optional description (space separated)
  - Line(s) with the whole sequence

```
>DNA_SEQUENCE_1
NNNNNNCTGGGGGACAGAACCCATGGTGGCCCCGGCTCCTCCCCAGTATCCAGTCCT
CCGTGAAGATGGAGGCCATTCC
60 chars
```

```
>DNA_SEQUENCE_1
NNNNNNCTGGGGGACAGAACCCATGGTGGCCCCGGCTCCTCCCCAGTATCCAGTCCT
CCGTGAAGATGGAGGCCATTCC
>DNA_SEQUENCE_2
GGGGGACAGAACCCATGGTGGCCCCGGCTCCTCCCCAGTATCCAGTCCT
>DNA_SEQUENCE_3
CTCCTCCCCAGTATCCAGTCCTGGGGGACAGAACCCATGGTGGCCCCGCCAGTATCCA
```

We can have multiple  
sequences in the same file  
(multifasta)

# FASTA

| IUPAC nucleotide code | Base                |
|-----------------------|---------------------|
| A                     | Adenine             |
| C                     | Cytosine            |
| G                     | Guanine             |
| T (or U)              | Thymine (or Uracil) |
| R                     | A or G              |
| Y                     | C or T              |
| S                     | G or C              |
| W                     | A or T              |
| K                     | G or T              |
| M                     | A or C              |
| B                     | C or G or T         |
| D                     | A or G or T         |
| H                     | A or C or T         |
| V                     | A or C or G         |
| N                     | any base            |
| . or -                | gap                 |

Nucleotide codes  
(IUPAC)

# FASTQ

- Typical extensions: .fq, .fastq
- Each read is composed by 4 lines:
  - "@" Read name and optional description (space separated)
  - Sequence
  - "+" (optionally: repeat the read name)
  - Base Quality Score

The diagram illustrates the structure of a FASTQ read. It consists of four lines of text, each preceded by a label in a box with an arrow pointing to its corresponding line.

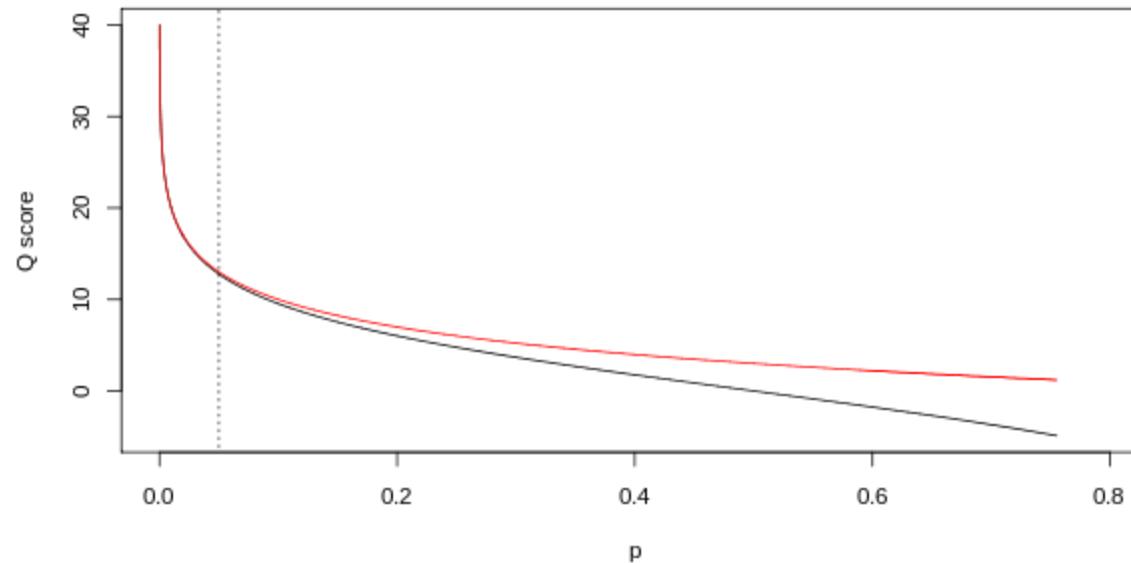
- readname**: @HWI-EAS209\_0006\_FC706VJ:5:58:5894:21141 ATCACG
- sequence**: TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTNNNNNNNNNTAGTTCTTGAGATTGTTGGGGGAGACATTTGTGATTGCCTTGAT
- comment**: +"#\$%&'()\*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^\_`abcdefghijklmnopqrstuvwxyz{|}JJJJJJ
- Quality**: (This line is empty in the provided example)

# Quality Score

- Phred quality scores (**QPhred**) are defined as a property which is **logarithmically related to the base-calling error probabilities  $p$**

$$\text{QPhred} = -10 \log_{10}(p)$$

- The higher the *QPhred* , the lower the probability that the base calling is erroneous



- The score is written as the character whose ASCII code is
  - QPhred + number (usually 33)

# ASCII code

| Dec | Hx     | Oct | Char                               | Dec | Hx     | Oct | Html  | Chr          | Dec | Hx     | Oct | Html  | Chr      | Dec | Hx     | Oct | Html   | Chr        |
|-----|--------|-----|------------------------------------|-----|--------|-----|-------|--------------|-----|--------|-----|-------|----------|-----|--------|-----|--------|------------|
| 0   | 0 000  | 000 | <b>NUL</b> (null)                  | 32  | 20 040 | 000 | &#32; | <b>Space</b> | 64  | 40 100 | 000 | &#64; | <b>Ø</b> | 96  | 60 140 | 000 | &#96;  | <b>`</b>   |
| 1   | 1 001  | 041 | <b>SOH</b> (start of heading)      | 33  | 21 041 | 041 | &#33; | <b>!</b>     | 65  | 41 101 | 065 | &#65; | <b>A</b> | 97  | 61 141 | 097 | &#97;  | <b>a</b>   |
| 2   | 2 002  | 042 | <b>STX</b> (start of text)         | 34  | 22 042 | 042 | &#34; | <b>"</b>     | 66  | 42 102 | 066 | &#66; | <b>B</b> | 98  | 62 142 | 098 | &#98;  | <b>b</b>   |
| 3   | 3 003  | 043 | <b>ETX</b> (end of text)           | 35  | 23 043 | 043 | &#35; | <b>#</b>     | 67  | 43 103 | 067 | &#67; | <b>C</b> | 99  | 63 143 | 099 | &#99;  | <b>c</b>   |
| 4   | 4 004  | 044 | <b>EOT</b> (end of transmission)   | 36  | 24 044 | 044 | &#36; | <b>\$</b>    | 68  | 44 104 | 068 | &#68; | <b>D</b> | 100 | 64 144 | 100 | &#100; | <b>d</b>   |
| 5   | 5 005  | 045 | <b>ENQ</b> (enquiry)               | 37  | 25 045 | 045 | &#37; | <b>%</b>     | 69  | 45 105 | 069 | &#69; | <b>E</b> | 101 | 65 145 | 101 | &#101; | <b>e</b>   |
| 6   | 6 006  | 046 | <b>ACK</b> (acknowledge)           | 38  | 26 046 | 046 | &#38; | <b>&amp;</b> | 70  | 46 106 | 070 | &#70; | <b>F</b> | 102 | 66 146 | 102 | &#102; | <b>f</b>   |
| 7   | 7 007  | 047 | <b>BEL</b> (bell)                  | 39  | 27 047 | 047 | &#39; | <b>'</b>     | 71  | 47 107 | 071 | &#71; | <b>G</b> | 103 | 67 147 | 103 | &#103; | <b>g</b>   |
| 8   | 8 010  | 050 | <b>BS</b> (backspace)              | 40  | 28 050 | 050 | &#40; | <b>(</b>     | 72  | 48 110 | 072 | &#72; | <b>H</b> | 104 | 68 150 | 104 | &#104; | <b>h</b>   |
| 9   | 9 011  | 051 | <b>TAB</b> (horizontal tab)        | 41  | 29 051 | 051 | &#41; | <b>)</b>     | 73  | 49 111 | 073 | &#73; | <b>I</b> | 105 | 69 151 | 105 | &#105; | <b>i</b>   |
| 10  | A 012  | 052 | <b>LF</b> (NL line feed, new line) | 42  | 2A 052 | 052 | &#42; | <b>*</b>     | 74  | 4A 112 | 074 | &#74; | <b>J</b> | 106 | 6A 152 | 106 | &#106; | <b>j</b>   |
| 11  | B 013  | 053 | <b>VT</b> (vertical tab)           | 43  | 2B 053 | 053 | &#43; | <b>+</b>     | 75  | 4B 113 | 075 | &#75; | <b>K</b> | 107 | 6B 153 | 107 | &#107; | <b>k</b>   |
| 12  | C 014  | 054 | <b>FF</b> (NP form feed, new page) | 44  | 2C 054 | 054 | &#44; | <b>,</b>     | 76  | 4C 114 | 076 | &#76; | <b>L</b> | 108 | 6C 154 | 108 | &#108; | <b>l</b>   |
| 13  | D 015  | 055 | <b>CR</b> (carriage return)        | 45  | 2D 055 | 055 | &#45; | <b>-</b>     | 77  | 4D 115 | 077 | &#77; | <b>M</b> | 109 | 6D 155 | 109 | &#109; | <b>m</b>   |
| 14  | E 016  | 056 | <b>SO</b> (shift out)              | 46  | 2E 056 | 056 | &#46; | <b>.</b>     | 78  | 4E 116 | 078 | &#78; | <b>N</b> | 110 | 6E 156 | 110 | &#110; | <b>n</b>   |
| 15  | F 017  | 057 | <b>SI</b> (shift in)               | 47  | 2F 057 | 057 | &#47; | <b>/</b>     | 79  | 4F 117 | 079 | &#79; | <b>O</b> | 111 | 6F 157 | 111 | &#111; | <b>o</b>   |
| 16  | 10 020 | 060 | <b>DLE</b> (data link escape)      | 48  | 30 060 | 060 | &#48; | <b>0</b>     | 80  | 50 120 | 080 | &#80; | <b>P</b> | 112 | 70 160 | 112 | &#112; | <b>p</b>   |
| 17  | 11 021 | 061 | <b>DC1</b> (device control 1)      | 49  | 31 061 | 061 | &#49; | <b>1</b>     | 81  | 51 121 | 081 | &#81; | <b>Q</b> | 113 | 71 161 | 113 | &#113; | <b>q</b>   |
| 18  | 12 022 | 062 | <b>DC2</b> (device control 2)      | 50  | 32 062 | 062 | &#50; | <b>2</b>     | 82  | 52 122 | 082 | &#82; | <b>R</b> | 114 | 72 162 | 114 | &#114; | <b>r</b>   |
| 19  | 13 023 | 063 | <b>DC3</b> (device control 3)      | 51  | 33 063 | 063 | &#51; | <b>3</b>     | 83  | 53 123 | 083 | &#83; | <b>S</b> | 115 | 73 163 | 115 | &#115; | <b>s</b>   |
| 20  | 14 024 | 064 | <b>DC4</b> (device control 4)      | 52  | 34 064 | 064 | &#52; | <b>4</b>     | 84  | 54 124 | 084 | &#84; | <b>T</b> | 116 | 74 164 | 116 | &#116; | <b>t</b>   |
| 21  | 15 025 | 065 | <b>NAK</b> (negative acknowledge)  | 53  | 35 065 | 065 | &#53; | <b>5</b>     | 85  | 55 125 | 085 | &#85; | <b>U</b> | 117 | 75 165 | 117 | &#117; | <b>u</b>   |
| 22  | 16 026 | 066 | <b>SYN</b> (synchronous idle)      | 54  | 36 066 | 066 | &#54; | <b>6</b>     | 86  | 56 126 | 086 | &#86; | <b>V</b> | 118 | 76 166 | 118 | &#118; | <b>v</b>   |
| 23  | 17 027 | 067 | <b>ETB</b> (end of trans. block)   | 55  | 37 067 | 067 | &#55; | <b>7</b>     | 87  | 57 127 | 087 | &#87; | <b>W</b> | 119 | 77 167 | 119 | &#119; | <b>w</b>   |
| 24  | 18 030 | 070 | <b>CAN</b> (cancel)                | 56  | 38 070 | 070 | &#56; | <b>8</b>     | 88  | 58 130 | 088 | &#88; | <b>X</b> | 120 | 78 170 | 120 | &#120; | <b>x</b>   |
| 25  | 19 031 | 071 | <b>EM</b> (end of medium)          | 57  | 39 071 | 071 | &#57; | <b>9</b>     | 89  | 59 131 | 089 | &#89; | <b>Y</b> | 121 | 79 171 | 121 | &#121; | <b>y</b>   |
| 26  | 1A 032 | 072 | <b>SUB</b> (substitute)            | 58  | 3A 072 | 072 | &#58; | <b>:</b>     | 90  | 5A 132 | 090 | &#90; | <b>Z</b> | 122 | 7A 172 | 122 | &#122; | <b>z</b>   |
| 27  | 1B 033 | 073 | <b>ESC</b> (escape)                | 59  | 3B 073 | 073 | &#59; | <b>;</b>     | 91  | 5B 133 | 091 | &#91; | <b>[</b> | 123 | 7B 173 | 123 | &#123; | <b>{</b>   |
| 28  | 1C 034 | 074 | <b>FS</b> (file separator)         | 60  | 3C 074 | 074 | &#60; | <b>&lt;</b>  | 92  | 5C 134 | 092 | &#92; | <b>\</b> | 124 | 7C 174 | 124 | &#124; | <b> </b>   |
| 29  | 1D 035 | 075 | <b>GS</b> (group separator)        | 61  | 3D 075 | 075 | &#61; | <b>=</b>     | 93  | 5D 135 | 093 | &#93; | <b>]</b> | 125 | 7D 175 | 125 | &#125; | <b>)</b>   |
| 30  | 1E 036 | 076 | <b>RS</b> (record separator)       | 62  | 3E 076 | 076 | &#62; | <b>&gt;</b>  | 94  | 5E 136 | 094 | &#94; | <b>^</b> | 126 | 7E 176 | 126 | &#126; | <b>~</b>   |
| 31  | 1F 037 | 077 | <b>US</b> (unit separator)         | 63  | 3F 077 | 077 | &#63; | <b>?</b>     | 95  | 5F 137 | 095 | &#95; | <b>_</b> | 127 | 7F 177 | 127 | &#127; | <b>DEL</b> |

# FASTQ – Single-end/Paired-end

One unique sample can have 1 or 2 files:

- If **single-end** Seq -> 1 file (name ".fastq")
- If **paired-end** Seq -> 2 files (names "\_R1.fastq" "\_R2.fastq")

**\_R1**

```
@SN1083:379:H8VA1ADXX:2:1101:1248:2144 1:N:0:12
CCTAAATGGTGCATGCTAGGAGGCCGTGCCCTCTTGAAAAGTTGTATGTGAA
+
BBBBFFFFFBFFFIIIIIFI<FFIIIIIFIIIIIFBFIIIIIIIIFFIIIIIFI
```

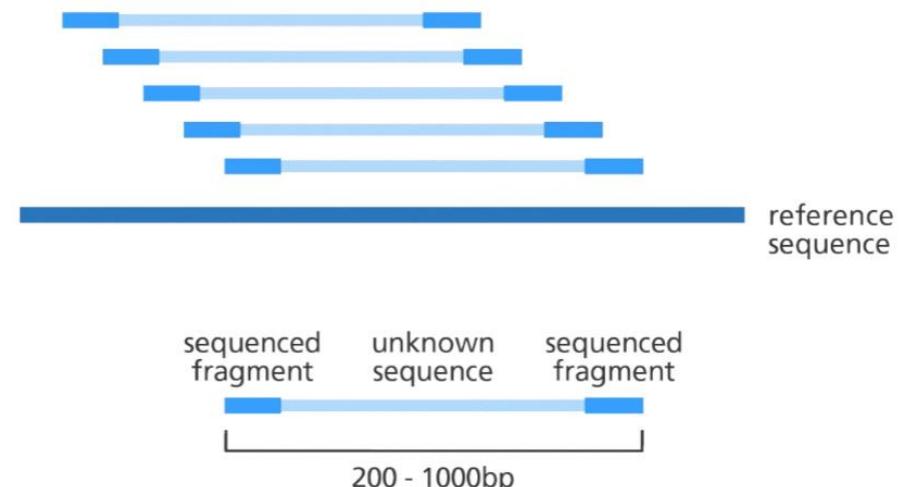
**\_R2**

```
@SN1083:379:H8VA1ADXX:2:1101:1248:2144 2:N:0:12
CATTTCGACGTTAATAAGCTCTGCGTACTTGCAAGCTATCTGCGCGAACG
+
BBBBFFFFFFFIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIFFF
```

Single-end reads



Paired-end reads



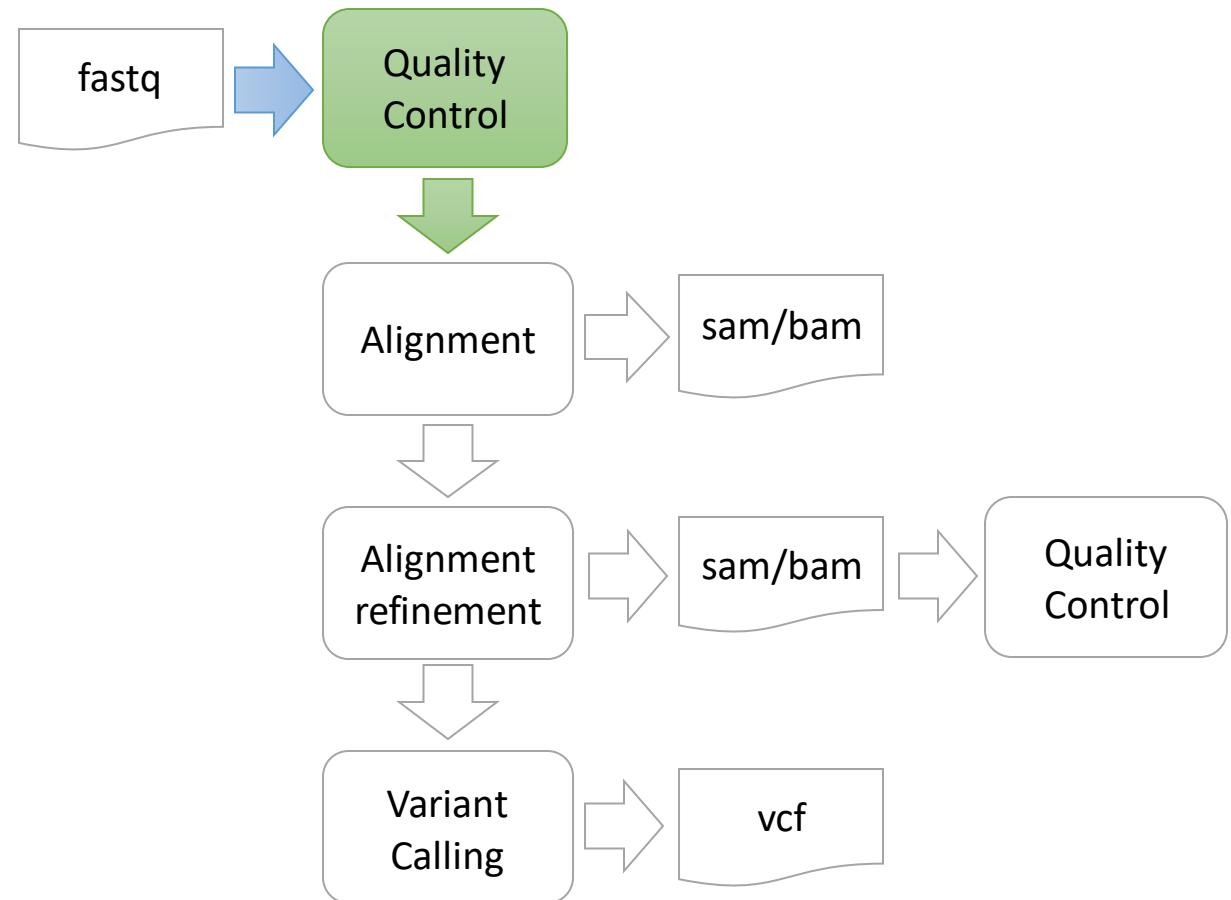
# Different NGS technologies have different capabilities

Table 1 | Main characteristics of current NGS technologies

| Technology | Run type   |            |           | Maximum read length  | Quality scores | Error rates | Refs     |
|------------|------------|------------|-----------|----------------------|----------------|-------------|----------|
|            | Single end | Paired end | Mate pair |                      |                |             |          |
| Illumina   | Yes        | Yes        | Yes       | 300 bp               | >30            | 0.0034–1%   | 59       |
| SOLiD      | Yes        | Yes        | Yes       | 75 bp                | >30            | 0.01–1%     | 60       |
| IonTorrent | Yes        | Yes        | No        | 400 bp               | ~20            | 1.78%       | 22       |
| 454        | Yes        | Yes        | No        | ~700 bp (up to 1 kb) | >20            | 1.07–1.7%   | 53,61    |
| Nanopore   | Yes        | No         | No        | 5.4–10 kb            | NA             | 10–40%      | 62–66    |
| PacBio     | Yes        | No         | No        | ~15 kb (up to 40 kb) | <10            | 5–10%       | 22,67–69 |

454, 454 pyrosequencing (Roche); NA, not applicable; Nanopore, Oxford Nanopore Technologies; NGS, next-generation sequencing; PacBio, Pacific Biosciences; SOLiD, sequencing by oligonucleotide ligation and detection (Thermo Fisher).

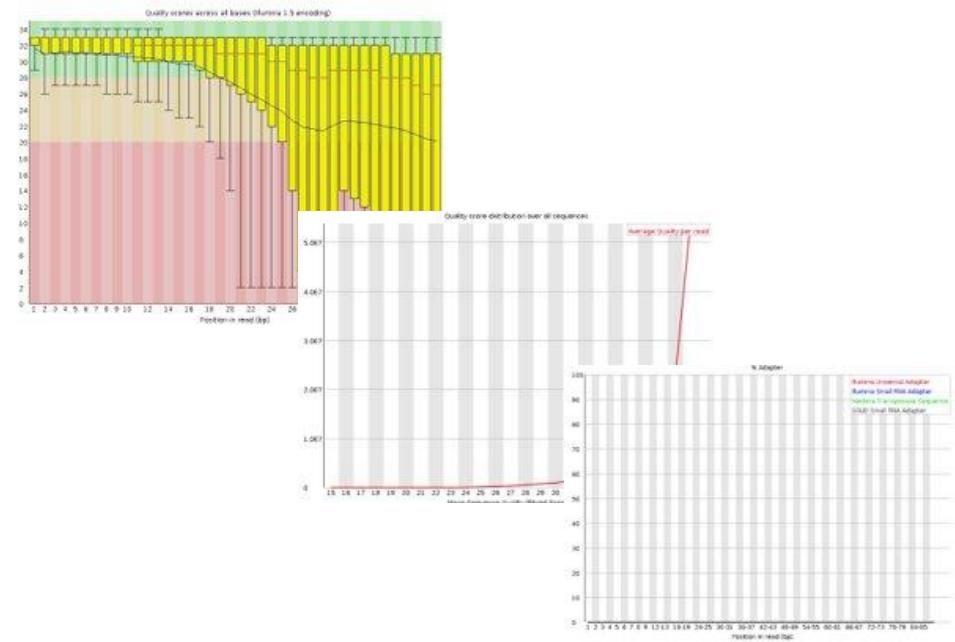
# Quality control



# Quality control

- Is the read number in the expected range?
- Have bases a good quality?
- Is there any contamination in the samples?

