

Secuenciación de Genomas Bacterianos: Herramientas y Aplicaciones

BU-ISCIII

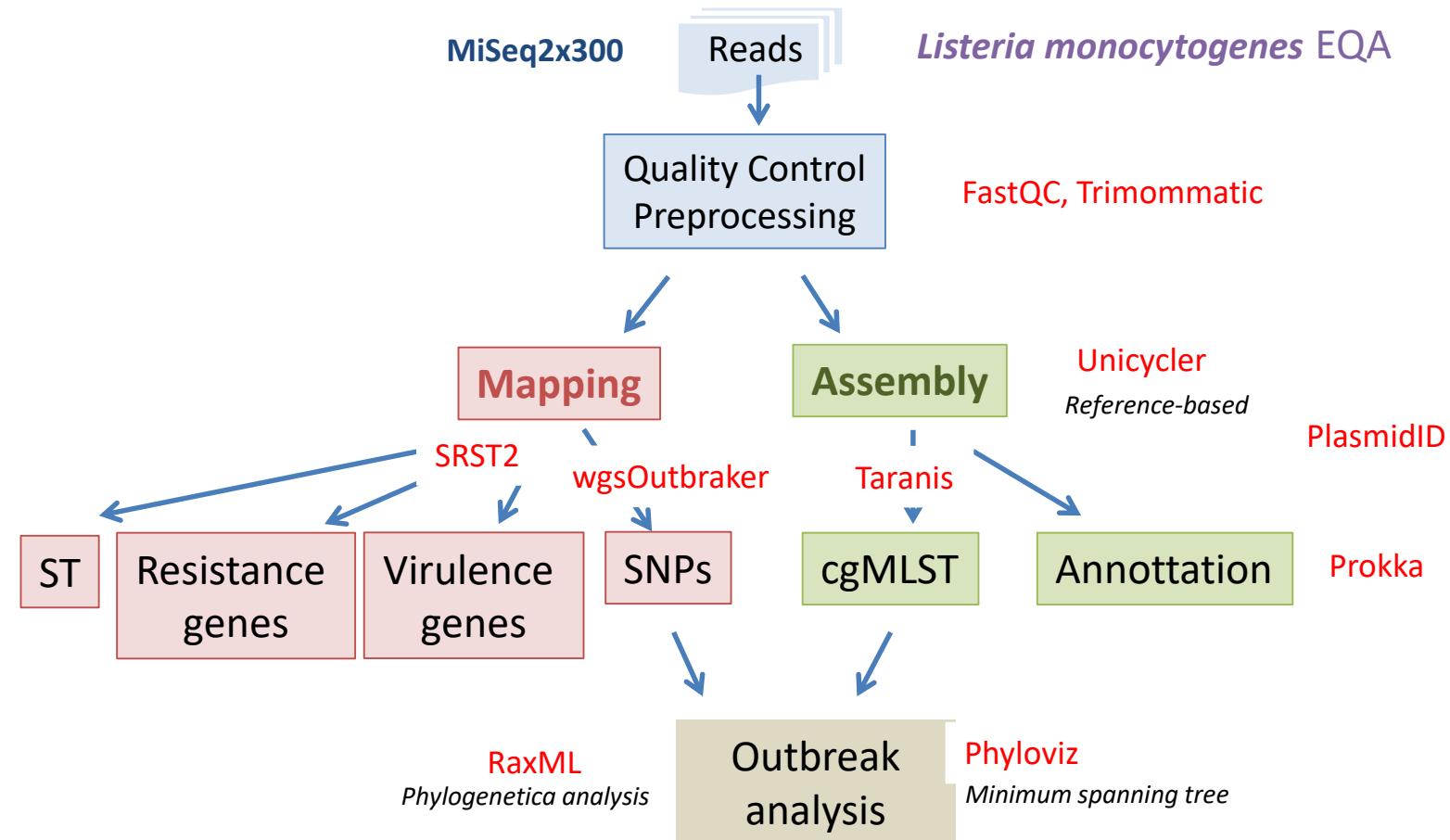
Unidades Centrales Científico Técnicas - SGSAFI-ISCIII

28 Junio - 2 Julio 2021, 3^a Edición
Programa Formación Continua, ISCIII

Learning aims and outcomes

- Understand some principles behind NGS and its applications to whole genome sequencing.
- Know the format files generated in NGS data analysis and the workflow analysis.
- Understand the uses of WGS in: specie, antimicrobial resistance genes and virulence factor genes identification, and for typing.
- Outbreak characterization based on SNPs or gene by gene approaches.

Training workflow



Teachers

- **Sara Monzón Fernández**, Biotecnóloga y Bioinformática (Analista de datos). Titulado Superior Especialista OPIS. Responsable técnico BU-ISCIII
- **Sarai Varona Fernández**, Bioquímica y Bioinformática (Analista de Datos). Contrato Titulado Superior asociado a proyecto (2021-2022)
- **Isabel Cuesta**, Dra Biología, Bioinformática (Científico de Datos). Científico Titular de OPIS. Coordinador BU-ISCIII

Session 1.1 - Secuenciación masiva de genomas bacterianos: situación actual

Isabel Cuesta

BU-ISCIII

Unidades Centrales Científico Técnicas - SGSAFI-ISCIII

28 Junio - 2 Julio 2021, 3^a Edición
Programa Formación Continua, ISCIII

Index

- BU-ISCI
- High throughput sequencing platforms update
- Bacterial genome sequencing, brief history
- Advantages of WGS
- Use of WGS in Europe

Qué es la Bioinformática?

PROBLEMAS
BIOLÓGICOS



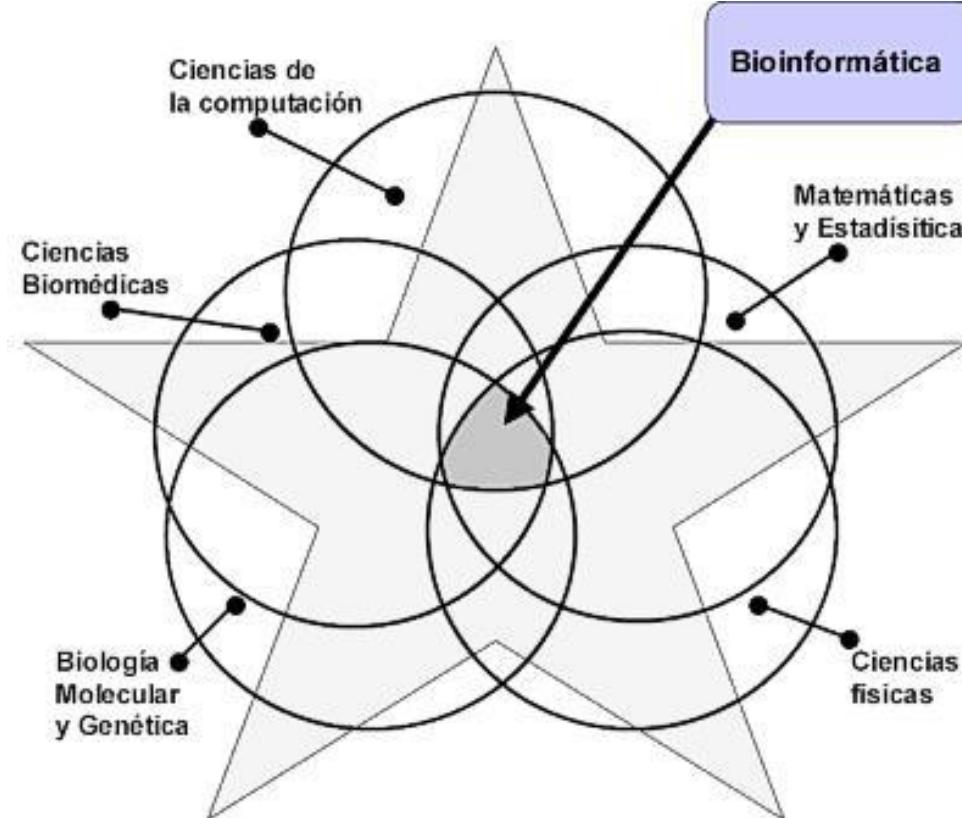
Procesamiento
de datos



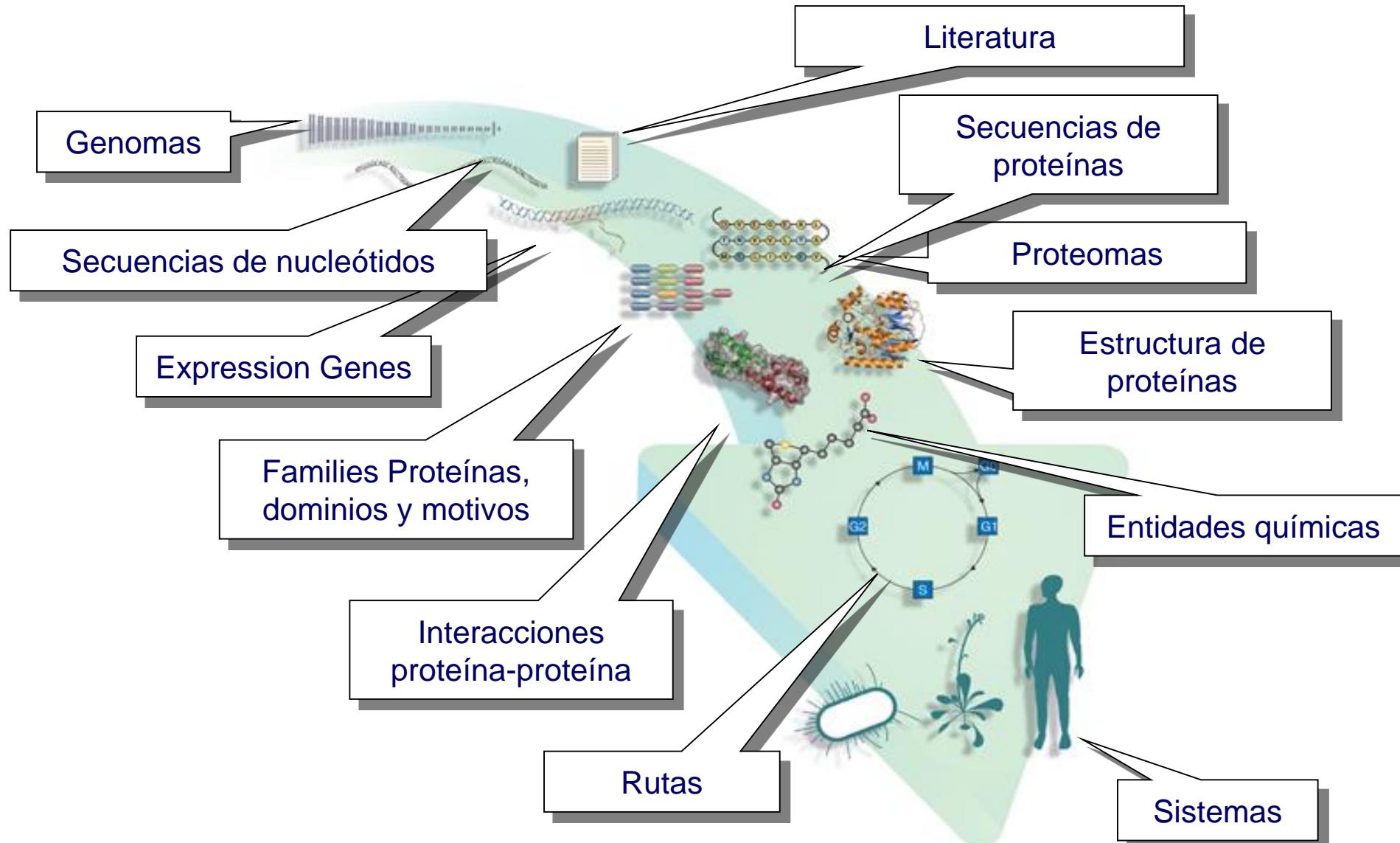
MÉTODOS
COMPUTACIONALES



Bioinformática es multidisciplinar



Tipos de datos dan idea de la dimensión de la Bioinformática



Why BU-ISCIII was founded

>_BU-ISCIII

Genomics Unit

2010



454



NextSeq500



MiSeq

2013



Minlon -
Nanopore

2019



NovaSeq
6000

2021

Bioinformatics Unit

2012



Service &
Support to
Researchers
on
HTS Data
Analysis



National
Microbiology
Centre (CNM)



Research Institute
for Rare Diseases
(IIER)



