# Session 3.1 – Mapping against reference genome and Variant Calling

**BU-ISCIII**

**Unidades Comunes Científico Técnicas – SGSAFI-ISCIII**

28-31 Octubre 2024, 6ª Edición
Programa Formación Continua, ISCIII

# Index

**<u>Mapping against reference genome and Variant Calling :</u>**

- Mapping vs Alignment
- What is mapping?
- How to choose a NGS mapper.
- SAM/BAM format
- Duplicate filter
- Variant Calling
- Source of error and mitigation strategies
- VCF and bed format
- GATK vs VARSCAN2
- High quality SNP selection

# Alignment

| Definition: |
|---|
| Arrange two or more nucleotide or aminoacid sequences to identify regions of similarity that may be a consequence of functional, structural, or evolutionary relationships. |

```
AAB24882    TYHMCQFHCRYVNNHSGEKLYECNERSKAFSCPSHLQCHKRRQIGEKTHEHNQCGKAFPT
AAB24881    --------------------YECNQCGKAFAQHSSLKCHYRTHIGEKPYECNQCGKAFSK

AAB24882    PSHLQYHERTHTGEKPYECHQCGQAFKKCSLLQRHKRTHTGEKPYE-CNQCGKAFAQ-
AAB24881    HSHLQCHKRTHTGEKPYECNQCGKAFSQHGLLQRHKRTHTGEKPYMNVINMVKPLHNS
```

# Multiple alignment (MSA)

## Definition:

A multiple alignment is a colection of three or more sequences partial or completely aligned.

# Mapping definition

| Definición: |
|---|
| Place a sequence inside a larger sequence. For example, determine the position of a read inside a reference genome. |

```
        Referencia/ genoma

...GTGGGCCGGCAATTCGATATCGCGCATATATTTCGGCGCATGCTTAGC...

Lecturas:

GCAATTCGATAT
GCGCATATATTT
TGGGCCGGCAAT
CGCATGCTTAGC
ATTCGATATCGC
GCCGGCAATTCG


        Mapeo

...GTGGGCCGGCAATTCGATATCGCGCATATATTTCGGCGCATGCTTAGC...
        GCAATTCGATAT                CGCATGCTTAGC
    TGGGCCGGCAAT        GCGCATATATTT
            ATTCGATATCGC

  GCCGGCAATTCG
```

# Alignment vs mapping

## Mapping:

- A mapping is regarded to be correct if it overlaps the true region.
- Each read maps independently
- From thousand to millions of sequences.

## Multiple alignment:

- An alignment is regarded to be correct only if each base is placed correctly.
- Minimizes differences among sequences
- From tens to hundred of sequences.

## Consideratiosn:

- An algorithm can be good at mapping but may not be good aligning.
- This is because the true alignment minimizes differences between reads, but the read mapper only sees the reference.

Hen Li. Mapping, Alignment and SNP Calling. MPG Next Gen Workshop 2011

# So in summary...

CTGACCTCATG<span style="color:red">TGATCCAC</span>CCGCCTTGGCC

Find best match for the read
in a reference sequence
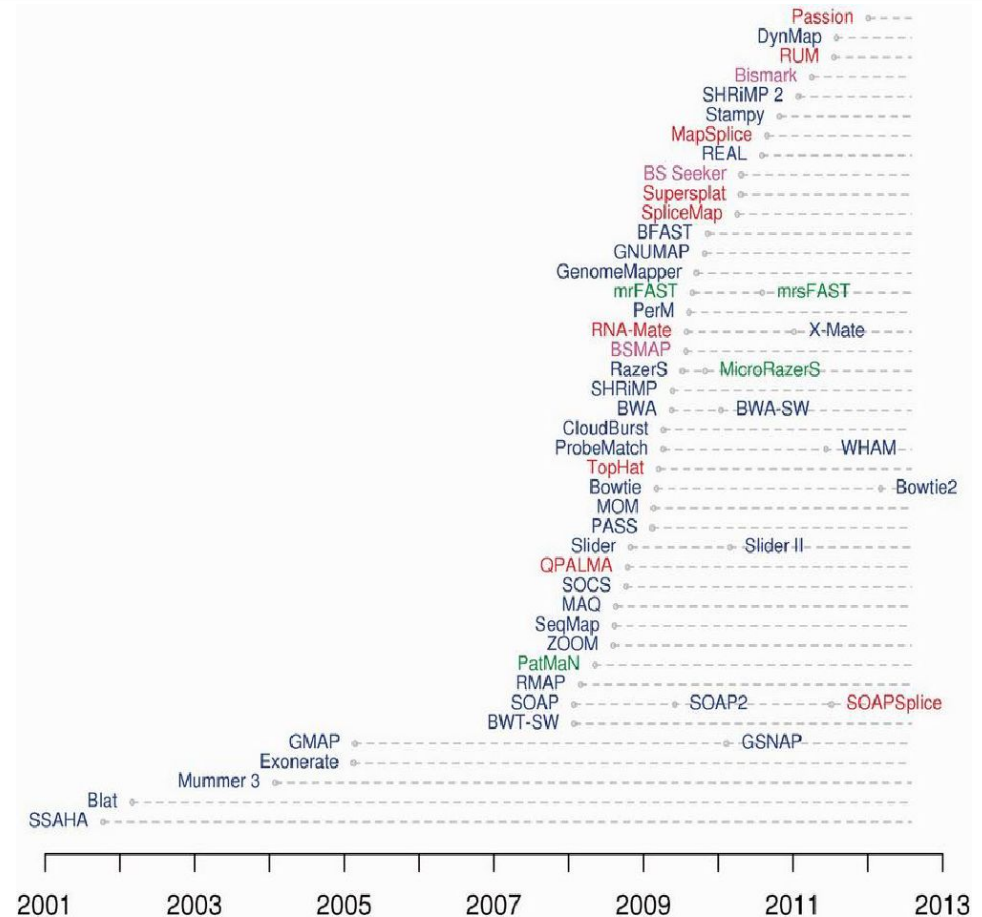
TGATCCAC

## Challenges

- Errors in reads
- Errors in libraries
- Repetitive regions (repeats, homologous regions)
- Homopolymers
- Individual polymorphisms

Pierre Lechat. Variants Calling lecture. Pasteur.fr

# What mapper should I use?

## Mappers:

- Más de 60 mappers available.
- Lots of papers reviewing its performamnce.

# What mapper should I use?

## Cosas a tener en cuenta:

- Computational resources vs sensibility
- Platform and type of experiment (Illumina/454/etc,paired-end,DNA/RNA/etc)
- Variation (indels allowance, mistmatch number,etc.)
- Repetitions (all regions, best match, random, user defined number...)

## Importante:

● Default options don't have to be the best:

"... there is no tool that outperforms all of the others in all the tests. Therefore, the end user should clearly specify his needs in order to choose the tool that provides the best results." - Hatem et al *BMC Bioinformatics* 2013, **14**:184

# End-to-end vs local alignment

End-to-end

Local

```
Read:          GACTGGGCGATCTCGACTTCG      Read:          ACGGTTGCGTTAATCCGCCACG
Reference: GACTGCGATCTCGACATCG         Reference: TAACTTGCGTTAAATCCGCCTGG
```

```
Alignment:                                             Alignment:
   Read:          GACTGGGCGATCTCGACTTCG         Read:          ACGGTTGCGTTAA-TCCGCCACG
                      | | | | |      | | | | | | | | | |   | | |                          | | | | | | | | |  | | | | | |
   Reference: GACTG--CGATCTCGACATCG         Reference: TAACTTGCGTTAAATCCGCCTGG
```

Bowtie2 manual.

# BWA MEM

Reference:
Read:

Seeding with SMEMs

'X' are mismtaches

Re-seeding with *k*-maximal inside long SMEMs

Chaining and chain filtering

Chains:

Filtered as too short

Seed extension

Ordered seed list:

(Successful)

(Skipped as contained in a found alignment)

(Unsuccessful extension)

## SMEM strategy

- Maximal exact match (MEM): an exact match that cannot be extended further in either direction
- Super-maximal exact match (SMEM): a MEM that is not contained in any other MEMs on the query coordinate (Li, 2012). At any query position, the longest exact match covering the position must be a SMEM.

## Seed-and-extend algorithm

## Local alignment

Hen LI. Aligning sequence reads, clone sequences and assembly con*gs with BWA-MEM. Poster. Broad Institute.

# BOWTIE2

**End-to-end alignment by default.**

**Three reporting modes:**
- Best alignment
- K alignments
- All alignments

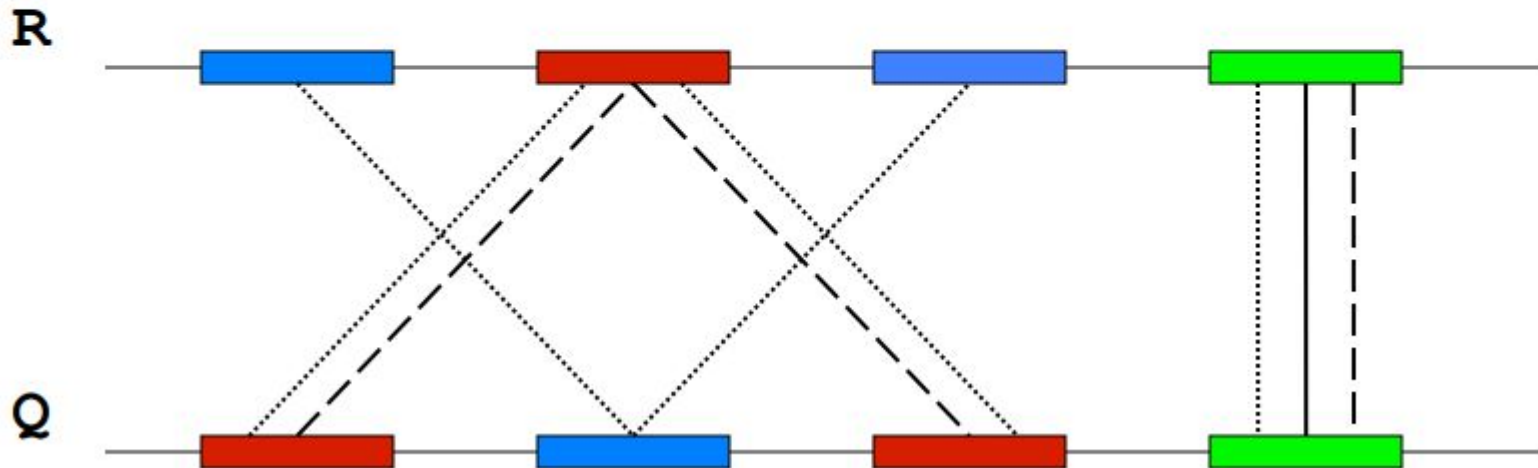**Lots of customizable parameters that change its performance.**