- Software:
  - FastQC
  - FastQ Screen



# FastQC

- **Input:** FastQ files. **Output report:** html (web-like) with plots
- **Raw files:** [https://fundacioncio-my.sharepoint.com/:u/g/personal/epineiro\\_cnio\\_es/EaWrloqPdb5GjXxak1O4TzkBUmlI6sqnupz8SLuJ2aT5-A?e=hv0Nss](https://fundacioncio-my.sharepoint.com/:u/g/personal/epineiro_cnio_es/EaWrloqPdb5GjXxak1O4TzkBUmlI6sqnupz8SLuJ2aT5-A?e=hv0Nss)



**Patient suffering ovarian cancer.**

**Whole-exome sequencing** data from two samples from the patient:

- Tumour sample.
- Matched normal sample (healthy tissue) from epithelium.

**Library protocol:** Agilent SureSelect V5 Human All Exons.

**Sequencing platform:** HiSeq 2000 (Illumina)

NOTE: This data was simulated and reduced in order to perform the computational analysis in class time.

# FastQC



```
> docker pull osvaldogc/ufv:2.0
```

1

```
> SAMPLES_LOCAL=/Path_directory_raw_data/  
> SAMPLES_DOCKER=/SAMPLES  
> RESULTS_LOCAL=/Path_directory_RESULTS  
> RESULTS_DOCKER=/RESULTS
```

2

```
> chmod 777 $RESULTS_LOCAL
```

3

```
> docker run --rm -v $SAMPLES_LOCAL:$SAMPLES_DOCKER -v $RESULTS_LOCAL:$RESULTS_DOCKER -it  
osvaldogc/ufv:2.0 /bin/bash
```

4

```
> fastqc /SAMPLES/fastq_file_name -o /RESULTS
```

FastQC execution

## Summary

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

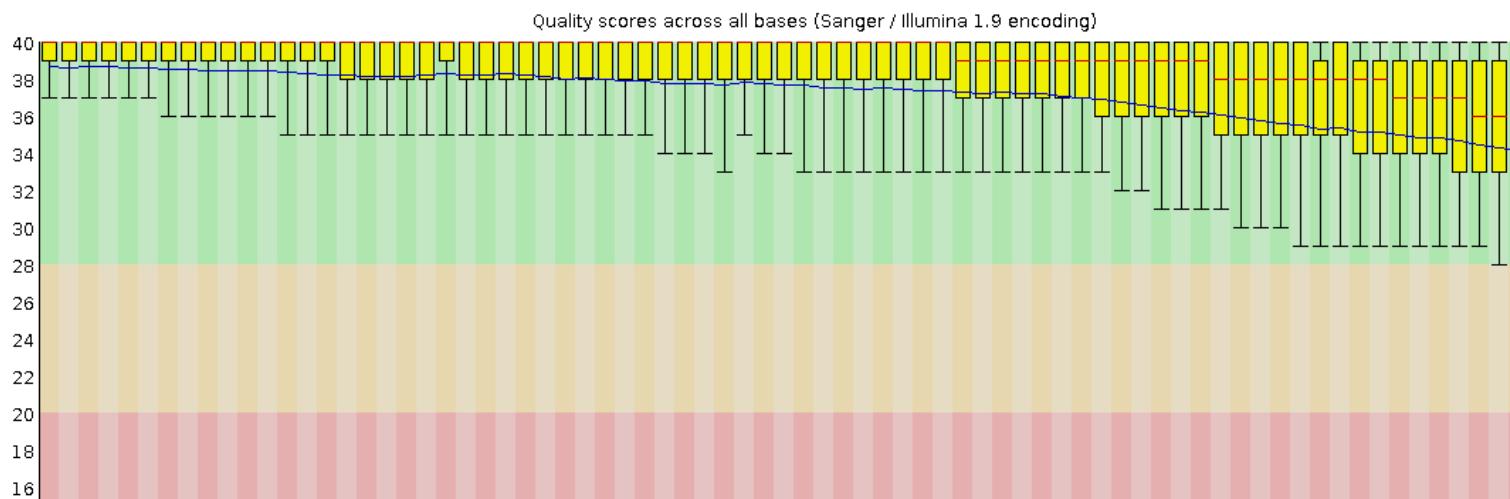
## Basic Statistics

| Measure                           | Value                   |
|-----------------------------------|-------------------------|
| Filename                          | WEx_Normal_R1.fastq     |
| File type                         | Conventional base calls |
| Encoding                          | Sanger / Illumina 1.9   |
| Total Sequences                   | 1400000                 |
| Sequences flagged as poor quality | 0                       |
| Sequence length                   | 75                      |
| %GC                               | 53                      |

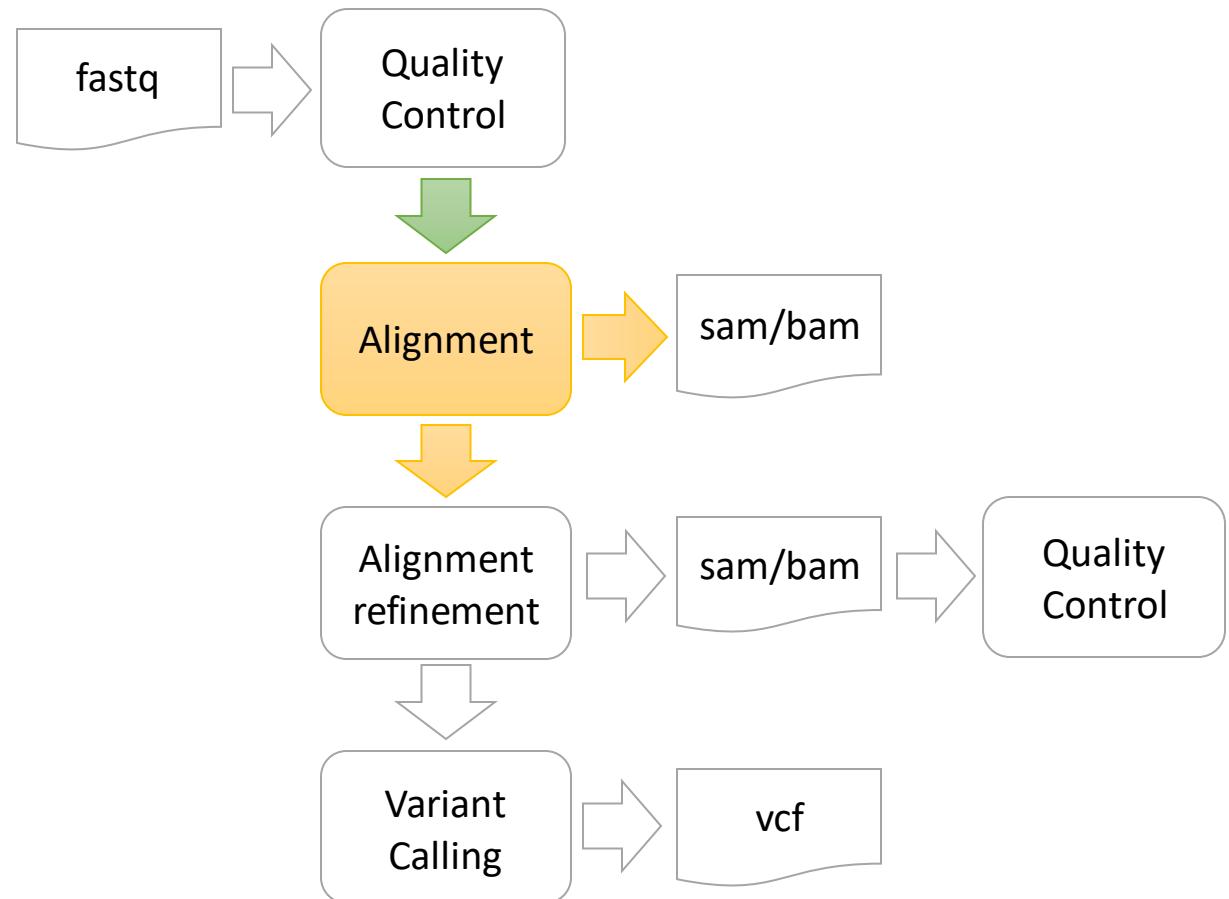
Phred+33 listed as Sanger/Illumina 1.9

Phred+64 listed as Illumina 1.5 (or lower)

## Per base sequence quality



# Alignment

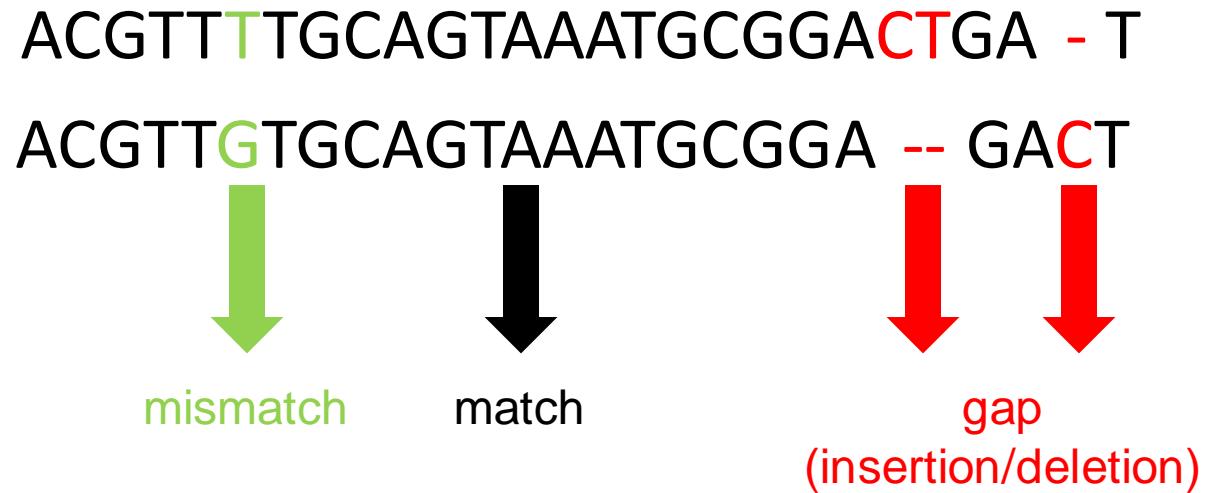


# What is an alignment?

ACGTCTTGA<sub>C</sub>TGG -TTAAAATAC  
AC - TCTTGA<sub>C</sub>TGGATTAA<sub>A</sub>CATAC

**Sequence alignment** is a way of arranging the sequences of DNA, RNA, or protein to **identify regions of similarity** that may be a consequence of functional, structural, or evolutionary relationships between the sequences.

# Elements



Alignment seeks to **reduce gaps and mismatches** and **maximize matches**.

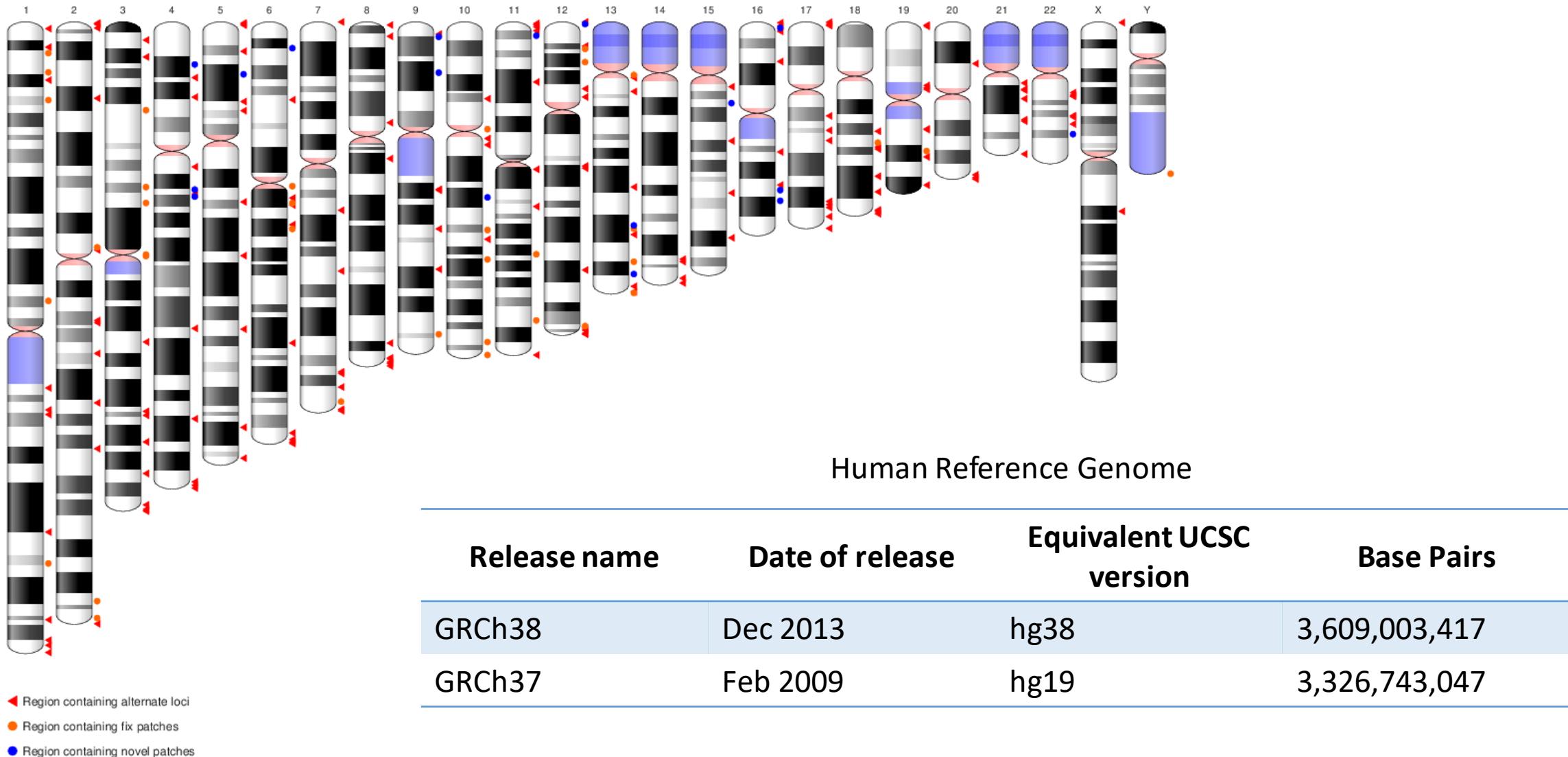
In the construction, each of these components has a penalty value associated. For gaps there is a penalty value for opening the gap and another for extending it.

ACGTT**T**TGCAGTAAATGCGGA**CTGAT**  
ACGTT**G**TGCAGTAAATGCGGA-**GACT**  
1 gap

ACGTT**T**TGCAGTAAATGCGGA**CTGAT**  
ACGTT**G**TGCAGTAAATGCGGA -- **GACT**  
1 extended gap

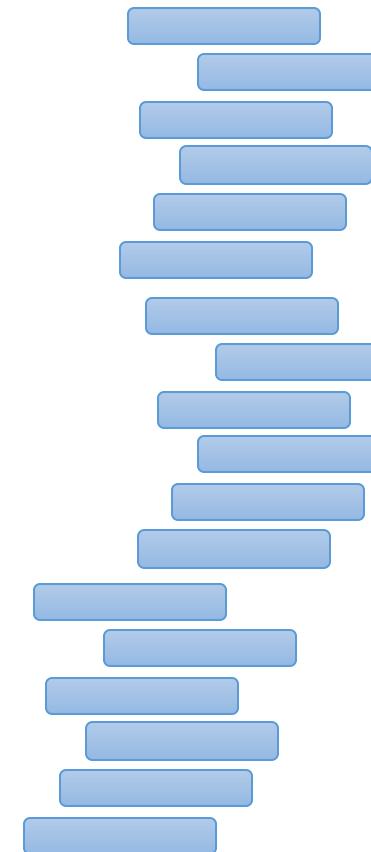
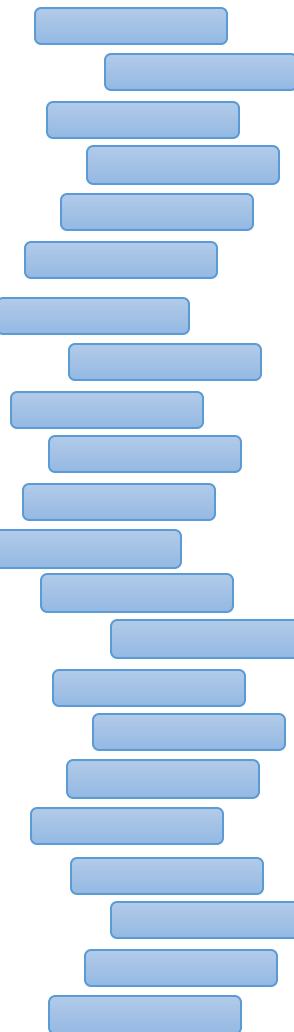
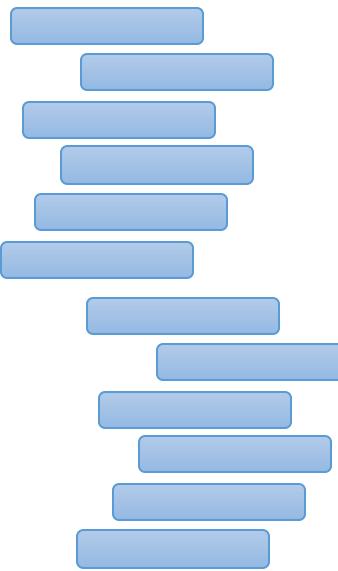
ACGTT**T**TGCAGTAAATGCGGA**CTGA** - T  
ACGTT**G**TGCAGTAAATGCGGA --**GACT**  
2 gaps

# Reference Genome





## REFERENCE GENOME

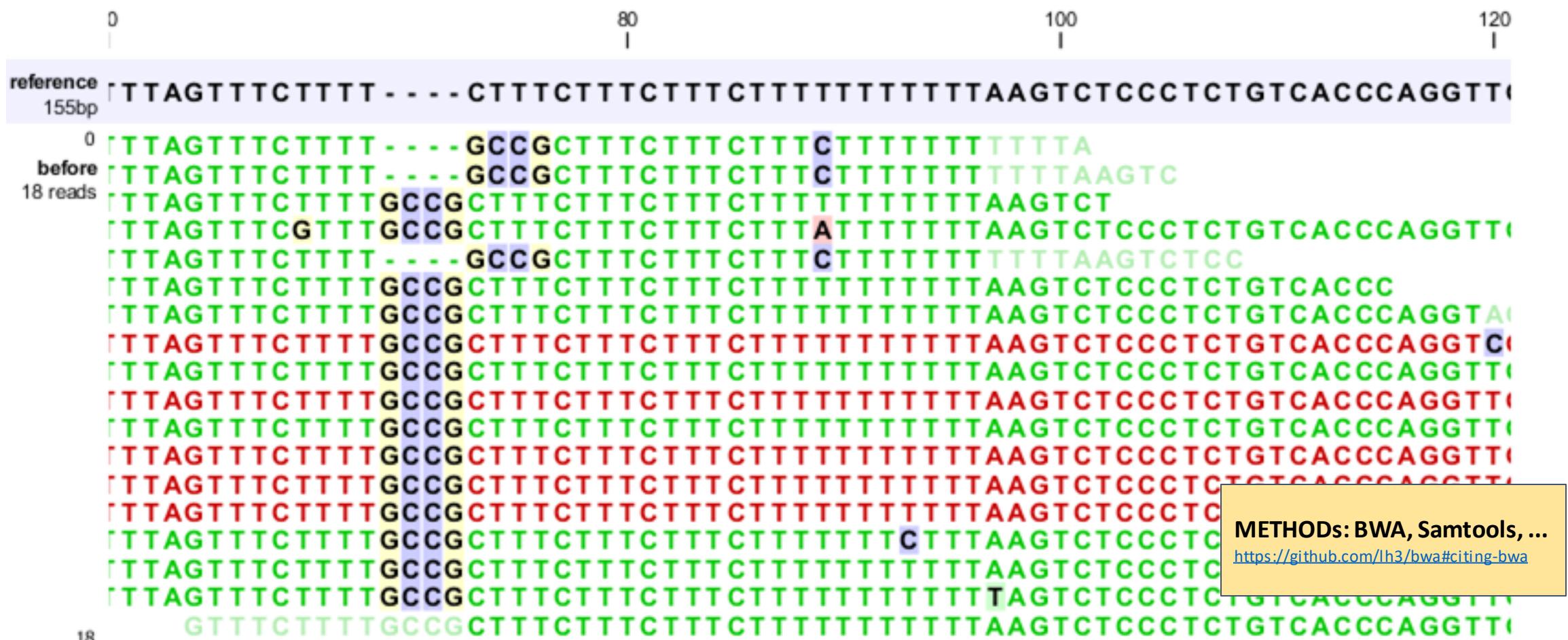


Reads

There is a problem  
with the dimension  
of the data



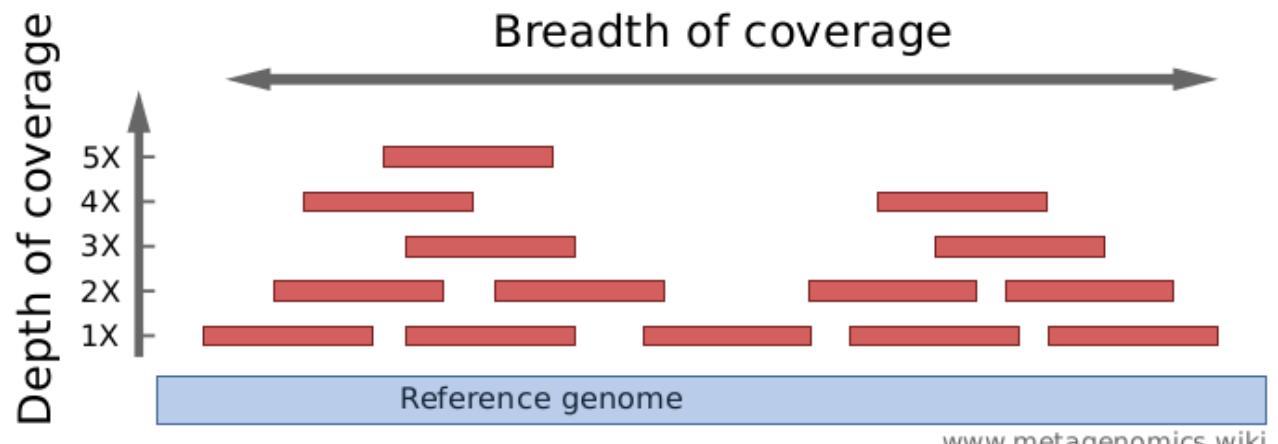
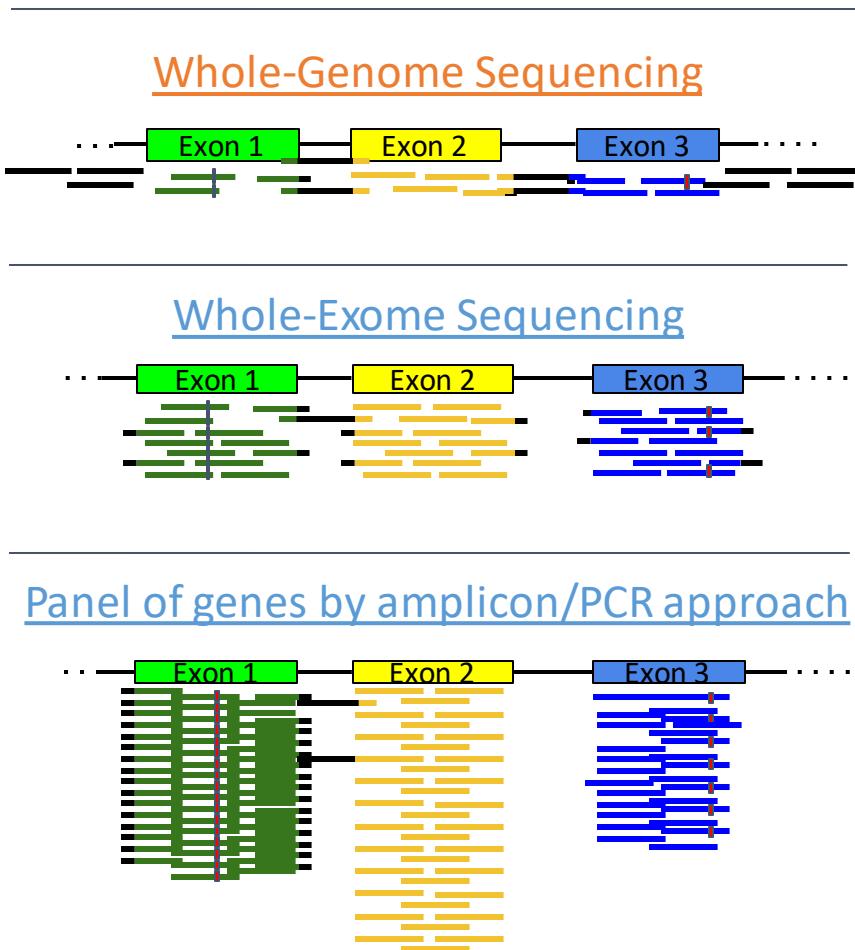
- Fast mapping on the reference genome by creating indexes. It is computationally intensive, but it is done only once.
- Search for candidate sites to align a given read by using seeds (fragments of a read).