Functional Unit
for Research in
Chronic Disease



Network of Biological
Alerts

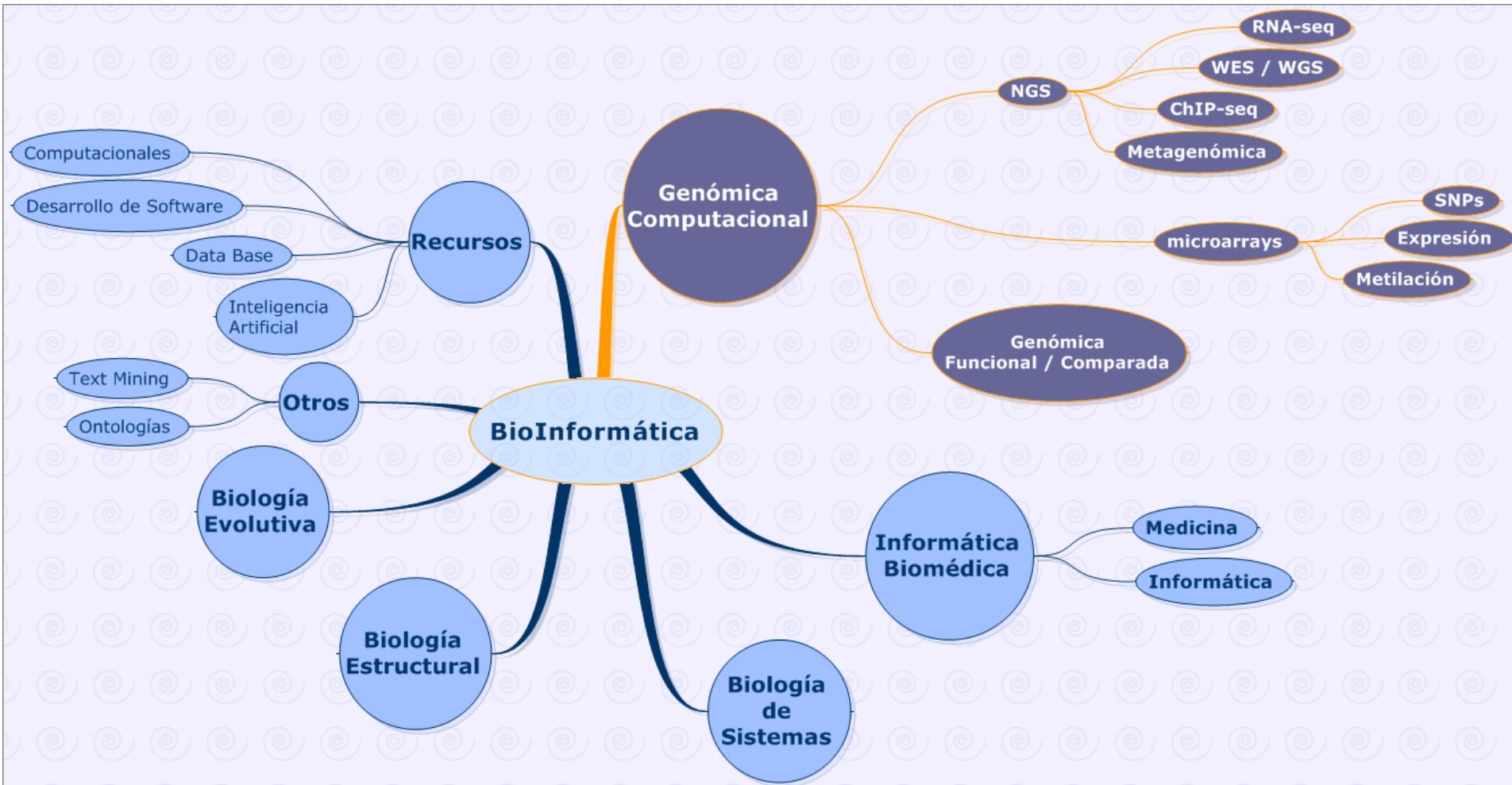


National Centre of
Tropical Medicine

National Environment
Health Centre



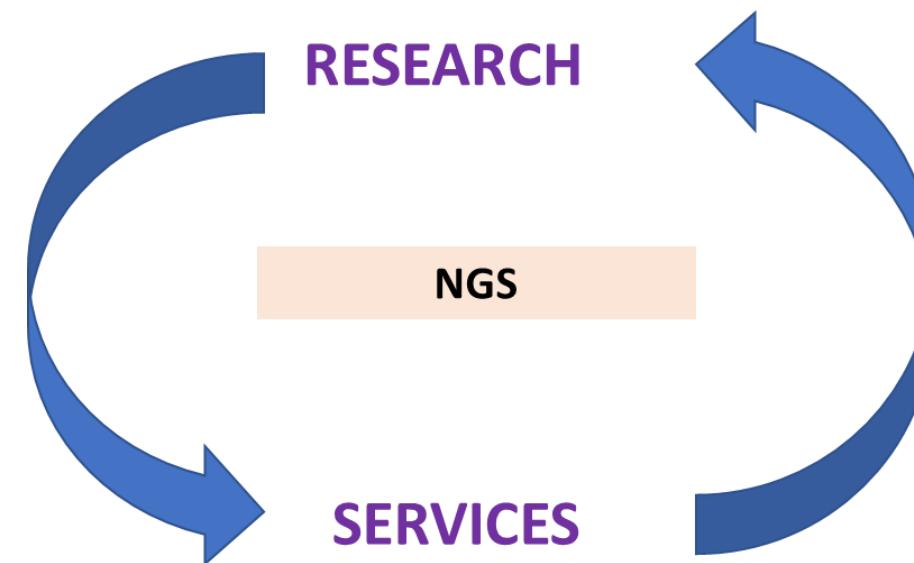
BU-ISCIII Mission - Activities



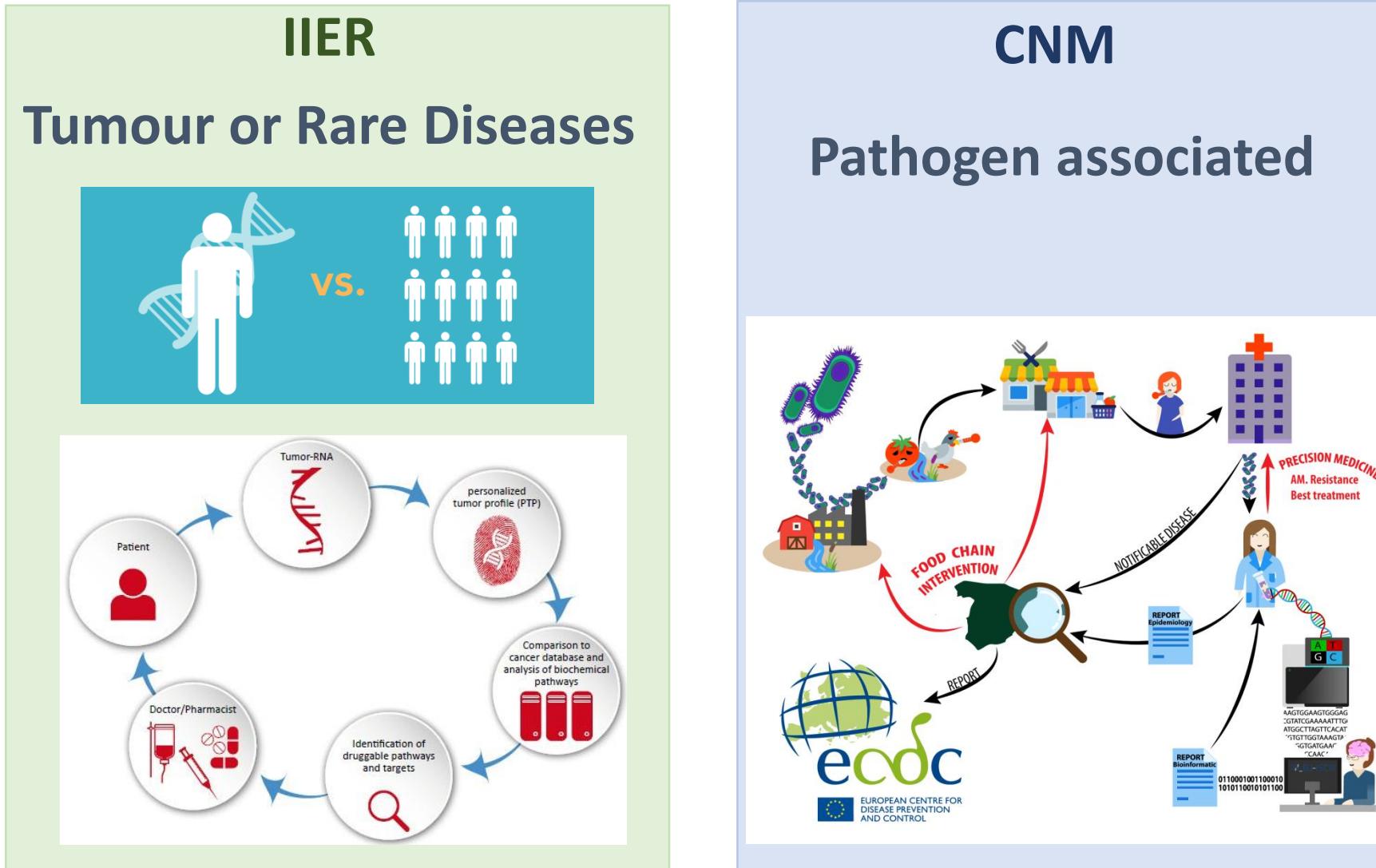
Bioinformatics Unit Activities

- Identify biological problems (PI / Groups) that could be target of NGS
- Early adopters: establish collaboration with.
- Be strategic providing transversal solutions → reusable tools

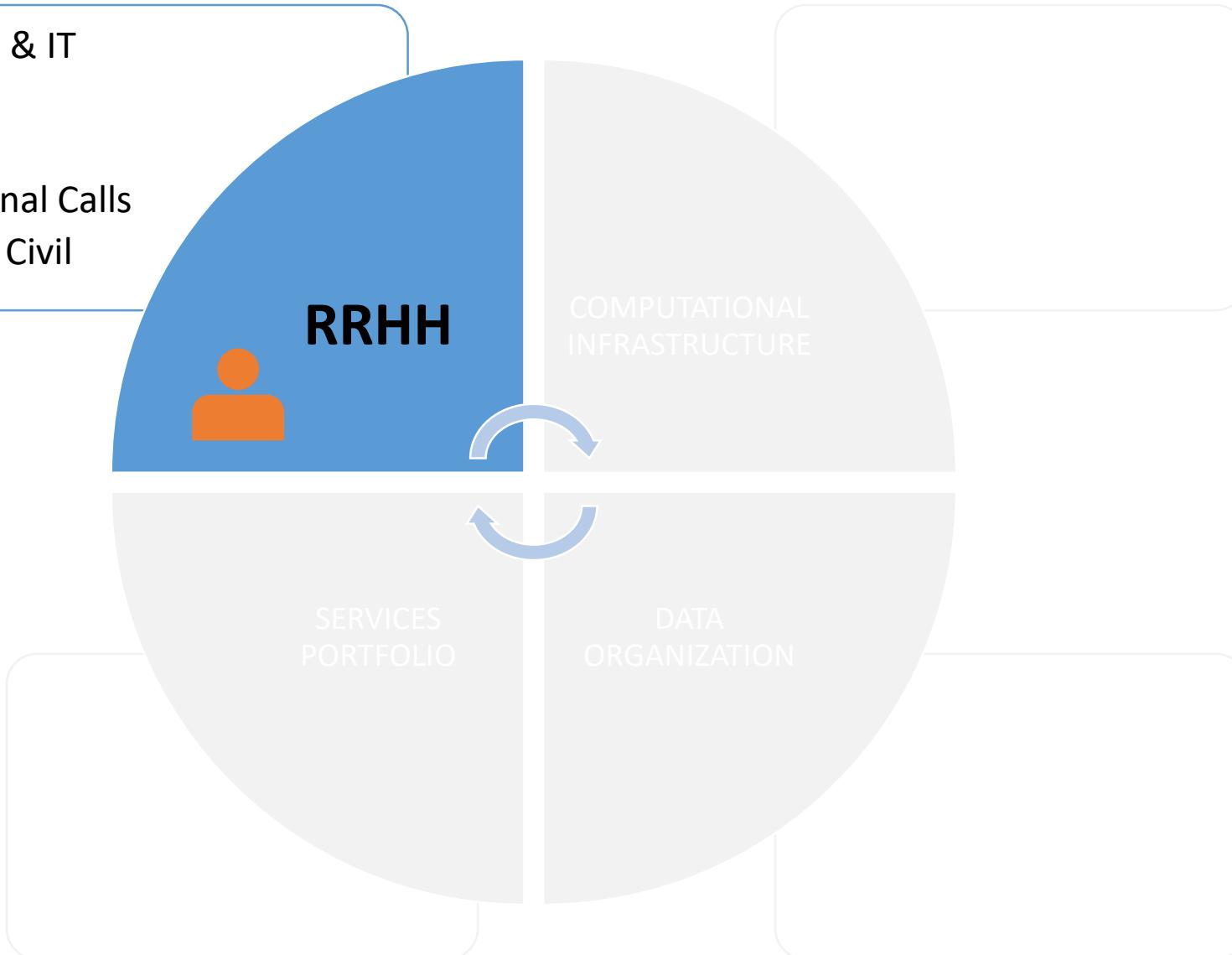
Provide scientific and technical solutions for using NGS in the diagnostic routine or research activity from different ISCIII labs



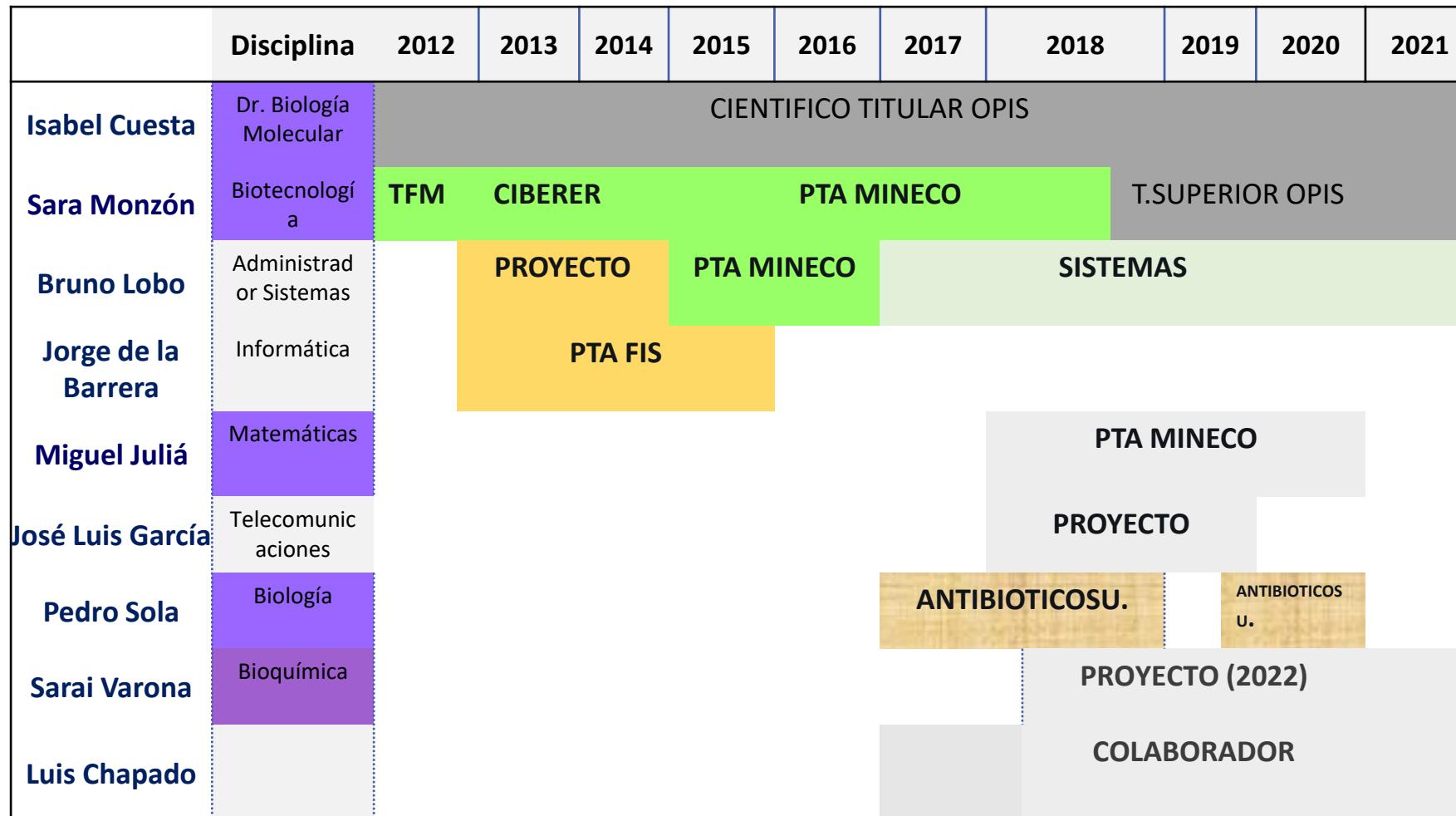
Clinical Bioinformatics - Precision Medicine



- TEAM: Bioinformatician & IT
- SOURCE OF FUNDING:
 - Research Project
 - National or International Calls
 - Permanent position – Civil Servant



Human resources



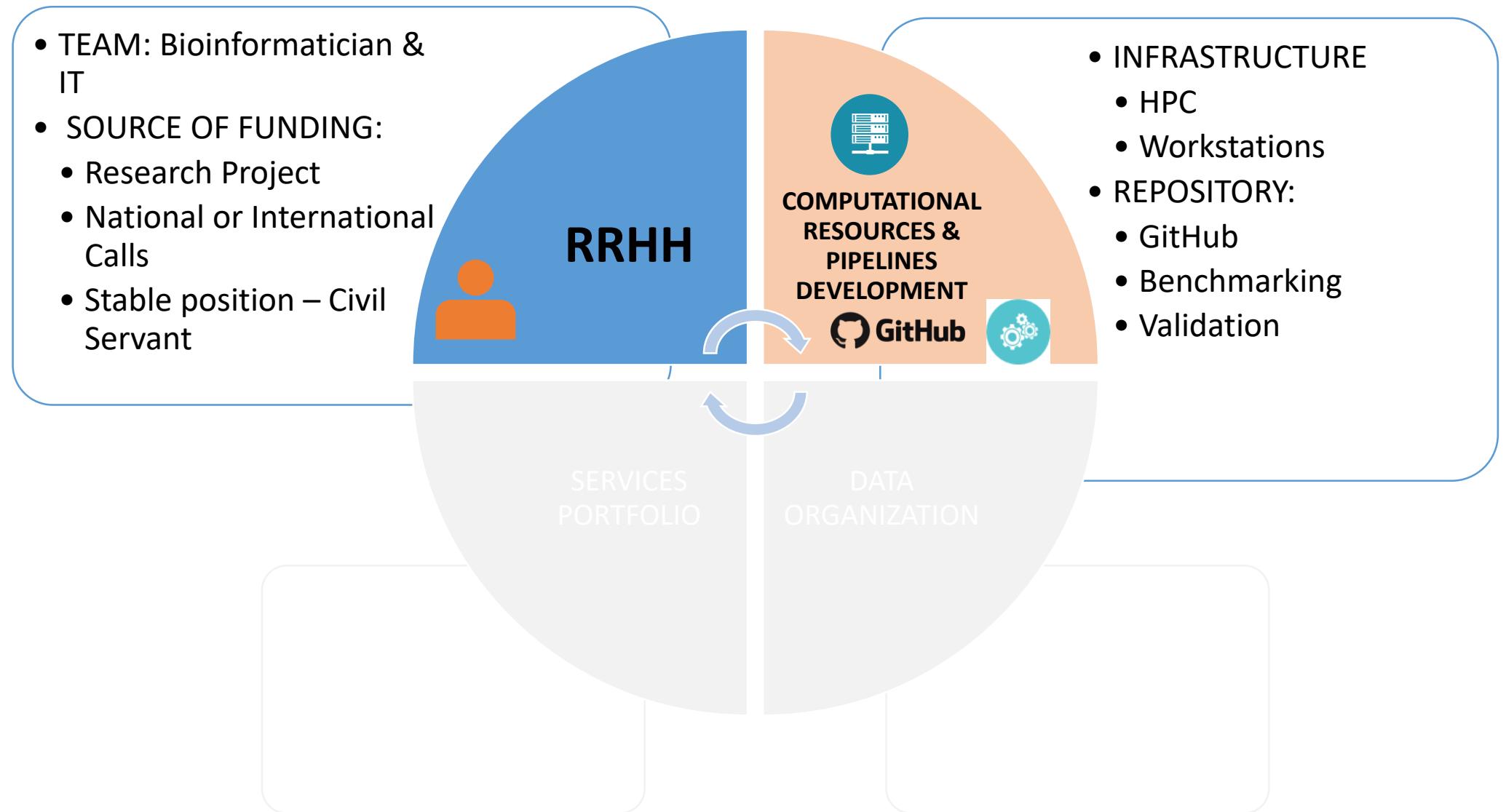
Master en
Bioinformática

CNM

IIER

BU-ISCIII

FUNCIONARIO



Computational Resources

- IT support: establish agreement with IT department including permission for using Linux.



