# Mapping against reference genome and Variant Calling

**BU-ISCIII**

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

2-10 Noviembre 2022
Programa Formación, AESAN

# 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. |

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# 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
```

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# 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
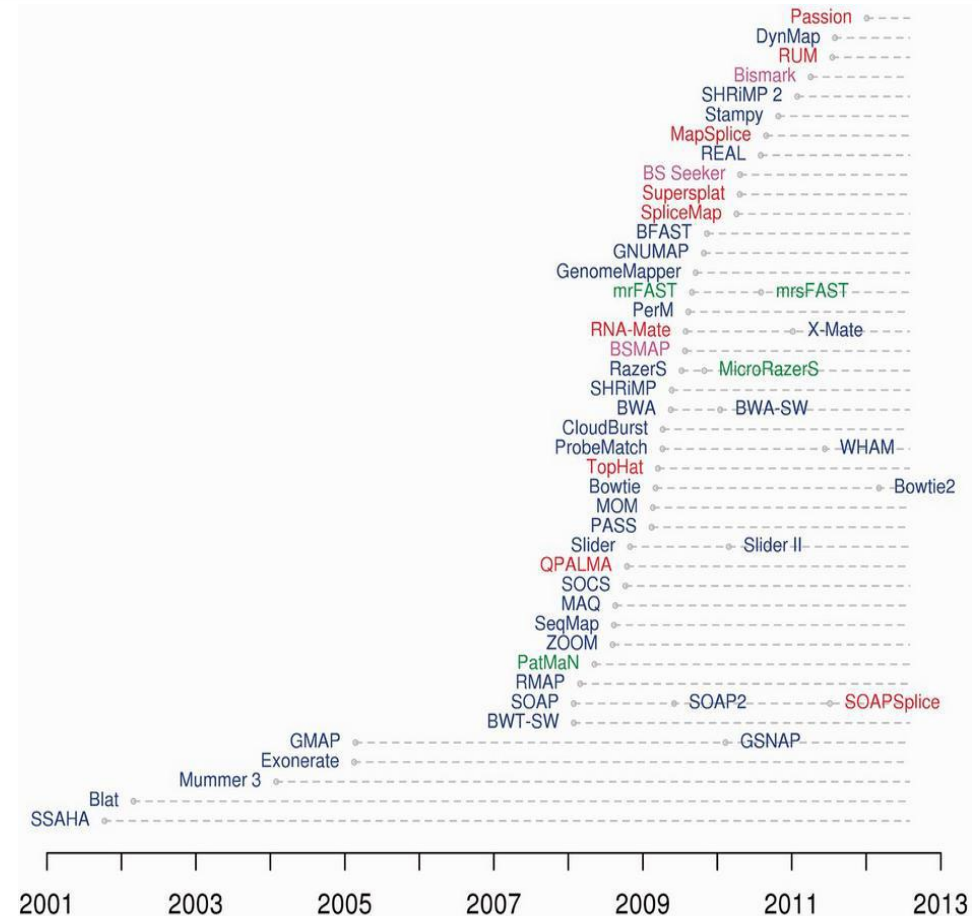
- Errors in reads

- 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. |

Secuenciación de genomas bacterianos: herramientas y aplicaciones
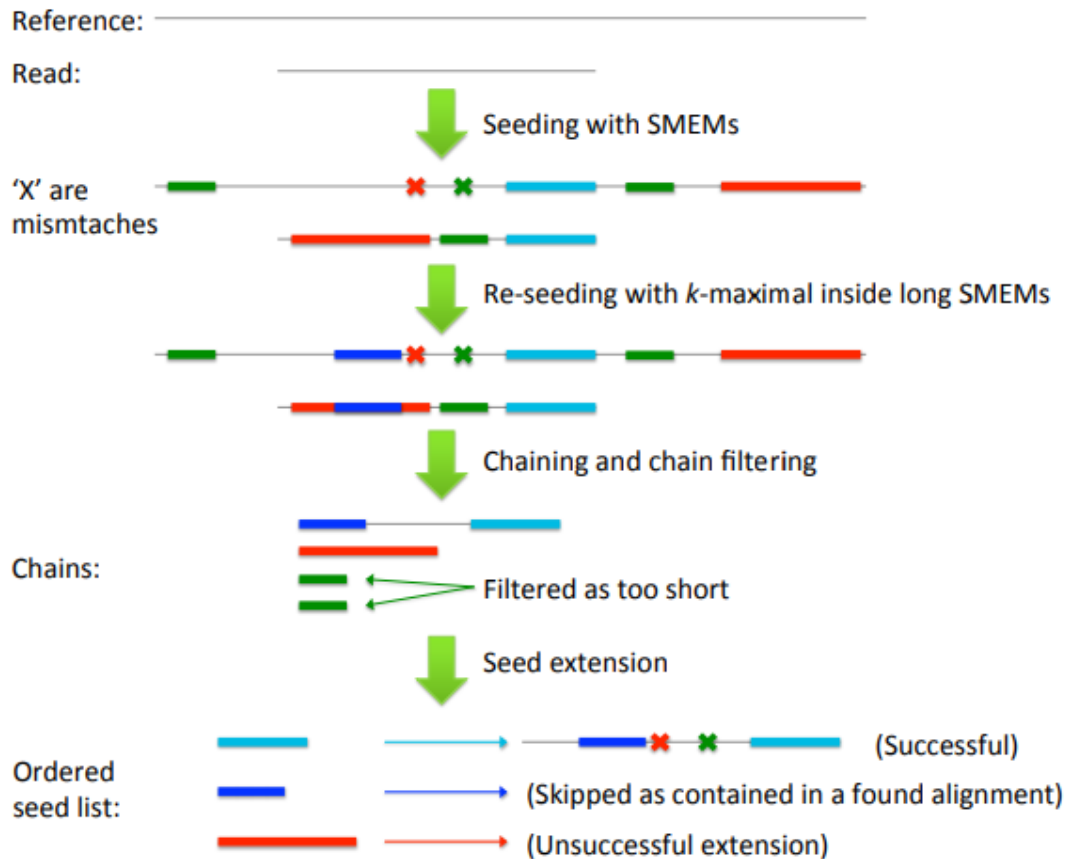
# What mapper should I use?

## 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

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# BWA MEM



## 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.

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# BOWTIE2

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

**Three reporting modes:**

- – Best alignment
- – K alignments
- – All alignments

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

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# 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.

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# Which aligner should I use for aligning reads agains a complete genome for variant calling?

Reference genome



Reads mapping uniquely

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

Secuenciación de genomas bacterianos:
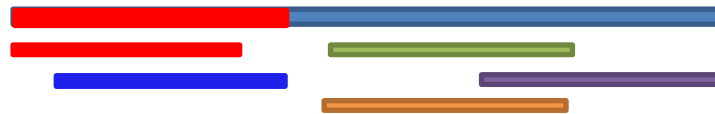herramientas y aplicaciones

# 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

Reads mapping uniquely only map in Allele 1. Which is the one more **covered**
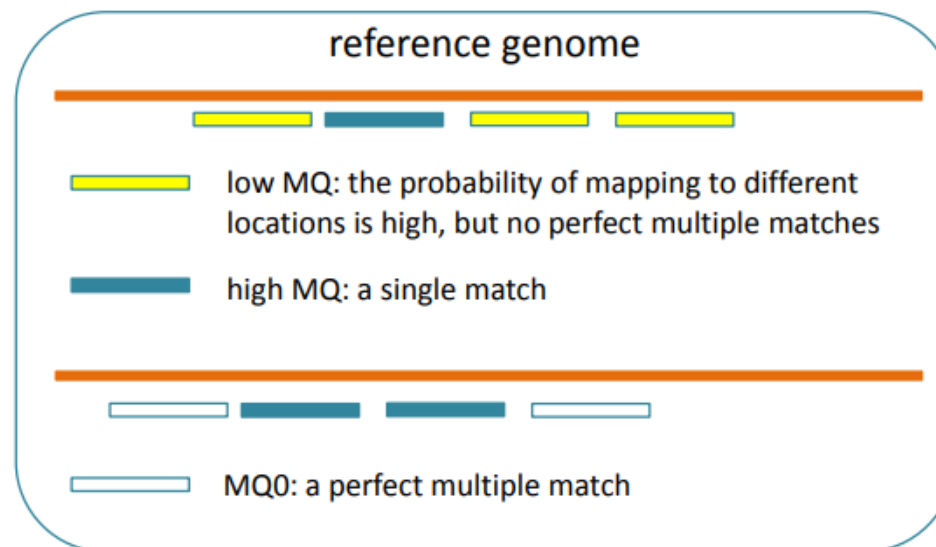
Resistance gene - Allele 2

Resistance gene - Allele 3

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# 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**



reference genome

low MQ: the probability of mapping to different locations is high, but no perfect multiple matches

high MQ: a single match

MQ0: a perfect multiple match

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

>&_BU-ISCIII

# 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
```