# Different types of variants detected by mapping reads



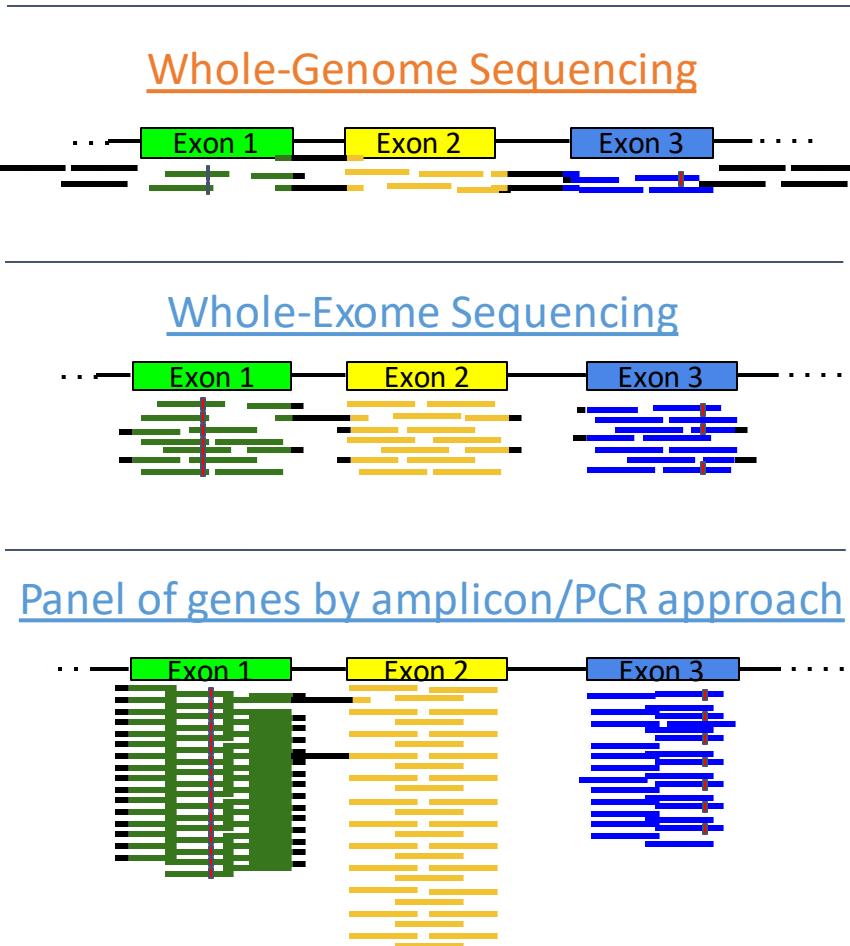
# Strategies on DNA-seq



**Depth:** times that a base is sequenced (in average)  
$$\frac{\text{Number of times a bases is covered by reads}}{\text{length of the sequenced genome}}$$

**Breadth:** percentage of the sequenced genome  
that is covered by the reads at a certain depth

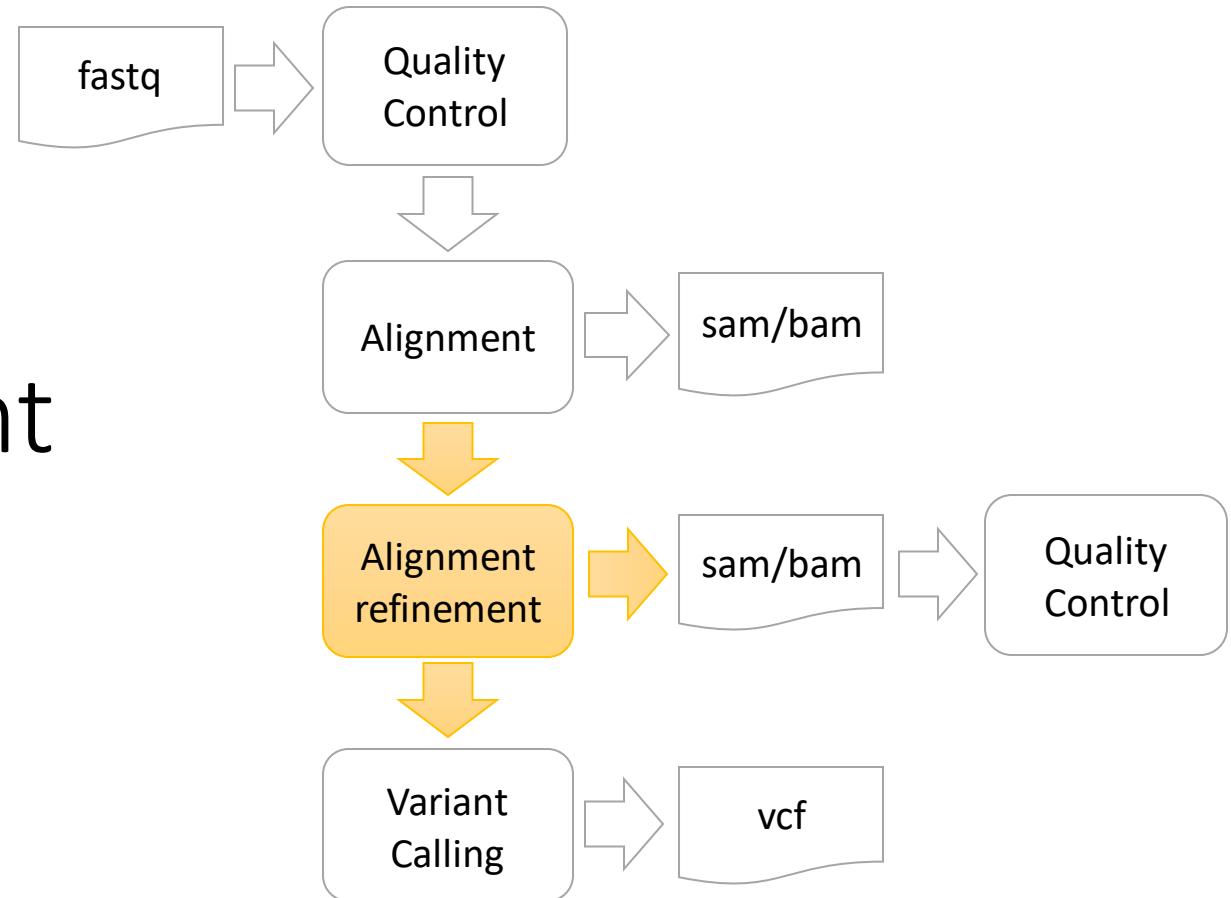
# Strategies on DNA-seq



| # bp pos seq    | Type of variants discovered   | Avg depth per pos | Cost |
|-----------------|---|-------------------|------|
| ~100 Gb         | - <b>coding variants*</b> , intronic and regulatory sites.<br>- <b>Structural variants</b><br>- <b>CNA</b><br>#Variants= 3M - 4M. | 10x               | High |
| ~32M b<br>50M b | - <b>coding variants*</b> .<br>- Some intronic and regulatory sites.<br>- <b>CNA (challenging)</b> .<br>#Variants= 20k - 60k.     | 20x - 80x         | Low  |
| ND              | Depends on the design<br>- Particular <b>coding variants*</b><br>- <b>CNA (challenging)</b><br># variants = ND                    | 1000x - 5000x     | Low  |

\***coding variants**: missense, stop gained, stop lost, frameshift, splice region...

# Alignment refinement



# Remove duplicates

- Duplicates derive from PCR amplification (library preparation): one fragment is sequenced multiple times.
- An error at the beginning of the PCR (first steps) is propagated.
- Therefore, duplicates are **worthless** for the analysis:  
*Duplicates are source of False Positives calls while only provide redundancy.*

**Solution: retrieve the best one, discard the duplicates:**

Duplicates share the  
**same alignment**  
properties : sequence,  
start and end positions



**METHOD:** by Picard-tools  
<http://broadinstitute.github.io/picard/>  
(alternatives : samtools)

Adapted from GATK

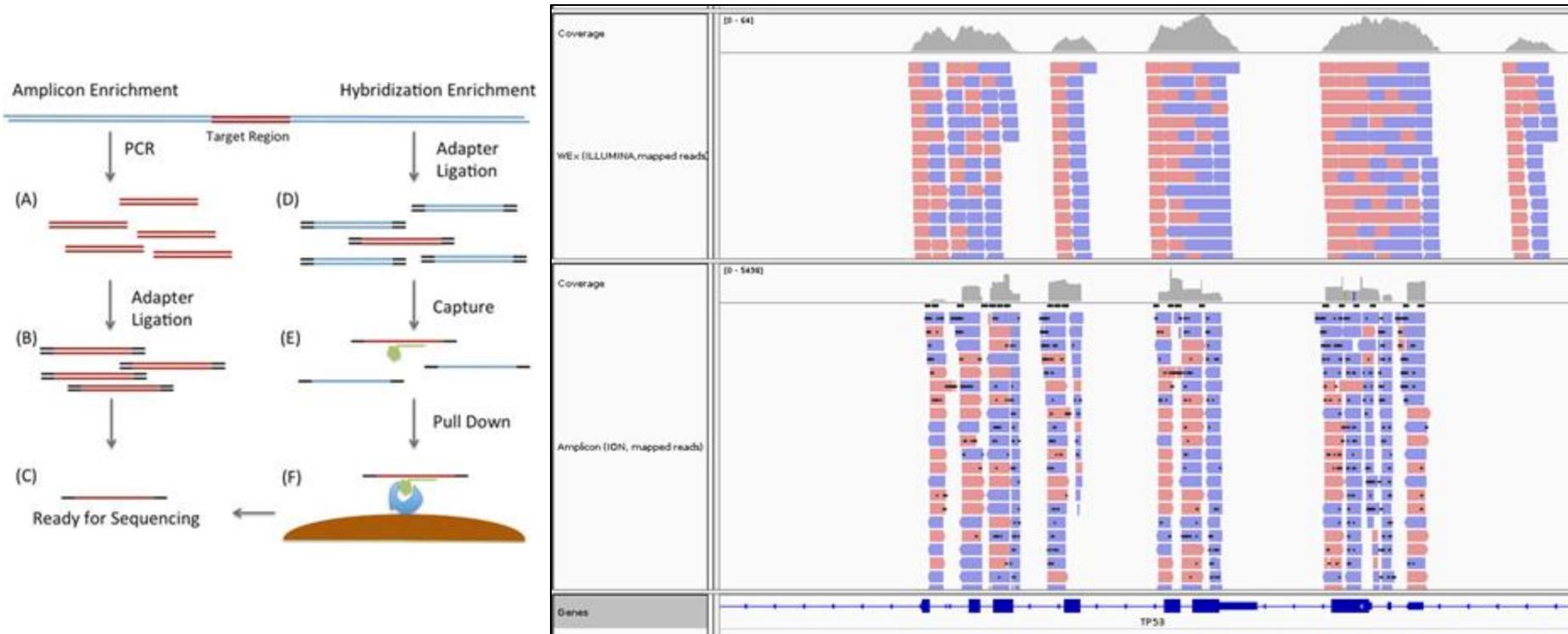


After marking duplicates, the variant caller will only see :



... and thus be more likely to make the right call

# Remove duplicates: Amplicon seq



**WARNING:** Do NOT remove duplicates in data derived from amplicon techniques (**Ion Torrent**).

More info.: <http://gatkforums.broadinstitute.org/discussion/5847/remove-duplicates-from-targetted-sequencing-using-amplicon-approach>

# Indel realignment

- Algorithms align reads very fast with high accuracy, but not perfectly.

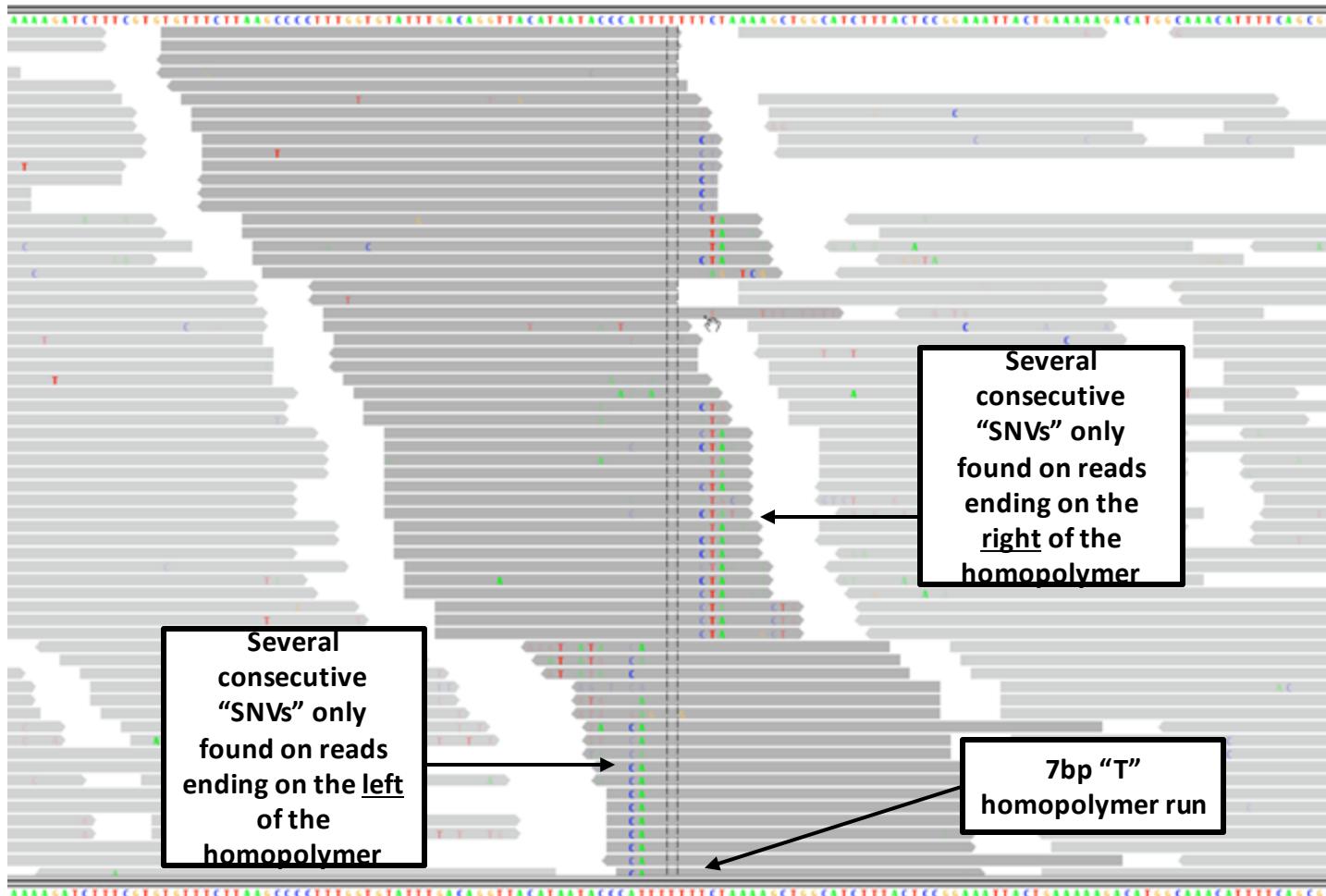
*During alignment, penalties on mismatches are much cheaper than gaps (indels).*

- Also, there are sometimes multiple solutions (alignments) for a given read. Aligners choose one randomly.
- Variant calling requires the most perfect alignment as possible to avoid False Positives.

**METHOD:** by GATK

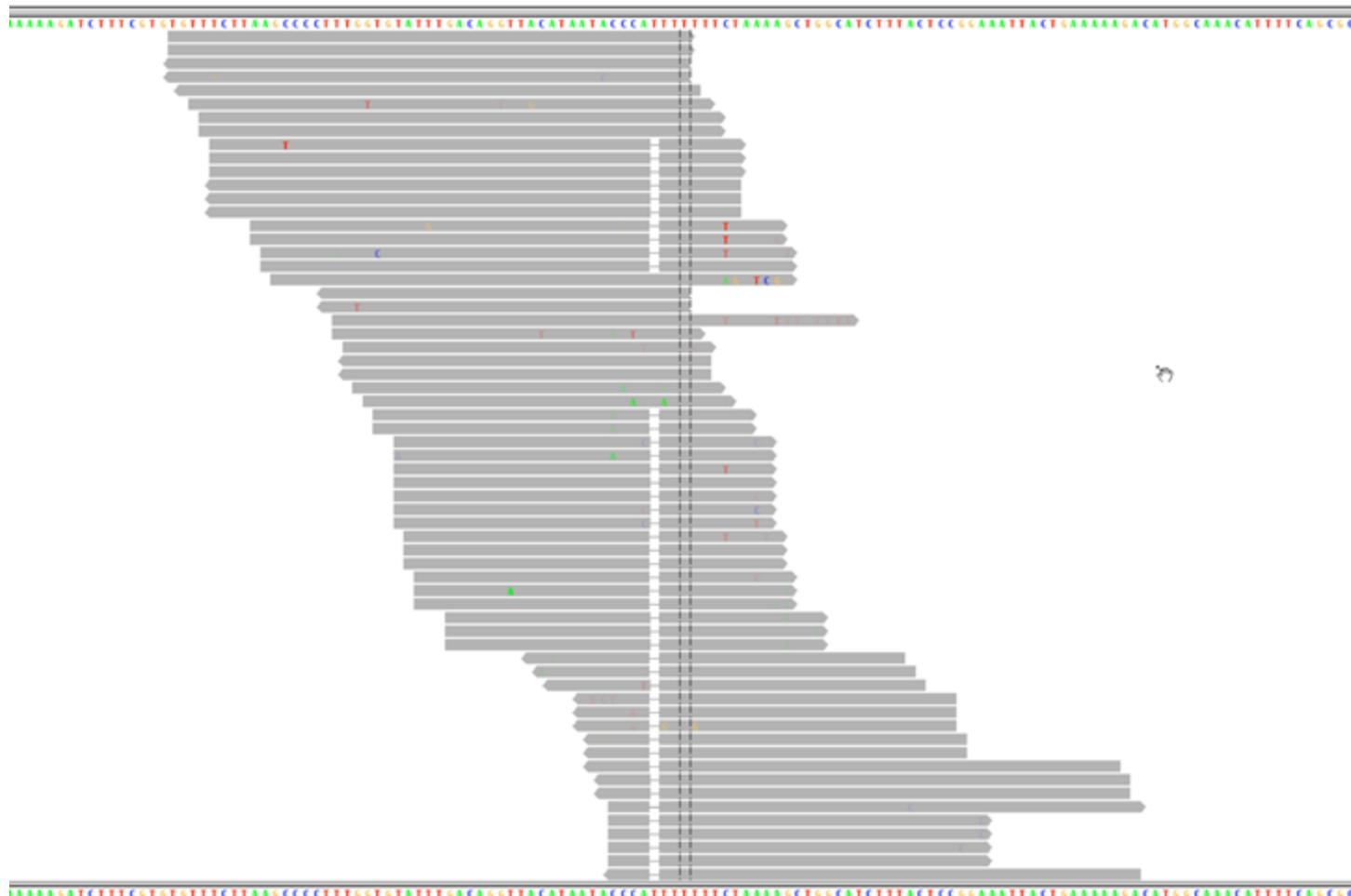
[https://www.broadinstitute.org/gatk/gatkdocs/org\\_broadinstitut\\_e\\_gatk\\_tools\\_walkers\\_indels\\_IndelRealigner.php](https://www.broadinstitute.org/gatk/gatkdocs/org_broadinstitut_e_gatk_tools_walkers_indels_IndelRealigner.php)

# Indel realignment



Taken from GATK team

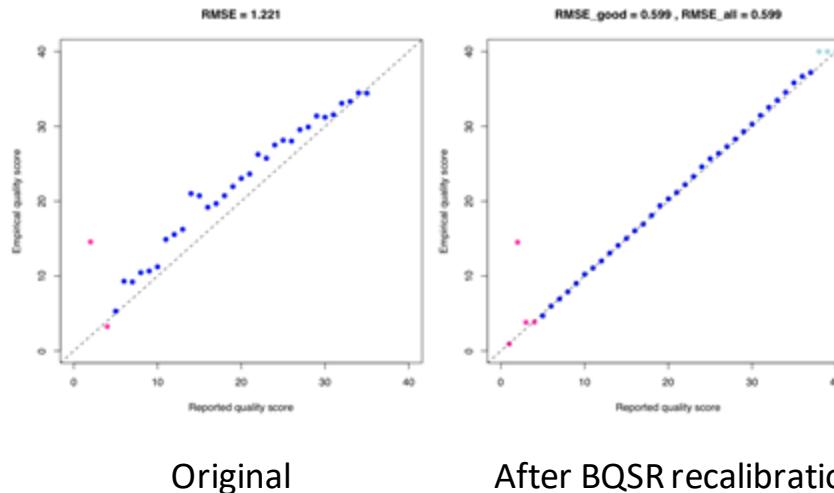
# Indel realignment



Taken from GATK team

# Base Quality Score Recalibration

- **Phred Quality score:** each position of the sequence has its particular **base Quality score**.
- The individual quality measures are crucial during Variant calling.
- Different NGS technologies have their particular **bias in Quality Score** depending on the context. Recalibration **correct empirically** these biases.