Workstations (5), 4cores, 64Gb, 8TB
Server, 4-quad, 120Gb, 16TB

Data Centre (CPD-ISCIII)



HPC 320 cores, 8TB RAM, 10Gbps.
2 flexible and scalable storages,
NetApp, 70 TB and 250TB

- Reproducibility of in-silico pipelines analysis

nextflow



Singularity containers
Admin support & environment independency
Sharing code easier

GitHub

<https://github.com/BU-ISCIII>

Bioinformatic Analysis: Software validation: ECDC EQAs

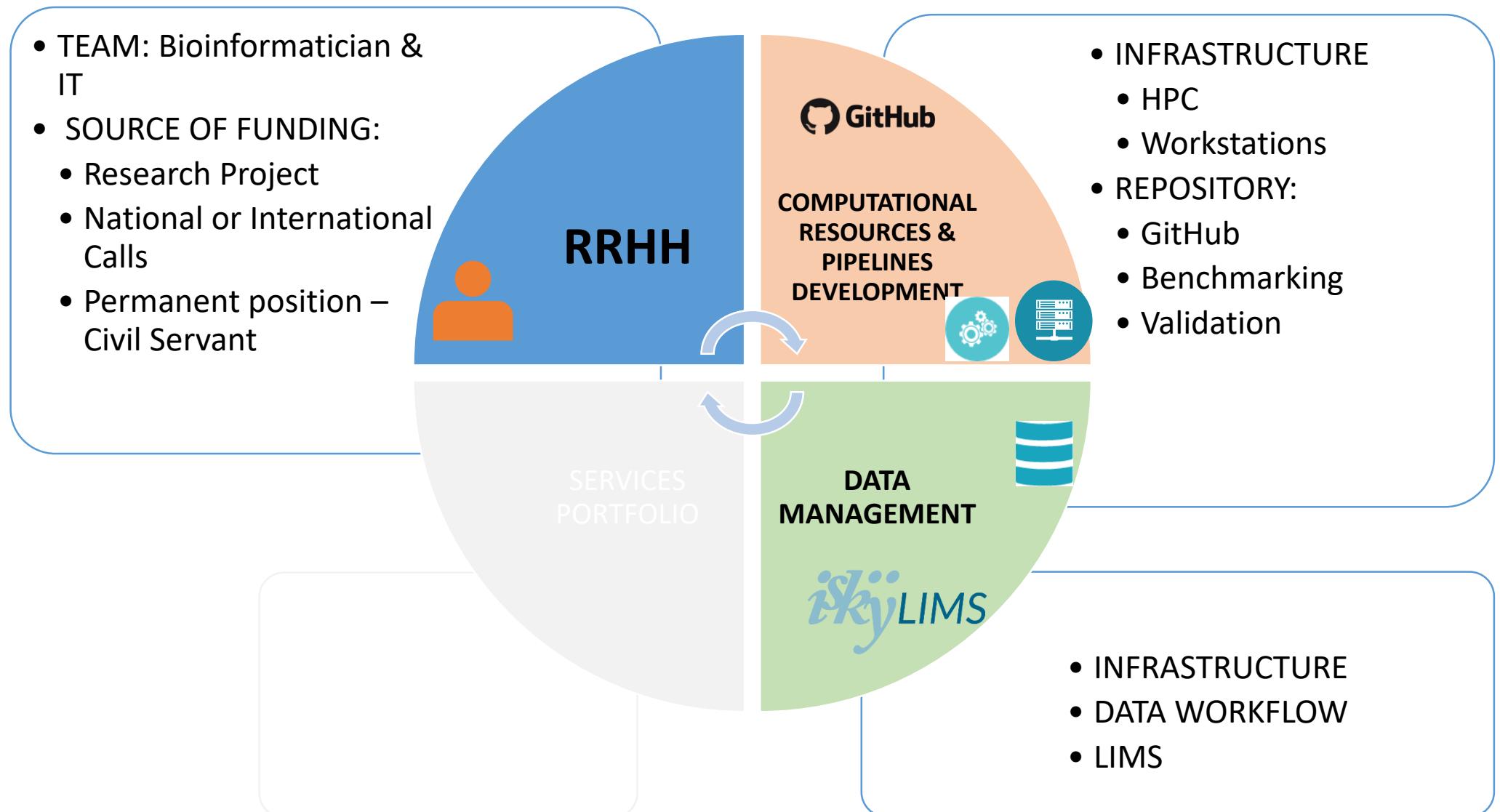
Table 5. Results of allele-based cluster analysis

| Lab ID | Approach | Allelic calling method | Allele based analysis | | | |
|--------------|-------------|-----------------------------|--------------------------|-------------------------------|---------------------------|----------------------------|
| | | | Assembler | Scheme | Difference within cluster | Difference outside cluster |
| EQA provider | BioNumerics | Assembly- and mapping-based | SPAdes | Applied Math (cgMLST/Pasteur) | 0-3 | 24-1112 |
| 19 | BioNumerics | Assembly- and mapping-based | SPAdes | Applied Math (cgMLST/Pasteur) | 0-3 | 25-1120 |
| 35 | SeqSphere | Assembly-based only | Velvet | Ruppitsch (cgMLST) | 0-2 | 16-1065 |
| 70 | SeqSphere | Assembly-based only | Velvet | Ruppitsch (cgMLST) | 0-2 | 16-1062 |
| 105* | SeqSphere | Assembly-based only | SPAdes v 3.80 | Ruppitsch (cgMLST) | 0-1* | 23-812 |
| 129 | SeqSphere | Assembly-based only | Velvet | | | |
| 135 | SeqSphere | Assembly-based only | CLC Genomic Workbench 10 | | | |
| 141 | SeqSphere | Assembly-based only | SPAdes 3.9.0 | | | |
| 142 | Inhouse | Assembly-based only | SPAdes | | | |
| 144 | SeqSphere | Assembly-based only | Velvet | | | |

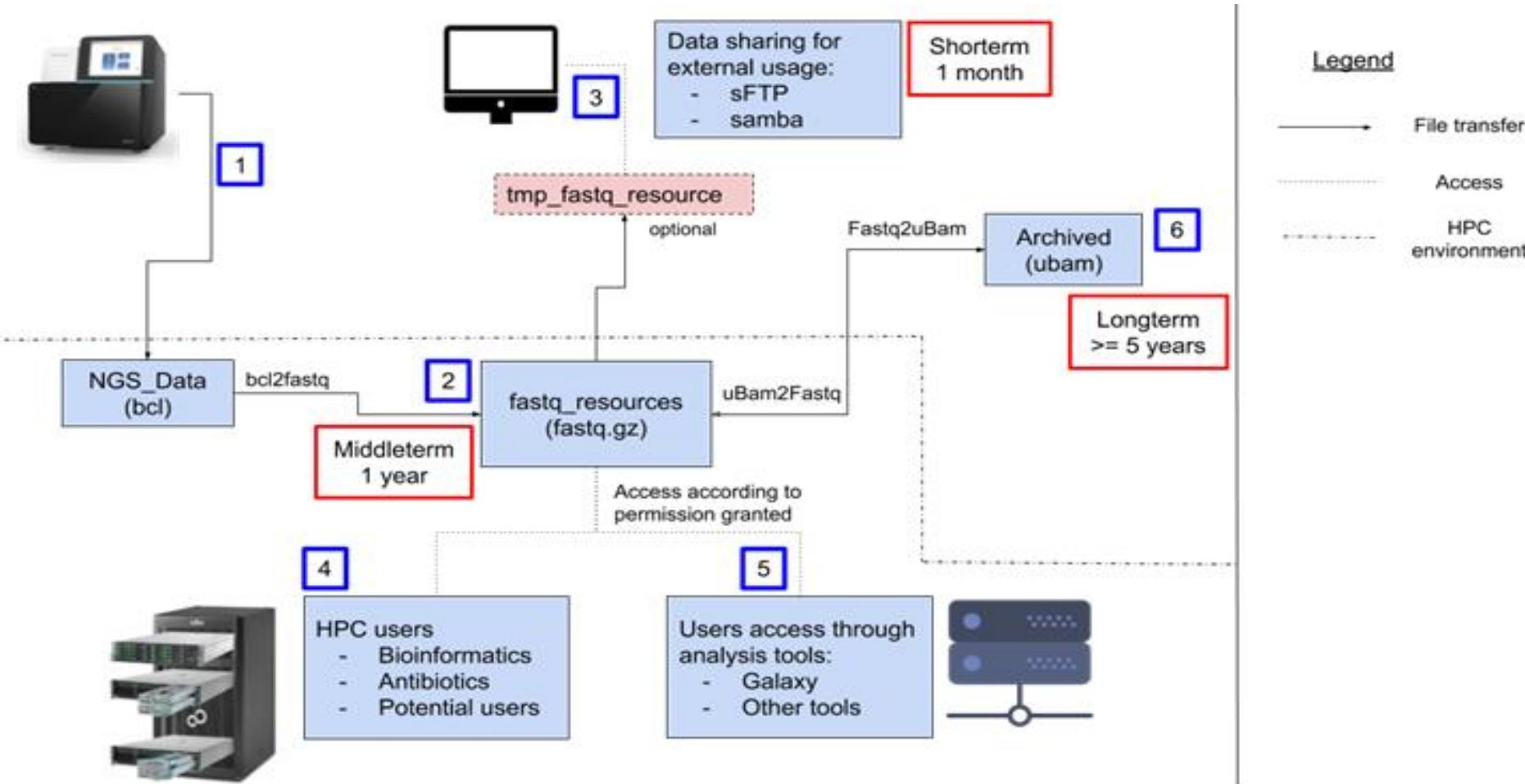
Table 4. Results of SNP-based cluster analysis

| Lab ID | SNP-based | | | | | | |
|----------|-----------------|--|--------------------------|--------------------------|-----------|-------------------------|--------------------------|
| | Approach | Reference | Read mapper | Variant caller | Assembler | Distance within cluster | Distance outside cluster |
| Provider | Reference-based | ST6 (REF4) | BWA | GATK | | 0-3 | 38-71 |
| 19* | Reference-based | ST6 ID 2362 | BWA | GATK | | 0-4 | 43-81 |
| 56 | Assembly-based | | | ksnp3 | SPAdes | 0-57* | 561-591 (6109) |
| 105 | Reference-based | ST6 J1817 | Bowtie2 | VARSCAN 2 | | 0-2* | 22-42 (1049) |
| 108 | Reference-based | In-house strain resp ST | CLC assembly cell v4.4.2 | CLC assembly cell v4.4.2 | | 0-2 | 37-72 |
| 142* | Reference-based | Listeria EGDe (cc9) | CLC Bio | CLC Bio | | 0-1219 | 1223-2814 (8138) |
| 146 | Reference-based | ST6 ref. CP006046 ST1 ref. F2365 ST213/ST382 no ref. | BWA | In-house | | 0-358 | |

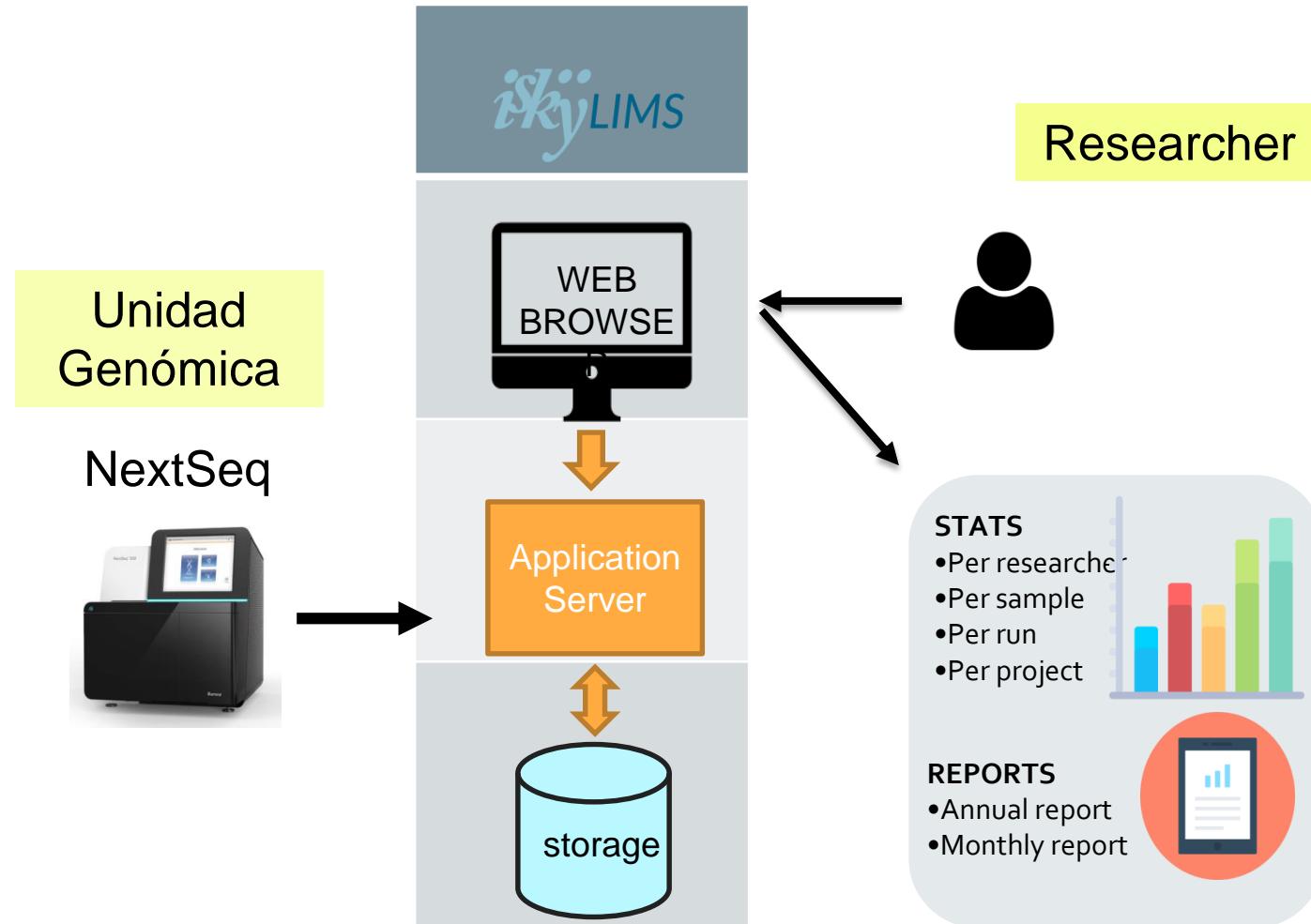
Fifth external quality assessment scheme for Listeria monocytogenes typing