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# 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

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# Flag explanation example 2

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# 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!!**

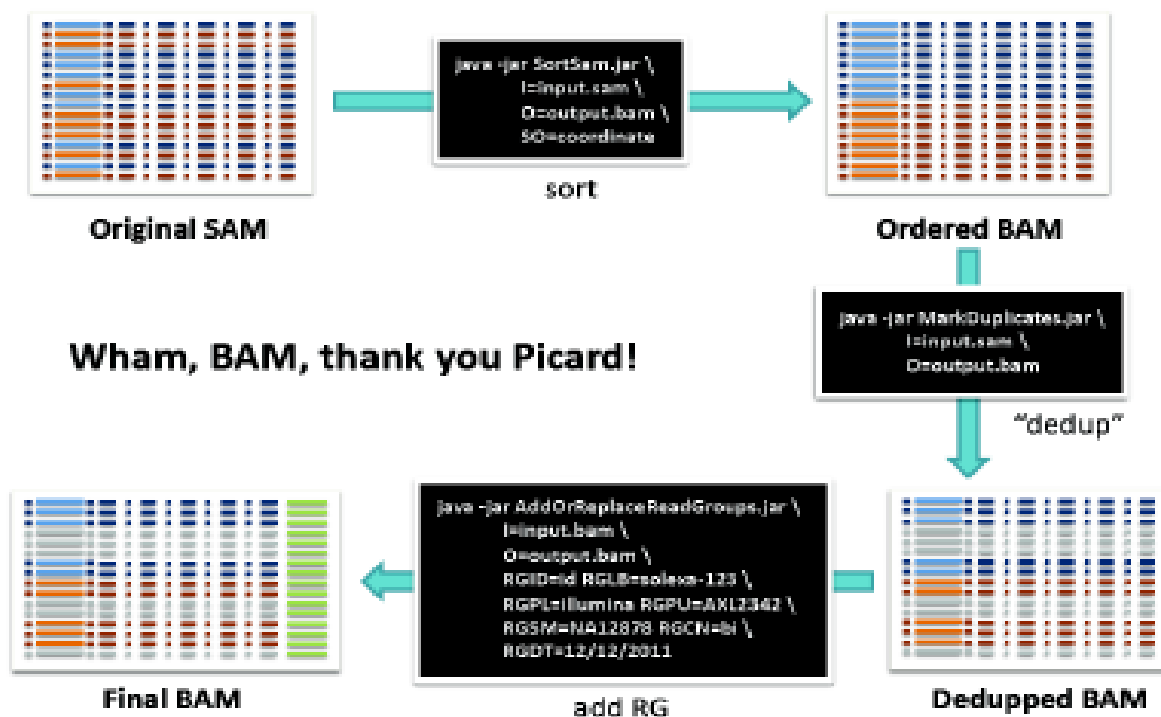Typical bam and sam format files weights from a S. grumpensis
SAM format file: 3.6 GB
BAM format file: 689 M

# 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

Secuenciación de genomas bacterianos:
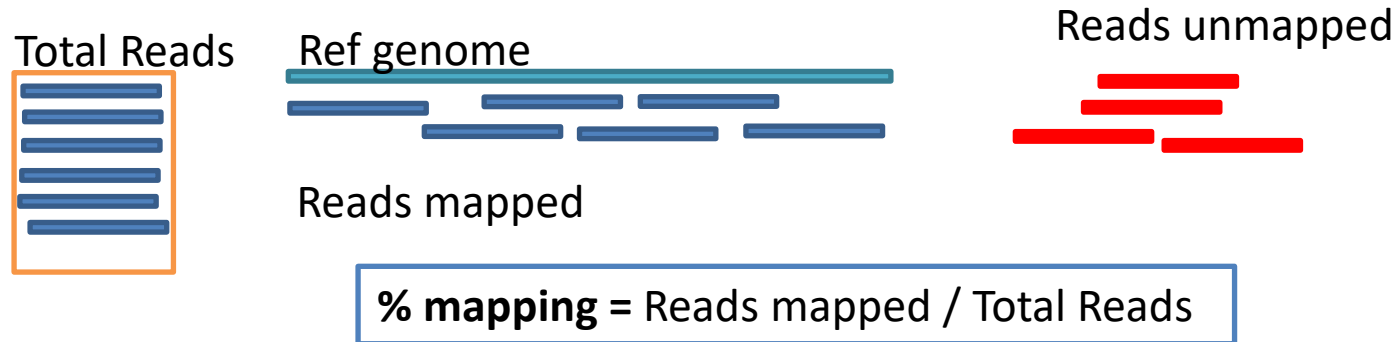herramientas y aplicaciones

# 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.

Análisis de Genomas Virales a través de la
plataforma Galaxy

# Mapping quality control

Picard
Samtools

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



Total Reads    Ref genome

Reads unmapped

Reads mapped

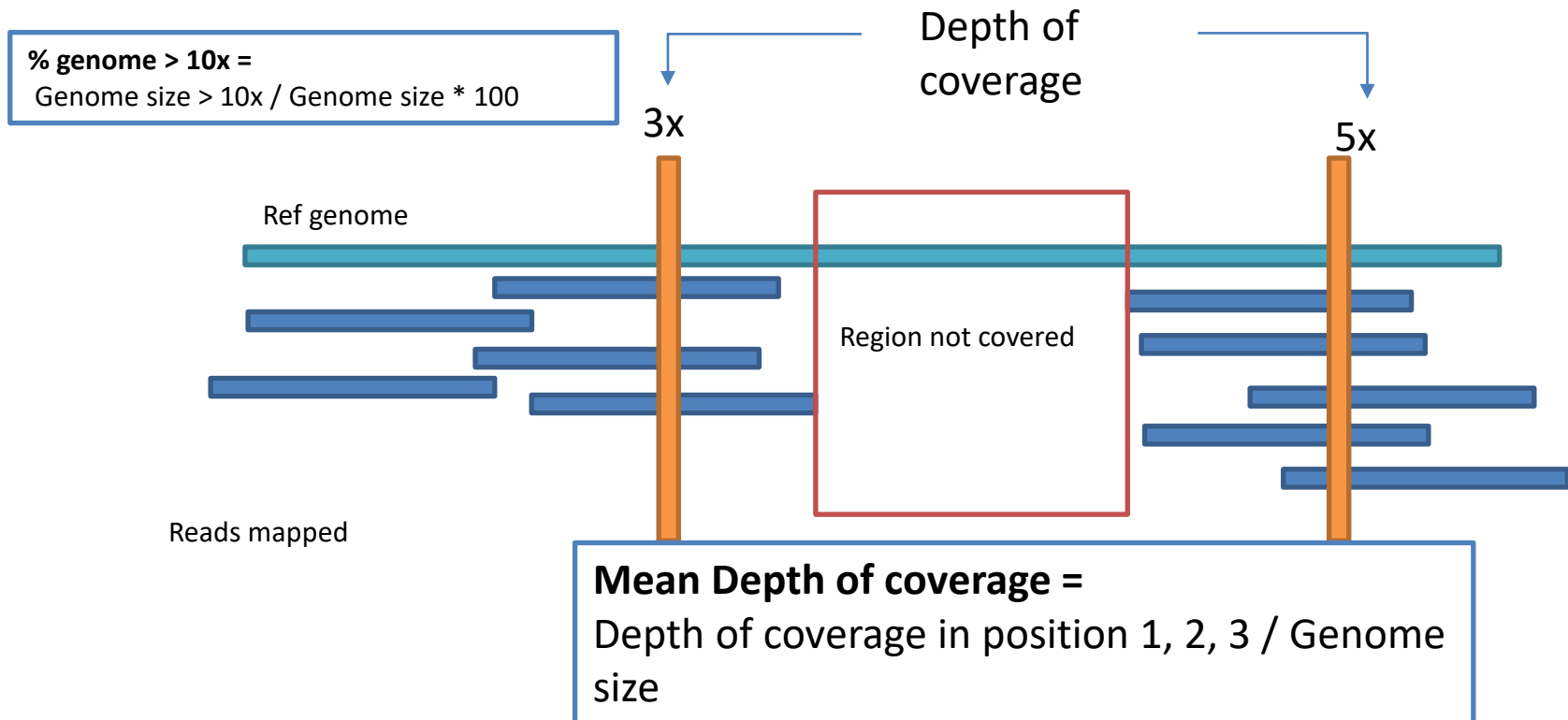**% mapping =** Reads mapped / Total Reads