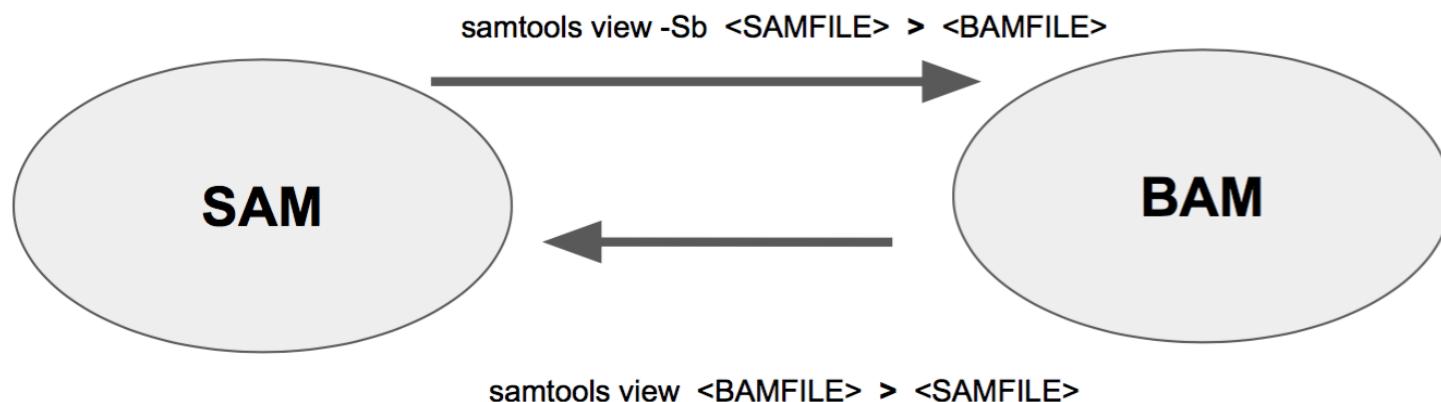
**METHOD:** by GATK  
<http://gatkforums.broadinstitute.org/discussion/44/base-quality-score-recalibration-bqsr>

Alignment  
file

sam & bam

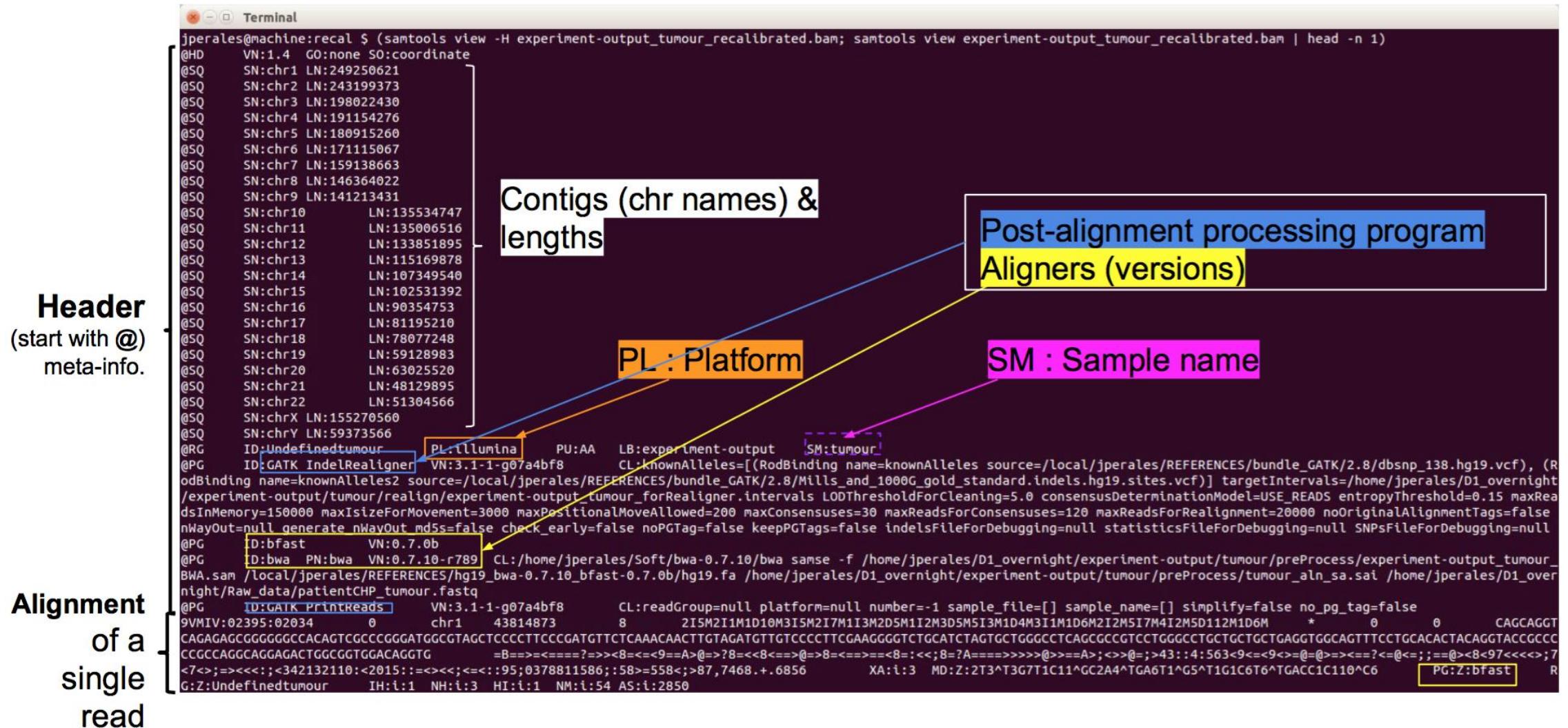
# Alignment - SAM/BAM

- **SAM** is the human readable text format (.sam extension)
- **BAM** is the binary, machine efficient format (.bam extension)
- Both contains exactly the same information and are interconvertible (samtools)



File specifications: <https://samtools.github.io/hts-specs/SAMv1.pdf>

# Alignment - SAM/BAM - Header



# Alignment - SAM/BAM - Alignments

| #1       | #2 | #3    | #4 | #5 | #6        | #7 | #8 | #9 | #10    | #11      | #12       |
|----------|----|-------|----|----|-----------|----|----|----|--------|----------|-----------|
| ReadName | 99 | chr10 | 2  | 30 | 3MD2M1I1M | =  | 14 | 20 | CATCTG | jjjjjjjj | z:Aligner |

If single-end:

7. reference sequence name of the alignment of the next read in sequence
8. position in the alignment of the next read in sequence
9. number of bases covered by reads from the same fragment. Plus/minus means the current read is the leftmost/rightmost read

## SAM FLAGS

| FIELDS |       |  |
|--------|-------|--|
| #Col   | Field | Description  |
| 1.     | QNAME | read name  |
| 2.     | FLAG  | bitwise FLAG* (unmapped, pair unmapped, properly mapped, ...)  |
| 3.     | RNAME | Reference sequence name (e.g. chr1).                           |
| 4.     | POS   | 1-based leftmost position.                                     |
| 5.     | MAPQ  | Mapping Quality (Phred-scaled). Scale 0 to 255.                |
| 6.     | CIGAR | extended CIGAR string.   |
| 7.     | MRNM  | Paired-end: Mate Reference sequence Name (= if same as RNAME). |
| 8.     | MPOS  | Paired-end: 1-based Mate position.                             |
| 9.     | TLEN  | Paired-end: Insert size  |
| 10.    | SEQ   | Read sequence  |
| 11.    | QUAL  | Base Quality Score from the Read sequence.                     |
| 12.    | OPT   | Optional Tags.   |

| Decimal | Description of read                       |
|---------|---|
| 1       | Read paired                               |
| 2       | Read mapped in proper pair                |
| 4       | Read unmapped                             |
| 8       | Mate unmapped                             |
| 16      | Read reverse strand                       |
| 32      | Mate reverse strand                       |
| 64      | First in pair                             |
| 128     | Second in pair                            |
| 256     | Not primary alignment                     |
| 512     | Read fails platform/vendor quality checks |
| 1024    | Read is PCR or optical duplicate          |
| 2048    | Supplementary alignment                   |

One of the reads is unmapped:  
[73](#), [133](#), [89](#), [121](#), [165](#), [181](#), [101](#), [117](#),  
[153](#), [185](#), [69](#), [137](#)

Both reads are unmapped:  
[77](#), [141](#)

Mapped within the insert size and in correct orientation:  
[99](#), [147](#), [83](#), [163](#)

Mapped within the insert size but in wrong orientation:  
[67](#), [131](#), [115](#), [179](#)

Mapped uniquely, but with wrong insert size:  
[81](#), [161](#), [97](#), [145](#), [65](#), [129](#), [113](#), [177](#)

# Alignment - SAM/BAM - CIGAR

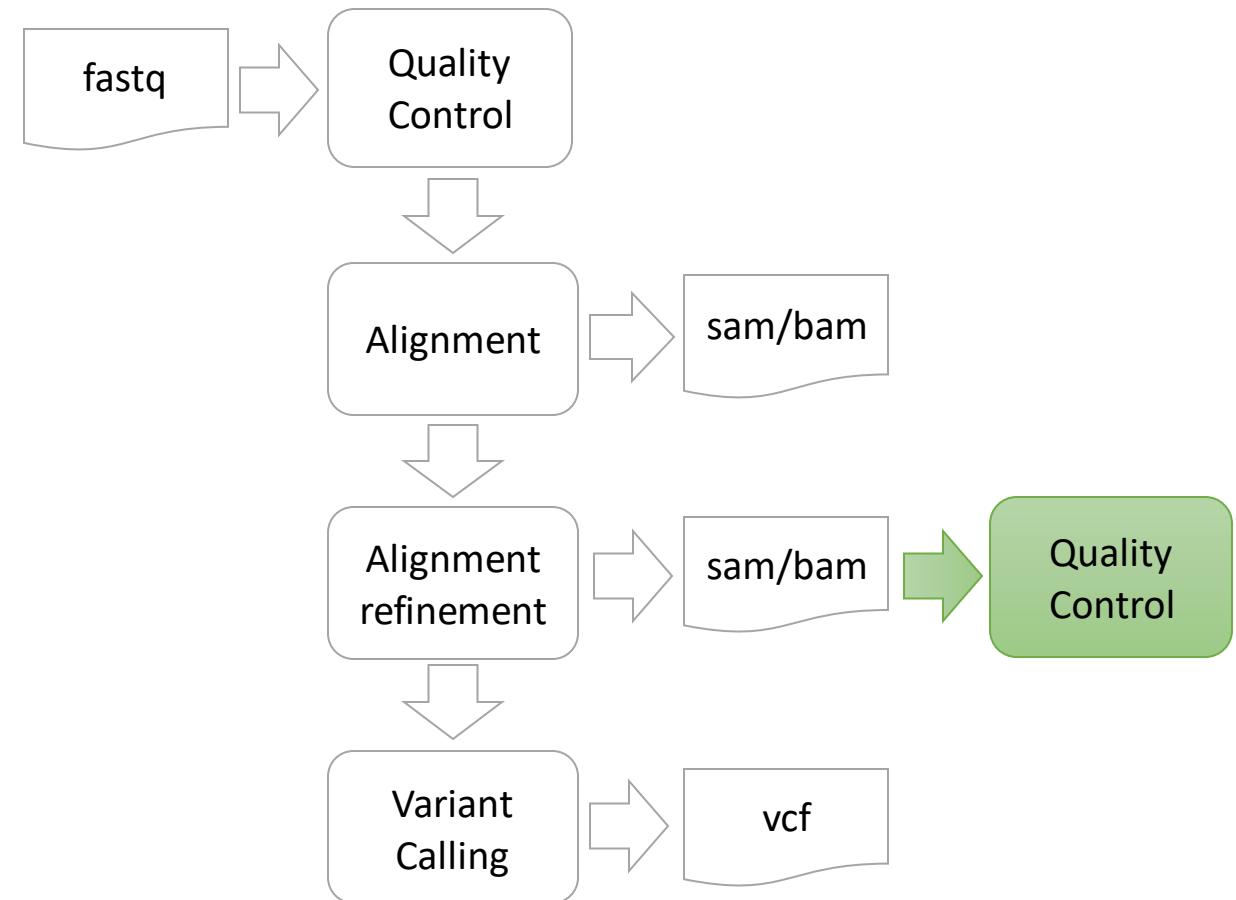
- Concise *Idiosyncratic Gapped Alignment Report*
- It is a compressed representation of an alignment
- Format:** A CIGAR string is made up of <integer><op> pairs
- Here, "op" is an operation specified as a single character, usually an upper-case letter (see table)

| RefPos:    | 1 2 3 4 5 6 7 8 9   |
|------------|---------------------|
| Reference: | C C A T A C T - G A |
| Read:      | C A T - C T A G     |
| POS:       | 2                   |
| CIGAR:     | 3M1D2M1I1M          |

| Op | Description   |
|----|---|
| M  | alignment match (can be a sequence match or mismatch) |
| I  | insertion to the reference                            |
| D  | deletion from the reference                           |
| N  | skipped region from the reference                     |
| S  | soft clipping (clipped sequences present in SEQ)      |
| H  | hard clipping (clipped sequences NOT present in SEQ)  |
| P  | padding (silent deletion from padded reference)       |
| =  | sequence match  |
| X  | sequence mismatch                                     |

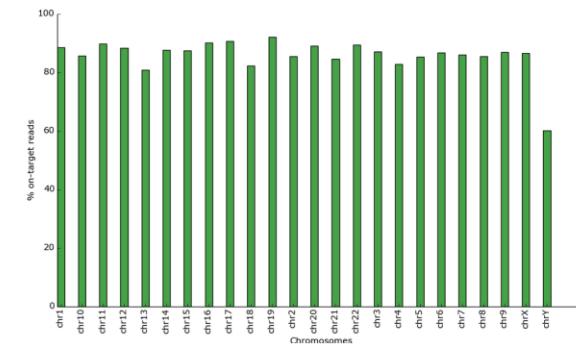
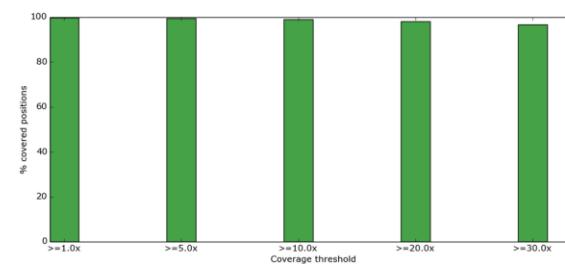
| Reference sequence with aligned reads   | CIGAR string | Explanation          |
|---|--------------|----------------------|
| C T G C A T G T T A G A T A A * * G A T A G C T G T G C T A<br>A A G G A T A * C T G<br>G A T A A * G G A T A | 1M2I4M1D3M   | Insertion & Deletion |
| T G T T A   | 5M1P1I4M     | Padding & Insertion  |
| a a a C A T G T T A G   | 5M15N5M      | Spliced read         |
| A A A C A T G T T A G   | 3S8M         | Soft clipping        |
|   | 3H8M         | Hard clipping        |

# Alignment Quality Control

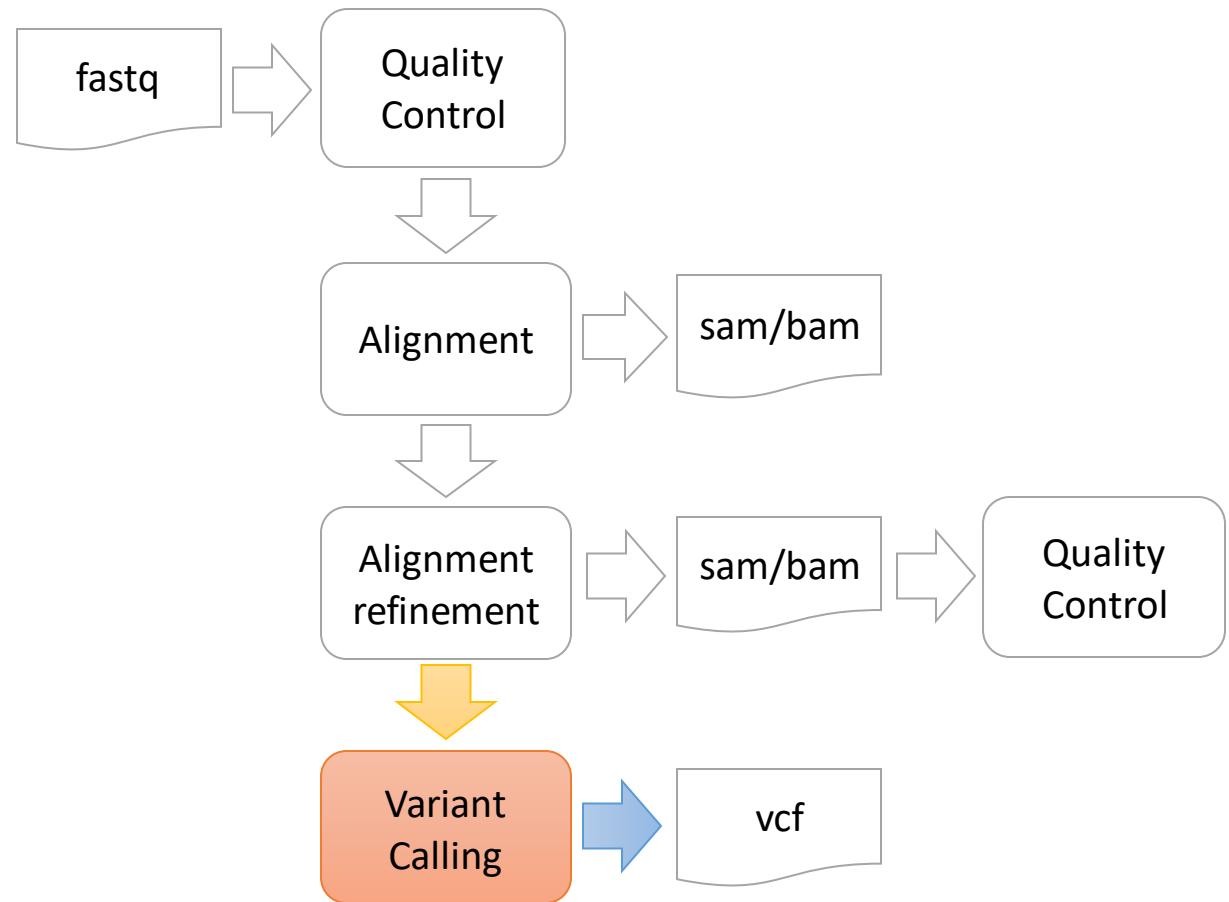


# Alignment Quality Control

- Mean sequencing depth
  - Is there enough coverage in regions of interest?
  - Are the reads on-target?
- 
- Software:
    - ngsCAT
    - QualiMap

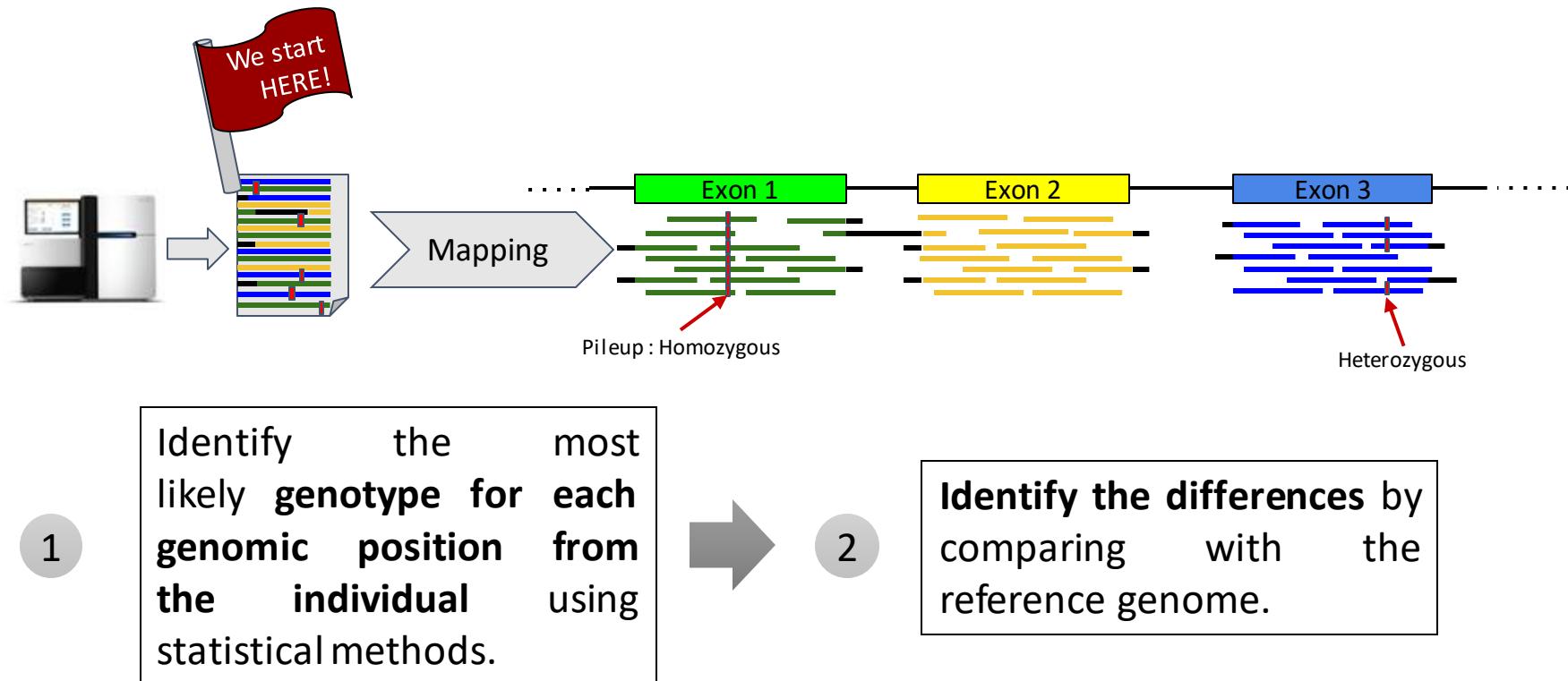


# Variant calling



# Fundamentals of Variant Calling

- How we detect genetic variants ?



# What is Crucial in Variant calling

- For clinical practices, the use of **gold standard methods** and **reproducible analysis** are mandatory.
- The analysis is based on the comparison against the **reference genome**:

*A single consensus sequence for the whole genome. It was built up from a high quality set of representative samples of the species (from different populations).  
It is the first-line comparison during analysis.*

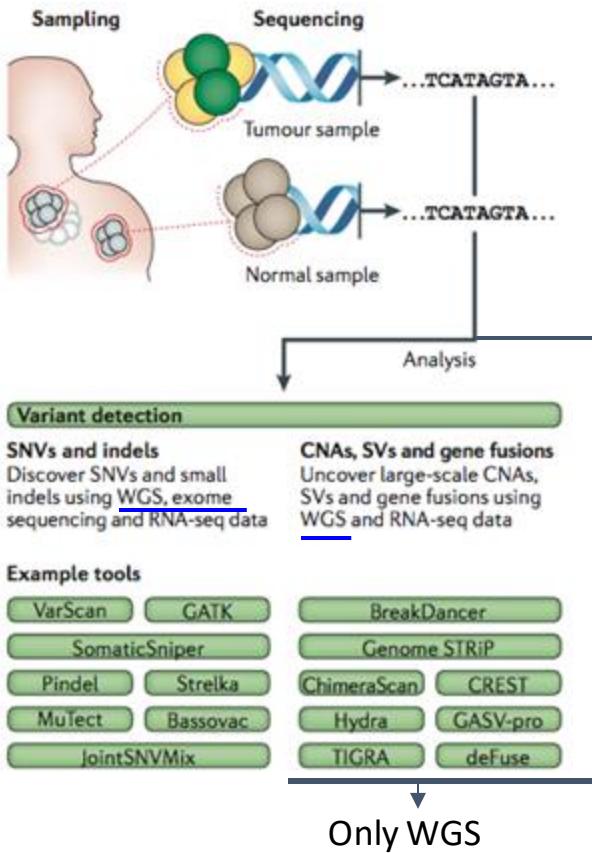
By Genome Reference Consortium (GRC) (<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>)

- Human assemblies (Versions):
  - + GRCh37/hg19 : former version. Released in 2012. It is still used for analysis.
  - + CRCh38/hg38 : current version (Sep. 2017). Released in 2014. More accurate, comprehensive (includes Haplotypes) and sophisticated.