# Example whole genome aligner: MUMMER

- **Maximal Unique Matcher (MUM)**
  - match <- exact match of a minimum length
  - maximal <- cannot be extended in either direction without a mismatch
  - unique
    - occurs only once in both sequences (MUM)
    - occurs only once in a single sequence (MAM)
    - occurs one or more times in either sequence (MEM)

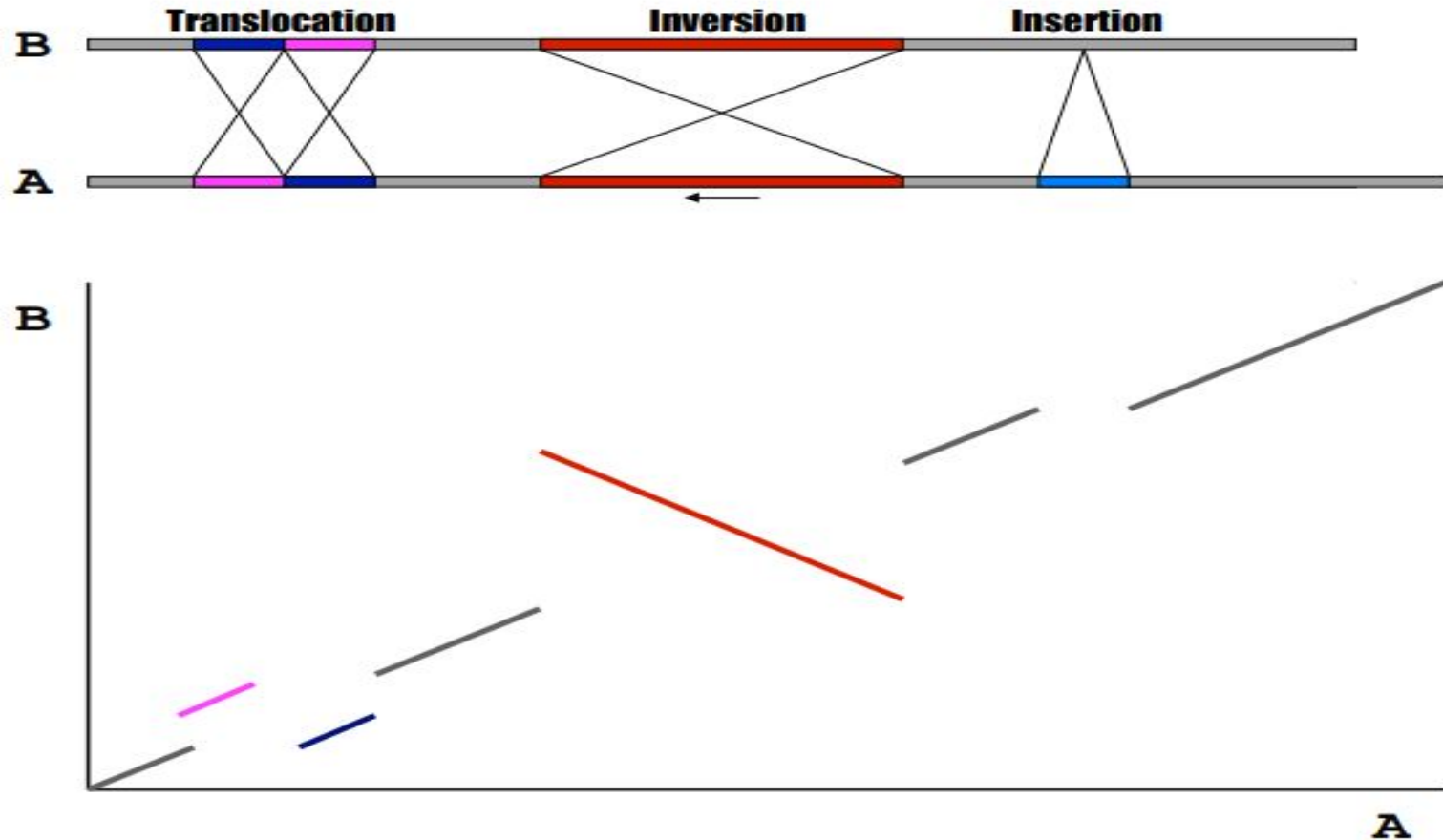Adam M. Phillippy. Whole Genome Alignment with MUMmer. Presentation.

# Example whole genome aligner: MUMMER

**MUM** : maximal unique match  ————————————

**MAM** : maximal almost-unique match  — — — — — —

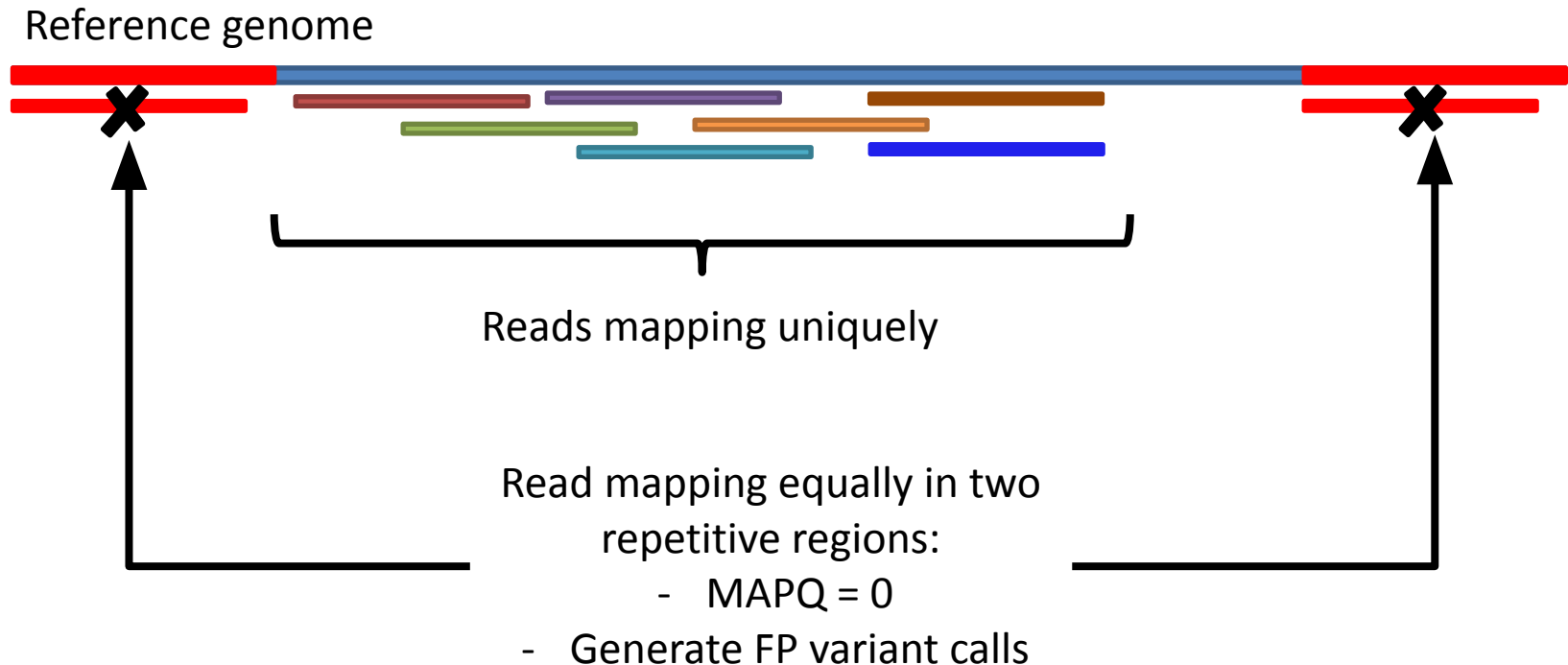**MEM** : maximal exact match  ············································

R

Q

Adam M. Phillippy. Whole Genome Alignment with MUMmer. Presentation.

# Example whole genome aligner: MUMMER



Adam M. Phillippy. Whole Genome Alignment with MUMmer. Lecture.

Reference genome

Reads mapping uniquely

Read mapping equally in two repetitive regions:
- MAPQ = 0
- Generate FP variant calls

# Which aligner should I use for aligning reads against a resistance gene database for determining with resistance genes I have in my sample?
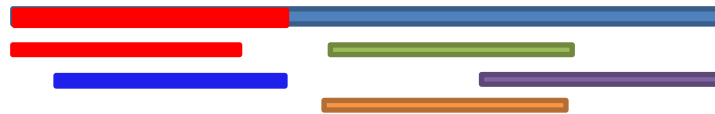
Homologus/repetitive region

Reads mapping to the repetitive/homologus region map against all alleles.
**We** allow one read to map to **several locations**.

Resistance gene - Allele 1

Resistance gene - Allele 2

Resistance gene - Allele 3

Reads mapping uniquely only map in Allele 1. Which is the one more **covered**
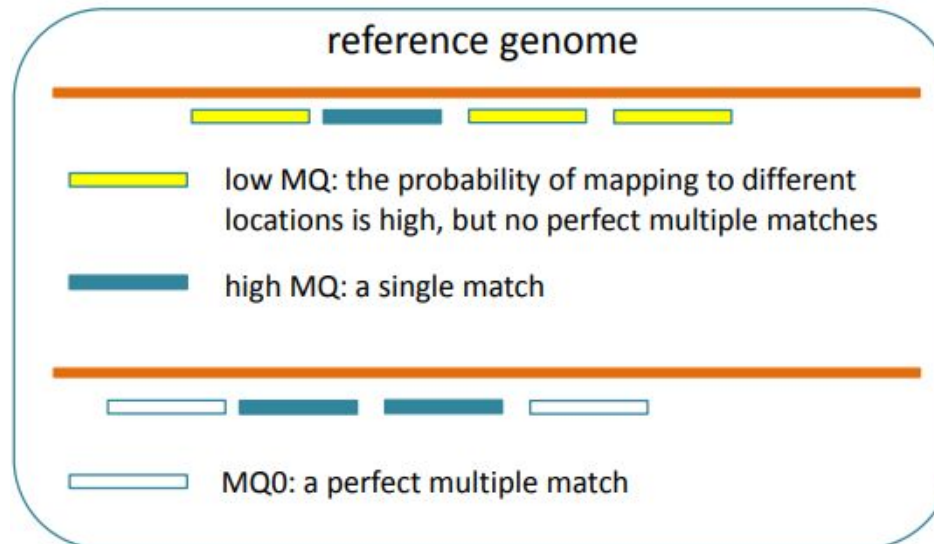
# MAPQ

- What if there are several possible places to align your sequencing read? This may be due to:
  - Repeated elements in the genome
  - Low complexity sequences
  - Reference errors and gaps

  **MQ is a phredScore of the quality of the alignment**

# MAPQ

**MAPQ is NOT comparable among mappers.**

**BWA:**

- MAPQ represents the probability of the read to be mapped correctly.
- MAPQ = 0 identifies unmapped reads and...