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.
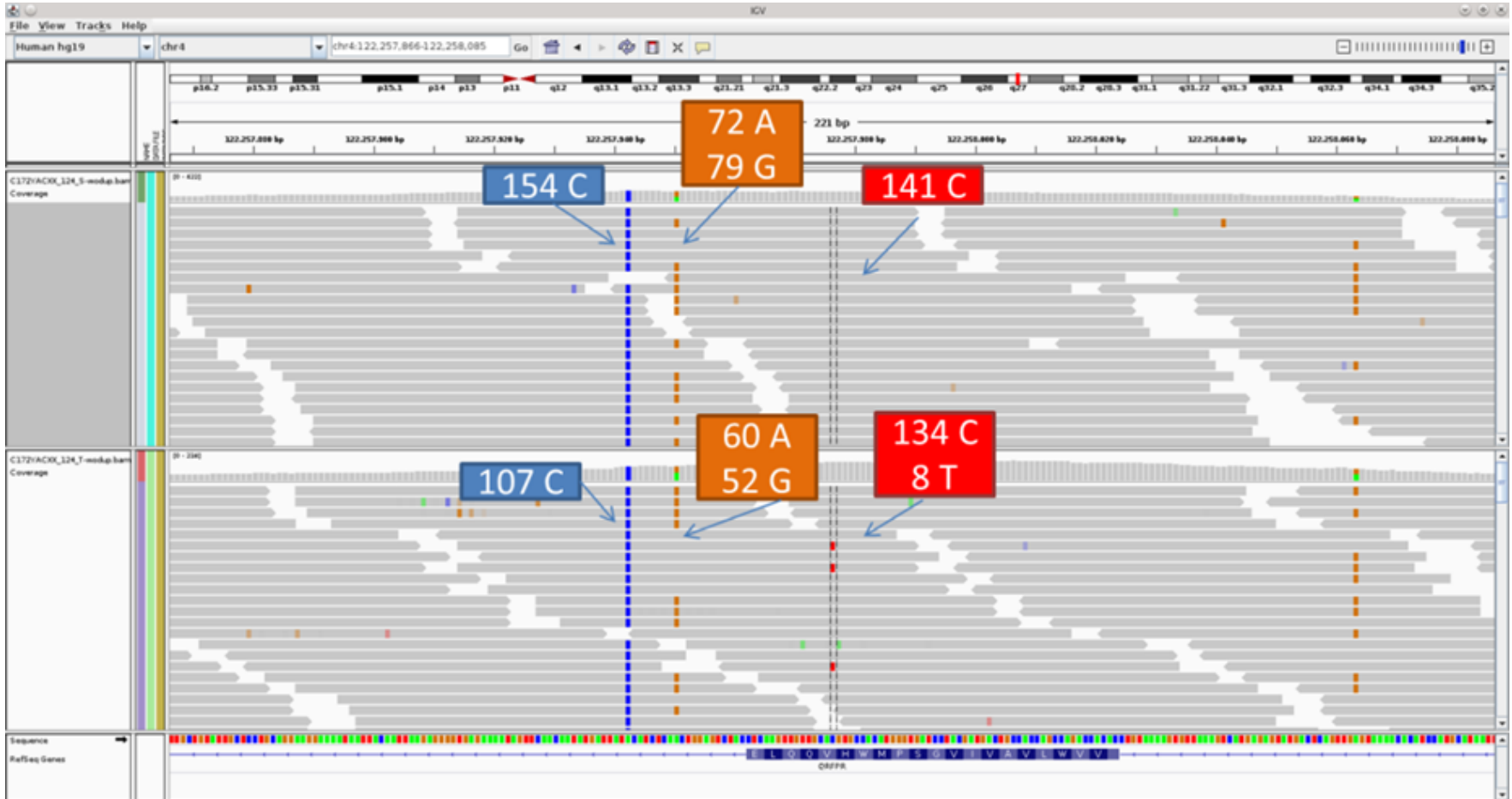
Picard
Samtools

**% 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

# 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.

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# 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

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# Sources of error and mitigation strategies

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# Reference selection

- Critial 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

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# 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!