We must keep consistency in the Genome Reference Version through the variant analysis.

- We must know what **regions along the genome were sequenced** in the experiment, that is, the sequencing library.

# Algorithms for Variant Calling



Several Methods have been published.

| Tool          | Year | Language      | Paired or pooled data | Segmentation | Feature   |
|---------------|------|---------------|-----------------------|--------------|---|
| ADTEX         | 2014 | Python, R     | Both                  | HMM          | Noise reduction<br>Ploidy estimation                |
| CONTRA        | 2012 | Python, R     | Both                  | CBS          | GC correction                                       |
| Control-FREEC | 2011 | C++, R        | Paired                | LASSO        | GC correction, mappability                          |
| EXCAVATOR     | 2013 | Perl, R       | Both                  | HSLM         | GC correction, mappability,<br>exon-size correction |
| ExomeCNV      | 2011 | R             | Paired                | CBS          | GC correction, mappability                          |
| Varscan2      | 2012 | Java, Perl, R | Paired                | CBS          | GC correction                                       |

Appropriate methods for Whole-Exome seq

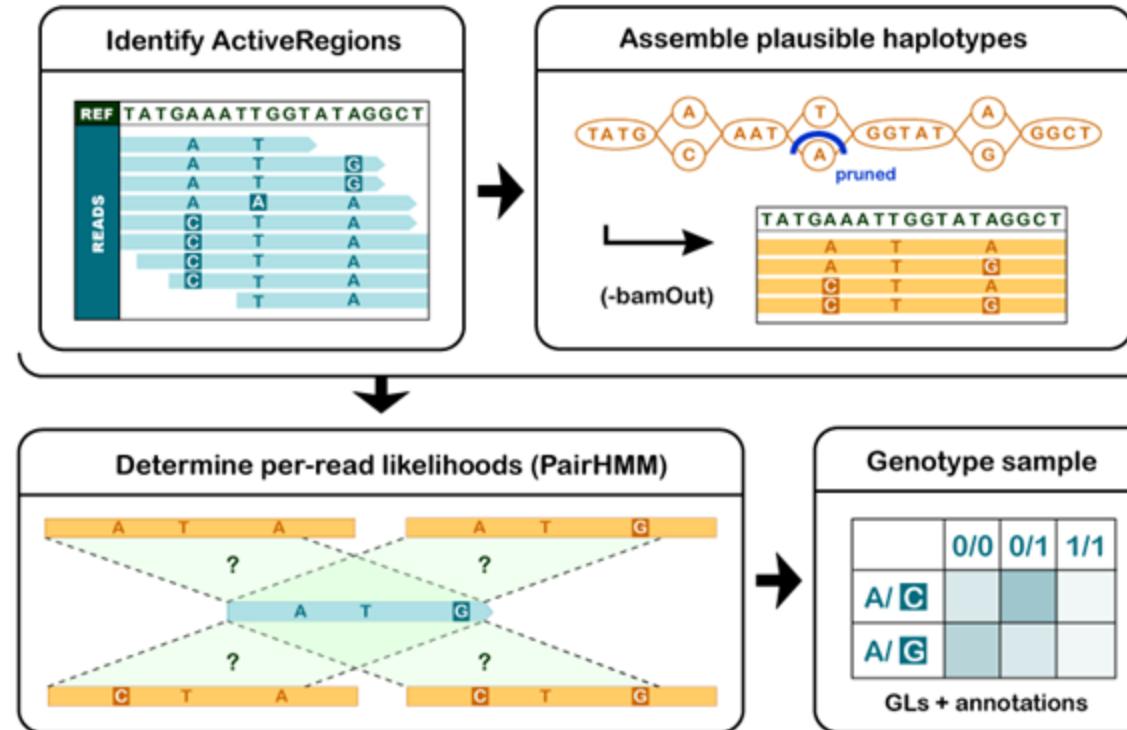
## Further reading:

Ding, L., Wendl, M. C., McMichael, J. F. & Raphael, B. J. Nat Rev Genet – (2014). doi:10.1038/nrg3767

Nam J.N. et al. Evaluation of somatic copy number estimation tools for whole-exome sequencing data. Brief. Bioinformatics (2015)

# Variant Calling for SNVs and Indels

**Haplotype Caller:** Variant calling based on the calculation of genotype likelihoods:



**Assumptions:** It bases the calling in the indicated ploidy (e.g. 2n)  
**Limited detection of low allele frequencies.**

Further reading:

<http://gatkforums.broadinstitute.org/discussion/4148/hc-overview-how-the-haplotypecaller-works>

HC steps 1-4: <https://software.broadinstitute.org/gatk/documentation/topic?name=methods>

# Variant Calling for somatic variants: MuTect

**SNV and Indel caller.**

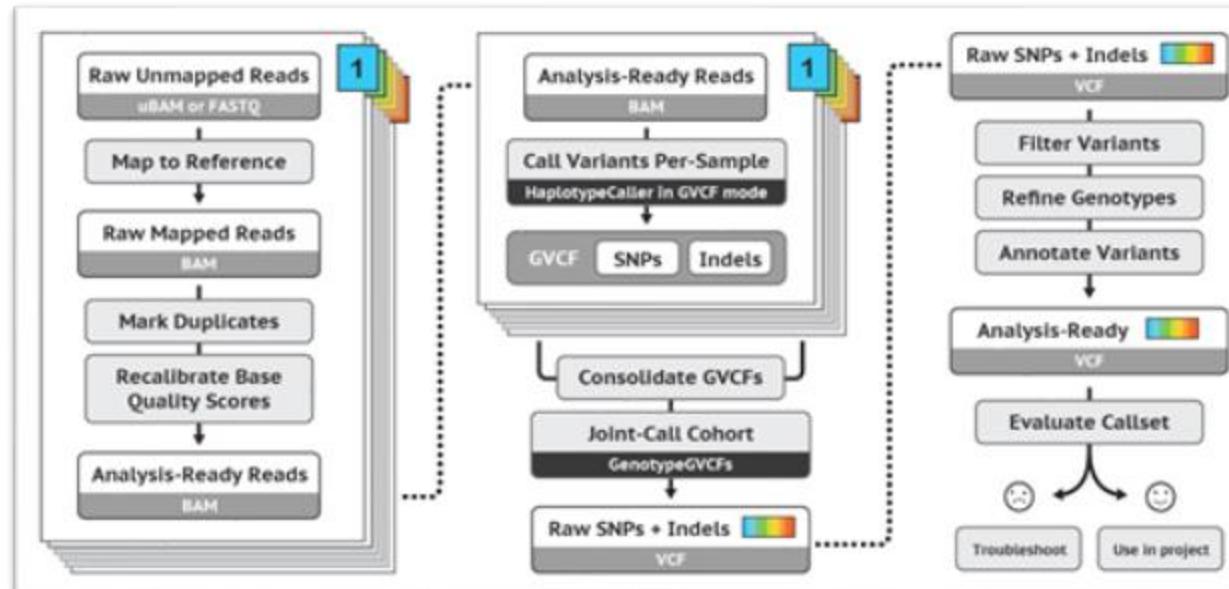
Similar logic to Haplotype Caller but:

- It allows variable allele frequencies.
- It includes logic to avoid germline variants.

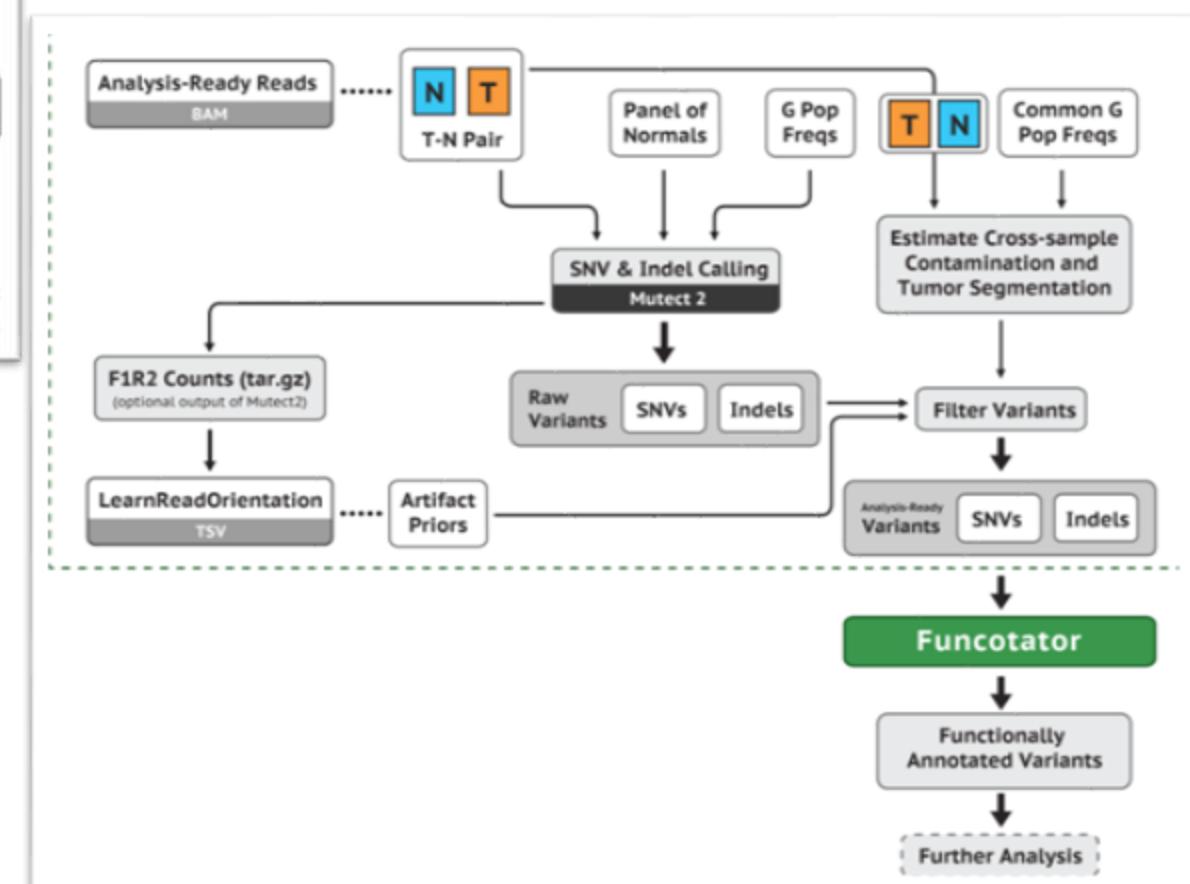
The screenshot shows two adjacent web pages. The left page is the CGA (Cancer Genome Analysis) homepage, featuring a sidebar with links to various tools like ABSOLUTE, BreakPoint, ChainFinder, Copynumber Pipeline, D-TrixO, dGanger, Firehose, GISTIC, HAPSEG, Indelocator, iMEx, and MuTect. The right page is the MuTect tool's documentation page, which includes a brief introduction, a 'How does it work?' section with a mathematical formula for calculating LOD scores, and a note about false positive rates. At the bottom, there is a table summarizing validation rates for different cancer types.

| publication                   | technology       | candidates | validated | no result | validation rate |
|-------------------------------|------------------|------------|-----------|-----------|-----------------|
| Multiple Myeloma <sup>a</sup> | Sequenom         | 97         | 92        | 5         | 94.85%          |
| Ovarian <sup>b</sup>          | Sequenom/PCR/454 | 1655       | 1483      | 172       | 89.41%          |
| Ovarian <sup>c</sup>          | Capture/Illumina | 6497       | 6232      | 265       | 95.82%          |
| Head and Neck <sup>d</sup>    | Sequenom         | 321        | 288       | 33        | 89.72%          |
| Breast <sup>e</sup>           | Sequenom/PCR/454 | 455        | 428       | 0         | 94.07%          |

# GATK Best Practices



Germline small-scale variants

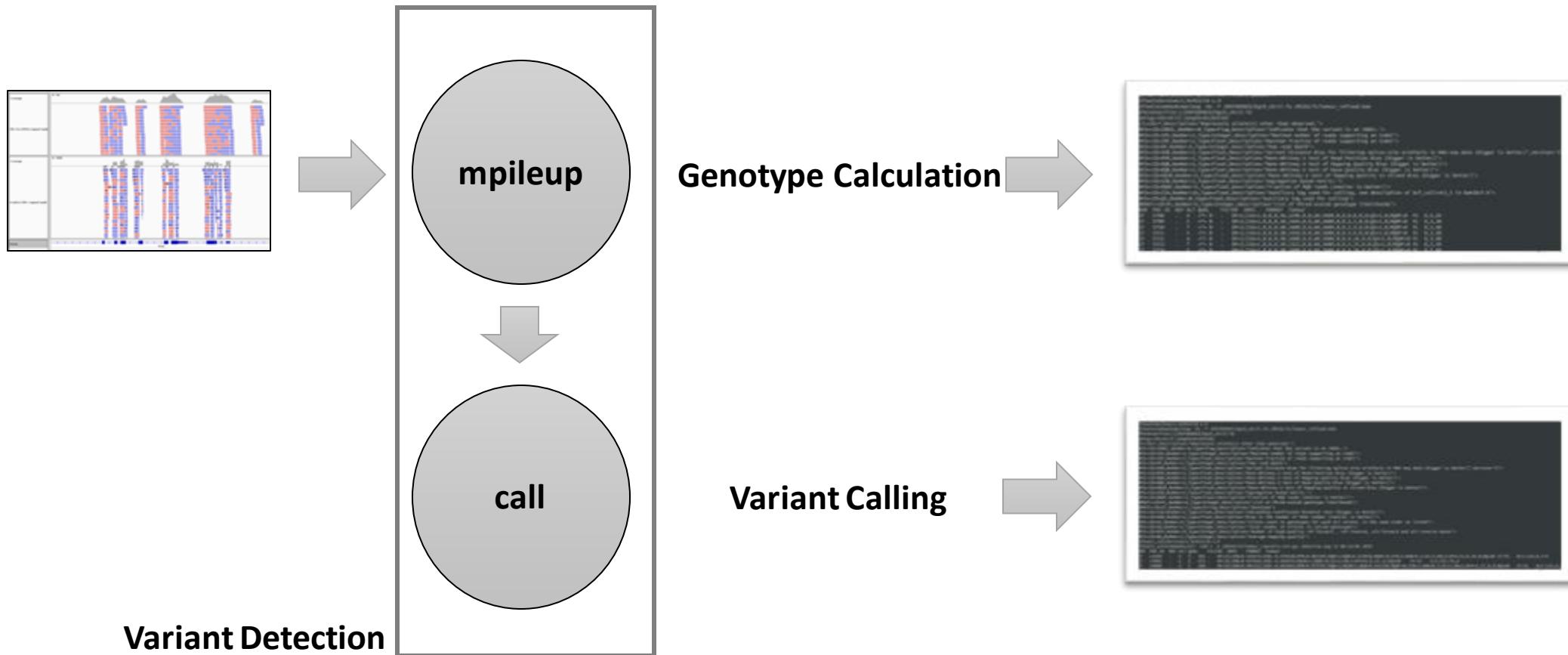


Somatic small-scale variants

# bcftools

<http://samtools.github.io/bcftools/>

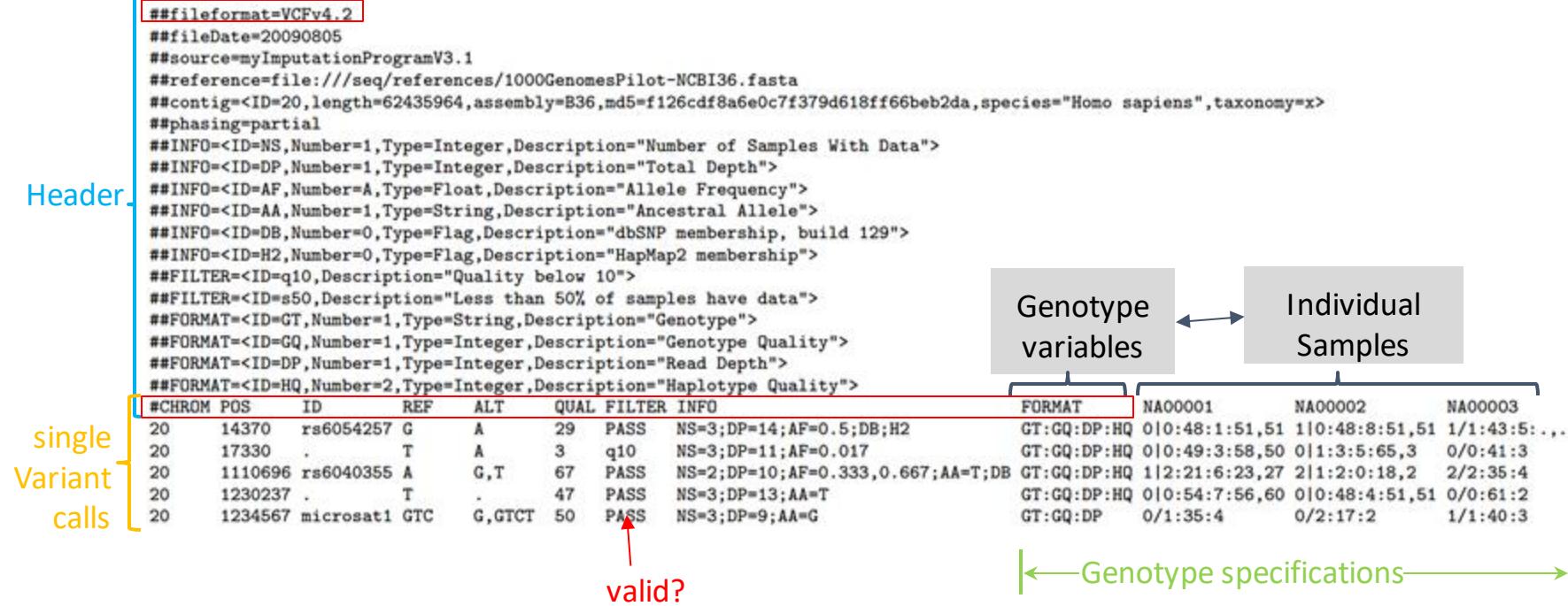
**Set of tools to call variants and manage VCF files.**



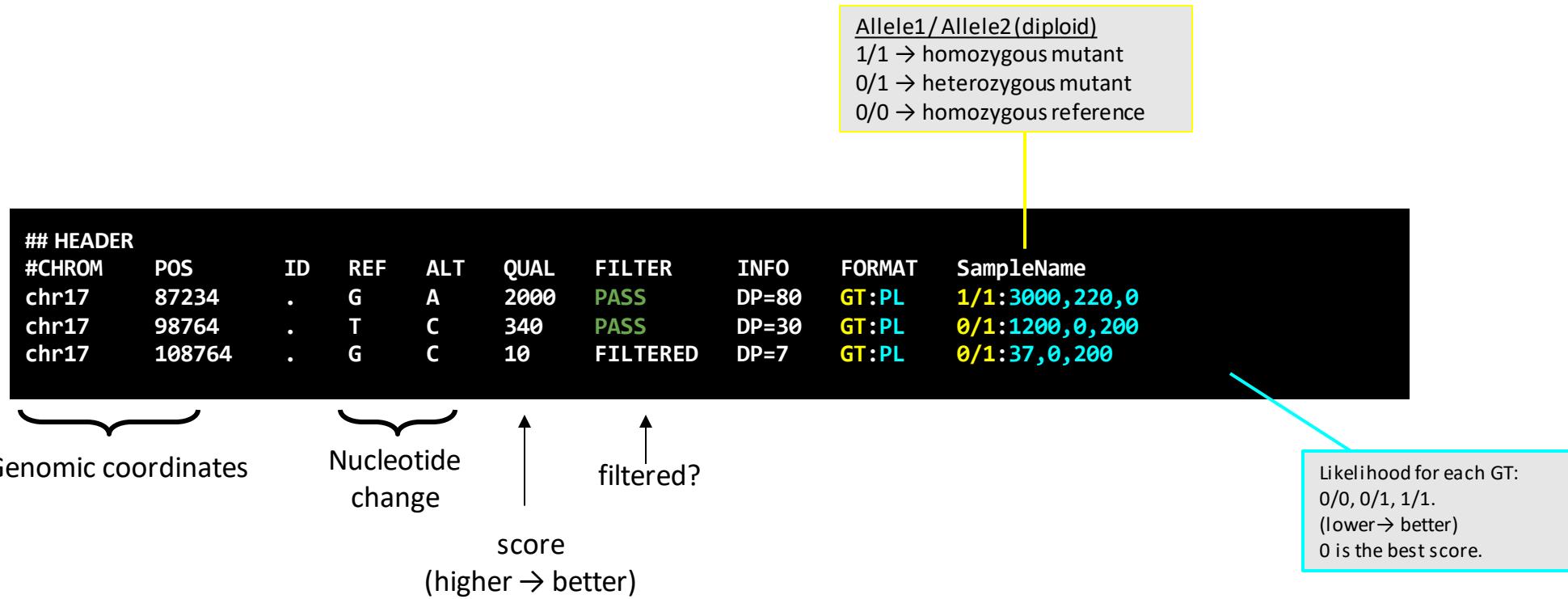
Variant's  
file

vcf

# Variant Calling Format - VCF



- Typical extension: .vcf (.bcf binary counterpart)
- Not all records in a VCF are true calls, the FILTER column specifies those which passed the calling
- QUAL is the score assigned to a given call. The greater QUAL is, the more reliable is. It is in log-scale



More info.:

<https://samtools.github.io/hts-specs/VCFv4.2.pdf>

<https://www.broadinstitute.org/gatk/guide/article?id=1268>

# Variant Detection

Hands on

# Data

- **Input** : FastQ files.
- **Output vcf**: variants in tumor sample, in normal sample and in both samples
- **Raw files:** [https://fundacioncnio-my.sharepoint.com/:u/g/personal/epineiro\\_cnio\\_es/EaWrloqPdb5GjXak1O4TzkBUmlI6sqnupz8SLuJ2aT5-A?e=hv0Ns](https://fundacioncnio-my.sharepoint.com/:u/g/personal/epineiro_cnio_es/EaWrloqPdb5GjXak1O4TzkBUmlI6sqnupz8SLuJ2aT5-A?e=hv0Ns)



**Patient suffering ovarian cancer.**

**Whole-exome sequencing** data from two samples from the patient:

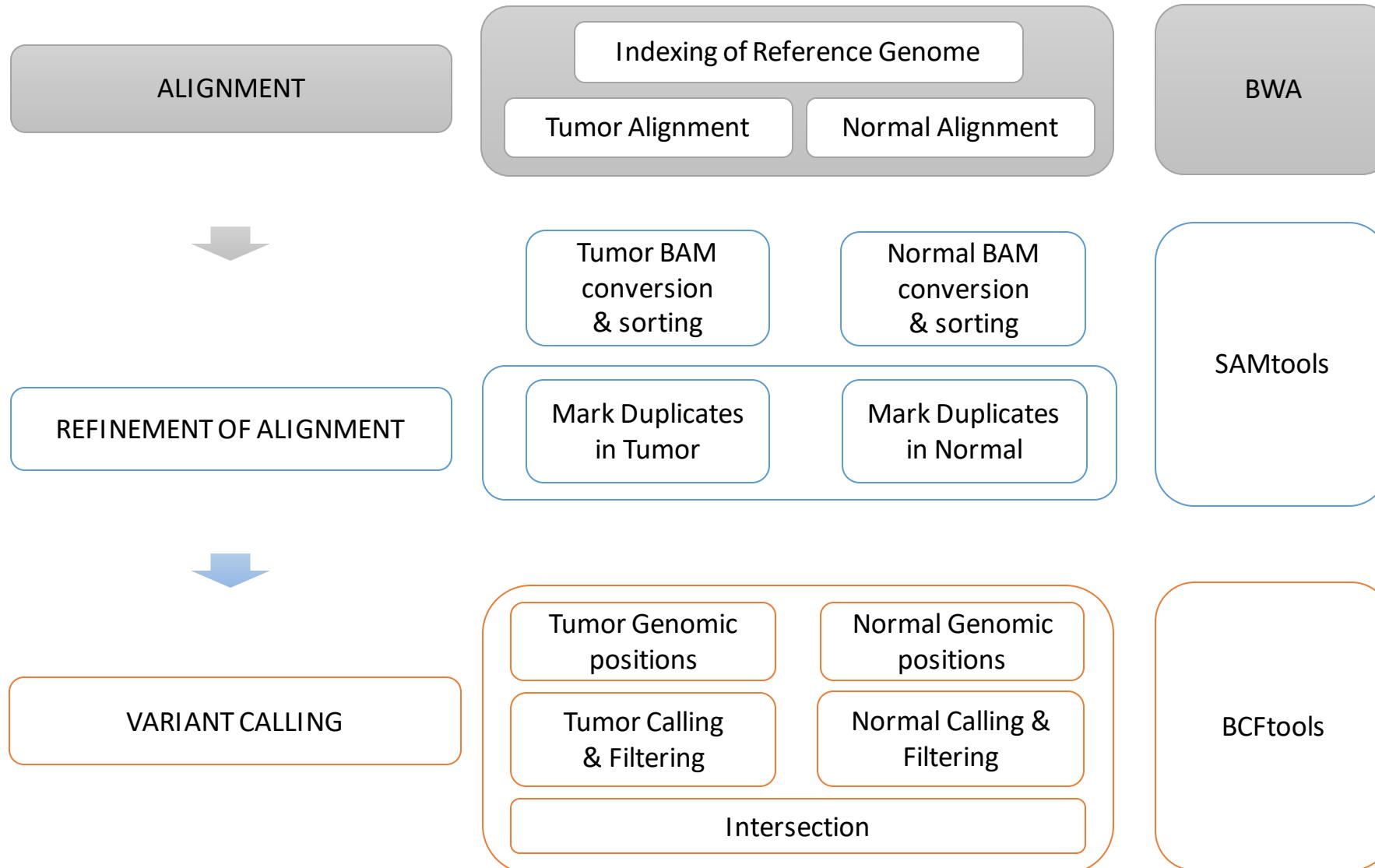
- Tumour sample.
- Matched normal sample (healthy tissue) from epithelium.