**Reads mapping to different locations!**

**BOWTIE2:**

- MAPQ represents the "uniqueness" of the read. A MAPQ < 10 indicates that there is at least a 1 in 10 chance that the read truly originated elsewhere
- MAPQ = 0 identifies unmapped reads

# SAM format

| Definición: |
|---|
| It's a specification that defines a generic format for storing nucleotide alignments. It describes a query alignment against a reference genome. |

```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
r001   99 ref  7 30 8M2I4M1D3M = 37  39 TTAGATAAAGGATACTG *
r002    0 ref  9 30 3S6M1P1I4M * 0   0 AAAAGATAAGGATA    *
r003    0 ref  9 30 5S6M       * 0   0 GCCTAAGCTAA       * SA:Z:ref,29,-,6H5M,17,0;
r004    0 ref 16 30 6M14N5M    * 0   0 ATAGCTTCAGC       *
r003 2064 ref 29 17 6H5M       * 0   0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001  147 ref 37 30 9M         = 7 -39 CAGCGGCAT         * NM:i:1
```

# SAM format

| Col | Field | Type | Regexp/Range | Brief description |
|-----|-------|------|--------------|-------------------|
| 1 | QNAME | String | [!-?A-~]{1,255} | Query template NAME |
| 2 | FLAG | Int | [0,$2^{16}-1$] | bitwise FLAG |
| 3 | RNAME | String | \*|[!-()+-<>-~][!-~]* | Reference sequence NAME |
| 4 | POS | Int | [0,$2^{31}-1$] | 1-based leftmost mapping POSition |
| 5 | MAPQ | Int | [0,$2^{8}-1$] | MAPping Quality |
| 6 | CIGAR | String | \*|([0-9]+[MIDNSHPX=])+ | CIGAR string |
| 7 | RNEXT | String | \*|=|[!-()+-<>-~][!-~]* | Ref. name of the mate/next read |
| 8 | PNEXT | Int | [0,$2^{31}-1$] | Position of the mate/next read |
| 9 | TLEN | Int | [$-2^{31}+1$,$2^{31}-1$] | observed Template LENgth |
| 10 | SEQ | String | \*|[A-Za-z=.]+ | segment SEQuence |
| 11 | QUAL | String | [!-~]+ | ASCII of Phred-scaled base QUALity+33 |

```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
r001   99 ref  7 30 8M2I4M1D3M = 37  39 TTAGATAAAGGATACTG *
r002    0 ref  9 30 3S6M1P1I4M *  0   0 AAAAGATAAGGATA    *
r003    0 ref  9 30 5S6M        *  0   0 GCCTAAGCTAA       * SA:Z:ref,29,-,6H5M,17,0;
r004    0 ref 16 30 6M14N5M     *  0   0 ATAGCTTCAGC       *
r003 2064 ref 29 17 6H5M        *  0   0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001  147 ref 37 30 9M          =  7 -39 CAGCGGCAT         * NM:i:1
```

# SAM format: flags

| Bit | Description |
|---|---|
| 0x1 | template having multiple segments in sequencing |
| 0x2 | each segment properly aligned according to the aligner |
| 0x4 | segment unmapped |
| 0x8 | next segment in the template unmapped |
| 0x10 | SEQ being reverse complemented |
| 0x20 | SEQ of the next segment in the template being reversed |
| 0x40 | the first segment in the template |
| 0x80 | the last segment in the template |
| 0x100 | secondary alignment |
| 0x200 | not passing quality controls |
| 0x400 | PCR or optical duplicate |
| 0x800 | supplementary alignment |

https://broadinstitute.github.io/picard/explain-flags.html

# Flag explanation example 1

# Flag explanation example 2

# SAM format: CIGAR string

| Op | BAM | Description |
| --- | --- | --- |
| M | 0 | alignment match (can be a sequence match or mismatch) |
| I | 1 | insertion to the reference |
| D | 2 | deletion from the reference |
| N | 3 | skipped region from the reference |
| S | 4 | soft clipping (clipped sequences present in SEQ) |
| H | 5 | hard clipping (clipped sequences NOT present in SEQ) |
| P | 6 | padding (silent deletion from padded reference) |
| = | 7 | sequence match |
| X | 8 | sequence mismatch |

# SAM vs BAM format

- SAM and BAM format are exactly the same.
  - SAM is a tabular plain text file.
  - BAM is its binary format. Binary meaning is in a compress format not human readable.
  - We **MUST** always use BAM format because it is optimized for computer-reading

**AND**

**BECAUSE IT SAVES A LOT OF DISK SPACE!!**

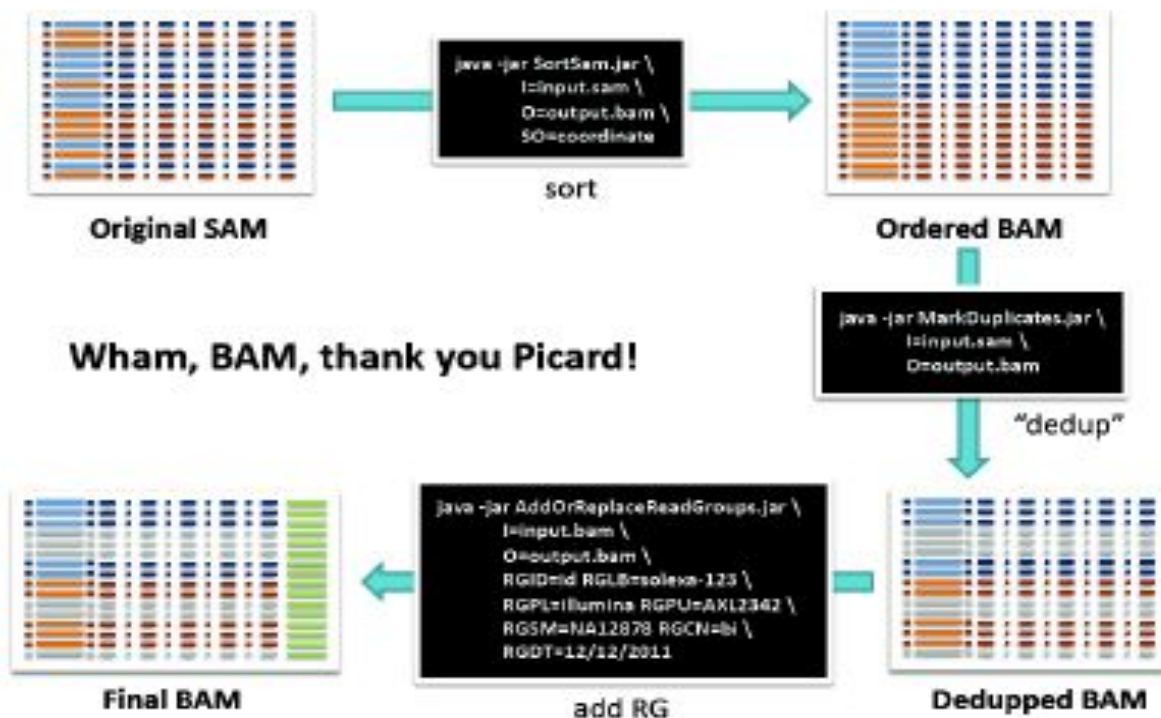Typical bam and sam format files weights from a S. grumpensis
SAM format file: 3.6 GB
BAM format file: 689 M

# Duplicate filter

- Duplicates are non-independent measurements of a sequence
  - Sampled from the exact same template of DNA
  - Violates assumptions of variant calling
- Errors in sample/library prep will get propagated to all the duplicates
- Just pick the "best" copy – mitigates the effects of errors
- **Definition**: sequences starting and finishing in the exact same coordinates. Both pairs if paired-end.
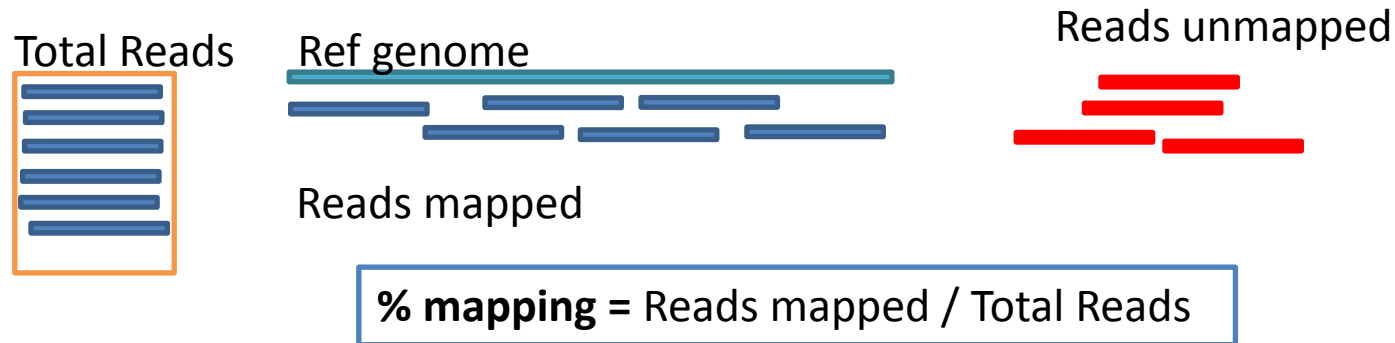
# Duplicate filter

# Mapping statistics

- % mapped: reads mapped/total reads

- % unmapped: reads unmapped/total reads

- % duplicates: reads belonging to same template/total

  reads

- Mean depth of coverage

- Coverage: % genome with at least one read mapped.

# Mapping quality control

Picard
Samtools

- **% mapping:** number of reads mapping againts reference genome.