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# VCF format



Secuenciación de genomas bacterianos: herramientas y aplicaciones

# 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
```
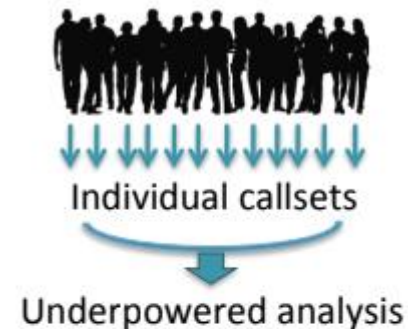
OBLIGATORIOS                                    OPCIONALES

# Pipelines for bacterial SNP-based analysis

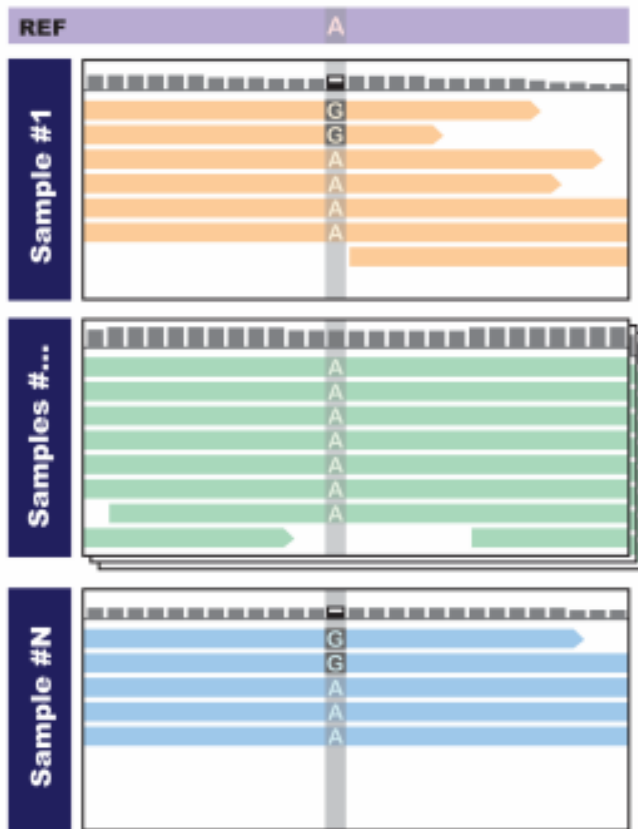| Software | Description | Other | References |
|---|---|---|---|
| CFSAN | VARSCAN variant calling | Terminal | Davis et al., 2015 |
| NASP | Variant calling with VarScan, solSNP,samtools and GATK | Terminal | Sahl et al., 2016 |
| Lyve-Set | VARSCAN variant calling | Terminal | Katz et al., 2017 |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

# 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

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# High Quality SNP selection

| CFSAN Filtering | | GATK |
|:---:|:---:|:---:|
| ✔ | PhredQ | ✔ |
| ✘ | Strand bias | ✔ |
| ✘ | MAPQ | ✔ |
| ✔ | AD filtering | ✘ |
| ✔ | SNP Cluster | ✔ |

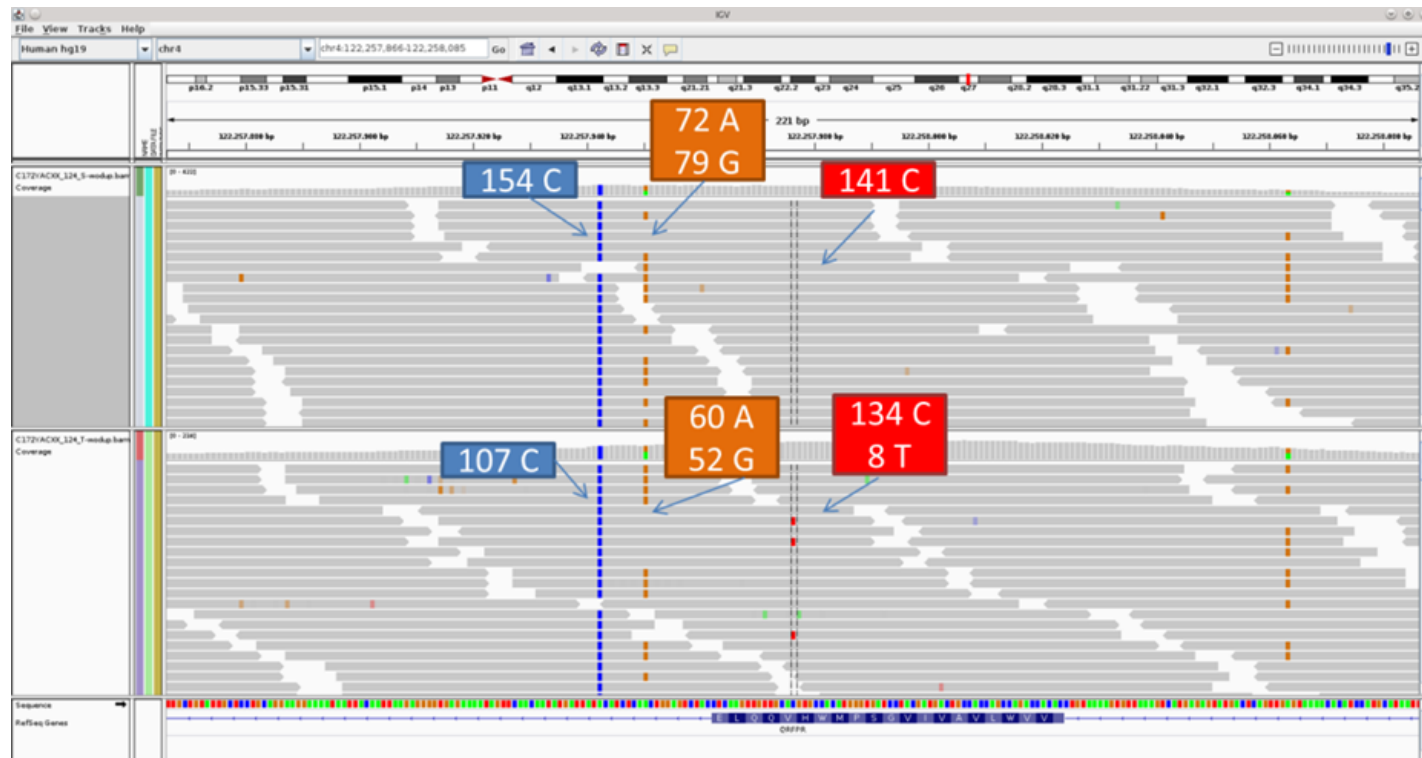Secuenciación de genomas bacterianos: herramientas y aplicaciones

# 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.

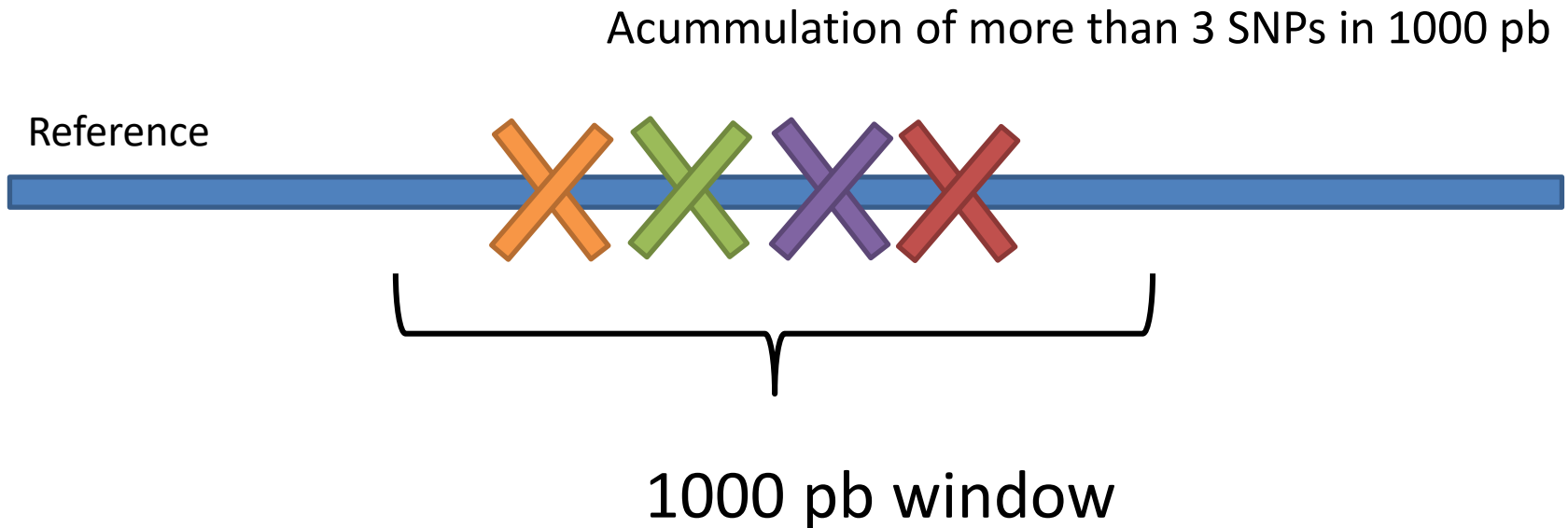Secuenciación de genomas bacterianos: herramientas y aplicaciones

# Population Allele frequency vs Sample Allele frequency

- **Alternate/Base allele frequency**: number of reads supporting the alternate allele vs total of reads.
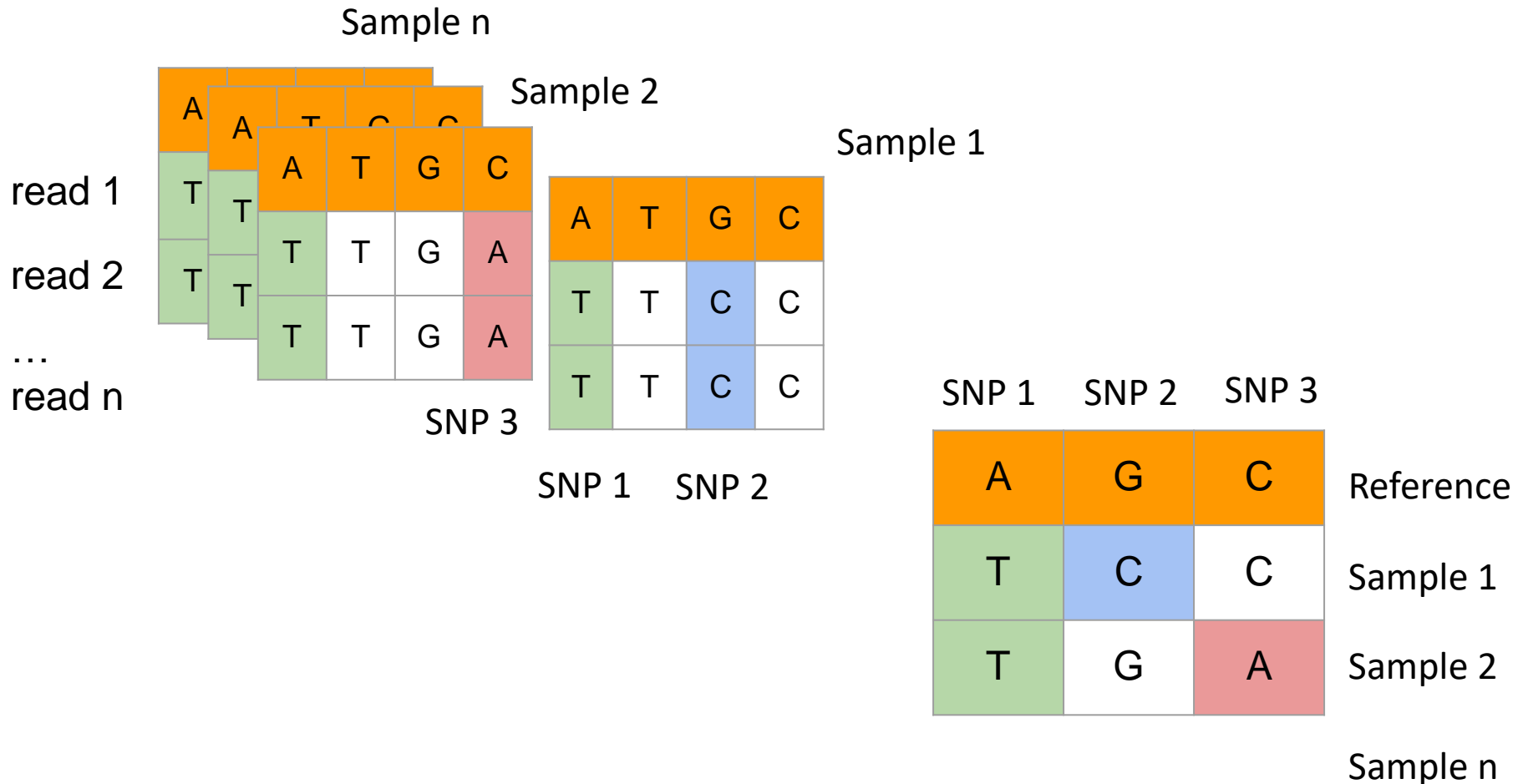
Secuenciación de genomas bacterianos: herramientas y aplicaciones

# SNP cluster filtering

Acummulation of more than 3 SNPs in 1000 pb

Reference

1000 pb window

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# What's next?

**SNP matrix creation**

**And**

**Phylogeny!**

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# Building a SNP matrix

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# Building a SNP matrix

- Once we have our multisample vcf:

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | RA-L2073 | RA-2805 |
|--------|-----|-----|-----|-----|---------|--------|--------|-----------|------------|-----------|
| NC_021827.1 | 276 | . | C | A | 291.68 | PASS | AC=1;... | GT:AD:DP:... | 0:13,0:13:... | 1:0,50:30 |