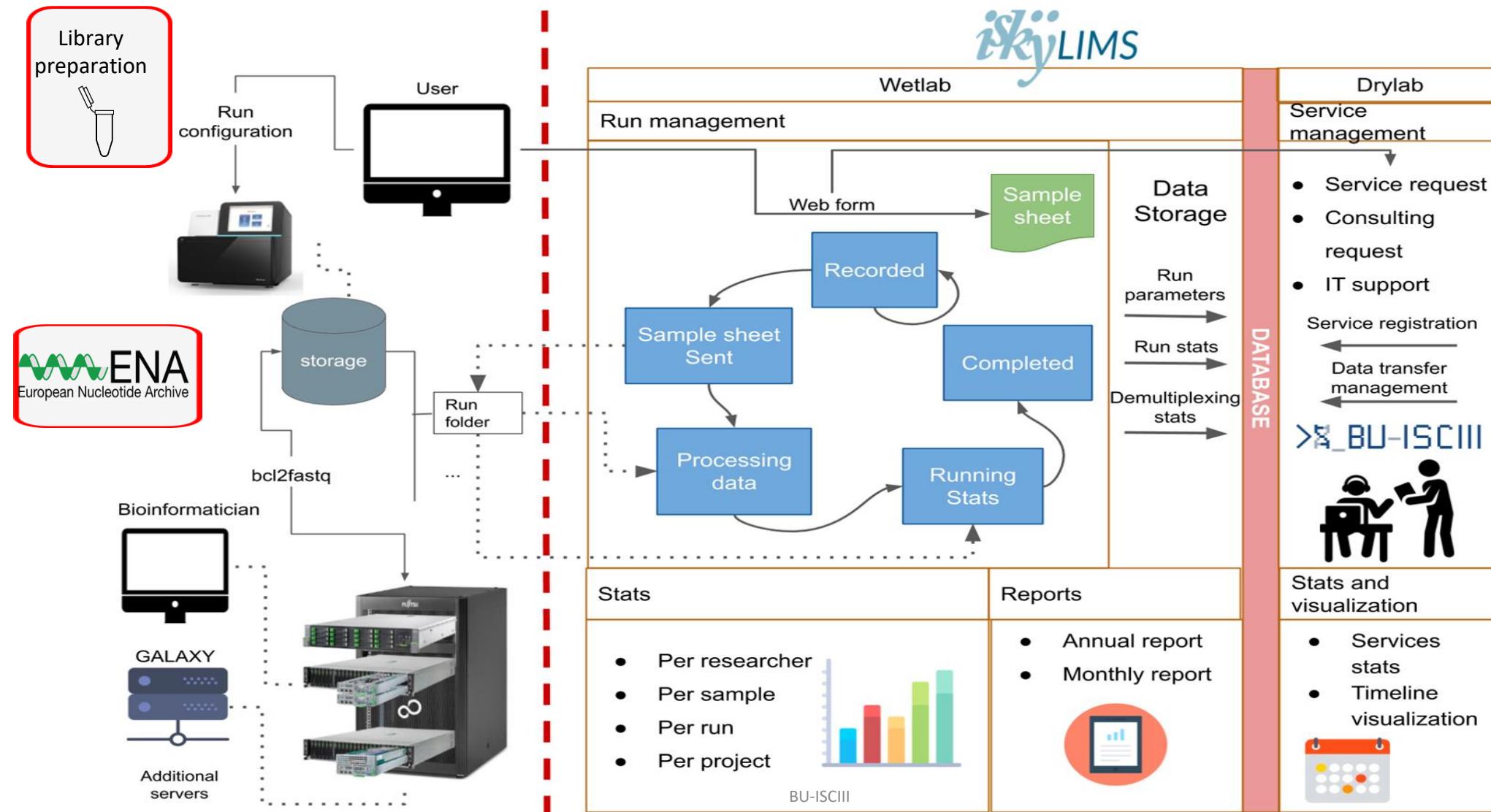
Infrastructure and data management



- Minimize I/O issues
- Maximize storage uses



Infrastructure and data management: LIMS



SERVICIOS DE LA UNIDAD

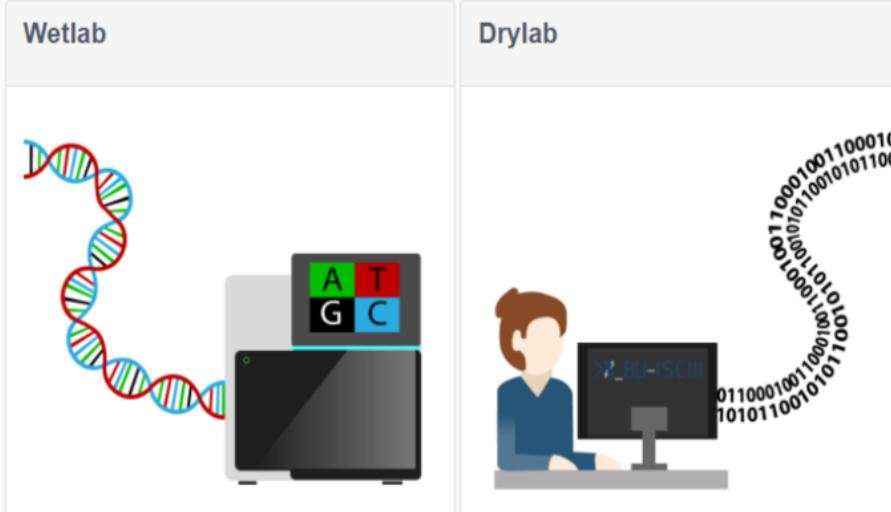


<https://iskylims.isciii.es/>



HOME ABOUT US TUTORIALS FAQS REGISTER CONTACT

 icuesta  [Login](#)



Logos



Connect



Links

- Contact
 - Getting started
 - FAQs

Sitemap

- iSkyLIMS home
 - Drylab page
 - Wetlab page

<https://iskylims.isciii.es/>

 smonzon [Logout](#) [My account](#)

BioInformatics

iSkyLIMS: DryLab

Welcome

This section will allow you to check BU-ISCIII service activity. Available processes are request new services, collaborations, counseling and infrastructure. You will be able to check the status of your ongoing services.



Services ongoing and queued

Under construction. This will be a table with services ongoing or queued

Timeline of services

Under construction. Kind of diagram with services dates.



Service Request Form

Form for requesting internal service to Bioinformatic Unit

Sequencing Data

User's projects*

BMartinez20161213
EXOMAS_ND_20170303
EXOMAS_ND_20170327_RE
EXOMAS_ND_20170228

Run specifications**File extension****Sequencing platform**

SERVICES REQUEST

[HOME](#)[SERVICES REQUEST](#)[COUNSELING REQUEST](#)[INFRASTRUCTURE REQUEST](#)

- Genomic Data Analysis
 - Download and quality analysis
 - Data download
 - Sequence quality analysis
 - Sequence pre-processing (quality filtering)
 - Next Generation Sequencing data analysis
 - DNAseq: Exome sequencing (WES) / Genome sequencing (WGS) / Target sequencing
 - Trio/family variant calling pipeline
 - Variant calling and annotation pipeline
 - Microbial: Whole genome outbreak analysis pipeline
 - Microbial: wgMLST
 - Microbial: MLST + virulence + AMR + plasmid analysis
 - Microbial: Assembly + automatic annotation
 - Microbial: plasmidID pipeline - strain plasmid characterization
 - RNAseq: Transcriptome sequencing
 - miRNA-Seq pipeline
 - mRNA-Seq pipeline
 - Amplicon sequencing (Deep sequencing)
 - Low frequency variant detection
 - Viral: assembly and minor variants detection
 - Metagenomics
 - 16S taxonomic profiling
 - Shotgun metagenomics profiling
 - Shotgun metagenomics - Virus genome reconstruction
 - CHIP-SEQ
 - Peak detection and annotation

SERVICES REQUEST



Service Description

Service description file*
 No file selected.

Service Notes*

COUNSELING REQUEST



Service selection

Available Services *

- Bioinformatics consulting and training
 - Bioinformatics analysis consulting
 - In-house and outer course organization
 - Student training in collaboration: Master thesis, research visit,...

Service Description

Service description file*

No file selected.

Service Notes*

INFRASTRUCTURE REQUEST

[HOME](#)[SERVICES REQUEST](#)[COUNSELING REQUEST](#)[INFRASTRUCTURE REQUEST](#)

Form for requesting Infrastructure service to Bioinformatic Unit

Service selection

Available Services *

- User support
 - Installation and support of bioinformatic software on Linux OS
 - Installation and access to Virtual machines in the Unit server containing bioinformatic software
 - Code snippets development
 - OT-2 robots

Service Description

Service description file

Ningún archivo seleccionado

Service Notes *

Infrastructure and data management

 [HOME](#) [RUN PREPARATION](#) [SEARCH](#) [STATISTICS](#) [REPORTS](#)

 bioinfoadm [Logout](#) [My account](#)

Statistics results for Investigator rabad

Projects using the sequencer NS500454 :

[Export Table To Excel](#)

| Project name | Date | Library Kit | Samples | Cluster PF | Yield Mb | % Q> 30 | Mean | Sequencer ID |
|---------------------------------|---------|--|---------|-------------|----------|---------|-------|--------------|
| NextSeq_CNM_191_20191004_RAbad | No Date | Nextera DNA CD Indexes (96 Indexes plated) | 48 | 149,441,968 | 45,876 | 89.98 | 33.70 | NS500454 |
| NextSeq_CNM_166_20190528b_Rabad | No Date | Nextera XT v2 Set B | 96 | 139,317,411 | 43,016 | 89.58 | 33.72 | NS500454 |
| NextSeq_CNM_166_20190528a_Rabad | No Date | Nextera XT v2 Set A | 82 | 102,267,350 | 31,623 | 89.26 | 33.65 | NS500454 |
| NextSeq_CNM_150_20190218B_RAbad | No Date | Nextera XT v2 set B | 20 | 17,335,577 | 5,352 | 86.77 | 33.17 | NS500454 |
| NextSeq_CNM_150_20190221A_RAbad | No Date | Nextera XT v2 Set A | 96 | 127,755,164 | 39,595 | 85.28 | 32.86 | NS500454 |
| NextSeq_CNM_166_20190528c_Rabad | No Date | Nextera XT v2 Set C | 96 | 152,945,860 | 47,264 | 89.38 | 33.68 | NS500454 |
| NextSeq_CNM_170_20190620_RAbad | No Date | IDT-ILMN Nextera UD Index Set A for Nextera DNA FI | 47 | 131,012,486 | 39,671 | 90.74 | 33.94 | NS500454 |
| NextSeq_CNM_171_20190624_RAbad | No Date | IDT-ILMN Nextera UD Index Set A for Nextera DNA FI | 47 | 140,488,964 | 42,597 | 89.61 | 33.72 | NS500454 |

- TEAM: Bioinformatician & IT
- SOURCE OF FUNDING:
 - Research Project
 - National or International Calls
 - Permanent position – Civil Servant

RRHH



COMPUTATIONAL RESOURCES & PIPELINES DEVELOPMENT



SERVICES PORTFOLIO & TRAINING



- COURSES, TFM, TFG
- DATA ANALYSIS
 - DNAseq
 - RNAseq
 - Metagenomics

- INFRASTRUCTURE
 - HPC
 - Workstations
- REPOSITORY:
 - GitHub
 - Benchmarking
 - Validation

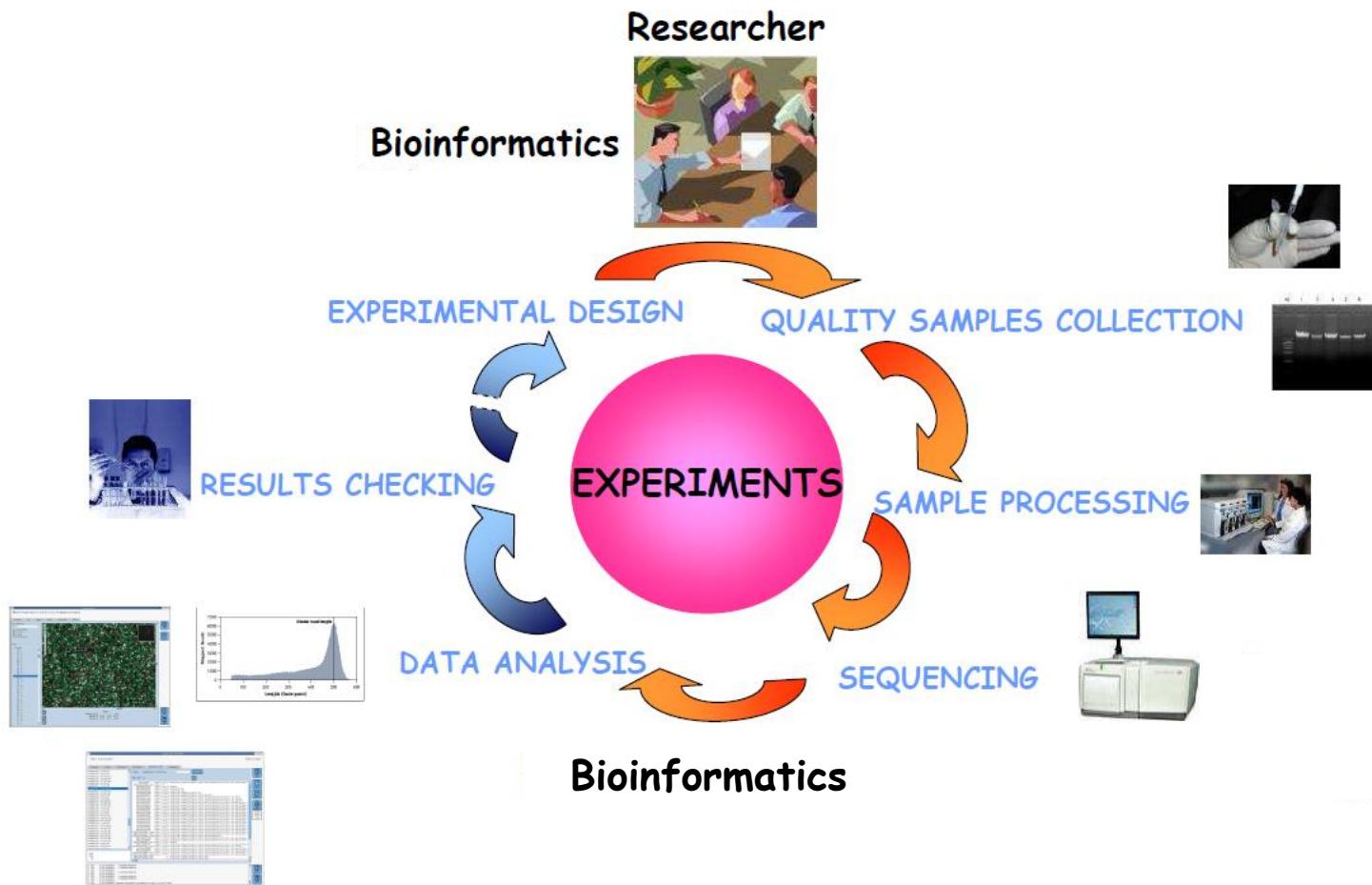
DATA MANAGEMENT



- INFRASTRUCTURE
- DATA WORKFLOW
- LIMS

- **GENÓMICA COMPUTACIONAL: ANÁLISIS DE DATOS MASIVOS**
Técnicas de secuenciación masiva (NGS)
- **ASESORIA Y FORMACIÓN EN BIOINFORMÁTICA**
Orientación en el análisis bioinformático
Organización de cursos internos y externos
- **SOPORTE A USUARIOS**
Generación y acceso a máquinas virtuales que contienen software bioinformático, ubicadas en los servidores de la Unidad

Workflow en NGS



Services Portfolio

| | | QC | Assembly | Reference based Mapping | Variant calling | Annotation | Pipelines |
|---------------------|------------------------------|------------------|---|-----------------------------------|----------------------------|------------------------------------|---|
| DNaseq | HUMAN | | | | | | |
| | WES Target -Panels | Report html | | (Bam file) | (Vcf file) | Desease model (Vcf file annotated) | .Trio / family .Tumor .Pampu caller |
| RNAseq | MICROBIAL | | | | | | |
| | WGS Amplicon | Report html | <i>De novo</i> / Reference (fasta file) | MLST, Resistance g, Virulence g | SNPs Phylogenetic analysis | Structural Functional | .WGSOutbraker .Plasmid ID |
| Metagenomics | mRNA | RSQC Report html | <i>De novo</i> (fasta file) | Transcripts coverage / expression | Variants (Vcf file) | Transcripts annotation | mRNA seq |
| | miRNA | | | | | | miRNA seq |
| | 16S taxonomic profile | Report html | <i>De novo</i> | Green genes DB | | species diversity | Qiime |
| | Shotgun | | | Genome RefSeq | | Pathogen / Genome coverage | PikaVirus |

Galaxy

172.23.2.60

Galaxy

Analyze Data Workflow Shared Data Visualization Help Login or Register Using 0 bytes

Tools search tools Get Data Send Data Collection Operations Text Manipulation Filter and Sort Join, Subtract and Group Convert Formats Extract Features Fetch Sequences Fetch Alignments Statistics Graph/Display Data MyTools IRMA NGS Data Quality Check Workflows All workflows

Welcome to our Galaxy platform!

This galaxy server has been built and is maintained by the Bioinformatics Unit of Instituto de Salud Carlos III in order to give an user friendly enviroment to run limited bioinformatic tools and data analysis. Contact us if you are interested in the service and want to take an introductory course to the use this platform.