Total Reads    Ref genome                                    Reads unmapped

Reads mapped
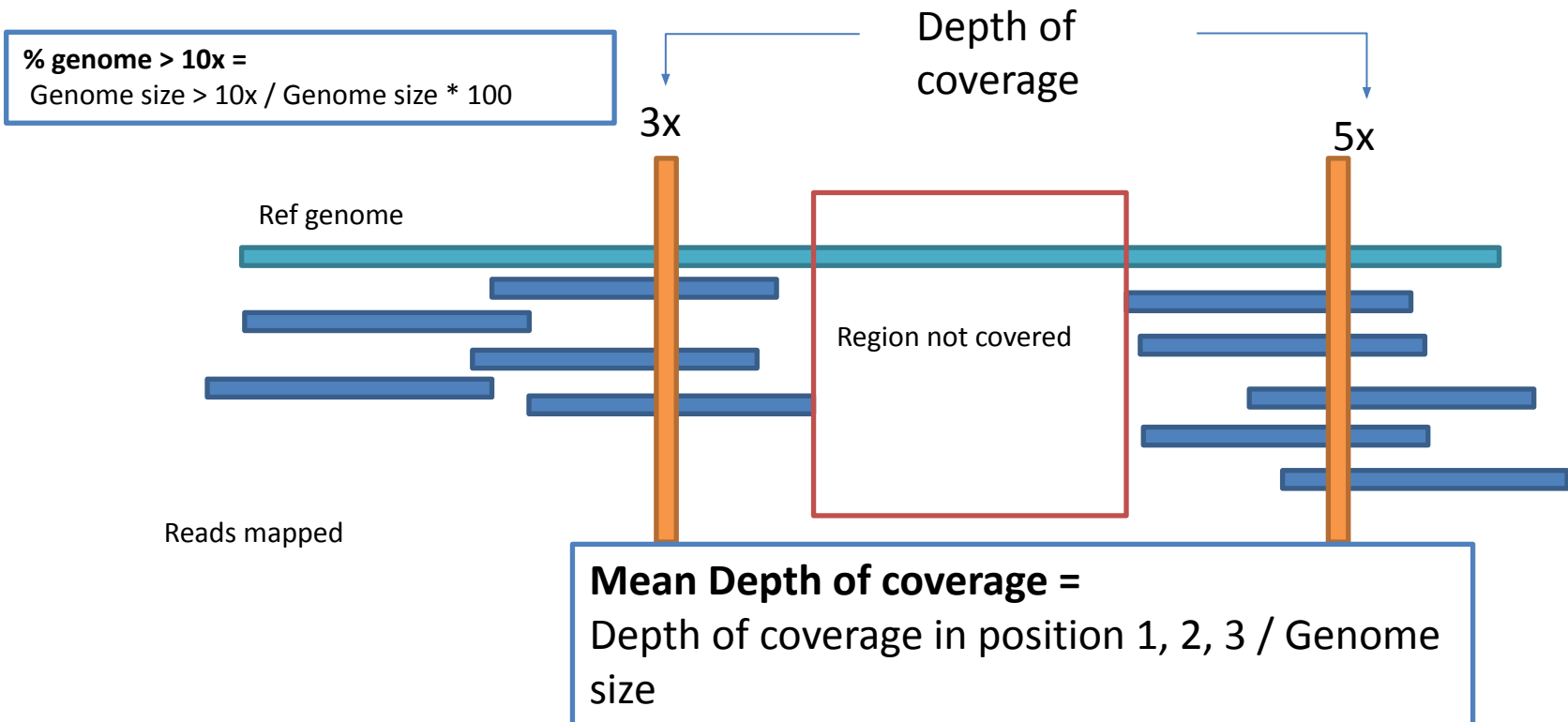
**% mapping =** Reads mapped / Total Reads

Mandatory parameter for microbial genomics!! It indicates us how many reads we have from our organism of interest. In human genomics this is almost always 99.99% unless something terrible happens. Not here!!!

# Mapping quality control

- **% genome > 10x:** percentage of genome covered with more than 10 reads.
- **Mean Depth of coverage:** mean of reads covering a genome position.

Picard
Samtools



**% genome > 10x =**
Genome size > 10x / Genome size * 100

Depth of coverage

3x

5x

Ref genome

Region not covered

Reads mapped

**Mean Depth of coverage =**
Depth of coverage in position 1, 2, 3 / Genome size

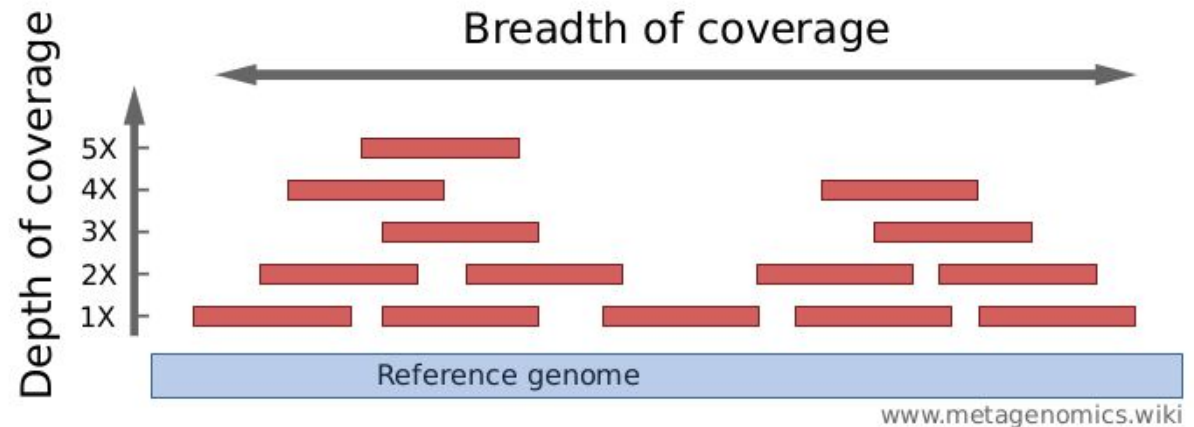# Depth of coverage vs coverage

## Breadth of coverage

How much of a genome is "covered" by short reads? Are there regions that are not covered, even not by a single read?

Breadth of coverage is the percentage of bases of a reference genome that are covered with a certain depth. For example: 90% of a genome is covered at 1X depth; and still 70% is covered at 5X depth.

## Depth of coverage

How strong is a genome "covered" by sequenced fragments (short reads)?

Per-base coverage is the average number of times a base of a genome is sequenced. The coverage depth of a genome is calculated as the number of bases of all short reads that match a genome divided by the length of this genome. It is often expressed as 1X, 2X, 3X,... (1, 2, or, 3 times coverage).
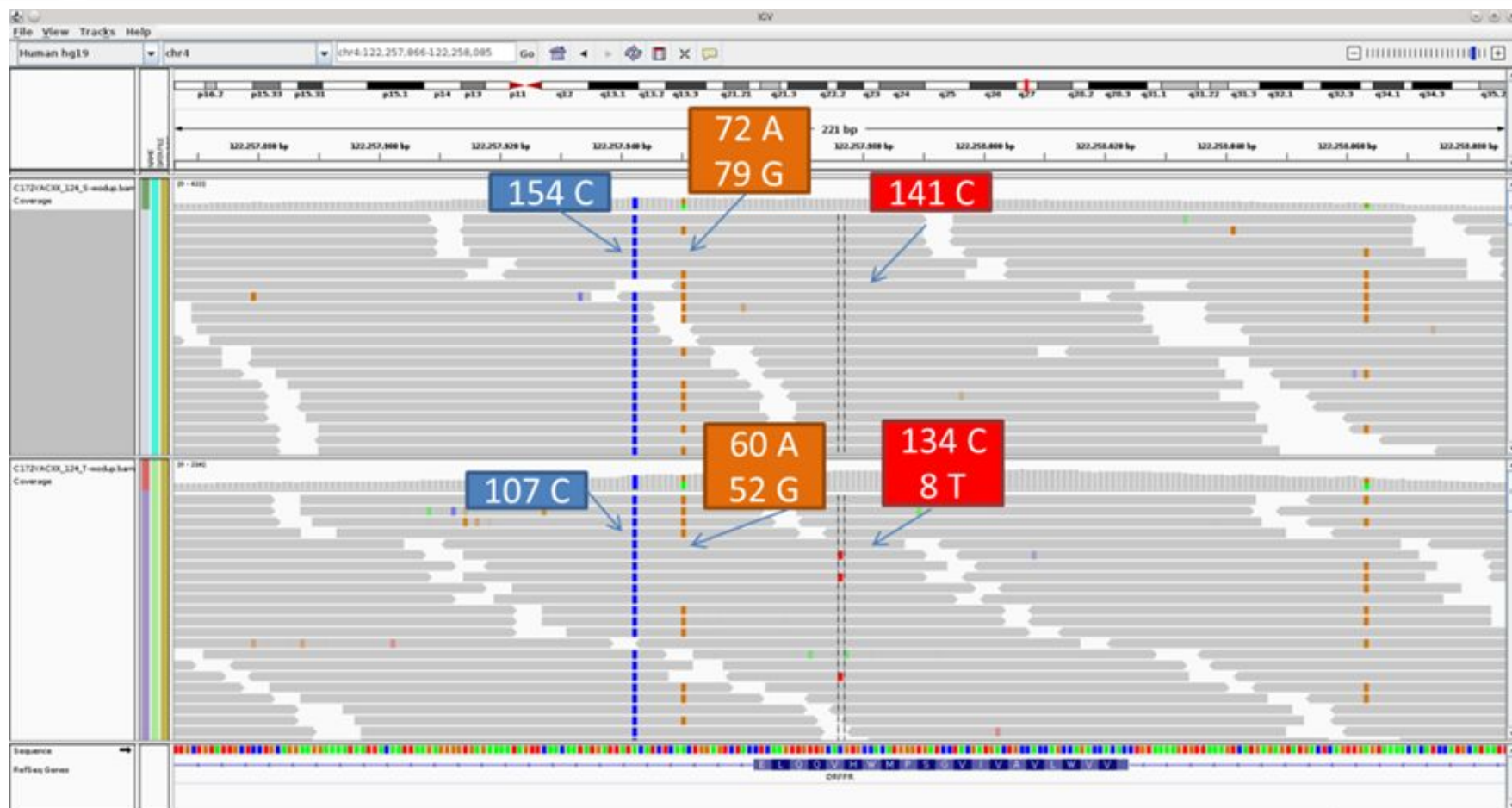