| NC_021827.1 | 731 | . | A | G | 2313.68 | PASS | AC=1;... | GT:AD:DP:... | 0:23,0:23:... | 1:0,10:10 |
| NC_021827.1 | 921 | . | C | T | 1841.68 | PASS | AC=1;... | GT:AD:DP:... | 0:53,0:53:... | 0:20,0:20 |

- We can generate the genotype for each sample

| #CHROM | POS | RA-L2073 | RA-2805 |
|--------|-----|----------|---------|
| NC_021827.1 | 276 | C | A |
| NC_021827.1 | 731 | A | G |
| NC_021827.1 | 921 | C | C |

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# Building a SNP matrix

- So… now we have a simple multifasta, where each nucleotide represents a SNP.

- This means that even the nucleotide positions are sequentially in the fasta, they don't have to be near each other in the genome!

- The SNP matrix file will look like this:

> RA-2073
CACGAAATTCCATTA
↑↑↑↑

>RA-2805
AGCTCATGCATATTA
↑↑↑↑

Each of this is a SNP:
First one is in position 276 in the genome
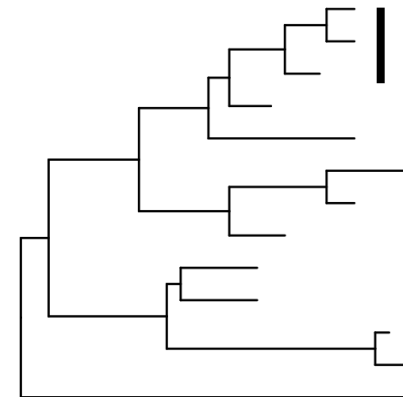Second one is in position 731 in the genome
Third one is in position  921 in the genome

Each SNP is ordered per SNP position. In this sample also first SNP is in position 276 in the genome
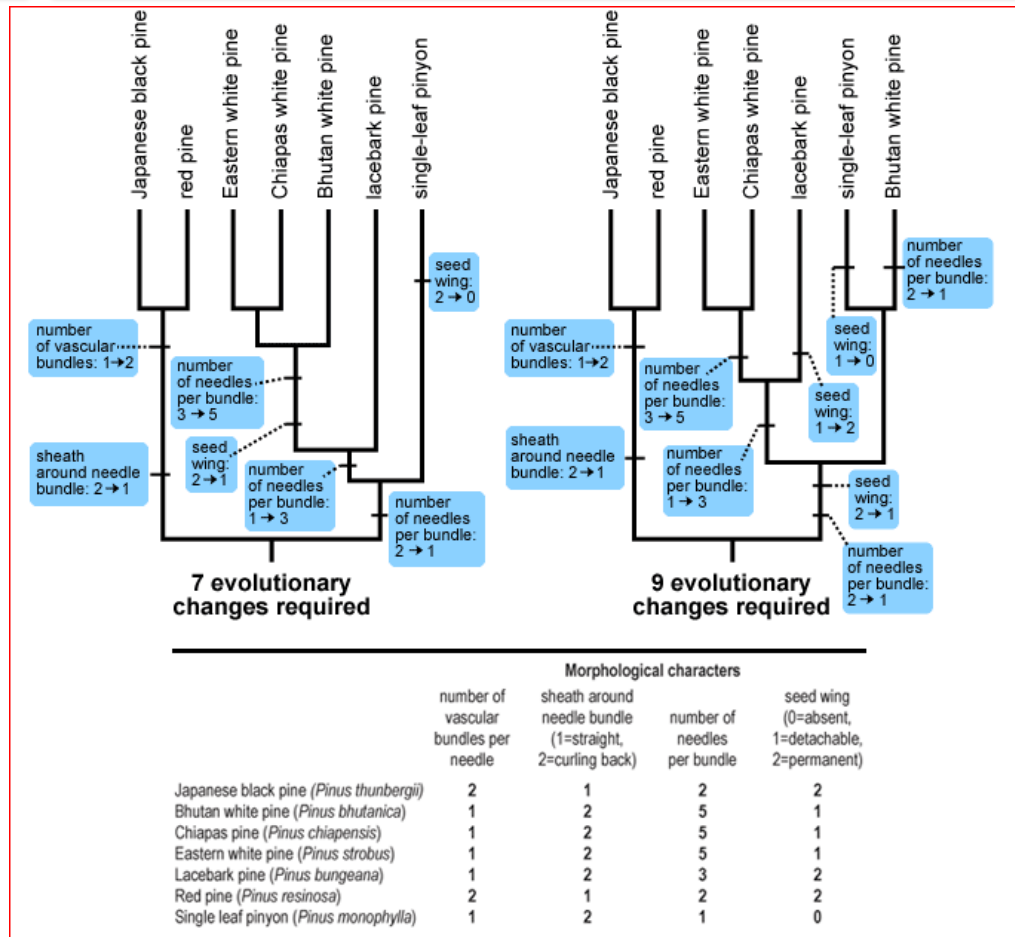
# Phylogeny

SNP matrix

| SNP 1 | SNP 2 | SNP 3 | |
|-------|-------|-------|-----------|
| A | G | C | Reference |
| T | C | C | Sample 1 |
| T | G | A | Sample 2 |

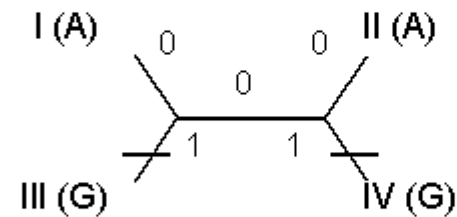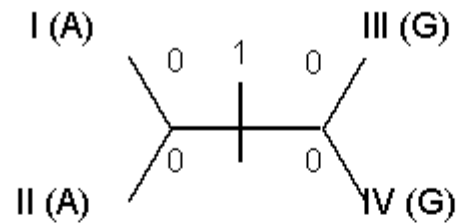Sample n

Phylogeny →

Outbreak!

# Maximum parsimony



- Search the most parsimonious tree
- The most simple hipothesis must be the correct.
- Search the tree that explains the relationships with the less changes possible.

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# Maximum Likelihood

- Searchs the most likely tree given the data and based in a evolutionary model.
- More sofisticated.
- Not prepared a priori for snp matrix.
- RAxML
    - Heterogenity rate disabled.
    - Branchs indicate the expected number of substitution per site.



- 0,1 = differences along that branch
- Which hypothesis is more likely, given that the change is rare?

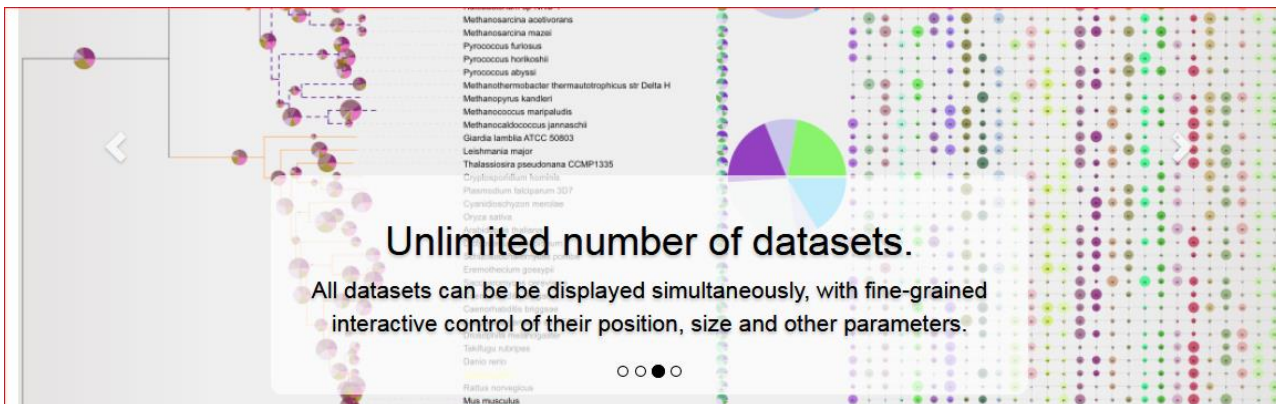Secuenciación de genomas bacterianos: herramientas y aplicaciones