**Library protocol:** Agilent SureSelect V5 Human All Exons.

**Sequencing platform:** HiSeq 2000 (Illumina)

NOTE: This data was simulated and reduced in order to perform the computational analysis in class time.

# ROAD MAP



# Step 1. Index the human reference genome

Reference genome: [https://fundacioncnio-my.sharepoint.com/:u/g/personal/epineiro\\_cnio\\_es/ETg0aiGbVzdChTX2VNv3PZ0BJMOmovsLtaDbMSbMoR6i8A?e=1ymXkA](https://fundacioncnio-my.sharepoint.com/:u/g/personal/epineiro_cnio_es/ETg0aiGbVzdChTX2VNv3PZ0BJMOmovsLtaDbMSbMoR6i8A?e=1ymXkA)

```
docker > docker pull osvaldogc/ufv:2.0 1
> SAMPLES_LOCAL=/Path_directory_raw_data/
> SAMPLES_DOCKER=/SAMPLES
> REFERENCE_LOCAL=/Path_to_reference_fasta/ 2
> REFERENCE_DOCKER=/REFERENCE
> RESULTS_LOCAL=/Path_directory_RESULTS
> RESULTS_DOCKER=/RESULTS

> chmod 777 $REFERENCE_LOCAL 3
> chmod 777 $RESULTS_LOCAL

> docker run --rm -v $SAMPLES_LOCAL:$SAMPLES_DOCKER -v $RESULTS_LOCAL:$RESULTS_DOCKER -v
$REFERENCE_LOCAL:$REFERENCE_DOCKER -it osvaldogc/ufv:2.0 /bin/bash 4
```

# Step 1. Index the human reference genome

```
> bwa index /REFERENCE/hg19_chr17.fa
```

Indexing the human reference genome

Note:

- The indexing of the reference genome must be done only the first time you are going to perform the alignment of any sample on it.
- During the indexing, several files are generated with the same prefix as the fasta file (hg19\_chr17.fa.\*)

Check that those files were generated:

```
> ls -l /REFERENCE
```

*hg19\_chr17.fa.amb  
hg19\_chr17.fa.ann  
hg19\_chr17.fa.bwt*

Check the creation of the index files

*hg19\_chr17.fa.pac  
hg19\_chr17.fa.sa*

# Step 2. Alignment

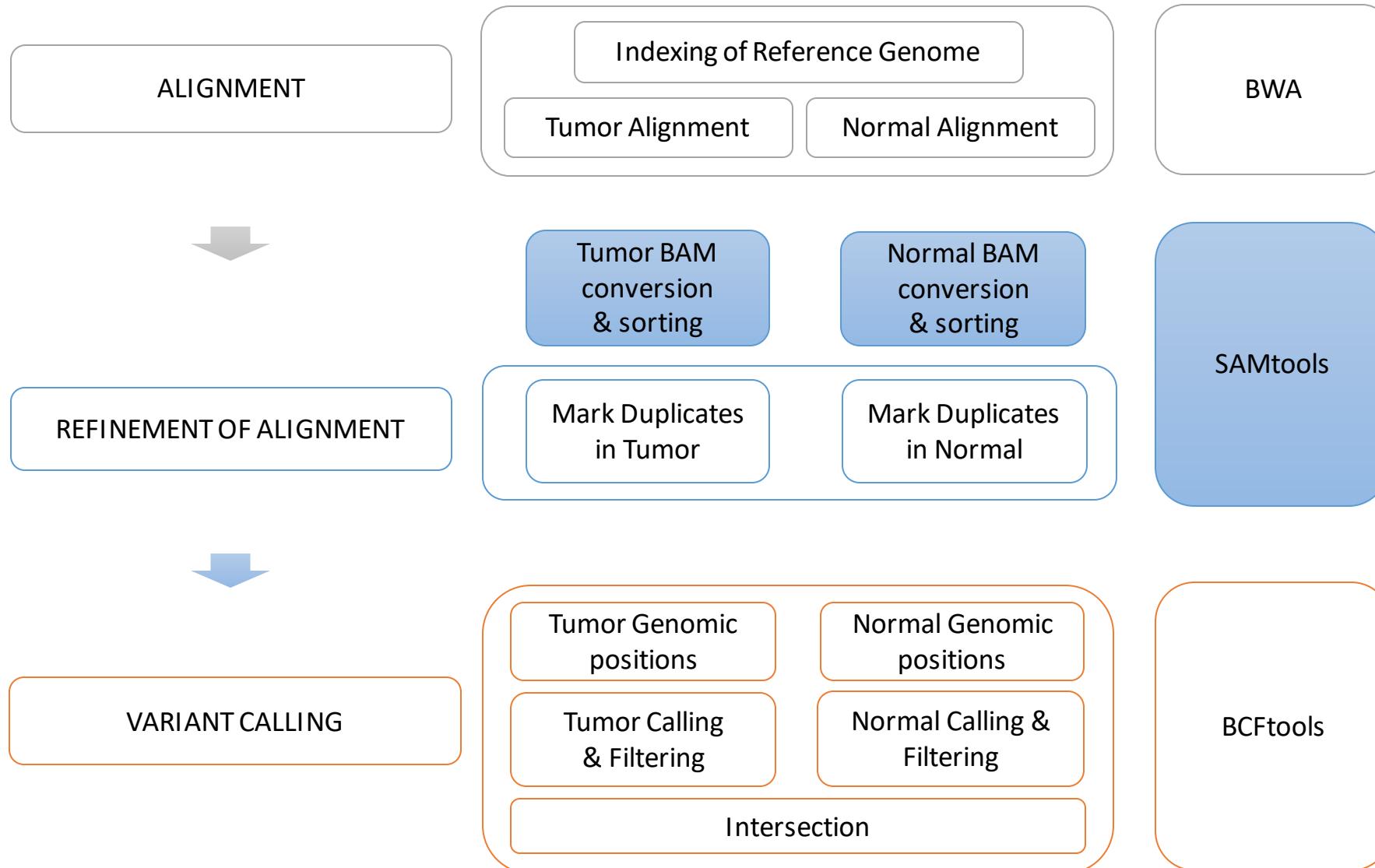
```
> bwa mem -R '@RG\@ID:OVCA\@SM:sample'  
/REFERENCE/hg19_chr17.fa  
/SAMPLES/WEx_sample_R1.fastq  
/SAMPLES/WEx_sample_R2.fastq > /RESULTS/sample.sam
```

Alignment of each of the 2 samples  
against the reference genome

```
[M::main_mem] read 133334 sequences (10000050 bp)...  
[M::mem_pestat] # candidate unique pairs for (FF, FR, RF, RR): (0, 60917, 0, 0)  
[M::mem_pestat] skip orientation FF as there are not enough pairs  
[M::mem_pestat] analyzing insert size distribution for orientation FR...  
[M::mem_pestat] (25, 50, 75) percentile: (333, 364, 400)  
[M::mem_pestat] low and high boundaries for computing mean and std.dev: (199, 534)  
[...]  
[M::worker2@0] performed mate-SW for 10461 reads  
[main] Version: 0.7.5a-r405  
[main] CMD: bwa mem -R @RG\@ID:OVCA\@SM:tumor /REFERENCE/hg19_chr17.fa /SAMPLES/WEx_Tumour_R1.fastq  
/SAMPLES/WEx_Tumour_R2.fastq  
[main] Real time: 169.040 sec; CPU: 129.535 sec
```

WARNING: You must replace **sample** with the ‘sample name’ (i.e. Tumour or Normal)

# ROAD MAP



# Step 3. Compression and sorting

```
> samtools view -S -b /RESULTS/sample.sam > /RESULTS/sample.bam
```

SAM to BAM conversion

BWA sometimes leaves unusual paired information on SAM records. We fix that using this function:

```
> samtools fixmate /RESULTS/sample.bam /RESULTS/sample_fixmate.bam
```

Fix paired information in files

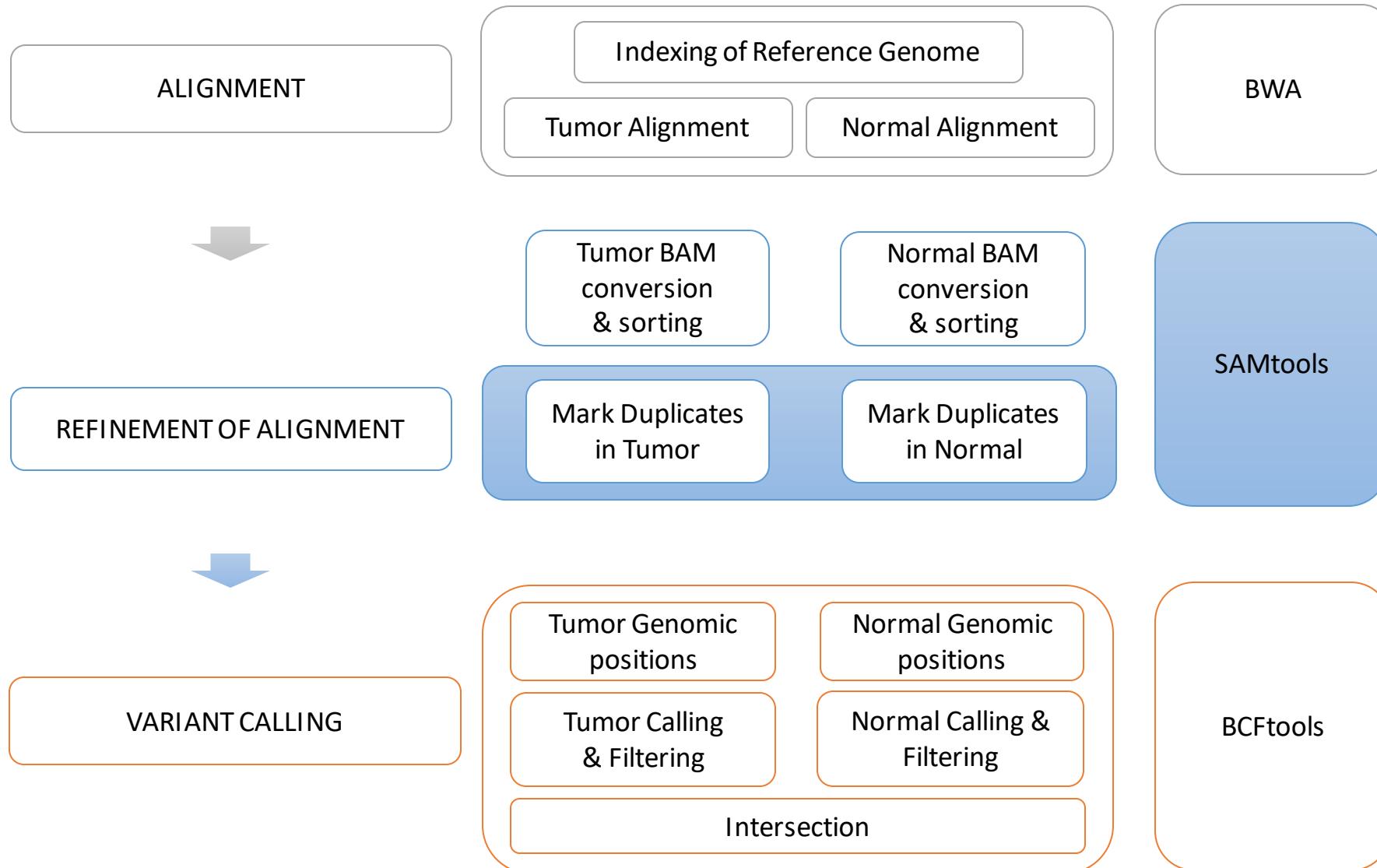
The variant caller requires that the alignment is sorted by genomic positions:

```
> samtools sort /RESULTS/sample_fixmate.bam /RESULTS/sample_sorted
```

Sort BAM files

WARNING: You must replace **sample** with the ‘sample name’ (i.e. Tumour or Normal)

# ROAD MAP



# Step 4. Refinement of the alignment

Next, we mark duplicates (from the PCR). So the Variant Caller will ignore them:

```
> samtools rmdup -S /RESULTS/sample_sorted.bam /RESULTS/sample_refined.bam
```

Mark duplicates

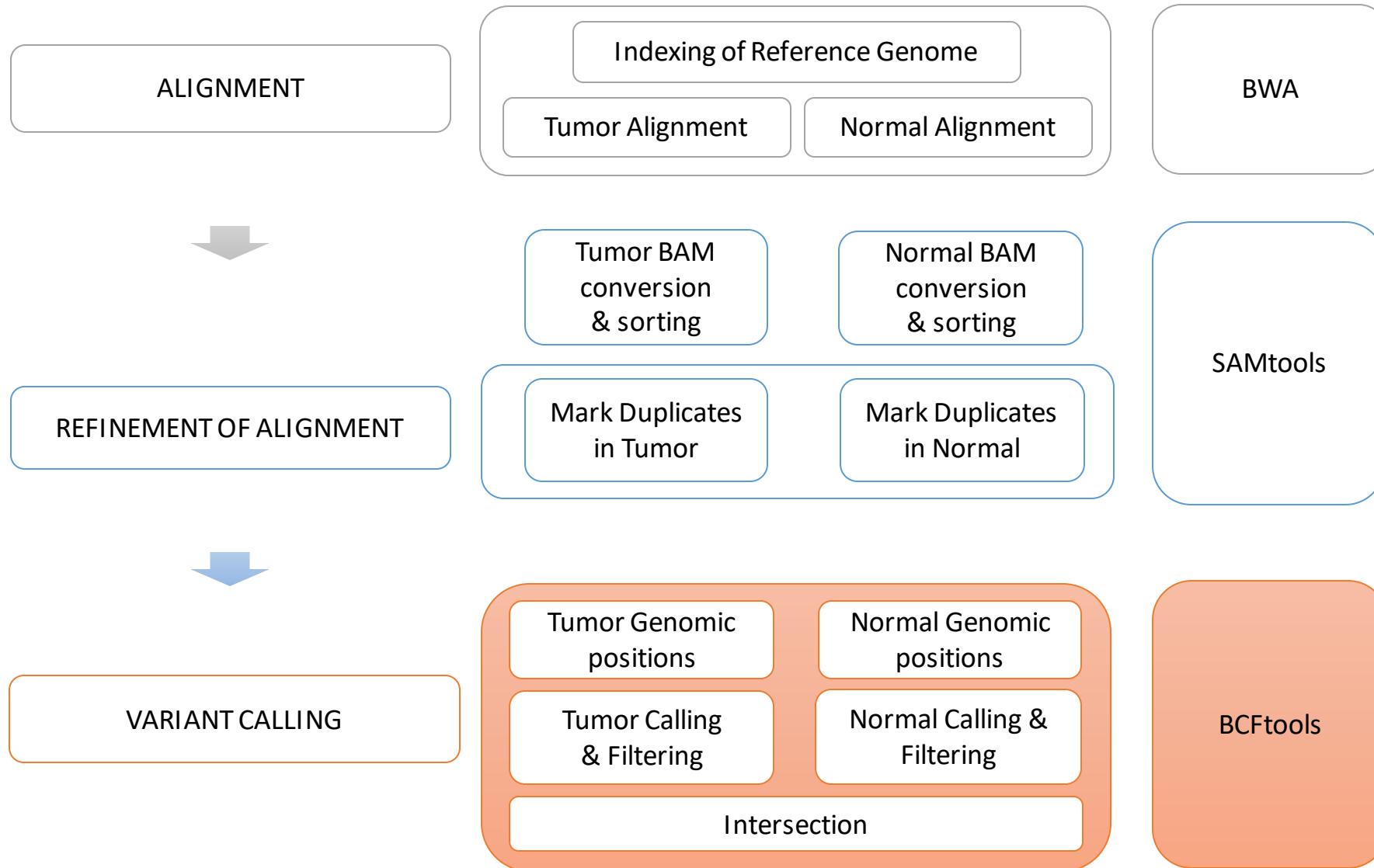
We index the BAM file:

```
> samtools index /RESULTS/sample_refined.bam
```

Index BAM files

WARNING: You must replace **sample** with the 'sample name' (i.e. Tumour or Normal)

# ROAD MAP



# Step 5. Variant calling

Identify "Active regions" & calculate the most likely genotype :

```
> bcftools mpileup -Ou -f /REFERENCE/hg19_chr17.fa /RESULTS/sample_refined.bam | bcftools call -vmO  
z -o /RESULTS/sample_rawcalls.vcf.gz
```

Variant calling

Create the index for the vcf file:

```
> bcftools index /RESULTS/sample_rawcalls.vcf.gz
```

Indexing

WARNING: You must replace **sample** with the 'sample name' (i.e. Tumour or Normal)

# Step 6. Variant calling filtering

Retrieve calls with high quality (include option):

```
> bcftools filter -i 'QUAL>10&&DP>10' /RESULTS/sample_rawcalls.vcf.gz
```

Filtering Option 1

Retrieve calls with high quality (exclude option):

```
> bcftools filter -e 'DP<10' /RESULTS/sample_rawcalls.vcf.gz
```

Filtering Option 2

WARNING: You must replace **sample** with the 'sample name' (i.e. Tumour or Normal)

# Step 7. Intersection. Discern between germline/somatic variants

Intersection of detected variants in the two samples from the patient:

```
> bcftools isec -i 'DP>10' /RESULTS/Tumour_rawcalls.vcf.gz /RESULTS/Normal_rawcalls.vcf.gz -p /RESULTS/
```

Intersection

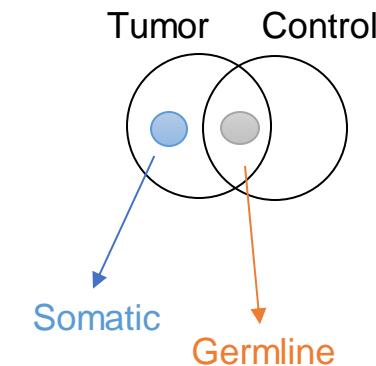
What set of variants are stored in each vcf:

```
> cat /RESULTS/README.txt
```

Check distribution of variants in vcf

Using the following file names:

/RESULTS//0000.vcf for records private to /RESULTS/Tumor\_rawcalls.vcf.gz  
/RESULTS//0001.vcf for records private to /RESULTS/Normal\_rawcalls.vcf.gz  
/RESULTS//0002.vcf for records from /RESULTS/Tumor\_rawcalls.vcf.gz shared by both /RESULTS/Tumor\_rawcalls.vcf.gz /RESULTS/Normal\_rawcalls.vcf.gz  
/RESULTS//0003.vcf for records from /RESULTS/Normal\_rawcalls.vcf.gz shared by both /RESULTS/Tumor\_rawcalls.vcf.gz /RESULTS/Normal\_rawcalls.vcf.gz  
root@f280f8d0f046:/SOFTWARE/bcftools-1.9# cat /RESULTS/README.txt



## *Case study :: Mutation results*

- Germline variants

There were detected █ germline variants in total:

- █ Single Nucleotide Variants.
- █ Indels.

- Somatic variants

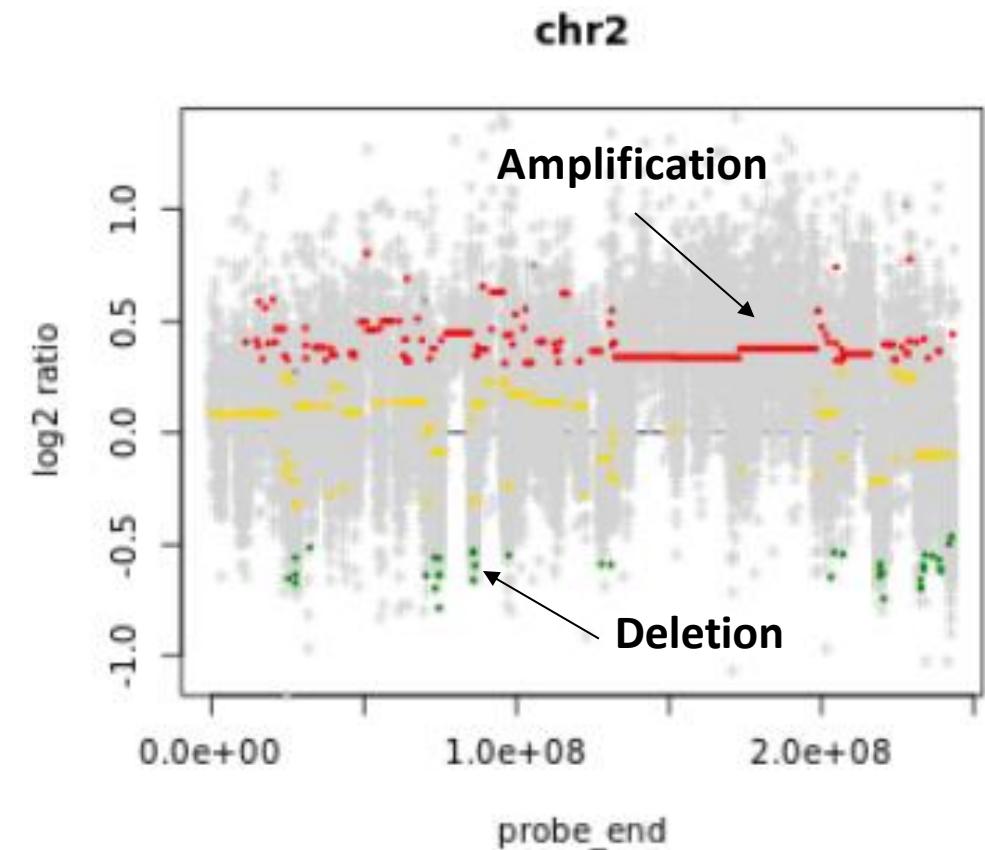
There were detected █ somatic variants in total:

- █ Single Nucleotide Variants.
- █ Indels.