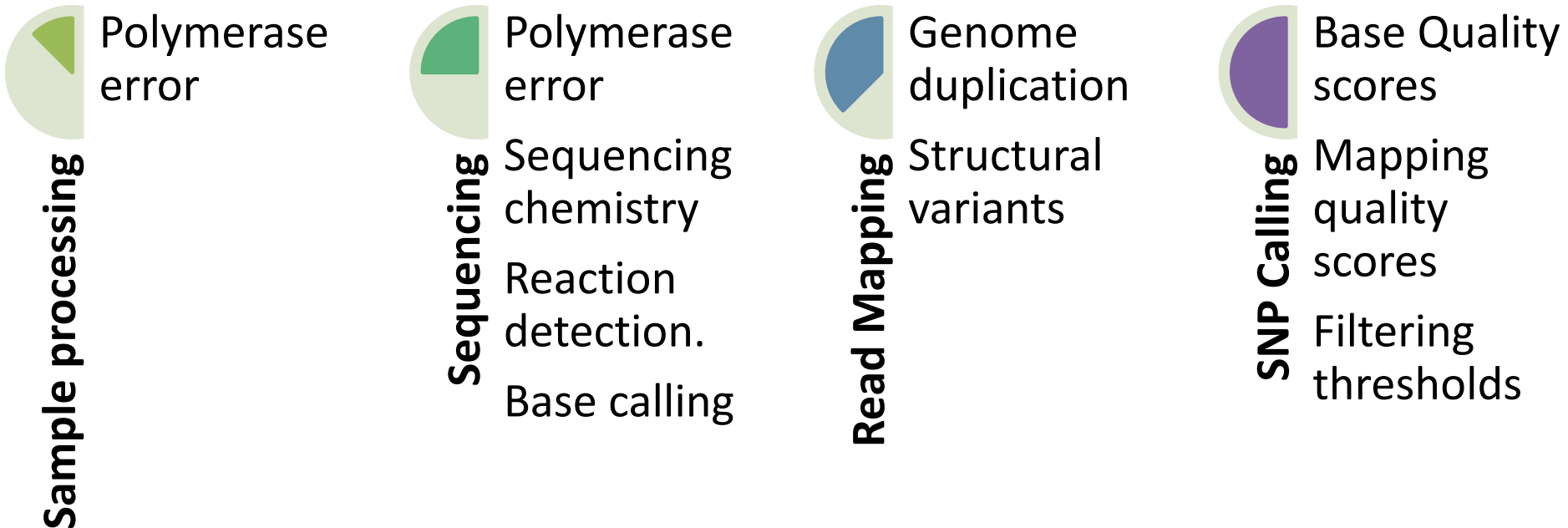www.metagenomics.wiki

# Variant Calling

- **<u>Variant calling concept is simple:</u>**

  **Find positions in our reads different from the reference.**

- We start with our secuences mapped against our reference genome, and we walk trough every column of the alignment counting the number of alleles found and comparing them against the reference.

# Sources of error and mitigation strategies

**Sample processing**
- Polymerase error

**Sequencing**
- Polymerase error
- Sequencing chemistry
- Reaction detection.
- Base calling

**Read Mapping**
- Genome duplication
- Structural variants

**SNP Calling**
- Base Quality scores
- Mapping quality scores
- Filtering thresholds

Adapted from Olson et al. Frontiers in Genetics. 2015

# Sources of error and mitigation strategies

- **<u>Sample processing errors.</u>**
  - Random errors.
  - Associated with polymerase errors . (1 in $10^{2-3}$ bases)
  - Homopolymers and tandem repeats experience higher indel error rates.
- **<u>Solutions:</u>**
  - Paired-end libraries.
  - Minimization of PCR cycles.

Adapted from Olson et al. Frontiers in Genetics. 2015

# Sources of error and mitigation strategies

- **Sequencing:**
  - Dependent on the platform.
  - Can be random and systematic.
  - 6% Illumina, 50% Roche (Ross et al.2013)
  - P.e Illumina commits error in the G/T channels.
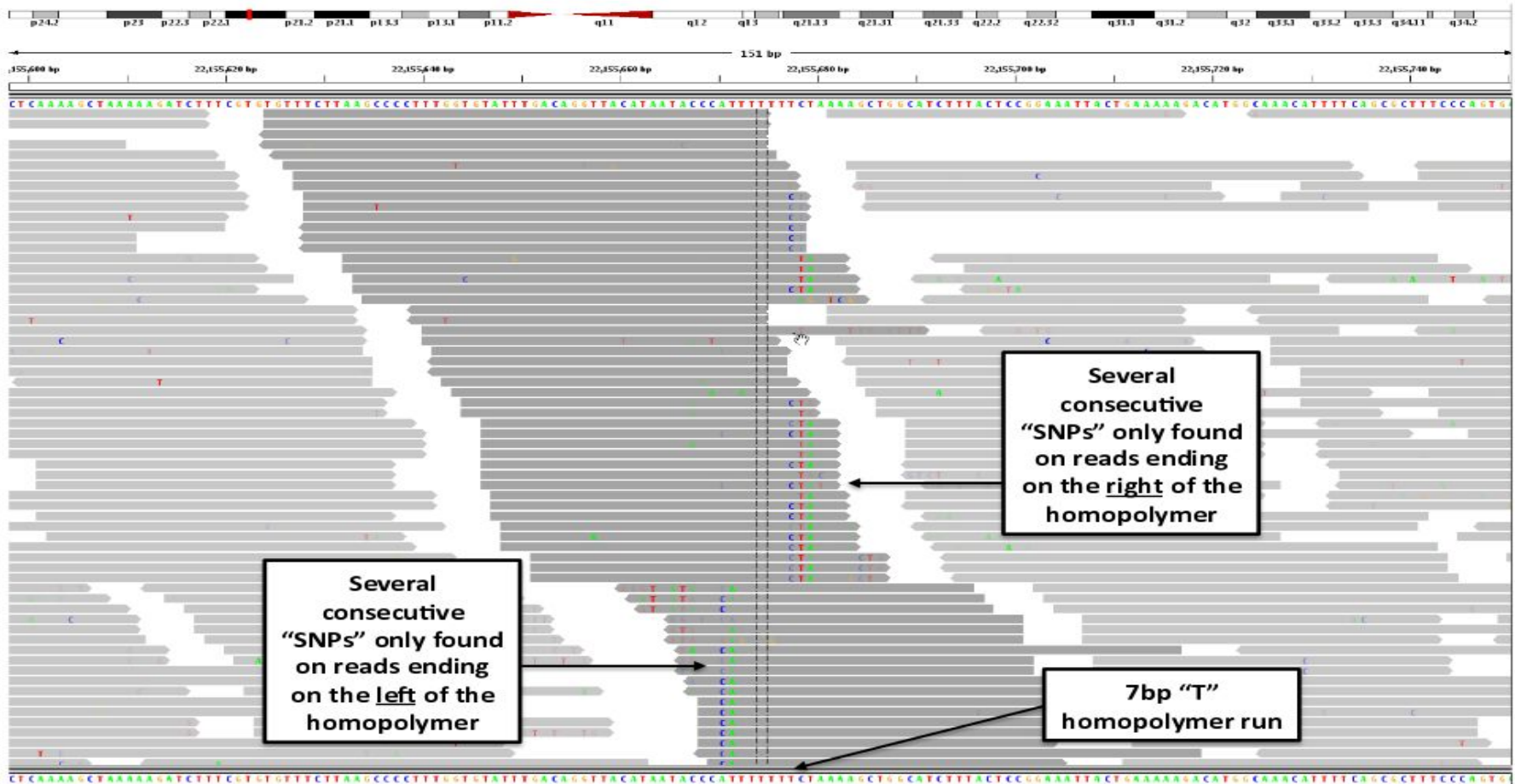
- **Solutions:**
  - Strand bias.

Adapted from Olson et al. Frontiers in Genetics. 2015

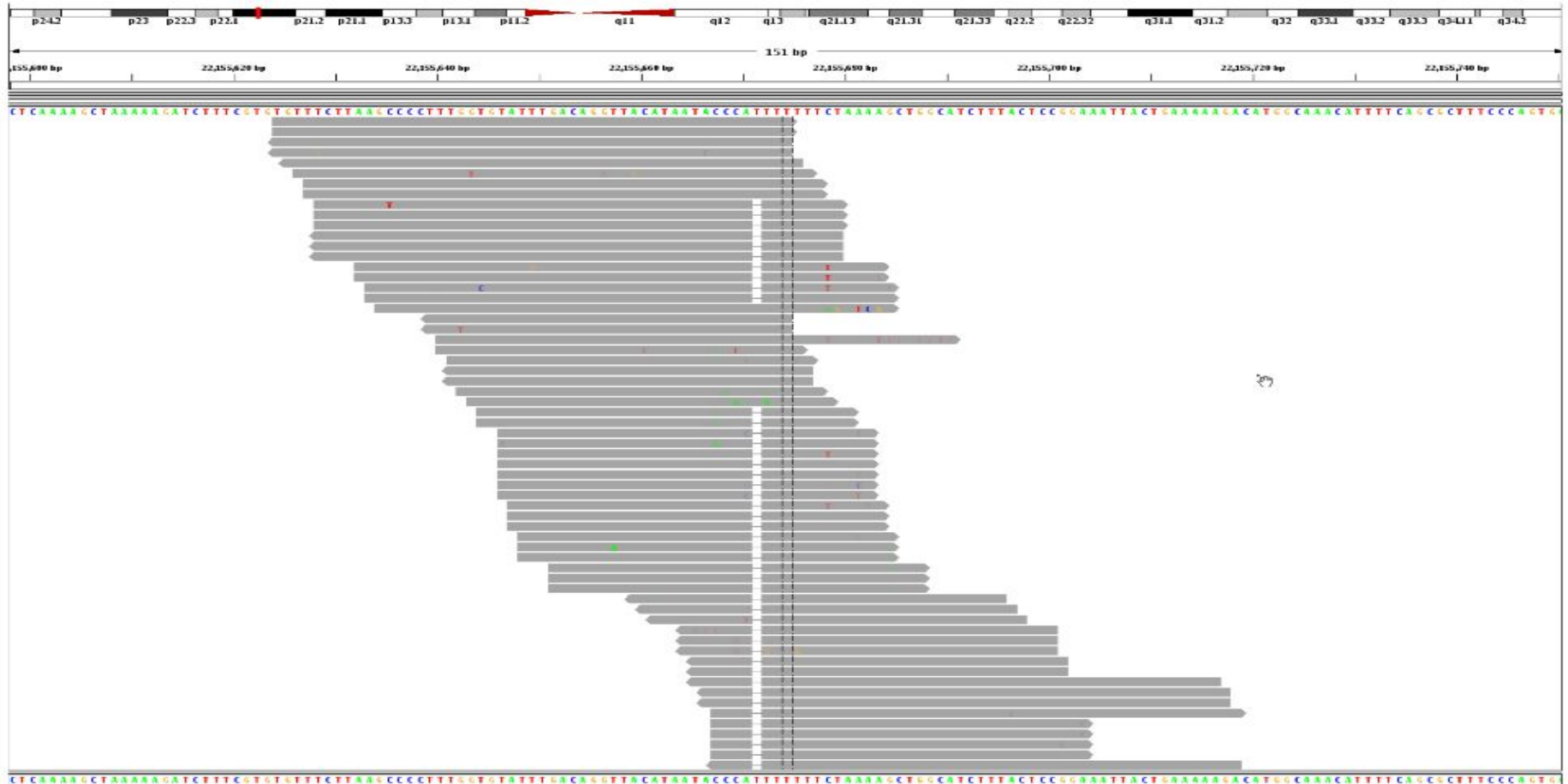# Sources of error and mitigation strategies

- **<u>Mapping errors:</u>**
  - Genomic duplication and structural variation.
  - High diverse areas.