> BU-ISCIII

Instituto de Salud Carlos III

THIS IS A PROTOTYPE. If you find any bugs please report them to mjuliam@isciii.es

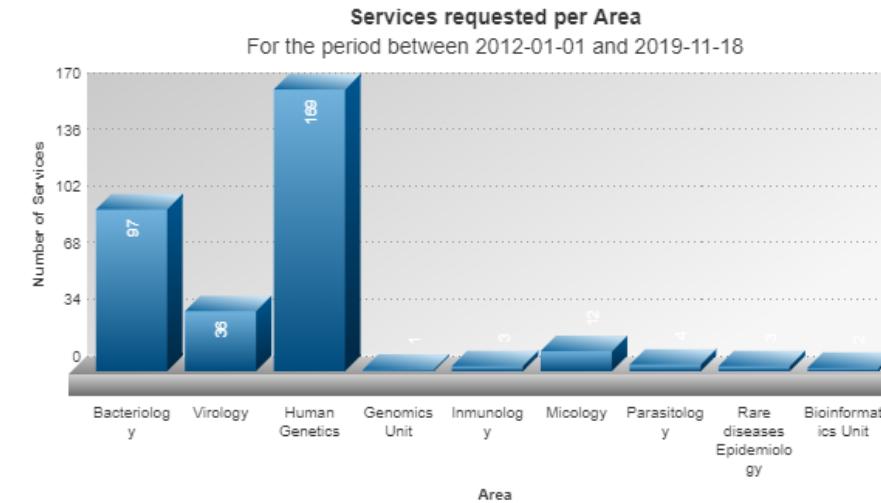
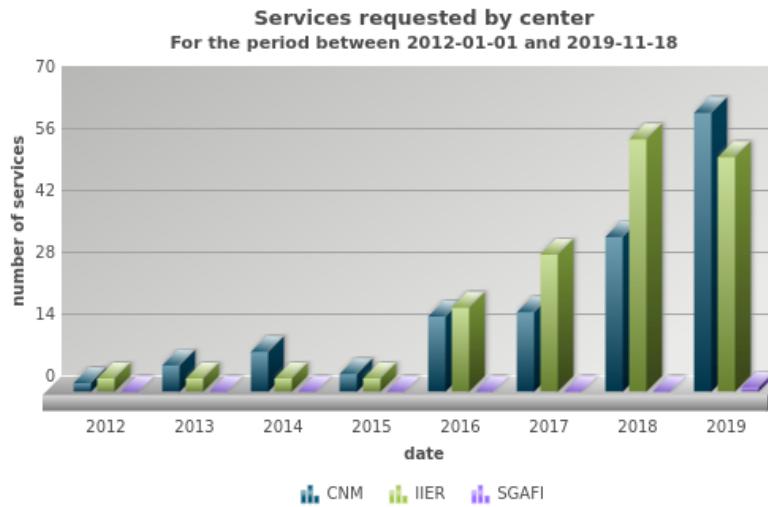
Take an interactive tour: Galaxy UI History Scratchbook

Galaxy is an open platform for supporting data intensive research. Galaxy is developed by [The Galaxy Team](#) with the support of [many contributors](#). If you use this platform to analyse your data, remember to cite both Galaxy Project and this server.

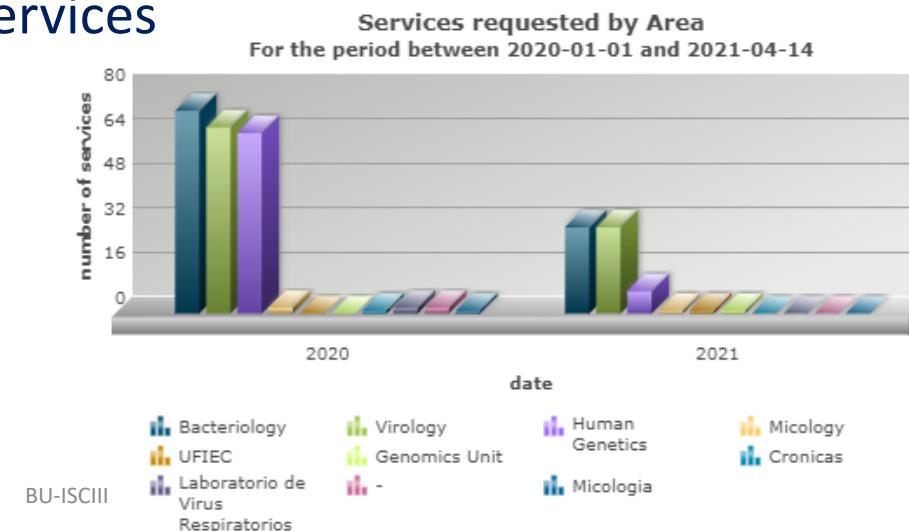
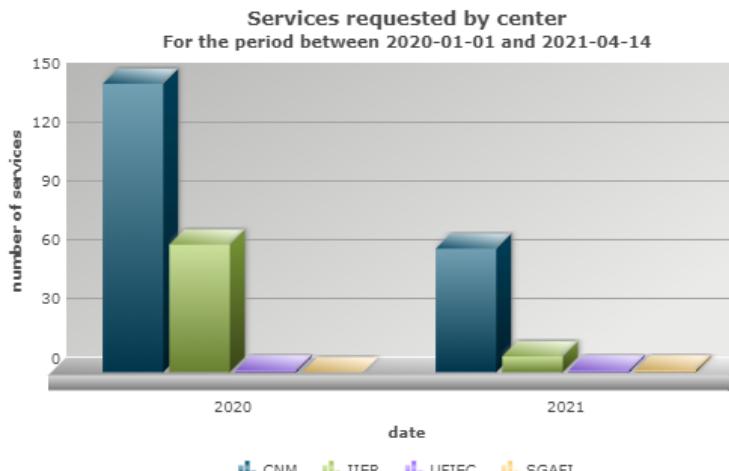
The [Galaxy Project](#) is supported in part by [NHGRI](#), [NSF](#), [The Huck Institutes of the Life Sciences](#), [The Institute for CyberScience at Penn State](#), and [Johns Hopkins University](#).

Number of services: 2012 – 2019

- 327 Services per center / CNM – IIER / 22 Researchers



2020-2021
200 Services



Training

Courses

ISCIII

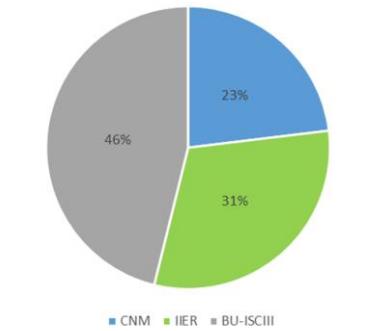
Introduction to massive sequencing data analysis, 2013-2021 (8 editions)

Secuenciación de genomas bacterianos: herramientas y aplicaciones, 2018-2021 (3 editions)

Análisis de genomas virales a través de la plataforma Galaxy, 2021 (1 edition)

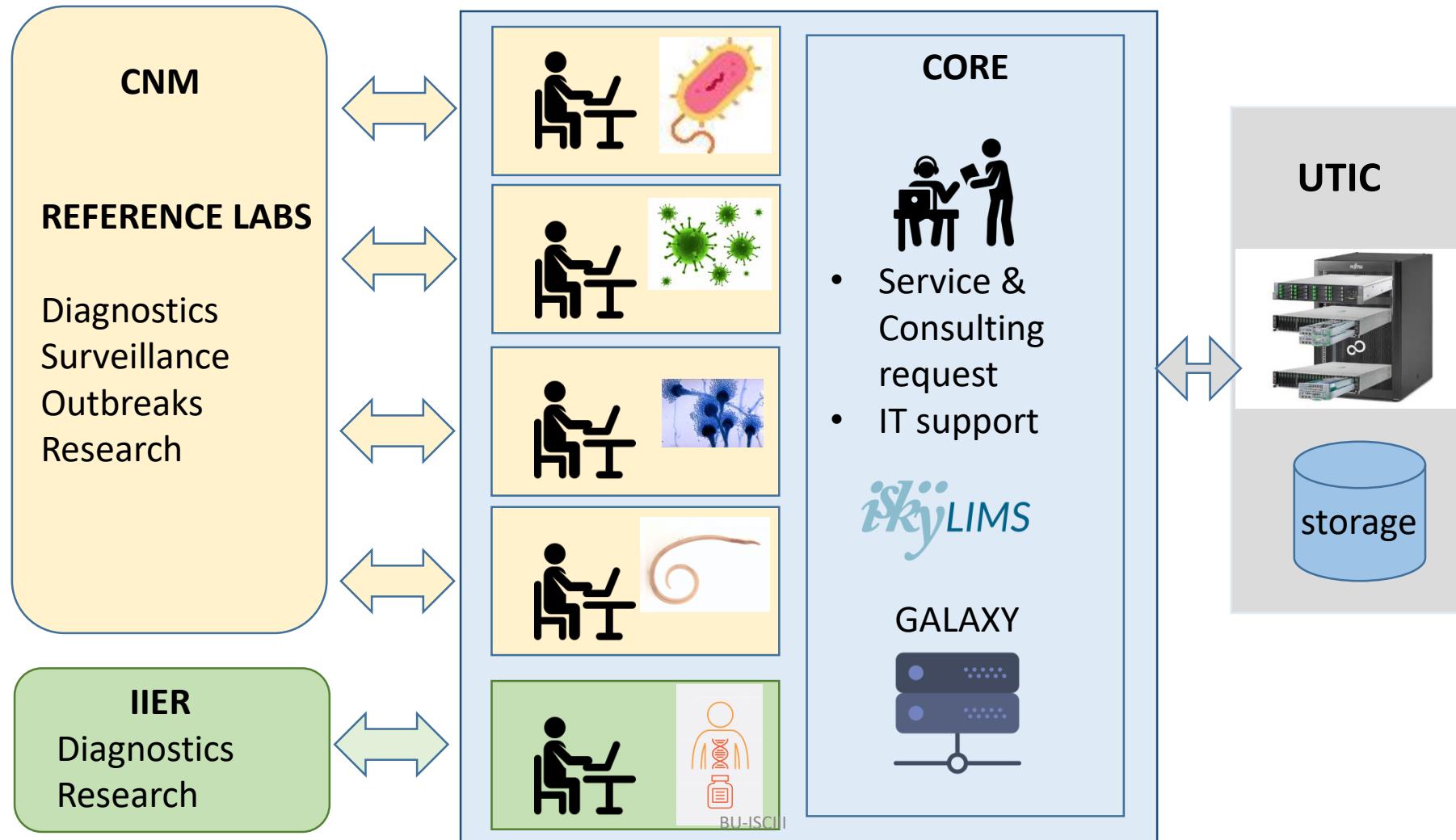
Master & Grade Students

- Bioinformática y Biología Computacional ENS-ISCIII
- Bioinformática UAM
- Genética y Biología Molecular UAM
- Microbiología aplicada a la salud pública e investigación en enfermedades infecciosas, U. Alcalá de Henares
- Sciences in Omics Data Analysis, Universidad de VIC, U. Central de Cataluña
- Complutense University



Hospitals Students

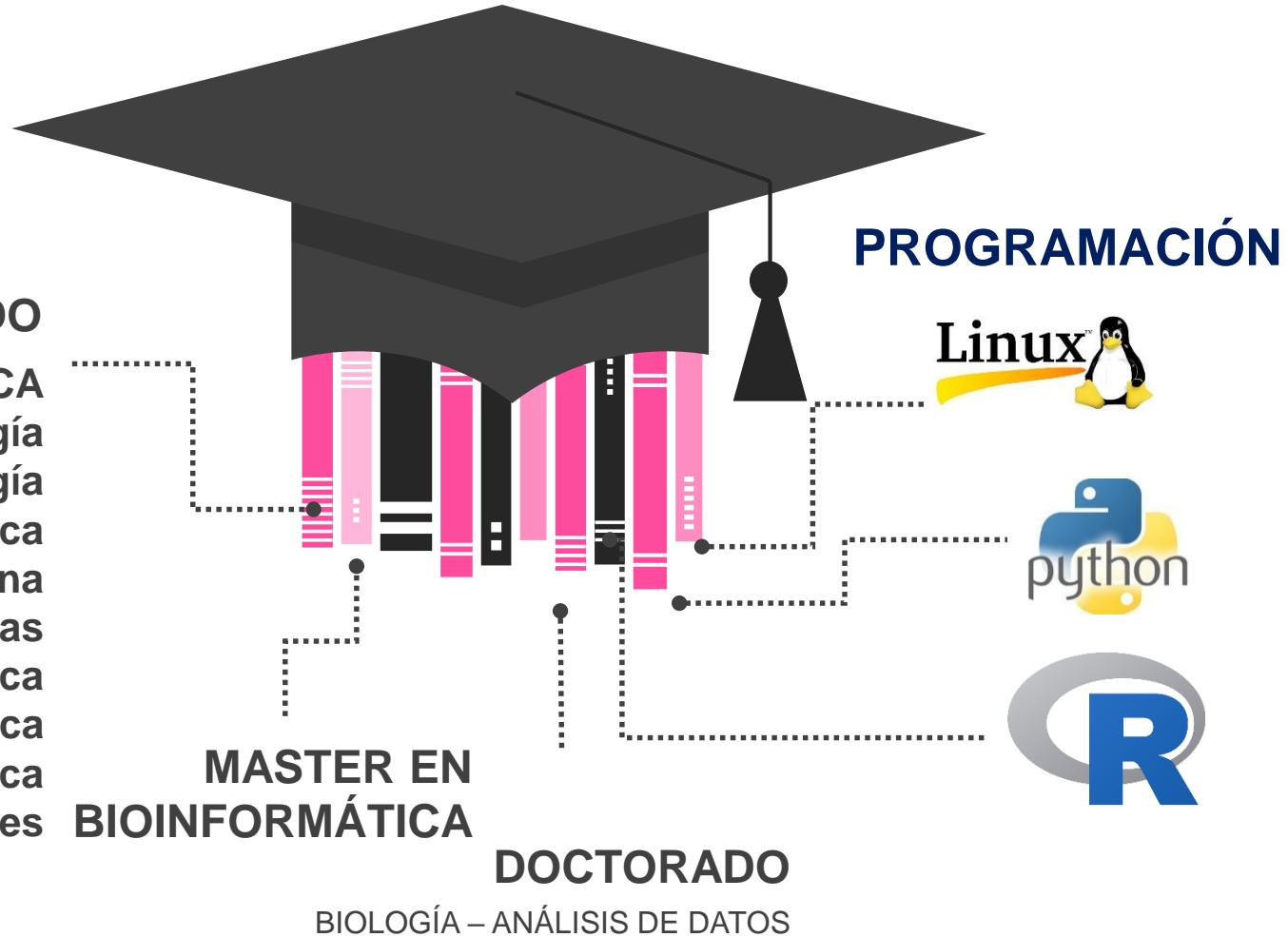
Roadmap: BU-ISCIII Model



FORMACIÓN EN BIOINFORMÁTICA

Universidad
Barcelona.

GRADO
BIOINFORMÁTICA
Biología
Biotecnología
Bioquímica
Medicina
Matemáticas
Química
Física
Informática
Telecomunicaciones



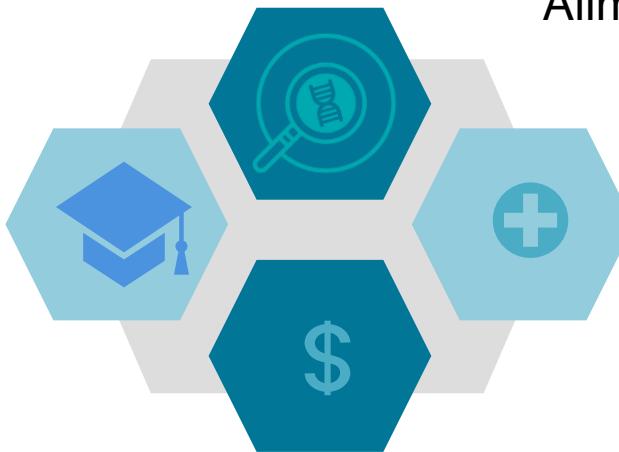
¿Dónde trabaja un Bioinformático?