# SNP distance

- SNP distance is calculated with N model. Simply the number of sites that differ between each pair of sequences.

- By default sites with at least one missing data is deleted for all sequences in R (complete deletion option in MEGA).

| dis_matrix.names | RA.L2073 | RA.L2281 | RA.L2327 | RA.L2391 | RA.L2450 | RA.L2677 | RA.L2701 |
|---|---|---|---|---|---|---|---|
| RA-L2073 | 0 | 9403 | 9028 | 80 | 46 | 46 | 49 |
| RA-L2281 | 9403 | 0 | 8777 | 9415 | 9397 | 9397 | 9402 |
| RA-L2327 | 9028 | 8777 | 0 | 9040 | 9022 | 9022 | 9027 |
| RA-L2391 | 80 | 9415 | 9040 | 0 | 74 | 74 | 79 |
| RA-L2450 | 46 | 9397 | 9022 | 74 | 0 | 38 | 45 |
| RA-L2677 | 46 | 9397 | 9022 | 74 | 38 | 0 | 45 |
| RA-L2701 | 49 | 9402 | 9027 | 79 | 45 | 45 | 0 |
| RA-L2782 | 9120 | 9183 | 4277 | 9132 | 9114 | 9114 | 9119 |
| RA-L2805 | 4 | 9403 | 9028 | 80 | 46 | 46 | 49 |
| RA-L2978 | 2 | 9401 | 9026 | 78 | 44 | 44 | 47 |

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# iTOL



https://itol.embl.de/

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# How to interpret our phylogeny

- **<u>Combination of:</u>**
  - SNP counts
  - Tree topologies
  - Bootstrap support

Pightling et al. Frontiers in Microbiology. 2018

# How to interpret our phylogeny

**TABLE 1** | Maximum pairwise SNPs measured during investigations into foodborne illness outbreaks and contamination events.

| Organism | Maximum SNP count (number) | Maximum SNP count (range) | | | Reference |
|---|---|---|---|---|---|
| | | <21 | 21–100 | >100 | |
| E. coli | 4 | X | | | Underwood et al., 2013 |
| E. coli | 15 | X | | | Eppinger et al., 2011 |
| L. monocytogenes | 9 | X | | | Chen et al., 2017c |
| L. monocytogenes | 12 | X | | | Chen et al., 2017a |
| L. monocytogenes | 18 | X | | | Li et al., 2017 |
| L. monocytogenes | 20 | X | | | Wang et al., 2015 |
| L. monocytogenes | 21 | | X | | Nielsen et al., 2017 |
| L. monocytogenes | 28 | | X | | Gilmour et al., 2010 |
| L. monocytogenes | 29 | | X | | Chen et al., 2017b |
| L. monocytogenes | 42 | | X | | Chen et al., 2016 |
| L. monocytogenes | 67 | | X | | Jackson et al., 2016 |
| S. enterica | 2 | X | | | Wuyts et al., 2015 |
| S. enterica | 3 | X | | | Allard et al., 2016 |
| S. enterica | 3 | X | | | Taylor et al., 2015 |