# Copy Number Variation (CNV)

# CNV

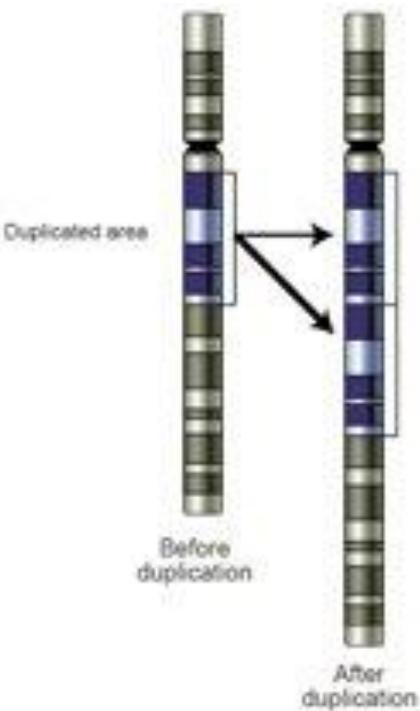
- Variations in the copy number of DNA regions.
- Longer than 1 kb (50 bp in some definitions)
- Can be deletions or amplifications
- Can change the expression of the genes
- Associated to phenotypes and diseases:
  - Cancer
  - Alzheimer
  - Autism
  - Diabetes



# CNV detection

Classical detection:

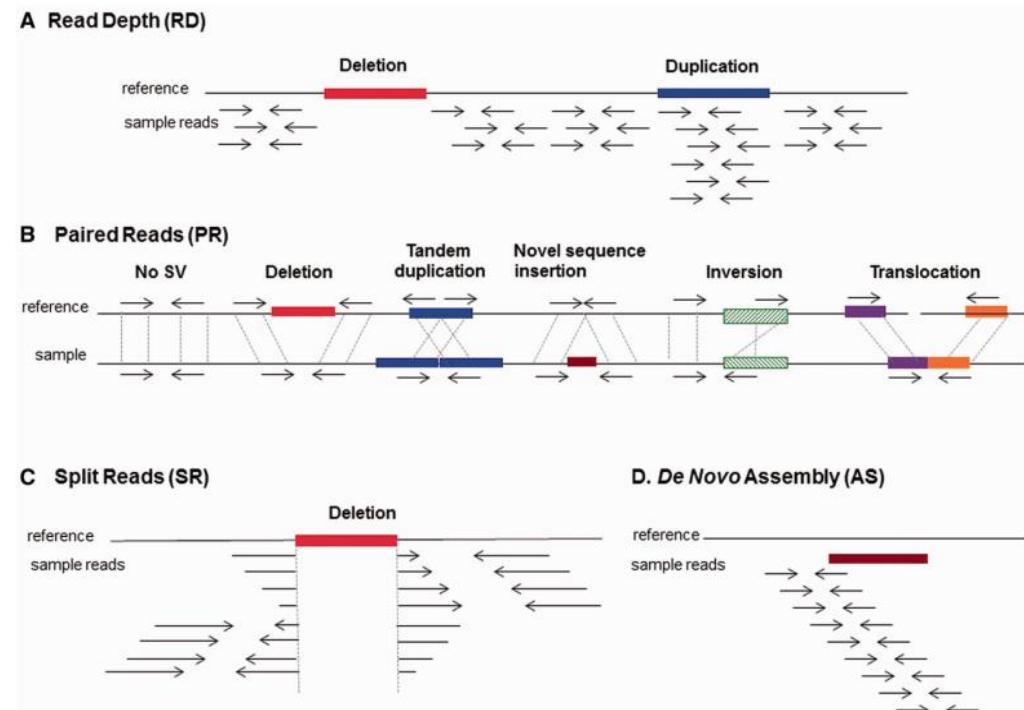
- FISH (Fluorescence In Situ Hybridization)
- aCGH (Comparative Genomic Hybridization)



NGS:

- Read Count
- Read pair
- Split reads
- Assembly

All of them rely in depth information, therefore, enough and uniform coverage is required for optimal results.



# Read Depth methodology

- Two main steps:
  - **Pre-processing:** reduce noise and sequencing bias
  - **Segmentation:** join by statistical methods regions with a similar number of overlapping reads
- Examples of segmentation methods:
  - CBS (Circular Binary Segmentation)
  - HMM (Hidden Markov Models)

# Algorithms for CNV in WGS

CNV detection methods on WGS data.

| Software       | Methods | Algorithm detail                         | Input data | Publish | Latest update | Accessibility | URL   | Programing Language | #Citations |
|----------------|---------|--|------------|---------|---------------|---------------|---|---------------------|------------|
| *Canvas        | RD      | Expectation-maximization (EM) clustering | BAM        | 2011    | 2018/3        | Y             | <a href="https://github.com/Illumina/canvas">https://github.com/Illumina/canvas</a>   | C#                  | 29         |
| *cn.MOPS       | RD      | Mixture Poisson model                    | BAM        | 2012    | 2018/10       | Y             | <a href="http://www.bioinf.jku.at/software/cnmops/cnmops.html">http://www.bioinf.jku.at/software/cnmops/cnmops.html</a>                       | R                   | 226        |
| CNVeM          | RD      | Expectation-maximization (EM) algorithm  | CSV        | 2013    | NA            | Y             | <a href="https://omictools.com/cnvem-tool">https://omictools.com/cnvem-tool</a>   | C                   | 14         |
| CNVer          | RP      | Maximum-likelihood, Graphic flow         | BAM        | 2010    | 2011/5        | N             | NA  | C                   | 158        |
| *CNVnator      | RD      | Mean shift algorithm                     | BAM        | 2011    | 2016/11       | Y             | <a href="https://github.com/abyzovlab/CNVnator">https://github.com/abyzovlab/CNVnator</a>   | C++                 | 640        |
| CNVRd2         | RD      | Expectation-maximization (EM) algorithm  | BAM/SAM    | 2014    | 2015/11       | Y             | <a href="https://bioconductor.org/packages/release/bioc/html/CNVRd2.html">https://bioconductor.org/packages/release/bioc/html/CNVRd2.html</a> | R                   | 13         |
| *Control-FREEC | RD      | LASSO regression                         | BAM/SAM    | 2011    | 2018/8        | Y             | <a href="http://boevalab.com/FREEC/">http://boevalab.com/FREEC/</a>   | C++                 | 190        |
| *GROM-RD       | RD      | Quantile normalization                   | BAM        | 2015    | 2017/5        | Y             | <a href="http://grigoriev.rutgers.edu/software/">http://grigoriev.rutgers.edu/software/</a>   | C                   | 7          |
| *iCopyDAV      | RD      | DoC approaches                           | BAM        | 2018    | 2018/3        | Y             | <a href="https://github.com/vogetihrsh/icopydav">https://github.com/vogetihrsh/icopydav</a>   | R,C++               | 1          |
| JointSLM       | RD      | Population-based approach                | SAM/BAM    | 2011    | NA            | N             | NA  | R                   | 49         |
| *LUMPY         | RD, PEM | A probabilistic framework                | BAM/CRAM   | 2014    | 2016/3        | Y             | <a href="https://github.com/ark5x/lumpy-sv">https://github.com/ark5x/lumpy-sv</a>   | C++                 | 157        |
| mrCaNaVAR      | RD      | mrFAST                                   | SAM        | 2009    | 2013/9        | Y             | <a href="http://mrcanavar.sourceforge.net/">http://mrcanavar.sourceforge.net/</a>   | C                   | 685        |
| *RDXplorer     | RD      | Event-wise testing algorithm             | BAM        | 2009    | 2013/4        | Y             | <a href="https://sourceforge.net/projects/rdxplorer/">https://sourceforge.net/projects/rdxplorer/</a>   | Python              | 496        |
| *ReadDepth     | RD      | Circular binary segmentation algorithm   | Bed Files  | 2011    | 2014/8        | Y             | <a href="https://github.com/chrisamiller/readDepth">https://github.com/chrisamiller/readDepth</a>   | R                   | 150        |
| *RSICNV        | RD      | Negative binomial transformations        | BAM        | 2017    | 2017/7        | Y             | <a href="https://github.com/yhhu/rsicnv">https://github.com/yhhu/rsicnv</a>   | C++                 | 2          |

# Algorithms for CNV in WES

**Table 1**

Selected tools for the performance analysis of CNV detection tools using WES data

| Tool name                  | ADTEx   | CONTRA  | cn.MOPS   | ExomeCNV  | VarScan 2   |
|----------------------------|---|---|---|---|---|
| <b>Characteristics</b>     |   |   |   |   |   |
| Control set required       | Yes   | Yes   | No  | Yes   | No  |
| Prog. Language             | Python, S/R   | Python, R   | R   | R   | Java  |
| Input format               | BAM, BED  | BAM, SAM, BED   | BAM, Read count matrices  | BAM, Pileup, GTF  | BAM, Pileup   |
| Segmentation Algorithm     | HMM   | CBS   | CBS   | CBS   | NA <sup>a</sup>   |
| OS                         | GNU, Linux  | Linux, Mac OS   | Linux, Mac OS, windows  | Linux, Mac OS, windows  | Linux, Mac OS, windows  |
| Methodology characteristic | DWT <sup>c</sup> for de-noising, use BAF <sup>d</sup>                   | Base-level log-ratio  | Bayesian approach for de-noising  | Statistical test for analyzing BAF data   | CMDS <sup>b</sup> for generating read counts                                |
| Year                       | 2014  | 2012  | 2012  | 2011  | 2012  |
| URL                        | <a href="http://adtex.sourceforge.net">http://adtex.sourceforge.net</a> | <a href="https://sourceforge.net/projects/contra-cnv/">https://sourceforge.net/projects/contra-cnv/</a> | <a href="http://www.bioinf.jku.at/software/cnmops/cnmops.html">http://www.bioinf.jku.at/software/cnmops/cnmops.html</a> | <a href="https://secure.genome.ucla.edu/index.php/ExomeCNV_User_Guide">https://secure.genome.ucla.edu/index.php/ExomeCNV_User_Guide</a> | <a href="http://varscan.sourceforge.net">http://varscan.sourceforge.net</a> |

<sup>a</sup>Segmentation is not imbedded in the tool. CBS is recommended for segmentation

<sup>b</sup>Correlation Matrix Diagonal Segmentation

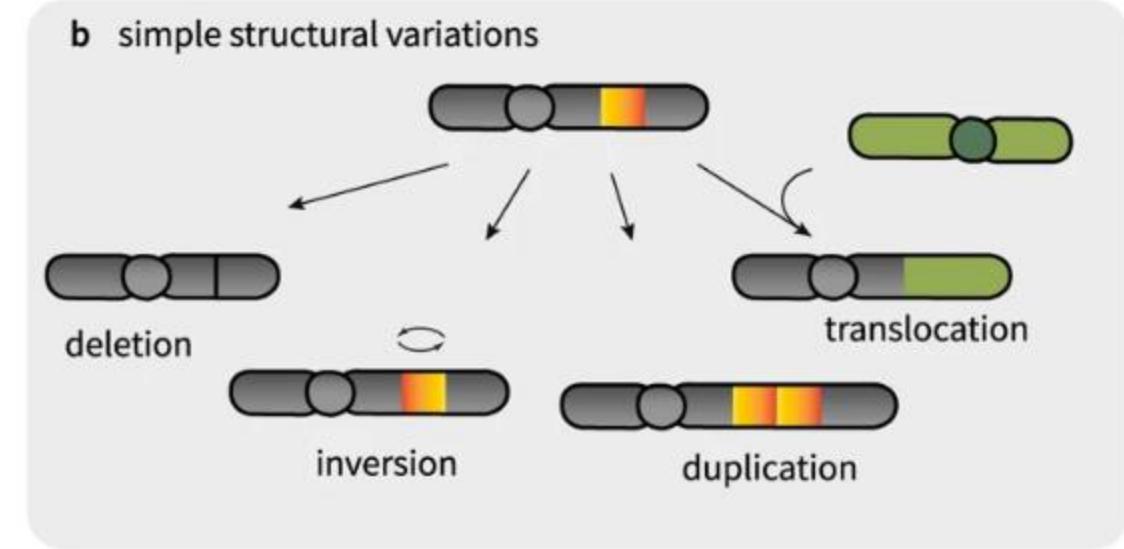
<sup>c</sup>Discrete wavelet transform

<sup>d</sup>B allele frequencies

# Structural Variants

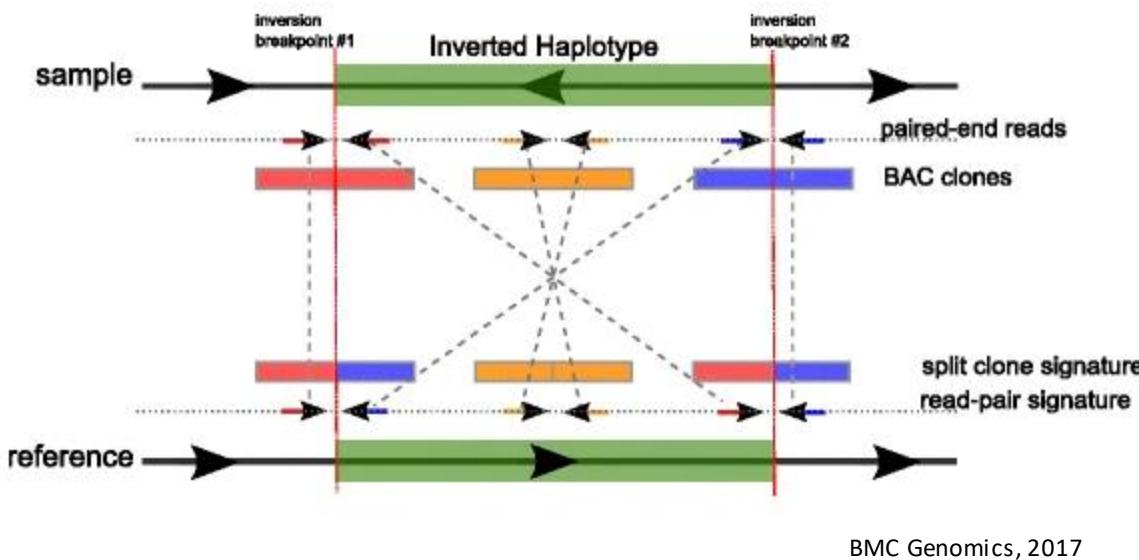
# Structural Variants

- Changes that affect to DNA structure
- Longer than 1 kb (> 50 bp in some definitions)
- CNV is a subclass of Structural Variants
- Another structural variants are inversions and translocations
- Associated to phenotypes and diseases

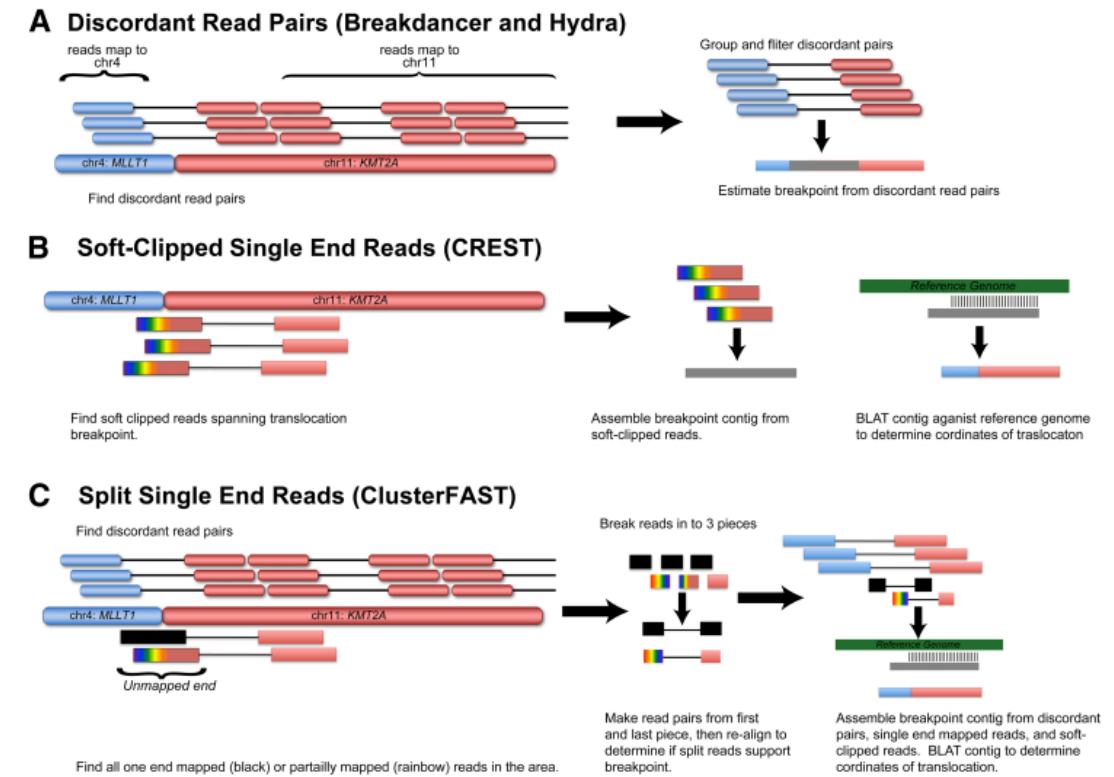


Experimental & Molecular Medicine. 2018

# Detection of inversions and translocations



Inversions and translocations are named balanced rearrangements



J Mol Diagn., 2014

# Algorithms for inversions and translocations

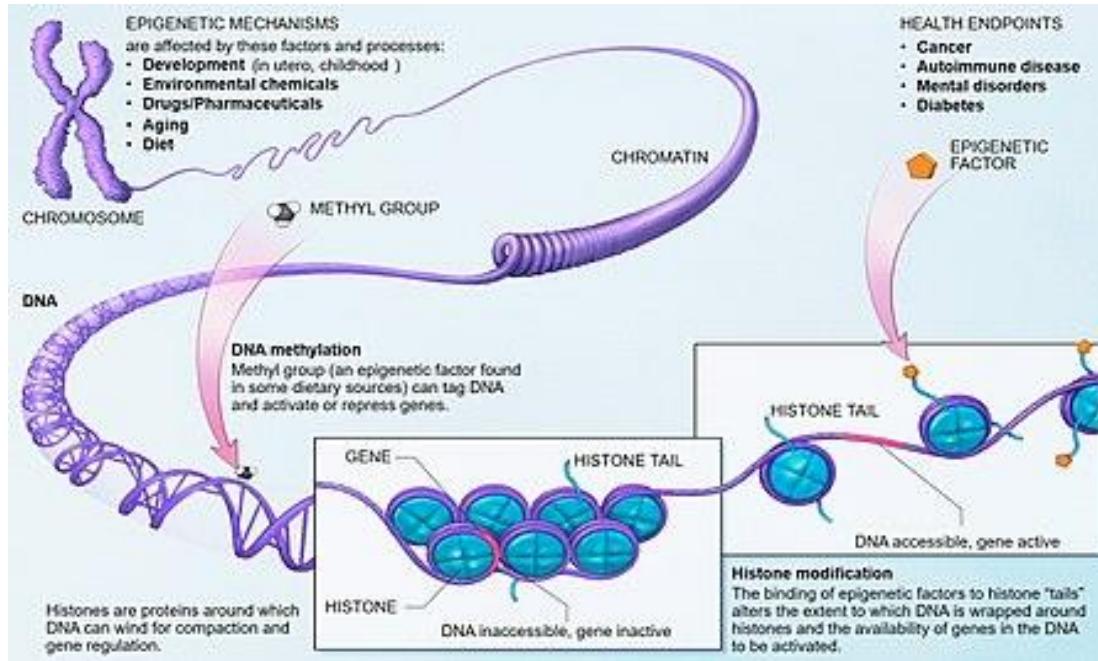
| Algorithm   | Inversions | Translocations | Discovery stage | Validation stage | Techniques | Algorithm   | Inversions | Translocations | Discovery stage | Validation stage | Techniques                  |
|-------------|------------|----------------|-----------------|------------------|------------|-------------|------------|----------------|-----------------|------------------|-----------------------------|
| LUMPY       | Yes        | Yes            | RD;SC;PR        | RD;SC;PR         | CL;SA      | PEMer       | Yes        | Yes            | PR              |                  | CL                          |
| MetaSV      | Yes        | Yes            | RD;SC;PR        | RD;SC;PR         | CL;SA;CA   | PeSV-Fisher | Yes        | Yes            | PR              | RD               | CL                          |
| Meerkat     | Yes        | Yes            | SC;PR;OEA       | SC;OEA           | CL;SA      | PRISM       | Yes        | Yes            | PR              | OEA              | CL;SA                       |
| SVMerge     | Yes        | Yes            | SC;PR;OEA       | SC;PR;OEA        | CL;SA;CA   | SVDetect    | Yes        | Yes            | PR              | PR               | CL                          |
| SoftSV      | Yes        | Yes            | SC;PR           | SC;PR            | CL;SA      | SVMiner     | Yes        |                | PR              | RD;PR            | CL;ST                       |
| BreakKmer   | Yes        | Yes            | SC;OEA          | PR               | CA         | Ulysses     | Yes        | Yes            | PR              | RD;PR            | CL;ST                       |
| ClipCrop    | Yes        | Yes            | SC              | SC               | CL;SA      | SLOPE       |            | Yes            | OEA             | OEA              | CL;SA                       |
| CREST       | Yes        | Yes            | SC              | SC               | CA         | SMUFIN      | Yes        | Yes            |                 |                  | CL;CA                       |
| Gustaf      | Yes        | Yes            | SC              | SC               | SA         |             |            |                |                 |                  | Adapted from Methods., 2016 |
| Socrates    | Yes        | Yes            | SC              | SC               | CL;SA      |             |            |                |                 |                  |                             |
| Bellerophon |            | Yes            | PR              | SC;PR            | CL;SA      |             |            |                |                 |                  |                             |
| BreakDancer | Yes        | Yes            | PR              |                  | CL;ST      |             |            |                |                 |                  |                             |
| DELLY       | Yes        | Yes            | PR              | SC;OEA           | CL;SA      |             |            |                |                 |                  |                             |
| FACTERA     |            | Yes            | PR              | SC               | CL;SA      |             |            |                |                 |                  |                             |
| GASV        | Yes        | Yes            | PR              |                  | CL         |             |            |                |                 |                  |                             |
| GASVPro     | Yes        | Yes            | PR              | RD;PR            | CL;ST      |             |            |                |                 |                  |                             |
| HYDRA       | Yes        | Yes            | PR              | PR               | CL         |             |            |                |                 |                  |                             |
| HYDRA-Multi | Yes        | Yes            | PR              | PR               | CL         |             |            |                |                 |                  |                             |
| inGAP-SV    | Yes        | Yes            | PR              | RD               |            |             |            |                |                 |                  |                             |

**RD:** Read Depth  
**SC:** Soft Clipped  
**PR:** Paired Reads  
**OEA:** One End Anchored

**CL:** Clustering  
**SA:** Split-reads Alignment  
**CA:** Contig Assembly  
**ST:** Statistical Testing

# Epigenetic variation

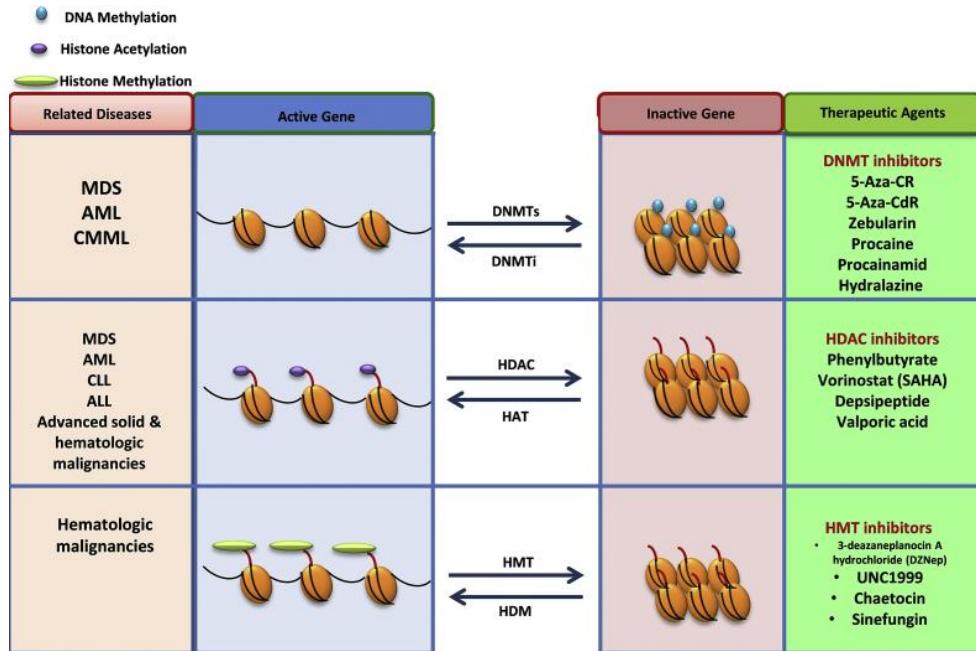
# Epigenome variants



- Genetic and epigenetic factors are interconnected.