UNIVERSIDAD

Biociencias
Informática

CENTRO DE INVESTIGACIÓN



EMPRESA

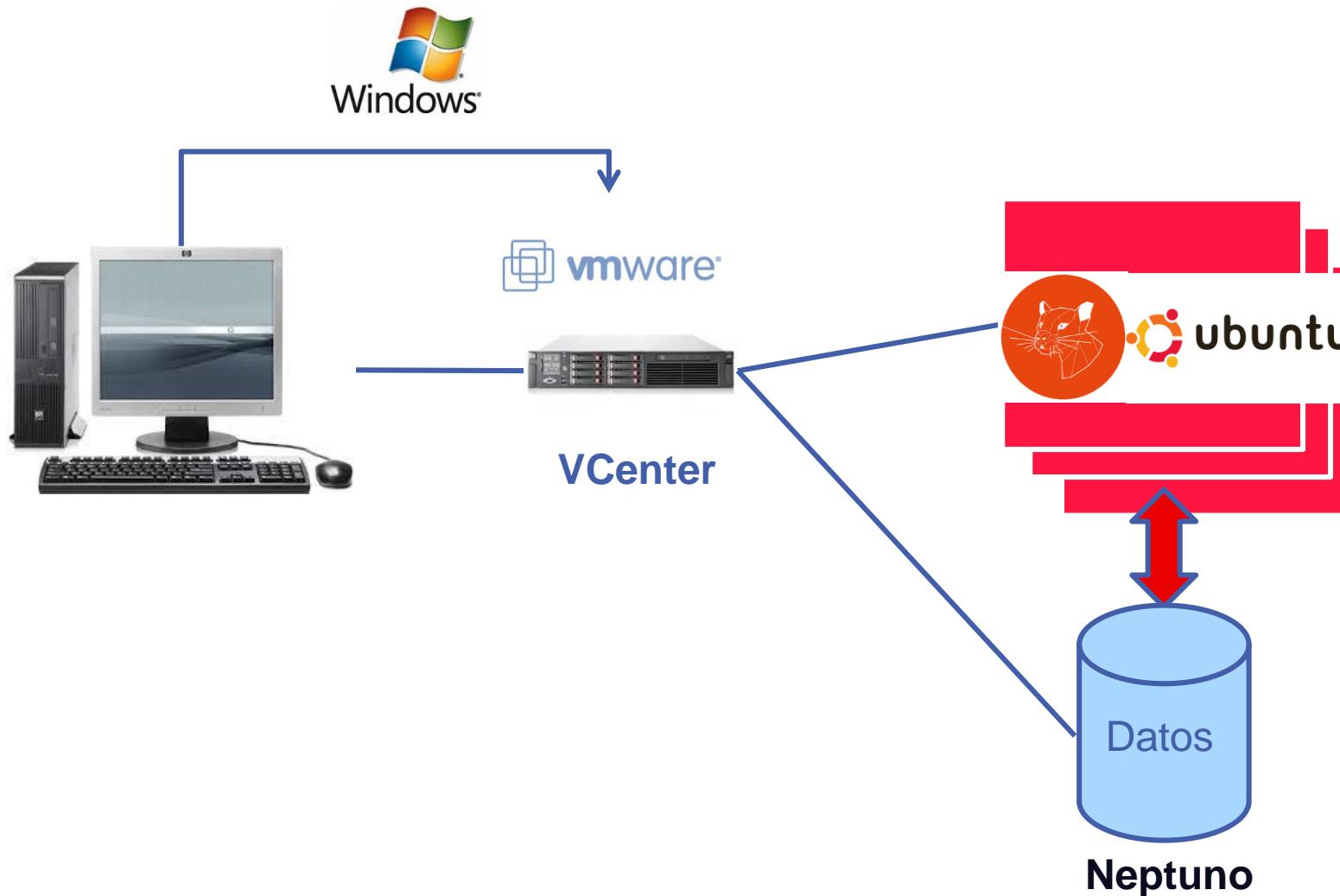
Bioinformática
Genética
Genómica

Biomedicina
Agricultura
Alimentación

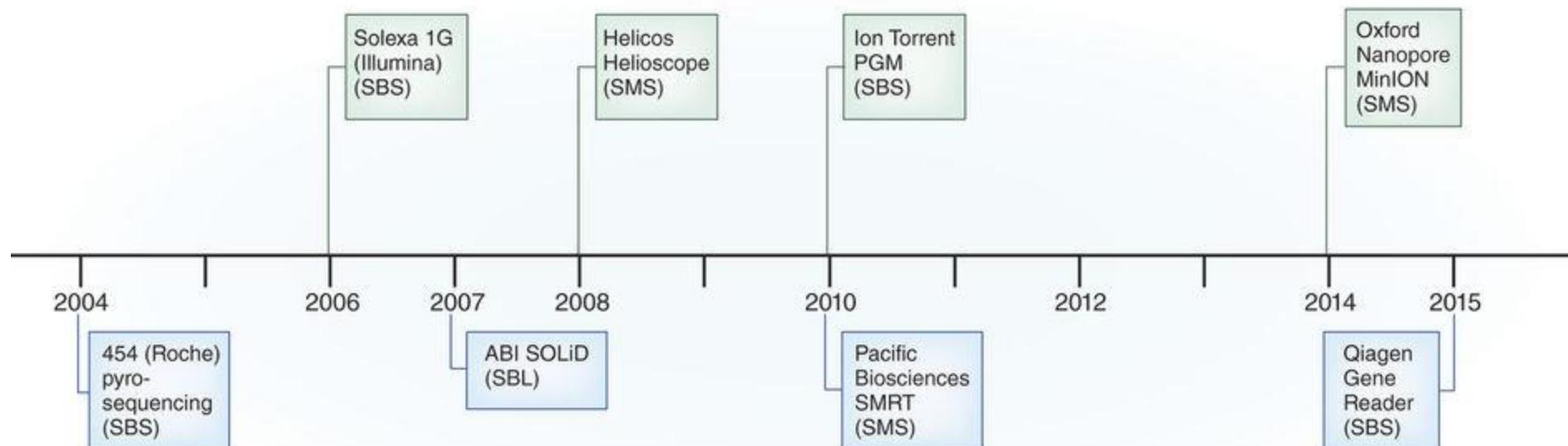
HOSPITAL

BIOINFORMÁTICO CLÍNICO
Genética
Oncología
Cardiología

Recursos Informáticos para el curso

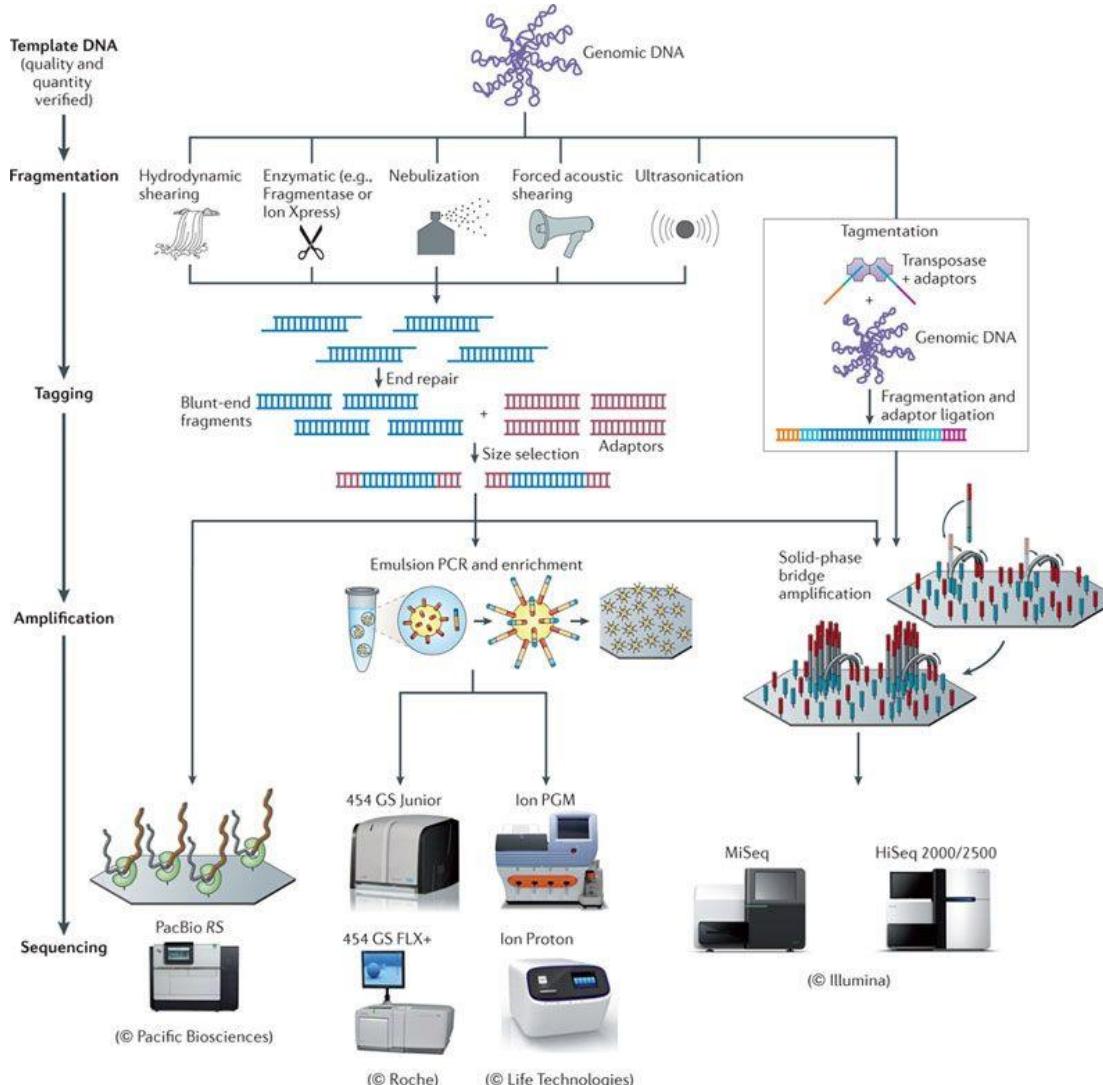


DNA sequencing technologies 2006-2016



Mardis, Nature Protocols 2017

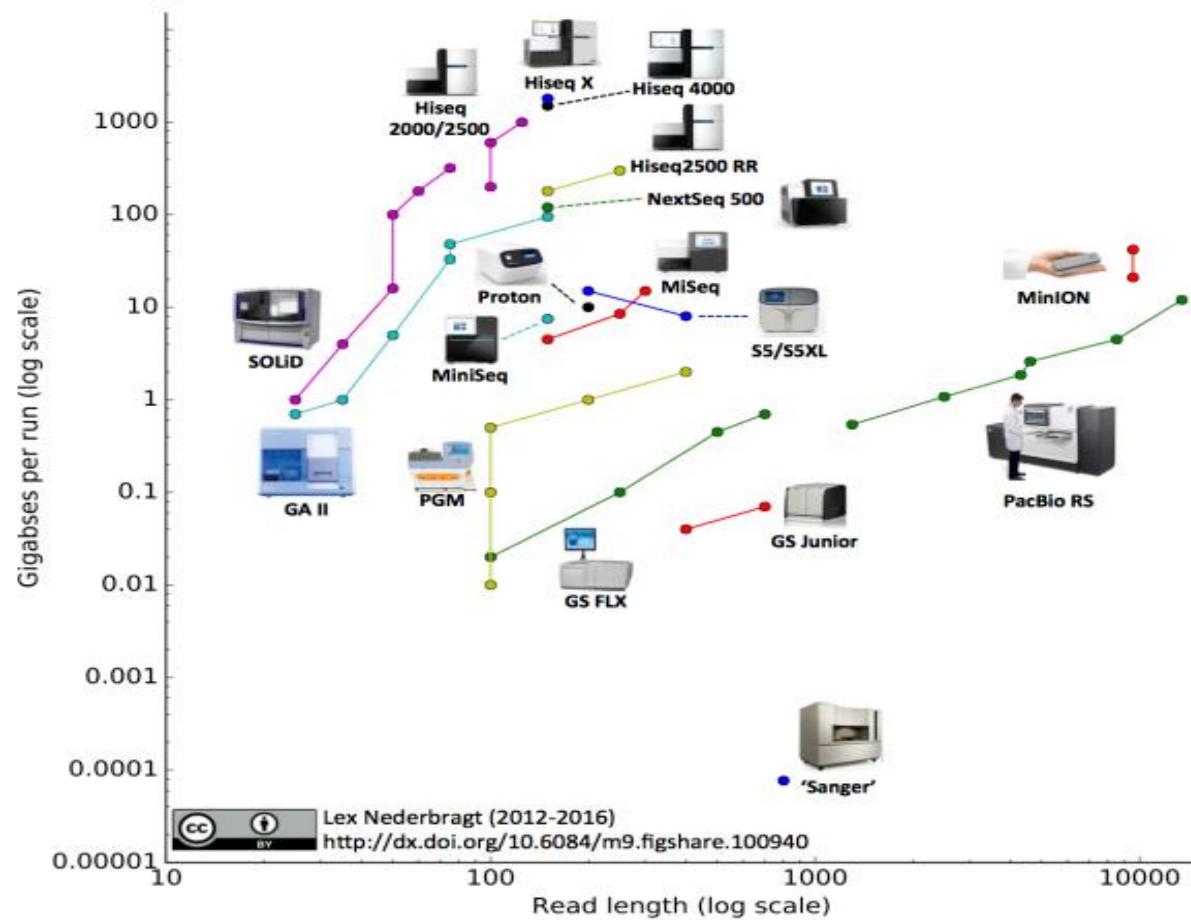
High-throughput sequencing platforms



Secuenciación de genomas bacterianos: herramientas y aplicaciones

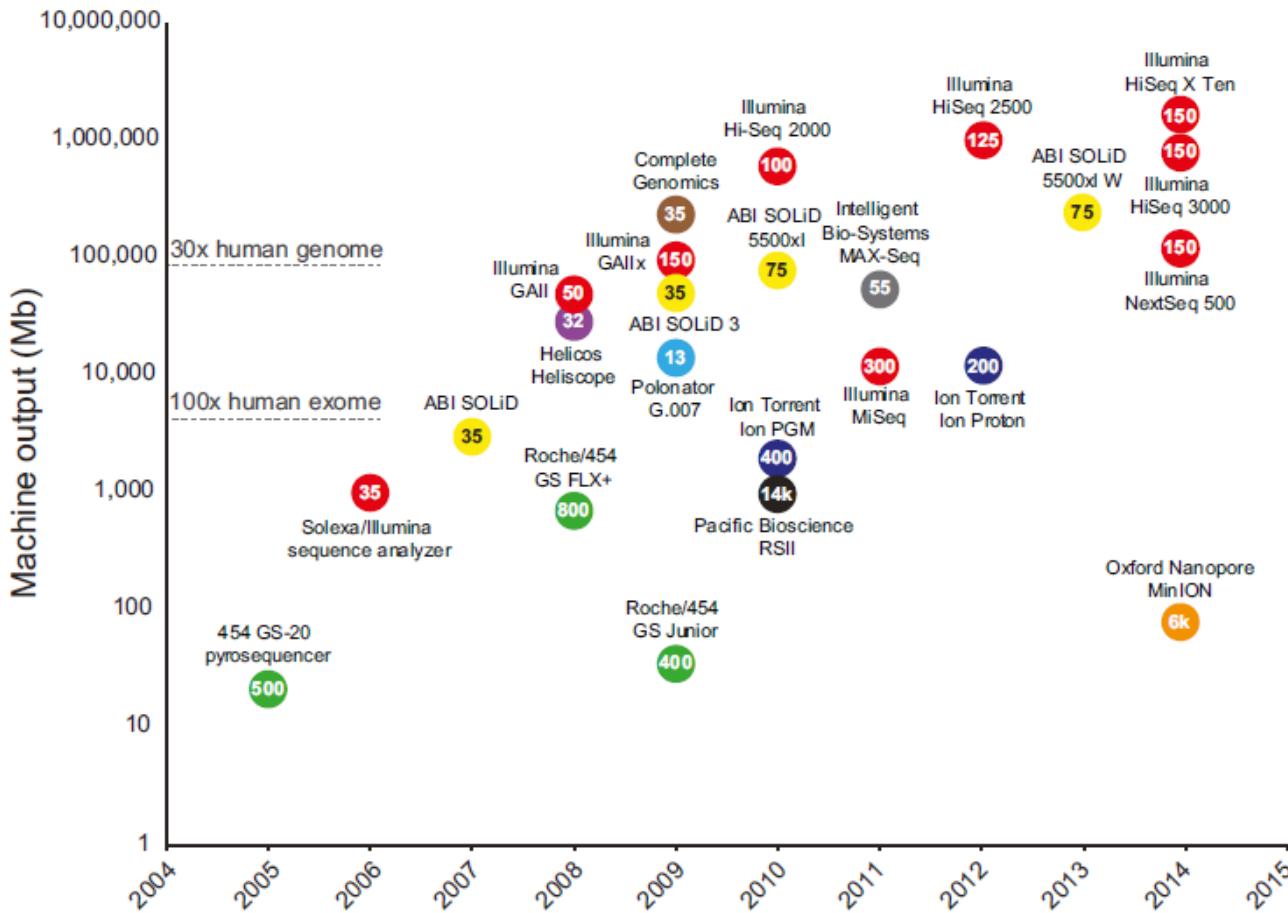
Nature Reviews | Microbiology Loman et al, 2012