| S. enterica | 6 | X | | | Hoffmann et al., 2016 |
| S. enterica | 12 | X | | | Octavia et al., 2015 |
| S. enterica | 30 | | X | | Leekitcharoenphon et al., 2014 |

*The maximum SNP counts for isolates that were traced back to the same source in the original study are presented. Whether the maximum SNP counts are less than 21 SNPs, 21 to 100 SNP, or greater than 100 SNPs is also indicated.*

Pightling et al. Frontiers in Microbiology. 2018
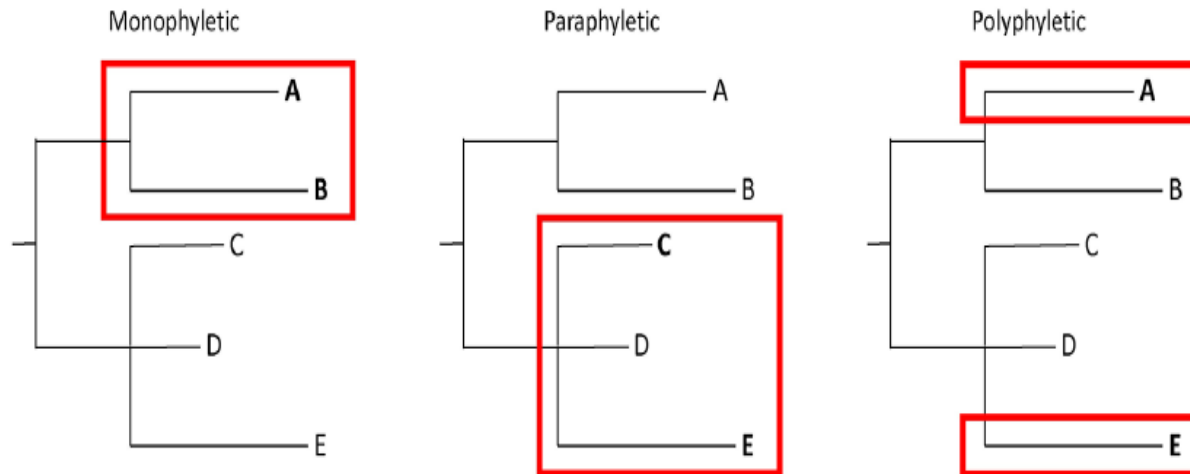
# How to interpret our phylogeny



FIGURE 2 | Illustration of monophyletic, paraphyletic, and polyphyletic groupings. A monophyletic topology exists when isolates of interest (e.g., A and B) group together to the exclusion of all others. A paraphyletic topology is one in which isolates of interest (e.g., C and E) group together but not to the exclusion of all others (e.g., D). A polyphyletic topology exists when isolates of interest do not form a group (e.g., A and E).

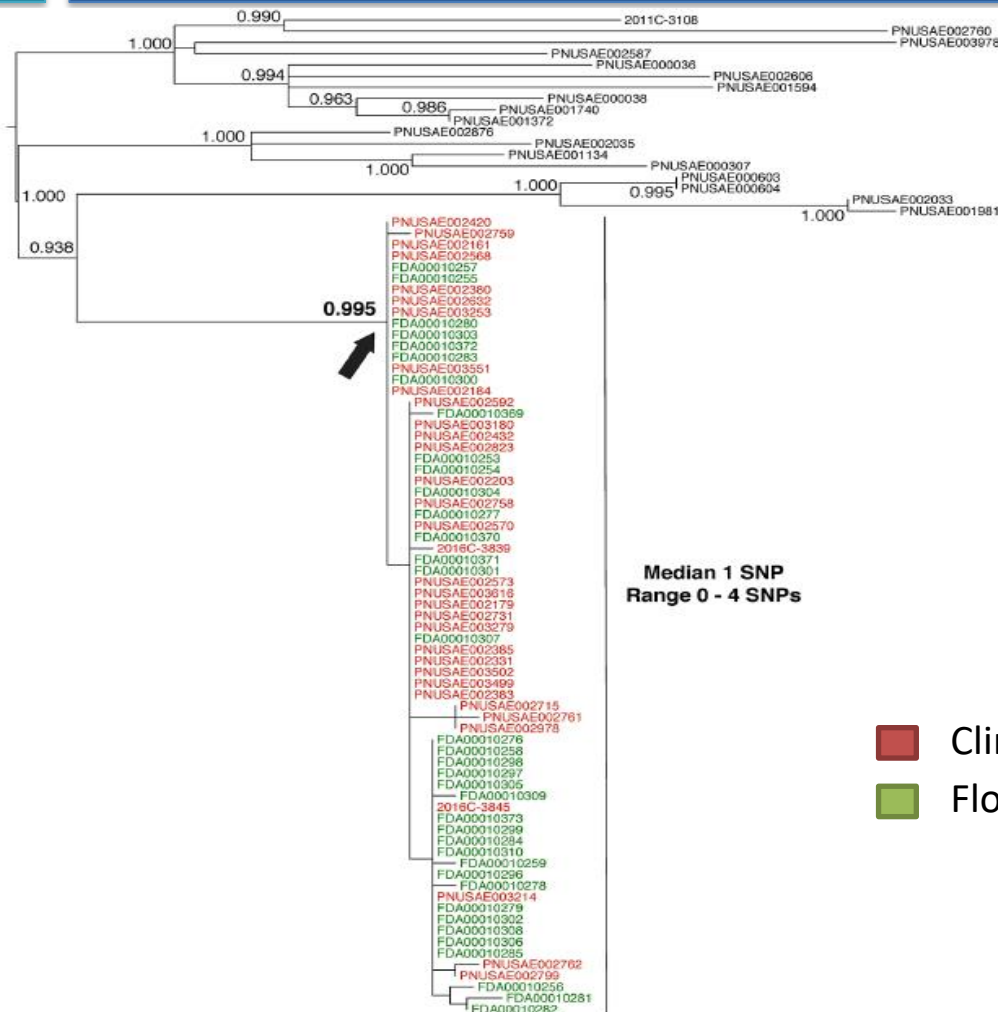Pightling et al. Frontiers in Microbiology. 2018

# How to interpret our phylogeny

**TABLE 2 |** Conditions used to determine whether whole-genome sequence analyses support a match between two or more genomes.

|  | Supports | Neutral | Does not support |
|---|---|---|---|
| SNP distance | <21 | 21–100 | >100 |
| Bootstrap support | >0.89 | 0.80–0.89 | <0.80 |
| Tree topology | Monophyletic | Paraphyletic | Polyphyletic |

Pightling et al. Frontiers in Microbiology. 2018

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# How to interpret our phylogeny



- Bootstrap support
- SNP count support
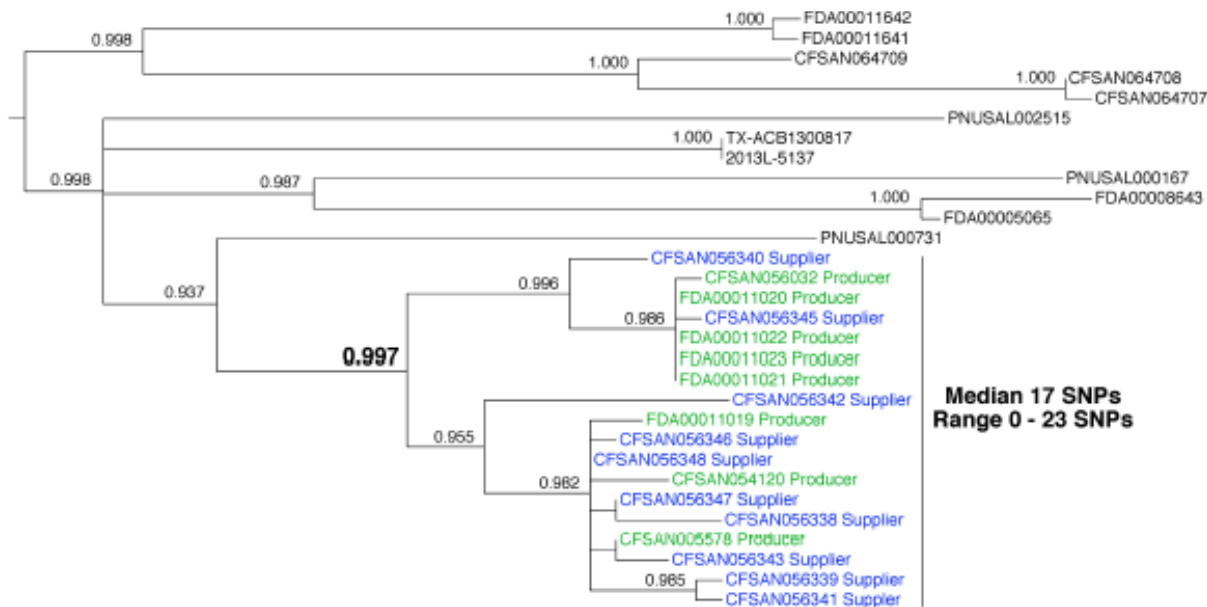- Topology support
- Epidemiology support

E. coli Clinical isolates - source

■ Clinical isolates
■ Flour isolates

Pightling et al. Frontiers in Microbiology. 2018

# How to interpret our phylogeny

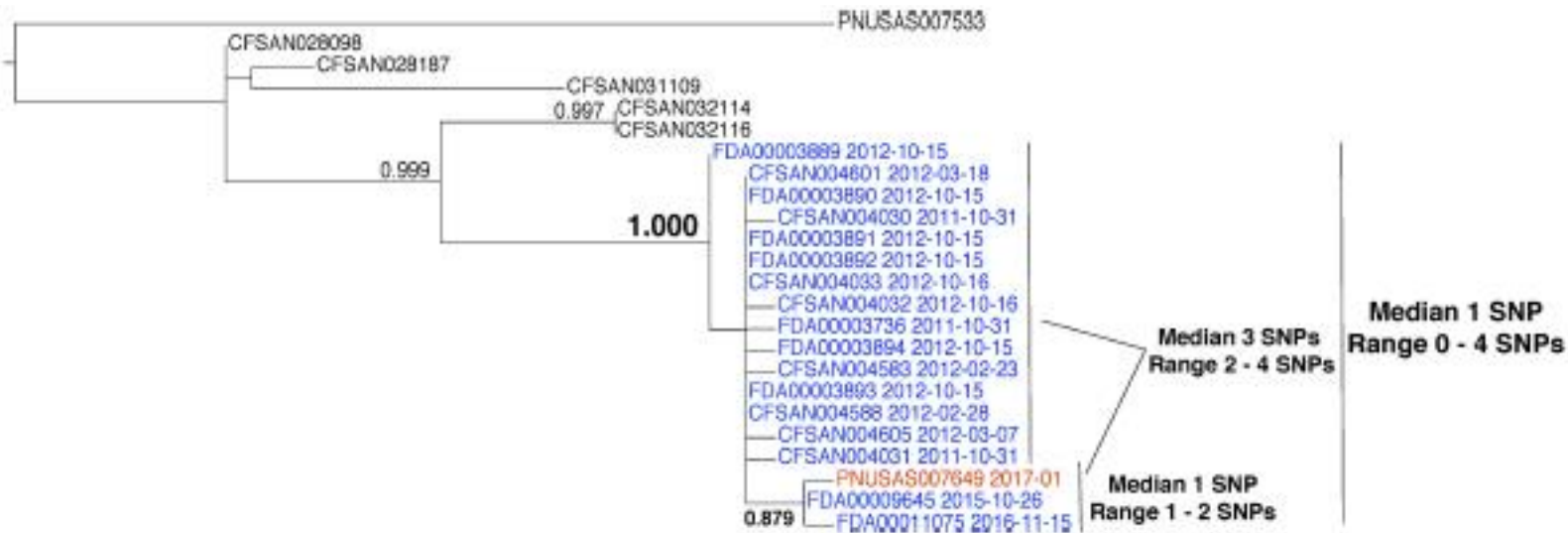L. monocytogenes ingredient supplier – Ice cream producer



- Bootstrap support
- SNP count support
- Topology support
- Epidemiology support

Pightling et al. Frontiers in Microbiology. 2018

# How to interpret our phylogeny

Resident pathogen



Pightling et al. Frontiers
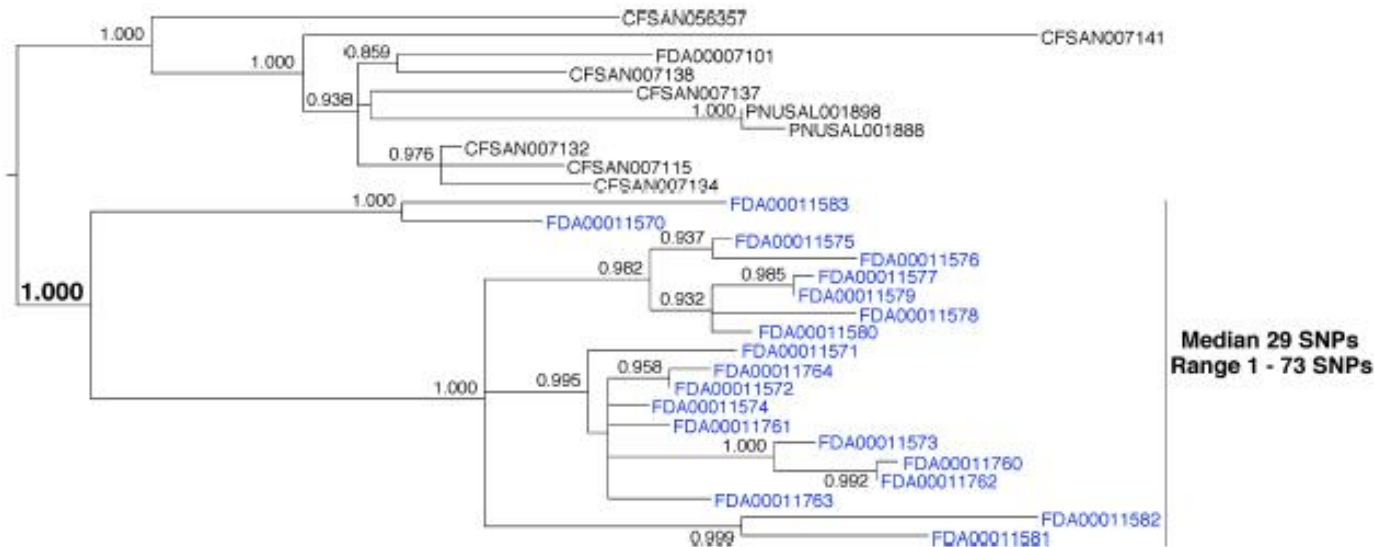in Microbiology. 2018