- **<u>Solutions</u>**
  - Paired-end libraries.
  - Long reads / fragments.
  - MAPQ
  - Realignment around indels.

Adapted from Olson et al. Frontiers in Genetics. 2015

# Sources of error and mitigation strategies

# Sources of error and mitigation strategies

# Sources of error and mitigation strategies

- **SNP calling step**
  - Errors may result in base calling errors.
  - FP and FN calls.
- **Solutions**
  - Strand bias
  - Base quality rank sum
  - MAPQ
  - Hard filters:
    - Depth of coverage
    - Minimun base call frequency.

Adapted from Olson et al. Frontiers in Genetics. 2015

# Reference selection

- Critical step <- Bias which SNPs are called.
- SNPs in genes not present in the reference **<u>WON'T</u>** be called.
- Less effect in clonal bacteria.
- Number of SNPs called vary **A LOT!**

- **Solutions:**
  - Kmerfinder

# Repetitive/Phage regions filtering

- **<u>PHASTER</u>**

- We can remove/mask phague/repetitive regions where reads won't map.
- This way those areas will be out of analysis.
- Problem: those areas could be important!

# GMI Proficiency test

1. **Proficiency Testing for bacterial WGS, 2012 an end-user survey of current capabilities, requirements and priorities**

2. **Proficiency Test Pilot, 2014 Wet lab and Dry lab**

   *Escherichia coli, Staphilococus aureus and Salmonella typhimurium*

3. **Full Proficiency Test, 2015**

   *Escherichia coli, Staphilococus aureus and Salmonella tiphimurium*

4. **Full Proficiency Test, 2016 Wet lab and Dry lab**

   *Campylobacter coli and C. jejuni, Listeria monocytogenes and klebsiella        pneumoniae*

- **Número de SNPs reportado por cada laboratorio parcipante**

| Lab | EC | SA | ST |
|-----|-----|-----|-----|
| GMI02 | 25731 | 1383 | 8968 |
| GMI04 | 25731 | 1383 | 8968 |
| GMI06 | 43264 | 6226 | 5822 |
| GMI10 | 13083 | 1797 | 12902 |
| GMI14 | 14687 | NA | 1431 |
| GMI26 | 92831 | 6164 | 31044 |
| GMI39 | 52590 | 2672 | 16034 |
| GMI42 | 9460 | NA | 12884 |
| GMI43 | 38532 | 4163 | 16562 |
| GMI46 | 63273 | 2341 | 9958 |
| GMI48 | 67034 | 2063 | 14080 |
| GMI58 | 79231 | NA | 19656 |
| GMI59 | 23561 | 2715 | 14199 |
| GMI13 | 9276 | 1628 | 8746 |
| GMI16 | 55473 | 2122 | 13630 |
| GMI21 | 5187829 | 2837196 | 5090636 |
| GMI22 | 33416 | 1597 | 13066 |
| GMI27 | 33664 | 2130 | 13297 |
| GMI30 | 607217 | 11881 | 12733 |
| GMI31 | NA | NA | 4141 |
| GMI32 | 14667 | 25949 | 28164 |
| GMI33 | 71822 | 5420 | 21668 |
| GMI35 | 6706 | 1334 | NA |
| GMI37 | 73355 | 2897 | 14294 |
| GMI40 | 45725 | 2033 | 11180 |
| GMI44 | 35039 | 1836 | 9446 |
| GMI45 | 5183821 | 2836332 | 5088344 |
| GMI47 | 20707 | 1805 | 12198 |
| GMI50 | 84 | NA | 1300 |
| GMI51 | 35521 | NA | 10042 |
| GMI55 | NA | 1644 | 9102 |
| GMI61 | NA | NA | 24 |
| GMI63 | NA | 2834703 | 5077509 |
| GMI7 | 21731 | 1673 | 9192 |
| GMI8 | 15972 | 1851 | 12979 |

# VCF format

# Bed format

chromosome  start  end  score  name  strand  thickstart  thickend  RGB

```
chr7    127471196    127472363    Pos1    0    +    127471196    127472363    255,0,0
chr7    127472363    127473530    Pos2    0    +    127472363    127473530    255,0,0
chr7    127473530    127474697    Pos3    0    +    127473530    127474697    255,0,0
chr7    127474697    127475864    Pos4    0    +    127474697    127475864    255,0,0
chr7    127475864    127477031    Neg1    0    -    127475864    127477031    0,0,255
chr7    127477031    127478198    Neg2    0    -    127477031    127478198    0,0,255
chr7    127478198    127479365    Neg3    0    -    127478198    127479365    0,0,255
chr7    127479365    127480532    Pos5    0    +    127479365    127480532    255,0,0
chr7    127480532    127481699    Neg4    0    -    127480532    127481699    0,0,255
```
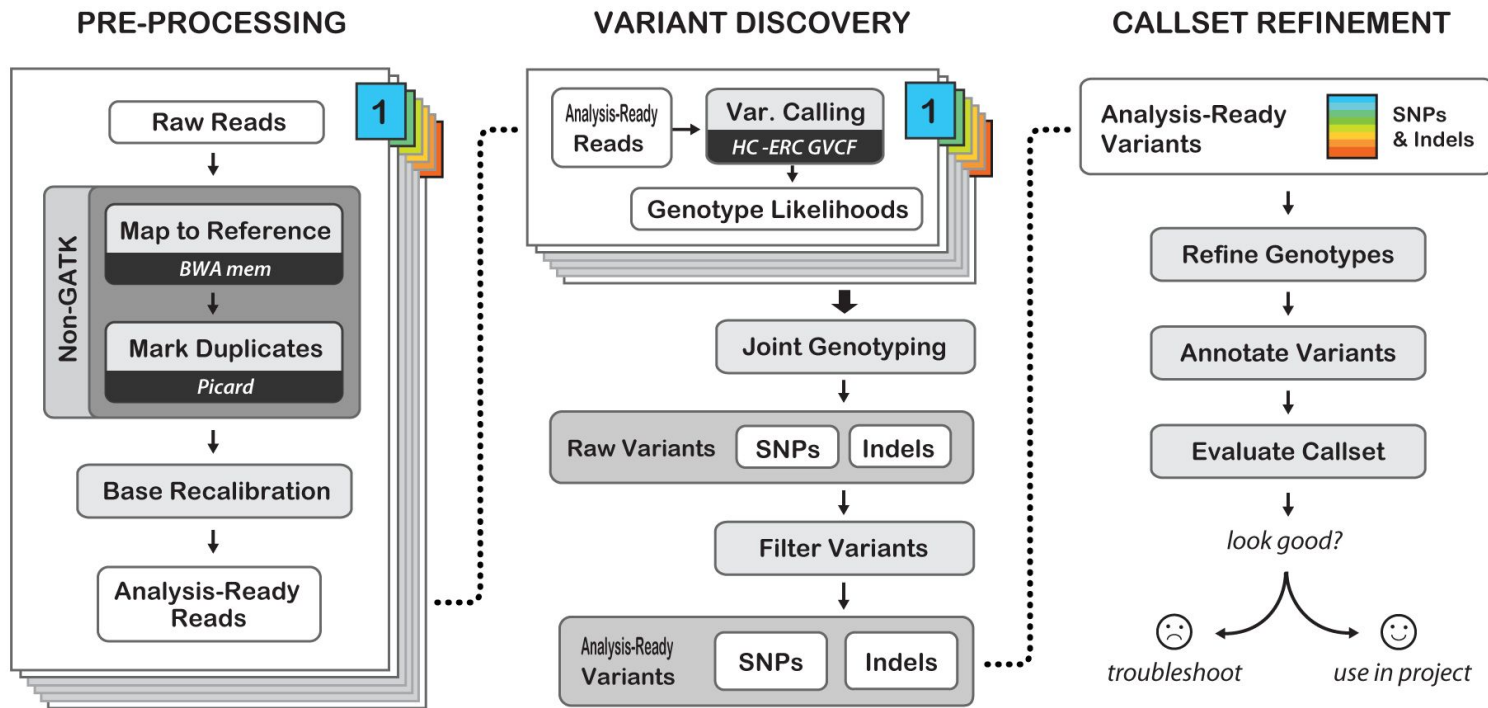
OBLIGATORIOS

OPCIONALES

# Pipelines for bacterial SNP-based analysis

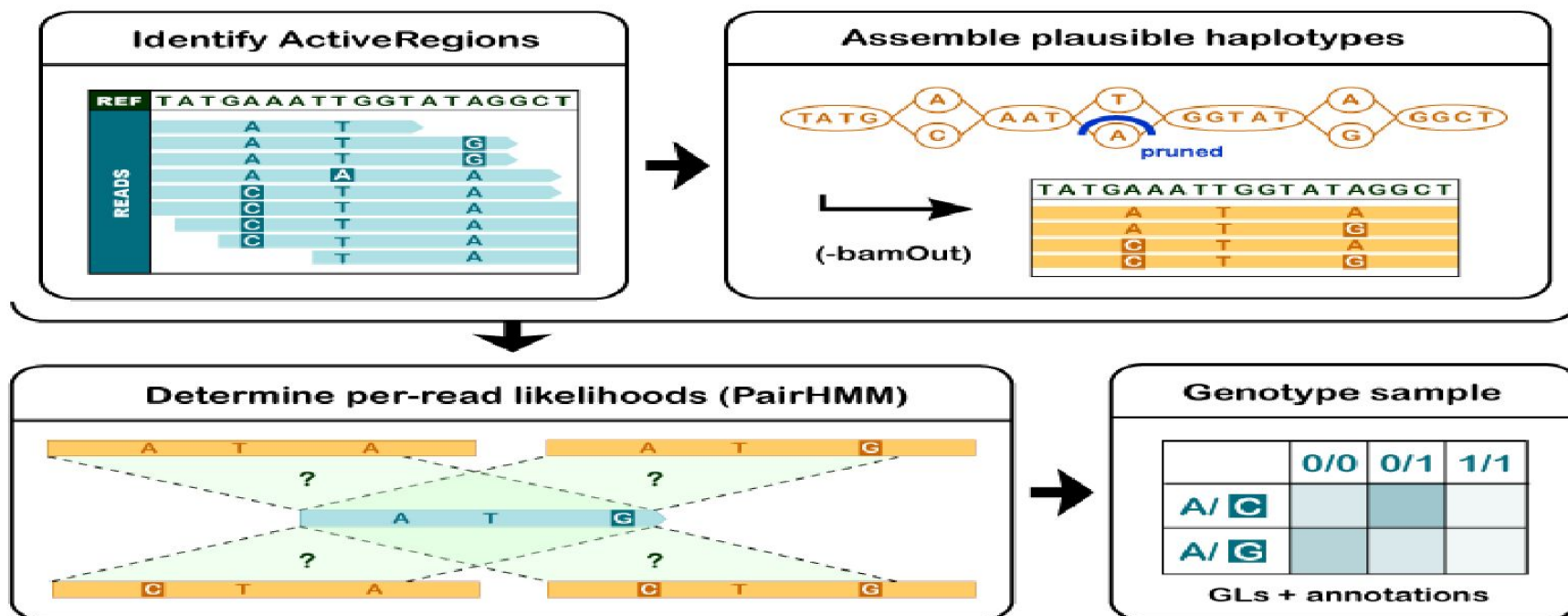| Software | Description | Other | References |
|---|---|---|---|
| **CFSAN** | VARSCAN variant calling | Terminal | Davis et al., 2015 |
| **NASP** | Variant calling with VarScan, solSNP,samtools and GATK | Terminal | Sahl et al., 2016 |
| **Lyve-Set** | VARSCAN variant calling | Terminal | Katz et al., 2017 |
| **KSNP** | Reference free variant calling. | Terminal | Gardner et al., 2015 |
| **SNVPhyl** | Variant calling with freebayes and samtools | Galaxy | Petkau et al., 2017 |
| **CSI phylogeny** | Variant calling with samtools. | Web | Kaas et al., 2014 |