- Changes not concerning nucleotide sequence. They can be reversible.
- Can be caused by external, environmental or biological factors.

# Epigenome variants



- Types

- DNA methylation
- Histone modification: covalent modification, methyl groups, acetyl groups, ...
- Binding of proteins to regulatory regions (silencers o enhancers)
- Non coding RNAs

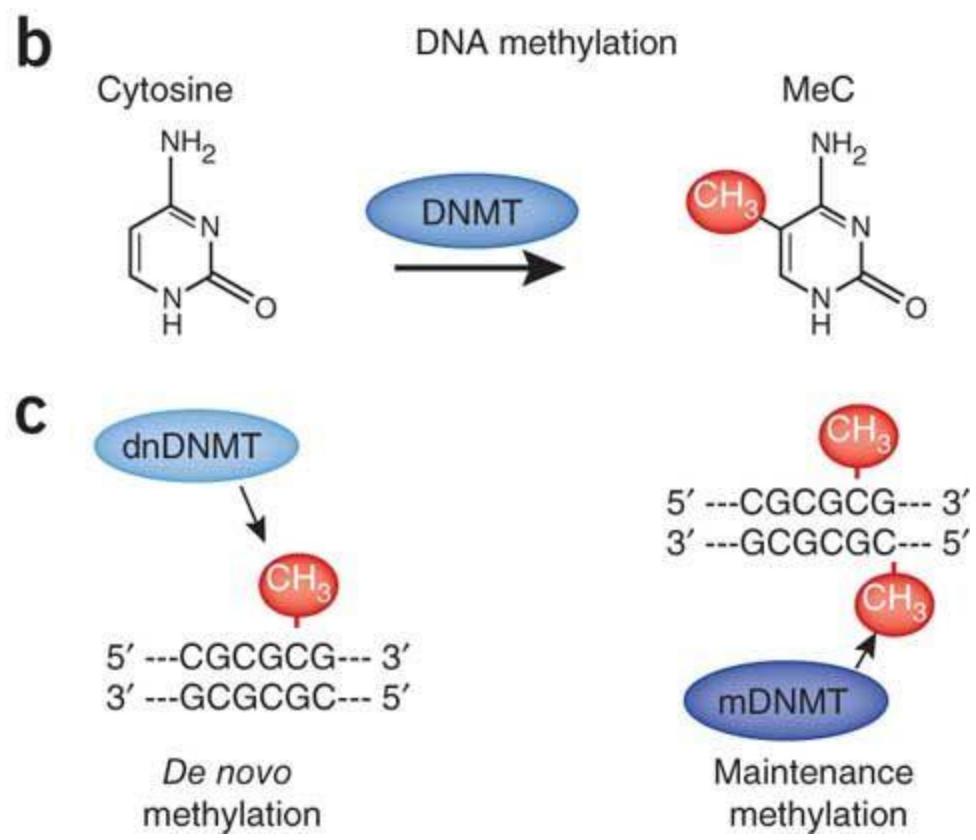
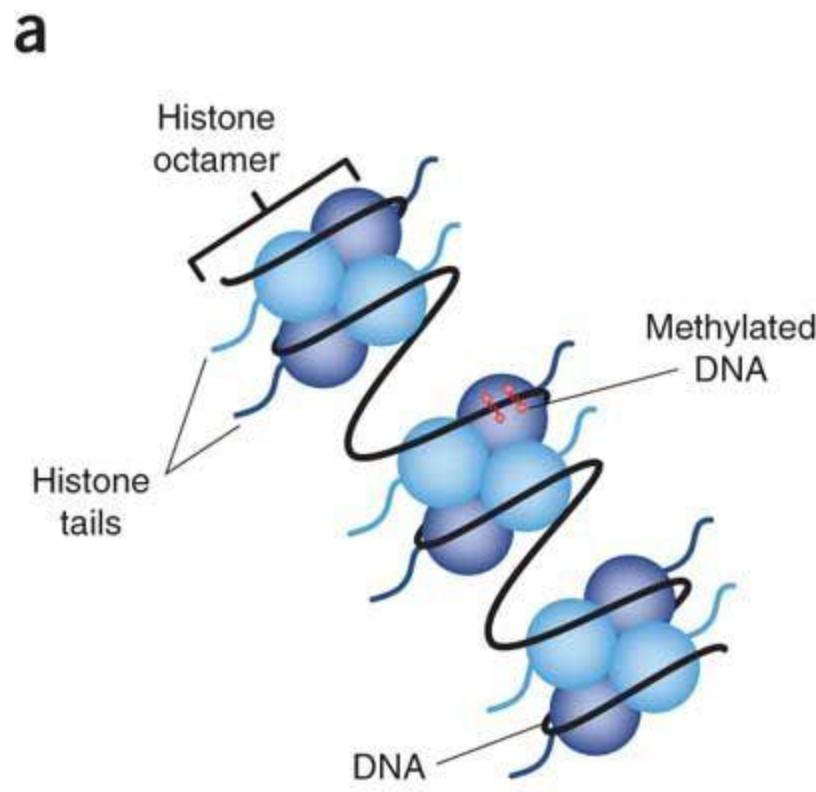
- Consequences

- Development: X chromosome inactivation, cellular lineage
- Phenotypic differences. I.e.: monozygotic twins
- Diseases: AS (UBE3A gene), cancer (IGF2)

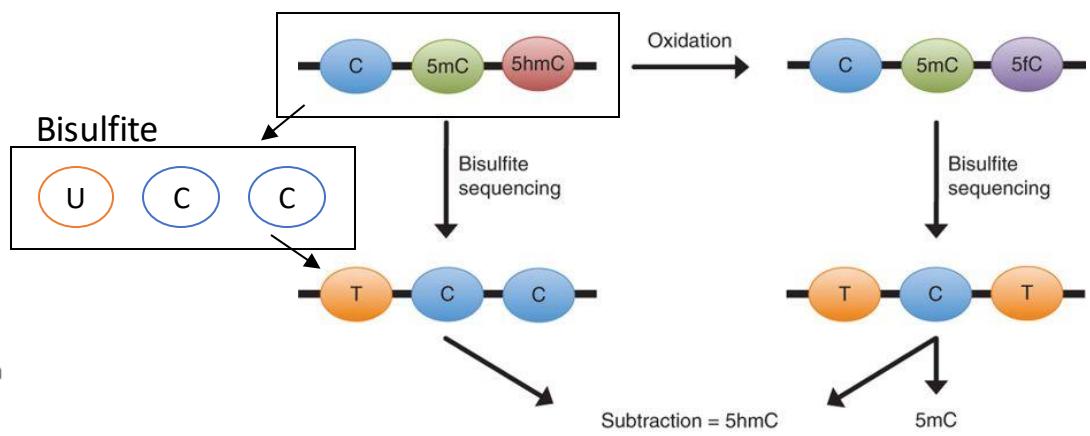
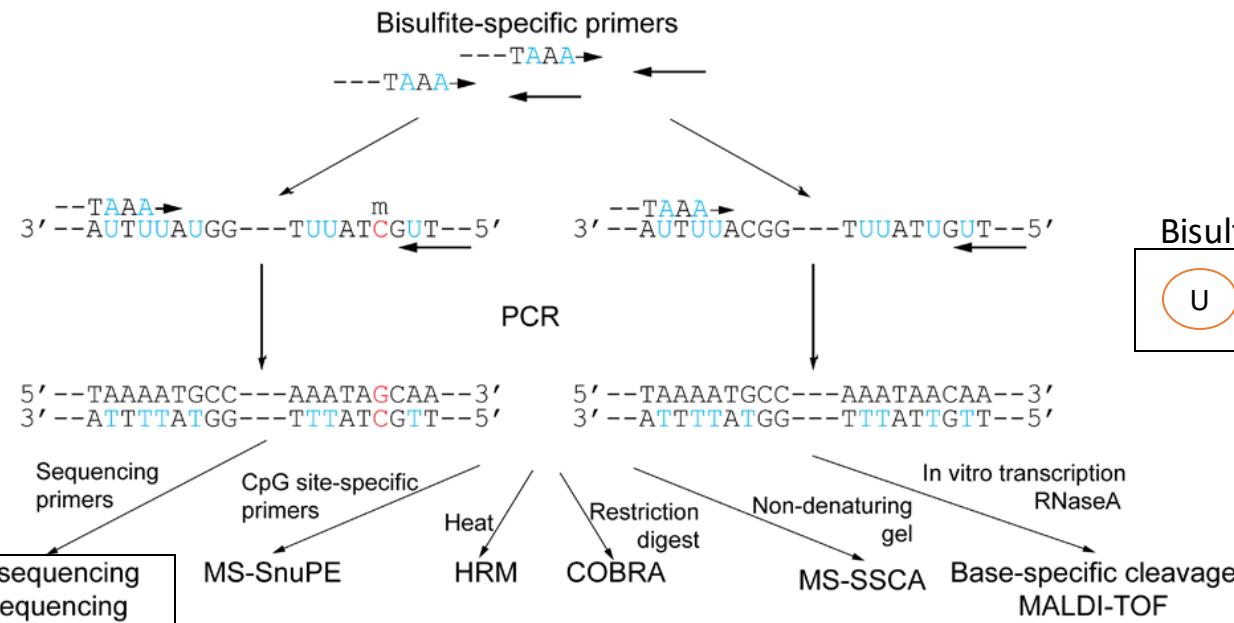
# NGS and epigenomics

- **Bisulfite sequencing:** DNA methylation
- **Chip-seq:** protein-DNA interactions, histone modifications

# DNA methylation



# Bisulfite sequencing



*Nature Protocols* volume8, pages1841–1851 (2013)

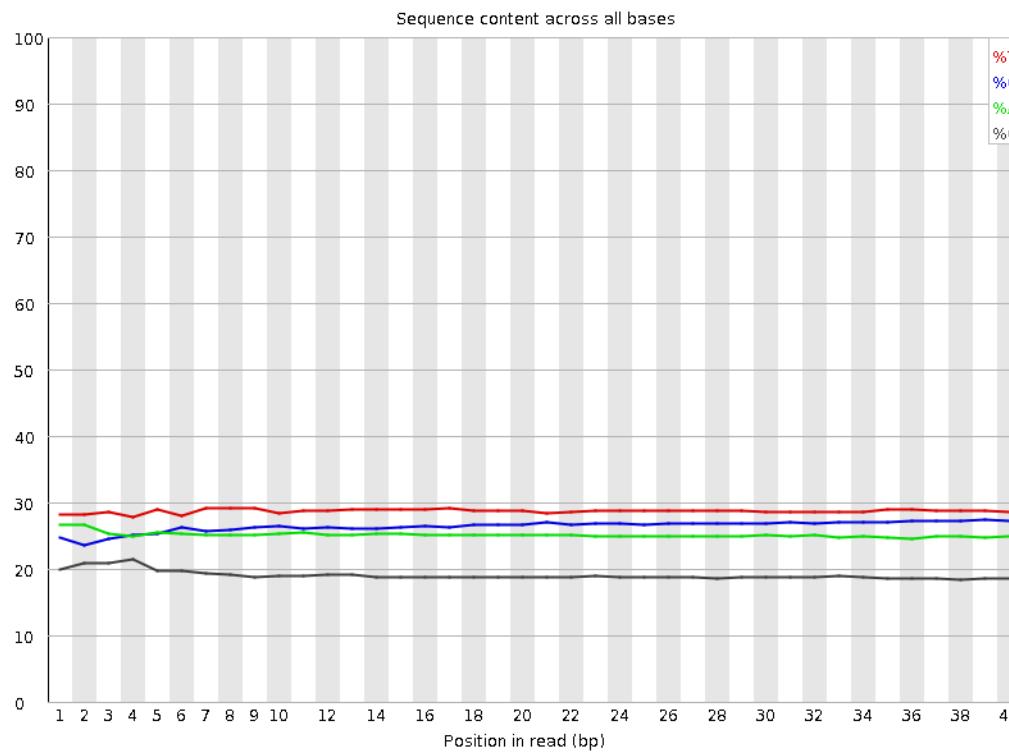
# Bisulfite sequencing - bioinformatic tools

| Step  | Bismark workflow | bwa-meth workflow     |
|---|------------------|-----------------------|
| Generate Reference Genome Index ( <i>optional</i> ) | Bismark          | bwa-meth              |
| Raw data QC   | FastQC           | FastQC                |
| Adapter sequence trimming                           | Trim Galore!     | Trim Galore!          |
| Align Reads   | Bismark          | bwa-meth              |
| Deduplicate Alignments                              | Bismark          | Picard MarkDuplicates |
| Extract methylation calls                           | Bismark          | MethylDackel          |
| Sample report                                       | Bismark          | -                     |
| Summary Report                                      | Bismark          | -                     |
| Alignment QC  | Qualimap         | Qualimap              |
| Sample complexity                                   | Preseq           | Preseq                |
| Project Report                                      | MultiQC          | MultiQC               |

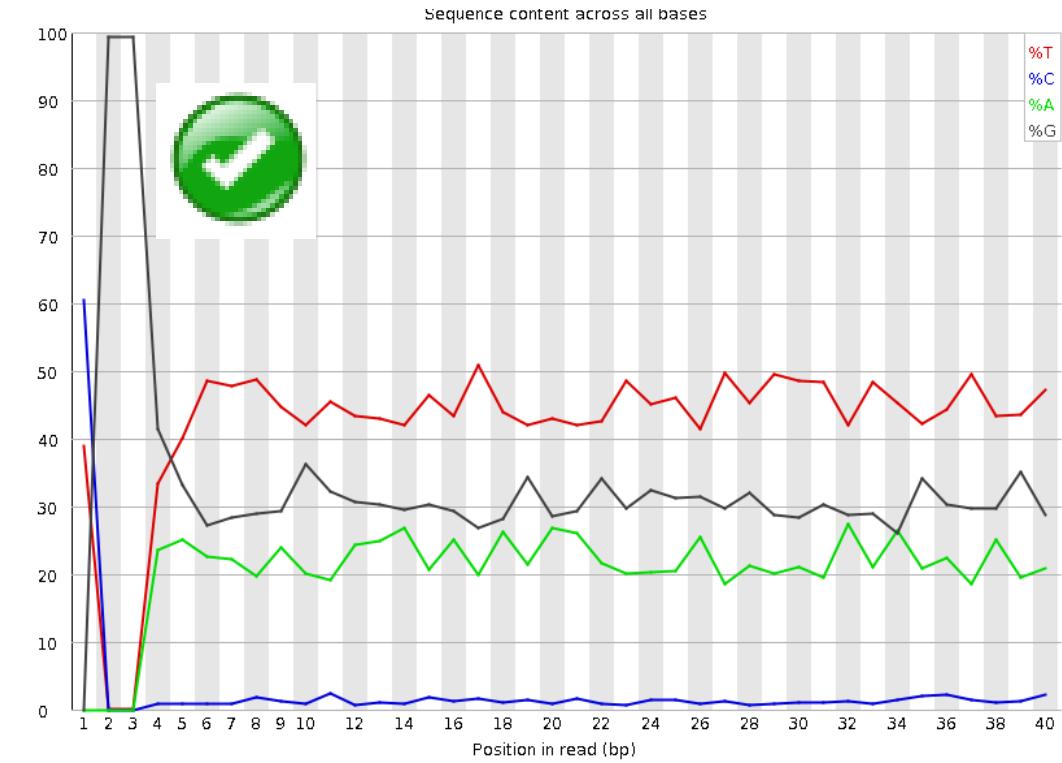
<https://github.com/nf-core/methylseq>

# Bisulfite sequencing - FastQC

## Per base sequence content



DNA sequencing

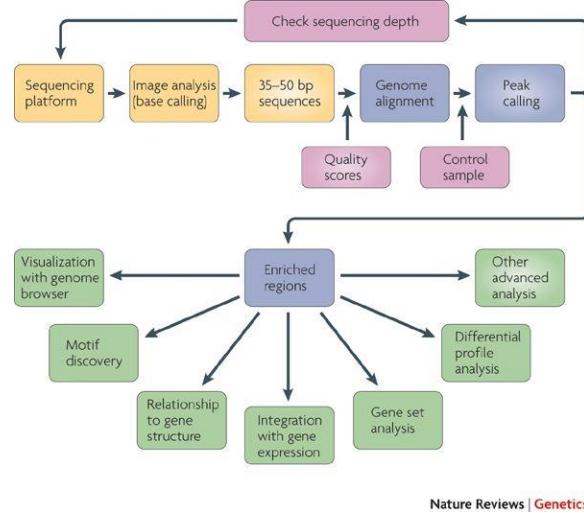


bisulfite sequencing

# Bisulfite sequencing - Alignment

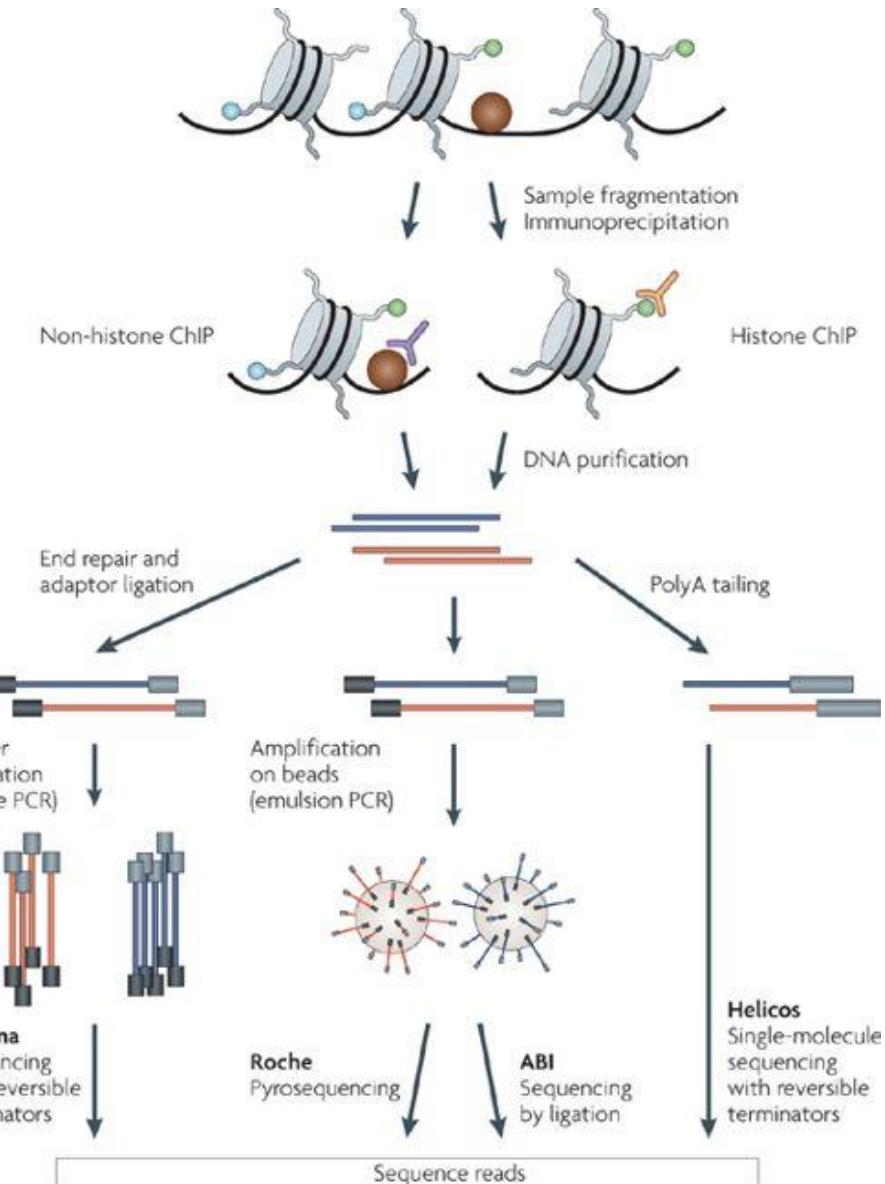
- More differences in bases, in relation to reference genome (C/T)
- Conventional aligners for DNA sequencing are not appropriate
- Examples of bisulfite sequencing specific aligners:
  - Bwa-meth
  - Bismark

# Chip-seq

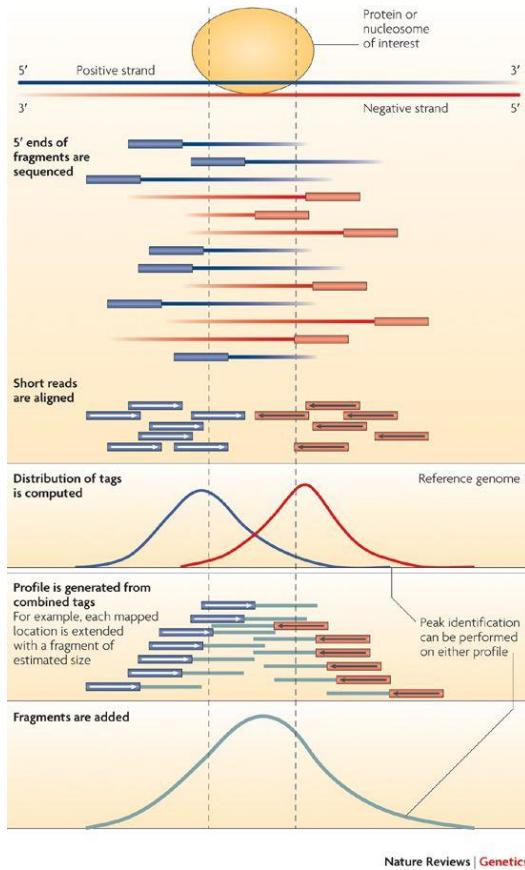


## Downstream analysis

- Protein binding sites motifs
- Annotate peak location



# Chip-seq



## Types of peaks

- Sharp: protein-DNA binding sites or histone modifications in regulatory elements
- Broad: histone modifications marking domains (transcribed or repressed regions)
- Mixed: RNA polymerase