- Bootstrap support
- SNP count support
- Topology support
- Epidemiology not support

🟥 Clinical isolates

🟦 Environmental isolates

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# How to interpret our phylogeny



Pightling et al. Frontiers in Microbiology. 2018

Median 29 SNPs
Range 1 - 73 SNPs

- Bootstrap support
- SNP count neutral
- Topology support

Clinical isolates

Environmental isolates

Environmental isolates from an inspection.

# How to interpret our phylogeny



- Bootstrap neutral
- SNP count support
- Topology neutral
- Epidemiology Not support

Median 3 SNPs
Range 0 - 7 SNPs

■ Clinical isolates

■ Food isolates

■ Environmental isolates

Pightling et al. Frontiers
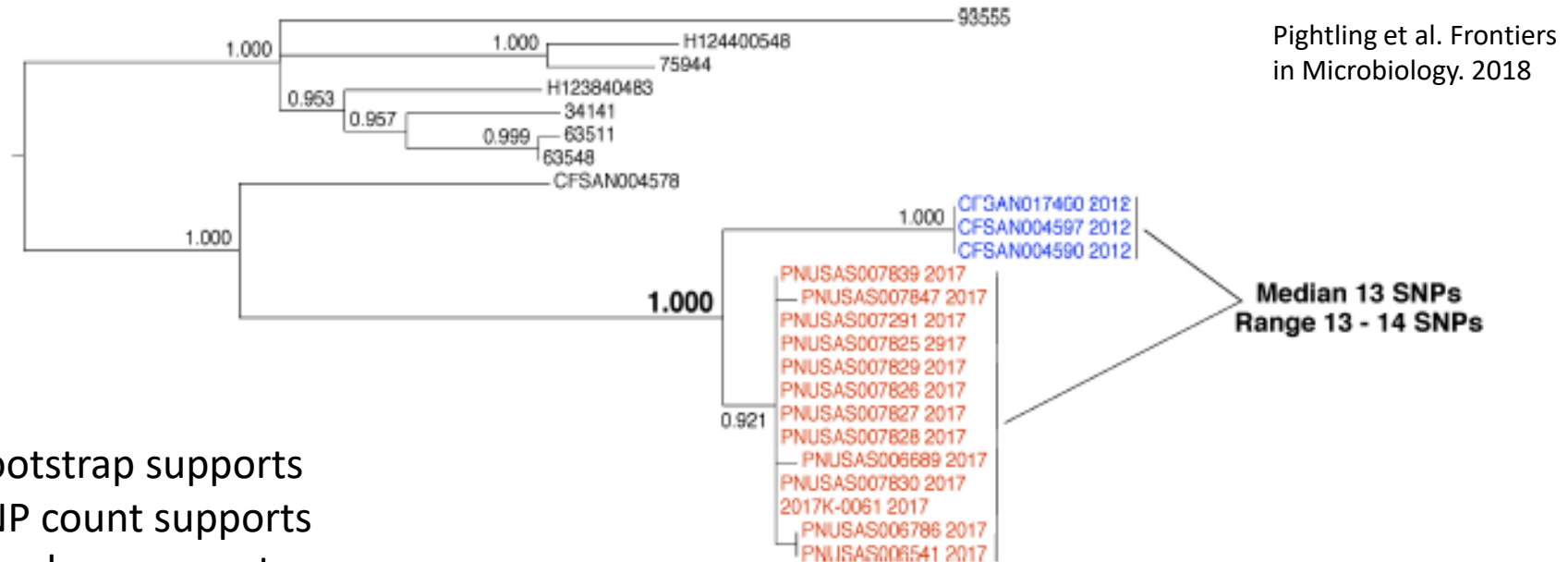in Microbiology. 2018

# How to interpret our phylogeny



Pightling et al. Frontiers in Microbiology. 2018

- Bootstrap supports
- SNP count supports
- Topology supports
- Epidemiology does not support

Clinical isolates
Environmental isolates

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# How to interpret our phylogeny

TABLE 3 | Characteristics of the examples presented in this paper.

| Example | SNP distance | Bootstrap support | Tree topology | Epidemiology, traceback, or compliance findings | Conclusion |
|---|---|---|---|---|---|
| Identifying the source of an *E. coli* outbreak | Supports | Supports | Supports | Supports | Match |
| Matching food isolates from one firm to environmental isolates from another firm | Supports | Supports | Supports | Supports | Match |
| Identifying a resident pathogen | Supports | Supports | Supports | Not applicable | Not applicable |
| Populations of environmental isolates can be very diverse | Neutral | Supports | Supports | Not applicable | Not applicable |
| Analyzing paraphyletic relationships | Supports | Neutral | Neutral | Does not support | No match |
| Evidence that isolates arose from the same source by WGS does not necessarily mean that they are linked | Supports | Supports | Supports | Does not support | No match |

Pightling et al. Frontiers in Microbiology. 2018

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# Thanks for your attention!

Secuenciación de genomas bacterianos:
herramientas y aplicaciones