# GATK



**Best Practices for Germline SNPs and Indels in Whole Genomes and Exomes - June 2016**

# GATK

- **1. Define active regions**
  - The program determines which regions of the genome it needs to operate on, based on the presence of significant evidence for variation.
- **2. Determine haplotypes by assembly of the active region**
  - For each ActiveRegion, the program builds a De Bruijn-like graph to reassemble the ActiveRegion, and identifies what are the possible haplotypes present in the data. The program then realigns each haplotype against the reference haplotype using the Smith-Waterman algorithm in order to identify potentially variant sites.
- **3. Determine likelihoods of the haplotypes given the read data**
  - For each ActiveRegion, the program performs a pairwise alignment of each read against each haplotype using the PairHMM algorithm. This produces a matrix of likelihoods of haplotypes given the read data. These likelihoods are then marginalized to obtain the likelihoods of alleles for each potentially variant site given the read data.
- **4. Assign sample genotypes**
  - For each potentially variant site, the program applies Bayes' rule, using the likelihoods of alleles given the read data to calculate the likelihoods of each genotype per sample given the read data observed for that sample. The most likely genotype is then assigned to the sample.
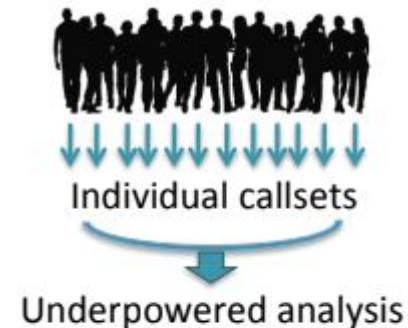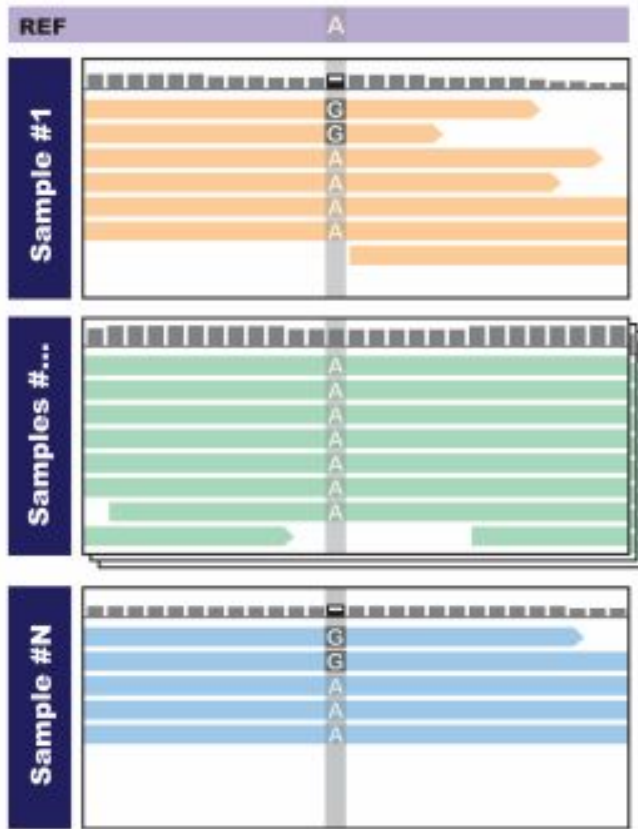
Best GATK practice guide.

# GATK

# Cohorts need to be analyzed together at variant calling step

- If we simply call variants on individual samples then merge lists of their variants, we miss a lot of important information

- Joint variant discovery rescues a lot of valuable information



Sequencing and variant calling pipelines MPG Primer @ Broad Institute Cambridge, 15 October, 2015

# Joint analysis empowers calls in difficult sites



- If we analyze Sample #1 or Sample #N alone we are not confident that the variant is real

- If we see both samples then we are more confident that there is real variation at this site in the cohort

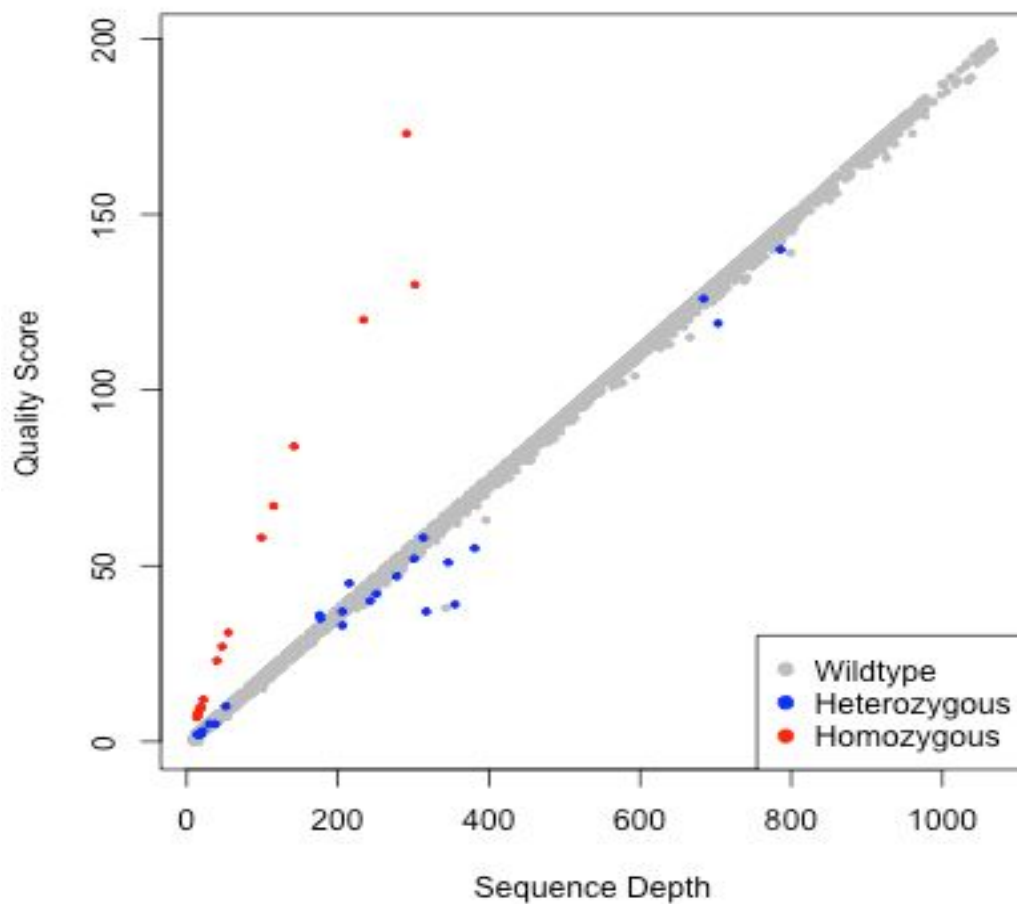Sequencing and variant calling pipelines MPG Primer @ Broad Institute Cambridge, 15 October, 2015

# GATK problems

- Haploid variant calling is a side project. GATK is mainly for diploid organism, and the development and improvement of the haploid algorithm is slow.
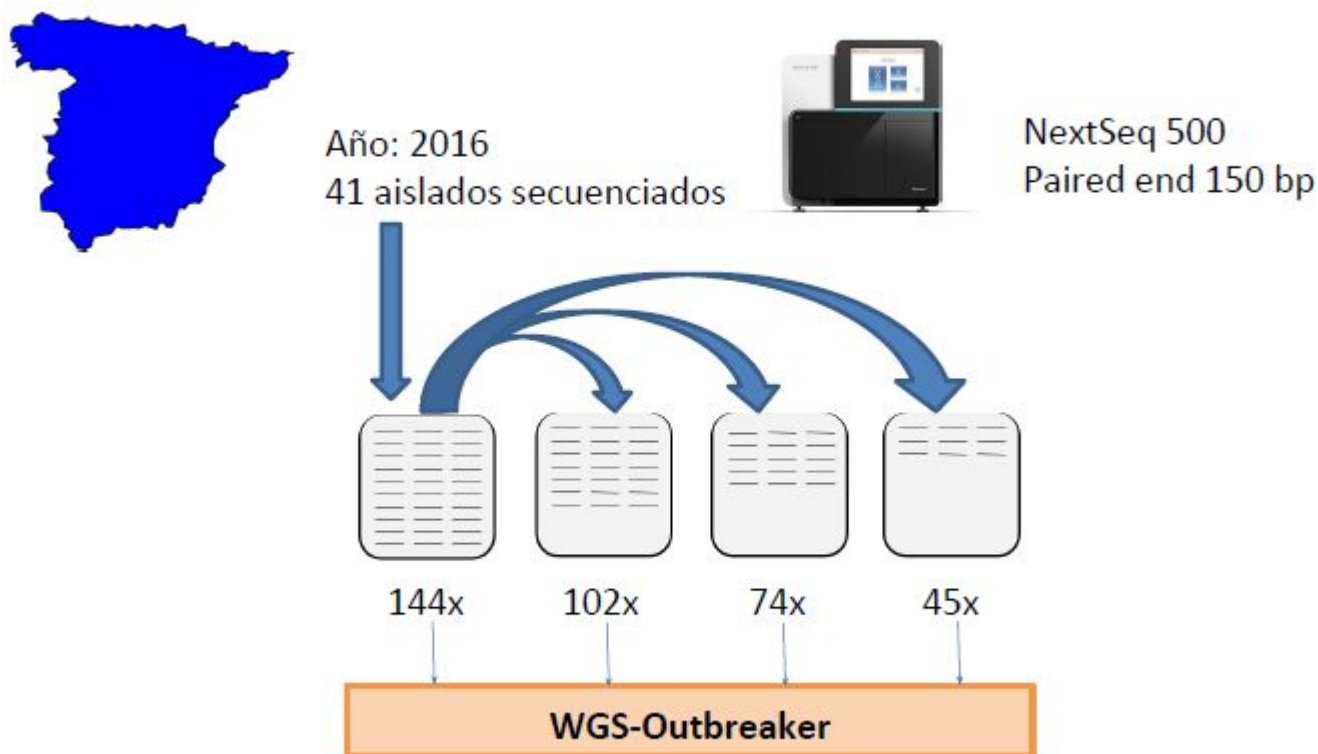- Manual filtering is needed for high quality SNPs selection.