High-Throughput Sequencing Technologies



<https://flxlexblog.wordpress.com/>

High-Throughput Sequencing Technologies

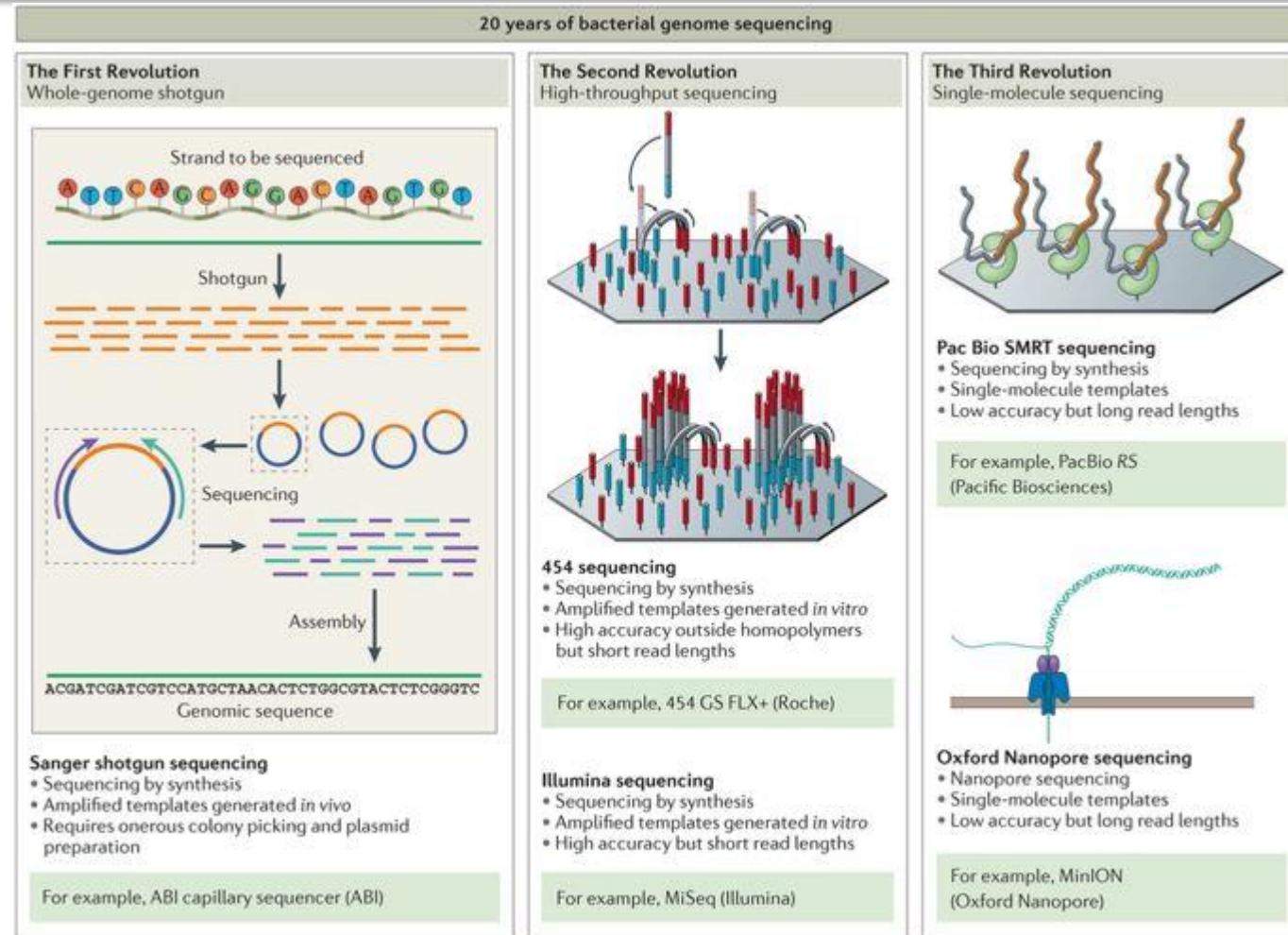


Numbers inside data points denote current read lengths.
Sequencing platforms are color coded.

Reuter et al., Mol Cell 2015

High-Throughput Sequencing Technologies

The three revolutions in sequencing technology that have transformed the landscape of bacterial genome sequencing

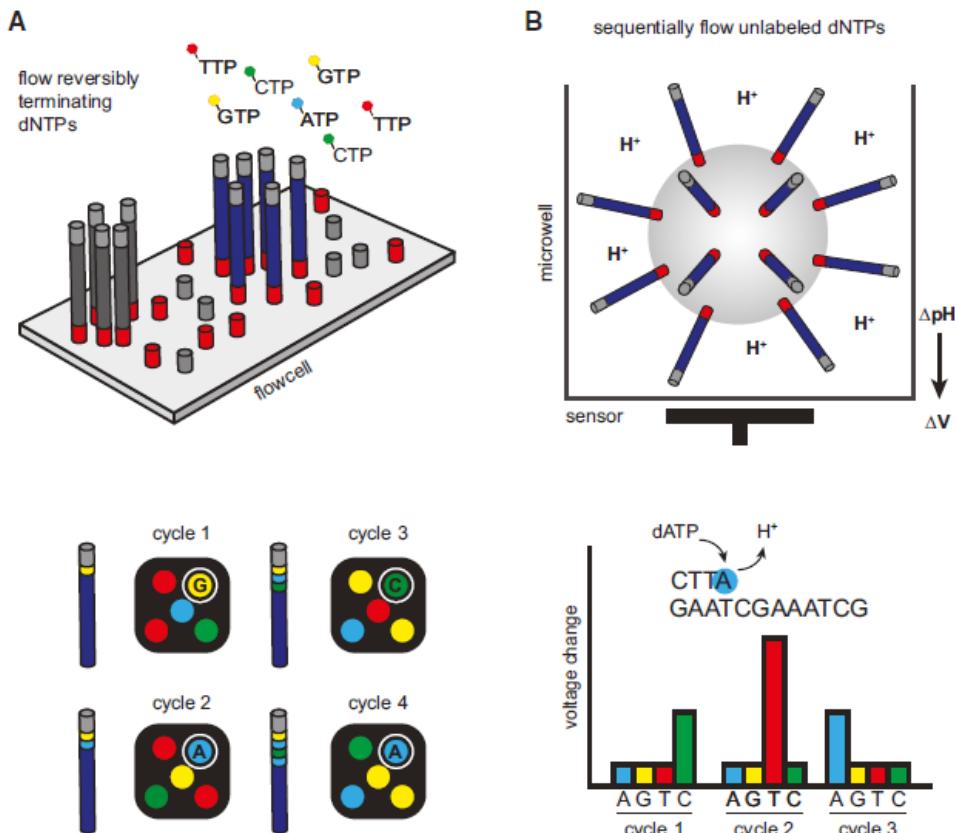


Nature Reviews | Microbiology

Secuenciación de genomas bacterianos: herramientas y aplicaciones

> BU-ISCIII

The Second-generation Sequencing Technologies



Clonal Amplification-Based Sequencing Platforms

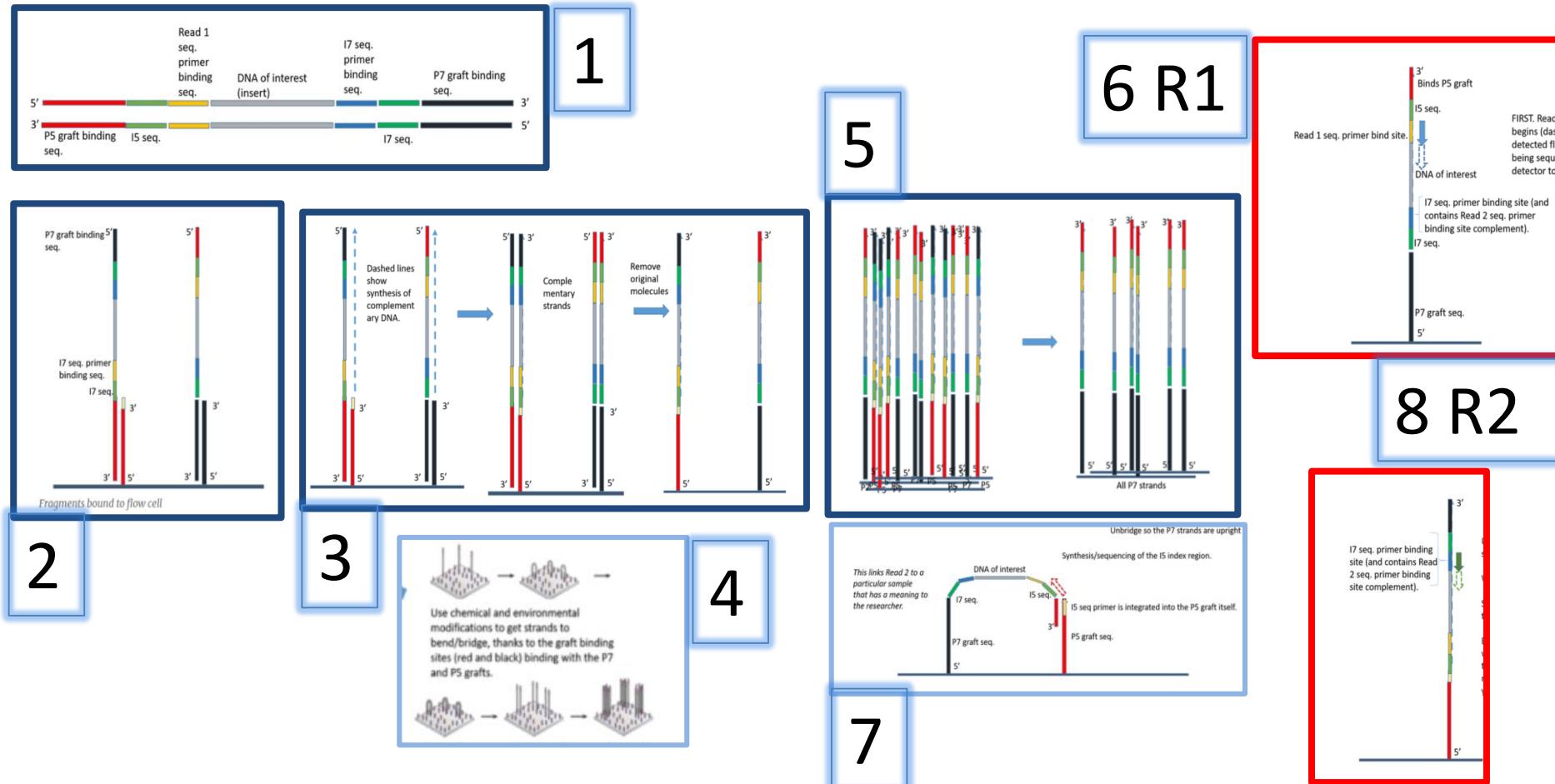
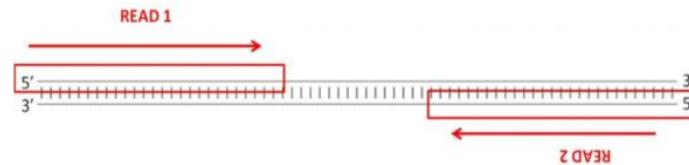
(A) Illumina's four-color reversible termination sequencing method.

(B) Ion Torrent's semiconductor sequencing method.

Reuter et al., Mol Cell 2015

Secuenciación de genomas bacterianos: herramientas y aplicaciones

Illumina sequencing



<https://kscbioinformatics.wordpress.com/2017/02/13/illumina-sequencing-for-dummies-samples-are-sequenced/>

Illumina Benchtop Sequencers



iSeq 100



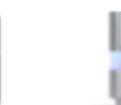
MiniSeq



MiSeq Series



NextSeq 550 Series



NextSeq 1000 & 2000

| Popular Applications & Methods | Key Application |
|---|-----------------|-----------------|-----------------|-----------------|-----------------|
| Large Whole-Genome Sequencing (human, plant, animal) | | | | | |
| Small Whole-Genome Sequencing (microbe, virus) | ● | ● | ● | ● | ● |
| Exome & Large Panel Sequencing (enrichment-based) | | | | ● | ● |
| Targeted Gene Sequencing (amplicon-based, gene panel) | ● | ● | ● | ● | ● |
| Single-Cell Profiling (scRNA-Seq, scDNA-Seq, oligo tagging assays) | | | | ● | ● |
| Transcriptome Sequencing (total RNA-Seq, mRNA-Seq, gene expression profiling) | | | | ● | ● |
| Targeted Gene Expression Profiling | ● | ● | ● | ● | ● |
| miRNA & Small RNA Analysis | ● | ● | ● | ● | ● |
| DNA-Protein Interaction Analysis (ChIP-Seq) | | | ● | ● | ● |
| Methylation Sequencing | | | | ● | ● |
| 16S Metagenomic Sequencing | | ● | ● | ● | ● |
| Metagenomic Profiling (shotgun metagenomics, metatranscriptomics) | | | | ● | ● |
| Cell-Free Sequencing & Liquid Biopsy Analysis | | | | ● | ● |

Benchtop Sequencer Sheds Light on Ebola Outbreak

Local scientists use the iSeq 100 Sequencing System to analyze transmission patterns and trace the origin of an Ebola outbreak in the Democratic Republic of the Congo.

[Read Article ▶](#)

<https://emea.illumina.com/systems/sequencing-platforms.html>

| Run Time | 9.5-19 hrs | 4-24 hours | 4-55 hours | 12-30 hours | 11-48 hours |
|-----------------------|------------|------------|-------------------------|-------------|--------------------------|
| Maximum Output | 1.2 Gb | 7.5 Gb | 15 Gb | 120 Gb | 330 Gb [*] |
| Maximum Reads Per Run | 4 million | 25 million | 25 million [†] | 400 million | 1.1 billion [*] |
| Maximum Read Length | 2 × 150 bp | 2 × 150 bp | 2 × 300 bp | 2 × 150 bp | 2 × 150 bp |

Illumina Production-Scale Sequencers



NextSeq 550 Series

NextSeq 1000 & 2000

NovaSeq 6000

| Popular Applications & Methods | Key Application | Key Application | Key Application |
|--|-----------------|-----------------|-----------------|
| Large Whole-Genome Sequencing (human, plant, animal) | | | ● |
| Small Whole-Genome Sequencing (microbe, virus) | ● | ● | ● |
| Exome & Large Panel Sequencing (enrichment-based) | ● | ● | ● |
| Targeted Gene Sequencing (amplicon-based, gene panel) | ● | ● | ● |
| Single-Cell Profiling (scRNA-Seq, scDNA-Seq, oligo tagging assays) | ● | ● | ● |
| Transcriptome Sequencing (total RNA-Seq, miRNA-Seq, gene expression profiling) | ● | ● | ● |
| Chromatin Analysis (ATAC-Seq, ChIP-Seq) | ● | ● | ● |
| Methylation Sequencing | ● | ● | ● |
| Metagenomic Profiling (shotgun metagenomics, metatranscriptomics) | ● | ● | ● |
| Cell-Free Sequencing & Liquid Biopsy Analysis | ● | ● | ● |

Optimized NGS Sample Tracking and Workflows

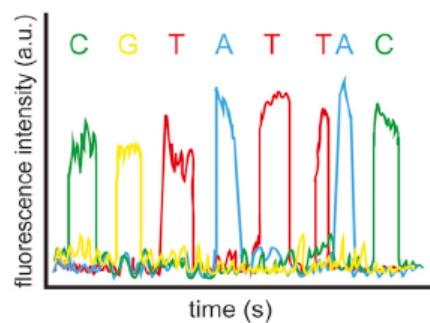
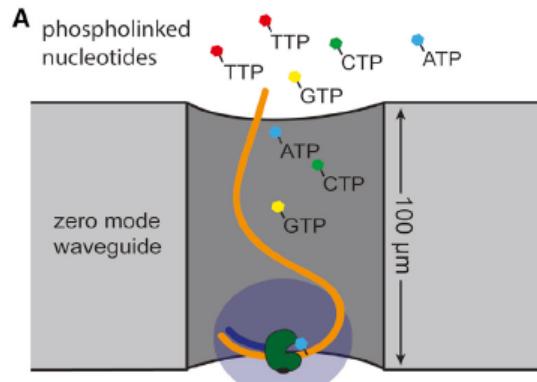
See how a Laboratory Information Management System (LIMS) enabled this large genomics lab to standardize lab procedures and cope with increasing sample volumes from diverse clients.