# VARSCAN2

- Uses a heuristic/statistic method instead of bayesian.
- Allows more flexibility and hard filters.
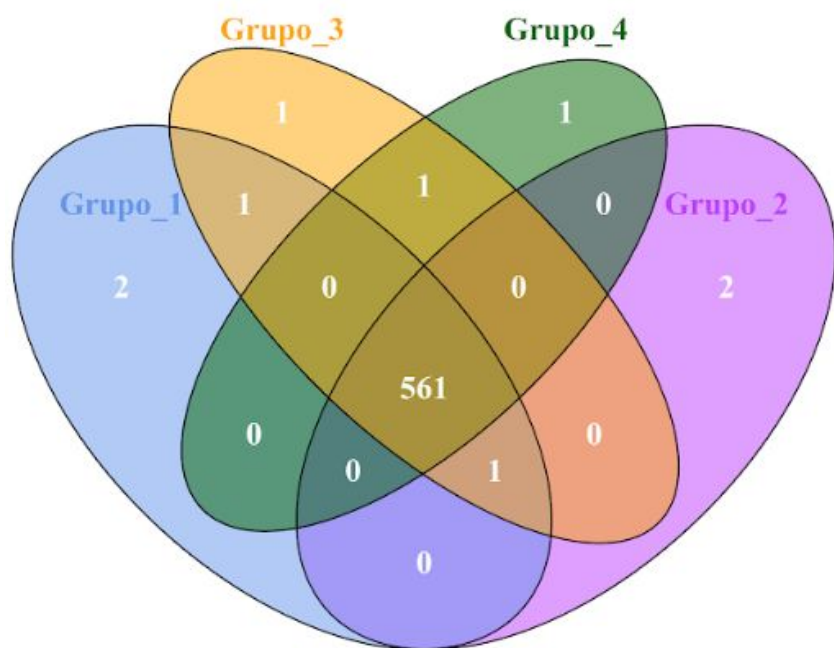- Used in many bacterial variant calling pipelines. P.e CFSAN snp-pipeline.

# VARSCAN2



Secuenciación de genomas bacterianos:
herramientas y aplicaciones

55

# Comparative VARSCAN - GATK

# Comparative VARSCAN - GATK



WGS-Outbreaker - GATK

CFSAN - VARSCAN

# Comparative VARSCAN - GATK

# High Quality SNP selection



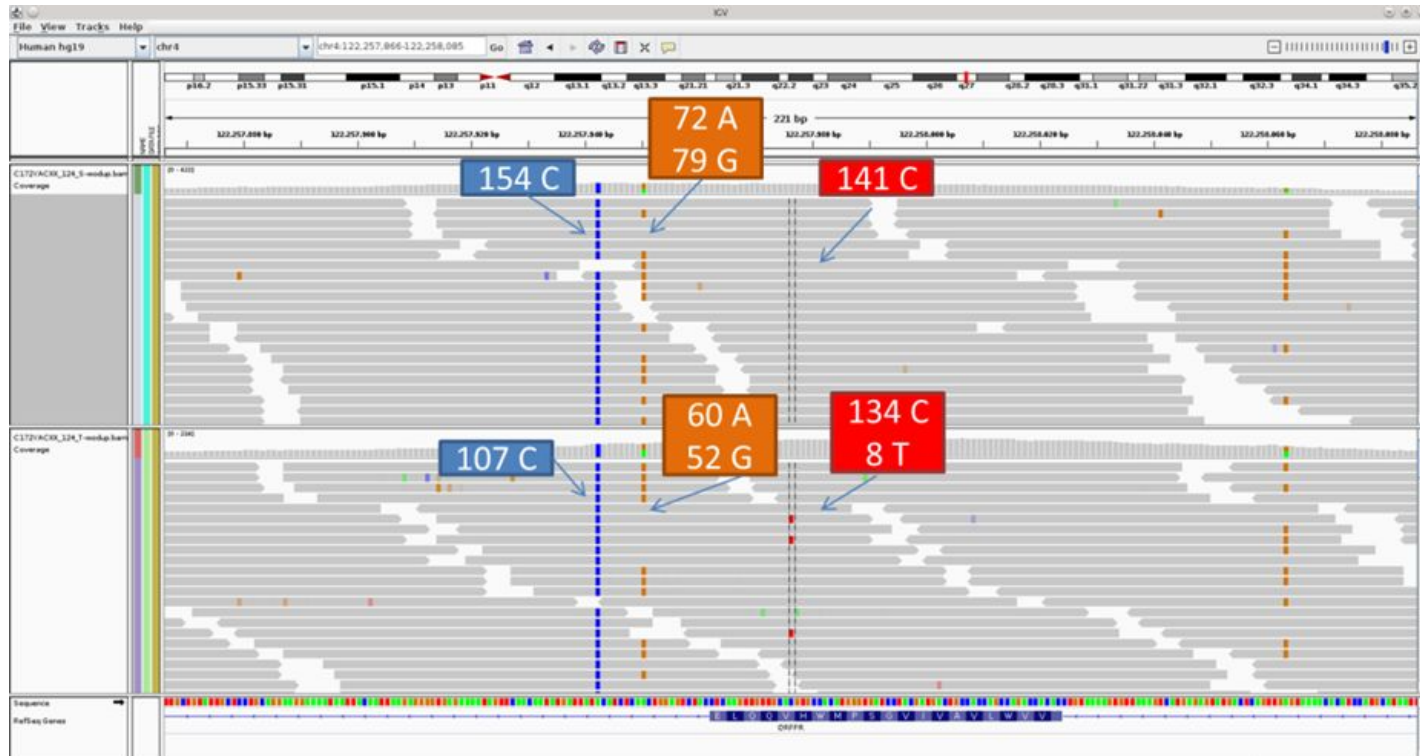| CFSAN Filtering | | GATK |
| :---: | :---: | :---: |
| ✓ | PhredQ | ✓ |
| ✗ | Strand bias | ✓ |
| ✗ | MAPQ | ✓ |
| ✓ | AD filtering | ✗ |
| ✓ | SNP Cluster | ✓ |

# Population Allele frequency vs Sample Allele frequency

- **Population allele frequency:** probability of finding an allele in the population. Number of individuals carrying an allele vs total of individuals in the population.
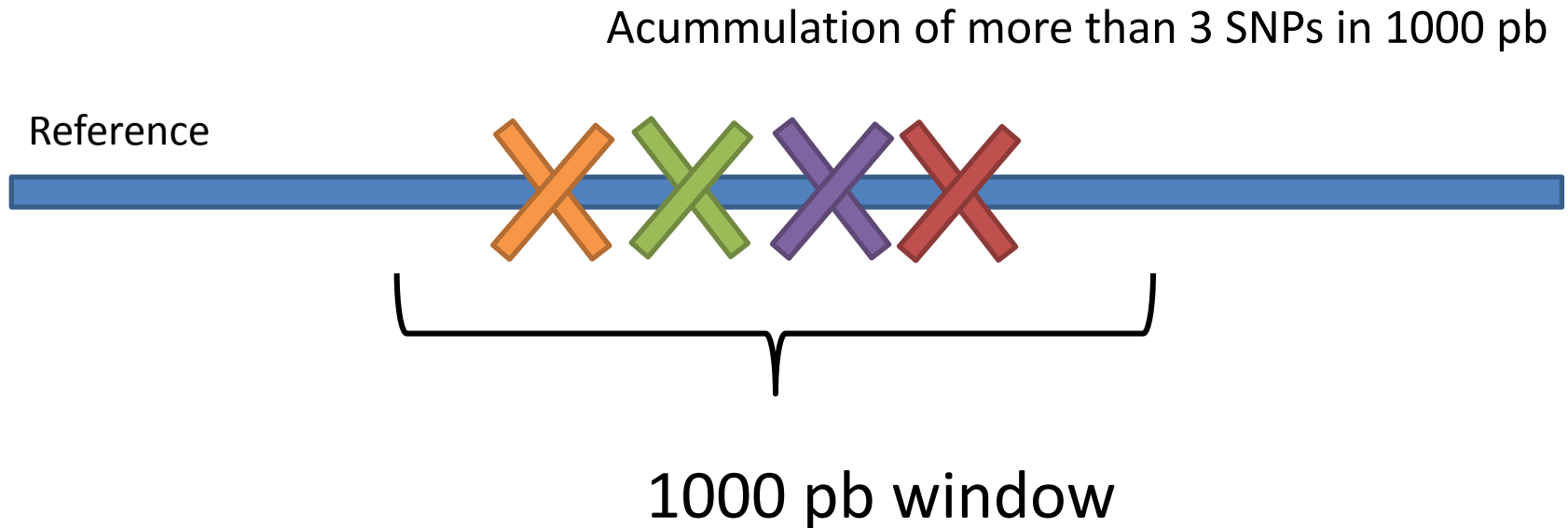
# Population Allele frequency vs Sample Allele frequency

- **Alternate/Base allele frequency**: number of reads supporting the alternate allele vs total of reads.

# SNP cluster filtering

Acummulation of more than 3 SNPs in 1000 pb

Reference



1000 pb window

# What's next?

**SNP matrix creation**

**And**

**Phylogeny!**

# Thanks for your attention!