[Read Case Study >](#)

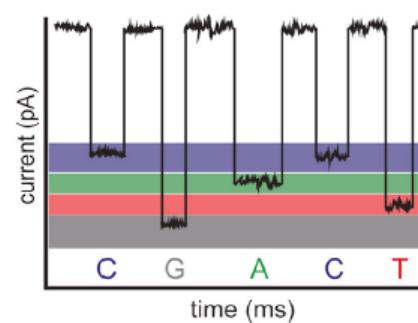
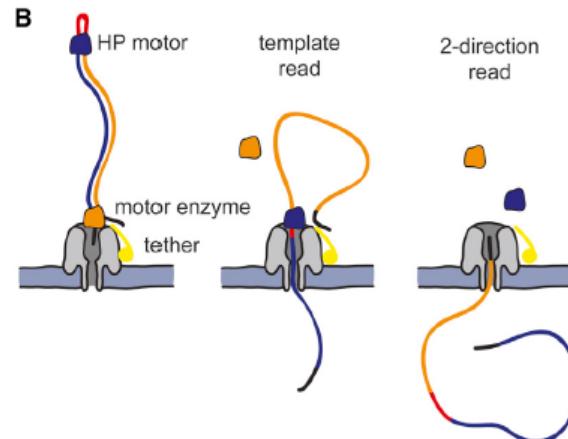
| | | | |
|-----------------------|-------------|--------------|---|
| Run Time | 12-30 hours | 11-48 hours | -13 - 38 hours (dual SP flow cells) -13-25 hours (dual S1 flow cells) -16-36 hours (dual S2 flow cells) -44 hours (dual S4 flow cells) |
| Maximum Output | 120 Gb | 330 Gb* | 6000 Gb |
| Maximum Reads Per Run | 400 million | 1.1 billion* | 20 billion |
| Maximum Read Length | 2 × 150 bp | 2 × 150 bp | 2 × 250** |

The Third-generation Sequencing Technologies

Single Molecule Sequencing Platforms



Pacific Bioscience's SMRT sequencing



Oxford Nanopore's sequencing strategy

Reuter et al., Mol Cell 2015

Secuenciación de genomas bacterianos: herramientas y aplicaciones

PacBio sequencing and its applications

Rhoads & Au, Gen Prot Bioinf 2015



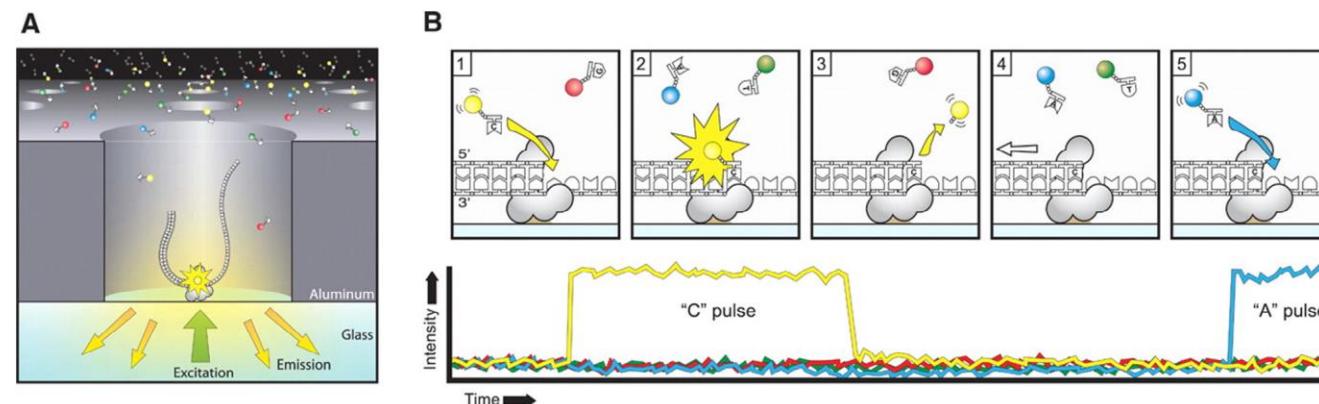
SMRTbell template: is a closed, single-stranded circular DNA that is created by ligating hairpin adaptors to both ends of a target dsDNA

Sequencing by light pulses: The replication processes in all ZMWs of a SMRTcell are recorder by a movie of light pulses, and the pulses corresponding to each ZMW can be interpreted to be a sequence of bases (**continuous long read, CLR**).

Both strands can be sequenced multiple times (passes) in a single CLR. CLR can be split to multiple reads (subreads) and CCS is the consensus sequence of multiple subreads



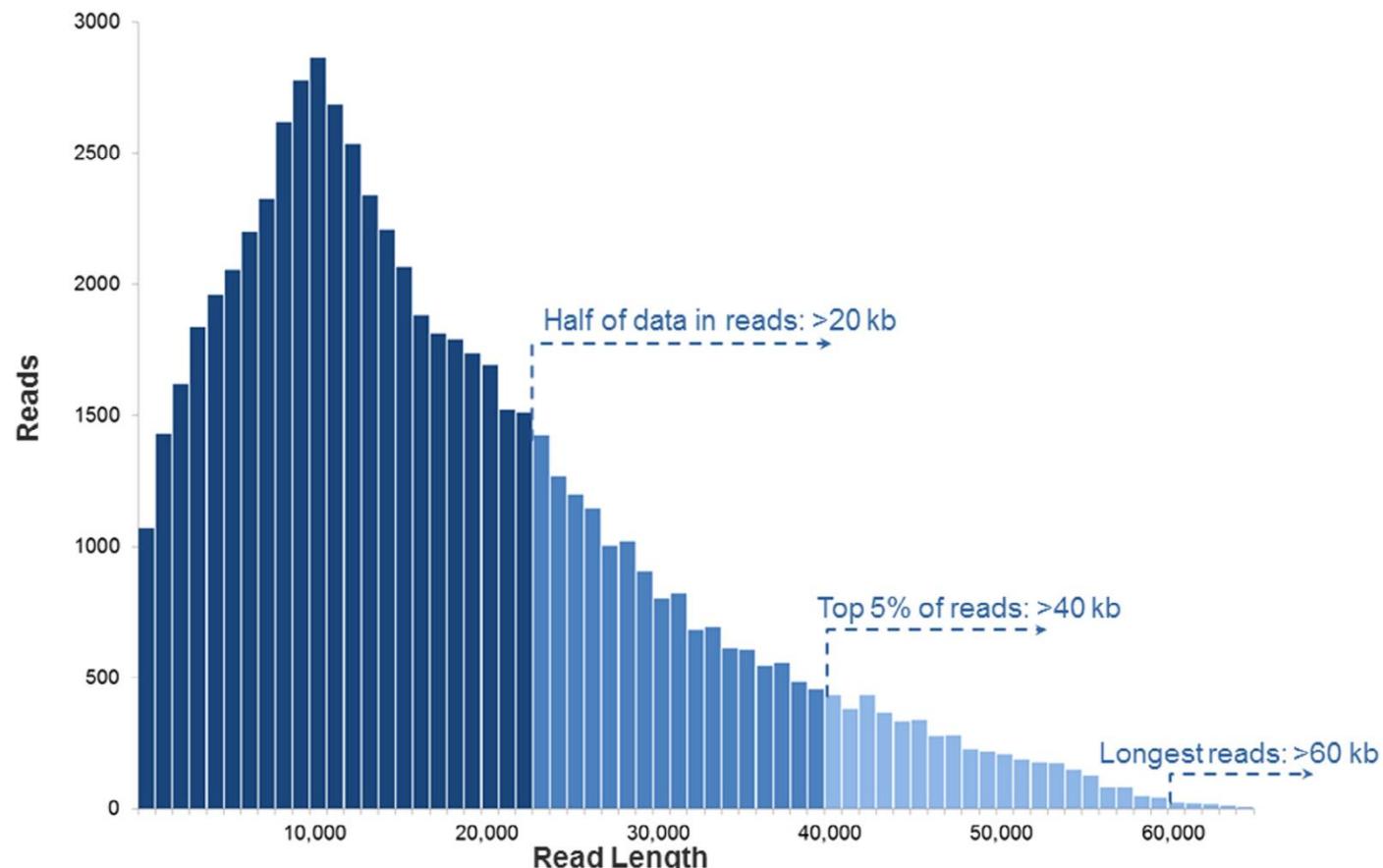
A single SMRT cell: this contains 150000 ZMWs (zero-mode waveguide). A SMRTbell diffuses into a ZMW. Approx 35000 -75000 ZMWs produce a read in a run lasting 0,5-4h resulting in 0,5-1Gb.



PacBio sequencing and its applications

Rhoads & Au, Gen Prot Bioinf 2015

PacBio RS II read length distribution using P6-C4 chemistry. Data are based on a 20kb size-selected E. coli library using a 4-h movie. A SMRTcell produces 0,5-1 billion bases.



PacBio sequencing and its applications

Rhoads & Au, Gen Prot Bioinf 2015

Table 2 *De novo* genome assemblies using hybrid sequencing or PacBio sequencing alone

| Species | Method | Tools | SMRT cells | Coverage | Contigs | Achievements | Ref. |
|--|--------|--------------------------------|--------------|---|-------------------------|---|------|
| <i>Clostridium autoethanogenum</i> | PacBio | HGAP | 2 | 179× | 1 | 21 fewer contigs than using SGS; no collapsed repeat regions (≥ 4 using SGS) | [7] |
| <i>Potentilla micrantha</i> (chloroplast) | PacBio | HGAP, Celera, minimus2, SeqMan | 26 | 320× | 1 | 6 fewer contigs than with Illumina; 100% coverage (Illumina: 90.59%); resolved 187 ambiguous nucleotides in Illumina assembly; unambiguously assigned small differences in two > 25 kb inverted repeats | [33] |
| <i>Escherichia coli</i> | PacBio | PBcR, MHAP, Celera, Quiver | 1 | 85× | 1 | 4.6 CPU hours for genome assembly (10× improvement over BLASR) | [31] |
| <i>Saccharomyces cerevisiae</i> | PacBio | PBcR, MHAP, Celera | 12 | 117× | 21 | 27 CPU hours for genome assembly (8× improvement over BLASR); improved current reference of telomeres | [31] |
| <i>Arabidopsis thaliana</i> | PacBio | PBcR, MHAP, Celera | 46 | 144× | 38 | 1896 CPU hours for genome assembly | [31] |
| <i>Drosophila melanogaster</i> | PacBio | PBcR, MHAP, Celera, Quiver | 42 | 121× | 132 | 1060 CPU hours for genome assembly (593× improvement over BLASR); improved current reference of telomeres | [31] |
| <i>Homo sapiens</i> (CHM1hert) | PacBio | PBcR, MHAP, Celera | 275 | 54× | 3434 | 262,240 CPU hours for genome assembly; potentially closed 51 gaps in GRCh38; assembled MHC in 2 contigs (60 contigs with Illumina); reconstructed repetitive heterochromatic sequences in telomeres | [31] |
| <i>Homo sapiens</i> (CHM1tert) | PacBio | BLASR, Celera, Quiver | 243 | 41× | N/A (local assembly) | Closed 50 gaps and extended into 40 additional gaps in GRCh37; added over 1 Mb of novel sequence to the genome; identified 26,079 indels at least 50 bp in length; cataloged 47,238 SV breakpoints | [32] |
| <i>Melopsittacus undulatus</i> | Hybrid | PBcR, Celera | 3 | 5.5× PacBio + 15.4× 454 = 3.83× corrected | 15,328 | 1st assembly of > 1 Gb parrot genome; N50 = 93,069 | [34] |
| <i>Vibrio cholerae</i> | Hybrid | BLASR, Bambus, AHA | 195 | 200× PacBio + 28× Illumina + 22× 454 | 2 | No N's in contigs; 99.99% consensus accuracy; N50 = 3.01 Mb | [30] |
| <i>Helicobacter pylori</i> | PacBio | HGAP, Quiver, PGAP | 8 per strain | 446.5× average among strains | 1 per strain | 1 complete contig for each of 8 strains; methylation analysis associated motifs with genotypes of virulence factors | [35] |

Note: N50, the contig length for which half of all bases are in contigs of this length or greater; MHC, major histocompatibility complex; SV, structural variation.

PacBio sequencing and its applications

Rhoads & Au, Gen Prot Bioinf 2015

Advantage

Closes gaps and completes genomes due to longer reads

Identifies non-SNP SVs

Achievements

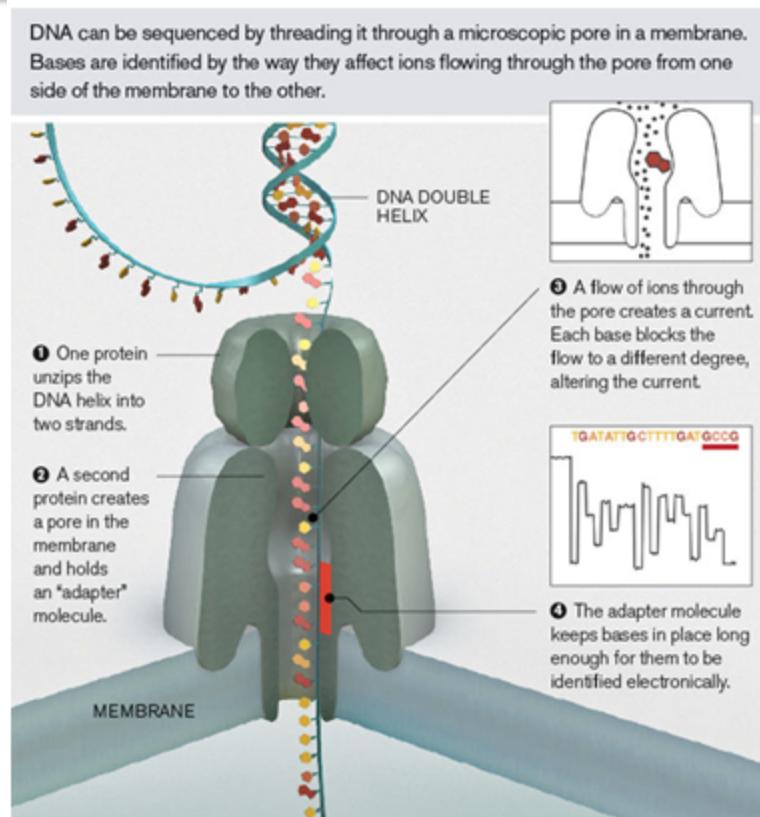
Produced highly-contiguous assemblies of bacterial and eukaryotic genomes

Discovered STRs (short tandem repeats)

Limitations

Both strands can be sequenced several times if the lifetime of the polymerase is long enough.

Nanopore-based fourth-generation DNA sequencing technology. ONT, Oxford Nanopore Technologies



'Strand sequencing' is a technique that passes intact DNA polymers through a protein nanopore, sequencing in real time as the DNA translocates the pore.

Nanopore sequencing also offers, for the first time, direct RNA sequencing, as well as PCR or PCR-free cDNA sequencing.

<https://nanoporetech.com/applications/dna-nanopore-sequencing>

Feng et al , Gen Prot Bioinf 2015

MinIon, OXFORD NANOPORE



<https://nanoporetech.com/news/movies#movie-24-nanopore-dna-sequencing>

Oxford Nanopore Technologies, MinION



The MinION is a portable sequencer; flow cells contain up to 512 nanopore sensors.

The Oxford Nanopore system processes the reads that are presented to it rather than generating read lengths. Sample-prep dependent, the longest read reported by a MinION user to date is >1 Mb.

Long reads confer many advantages, including simpler assembly and in the analysis of repetitive regions, phasing or CNVs.



For MinION / GridION
Flongle

Adapter to enable small, rapid nanopore sequencing tests, for mobile or desktop sequencers



MinION Mk1B

Your personal nanopore sequencer, putting you in control



MinION Mk1C

Your personal nanopore sequencer including compute and screen, putting you in control



GridION Mk1

Higher-throughput, on demand nanopore sequencing at the desktop, for you or as a service



PromethION 24/48

Ultra-high throughput, on-demand nanopore sequencing, for you or as a service

| | Flongle | MinION Mk1B | MinION Mk1C | GridION Mk1 | PromethION 24 | PromethION 48 |
|--|---|---|---|---|--|--|
| Number of channels per flow cell | 126 | 512 | 512 | 512 | 3000 | 3000 |
| Number of flow cells per device | 1 | 1 | 1 | 5 | 24 | 48 |
| Price per flow cell | \$90 | \$900 - \$475 | \$900 - \$475 | \$900 - \$475 | \$2000 - \$625 | \$2000 - \$625 |
| Run time | 1 min - 16 hours | 1 min - 72 hours | 1 min - 72 hours | 1 min - 72 hours | 1 min - 72 hours | 1 min - 72 hours |
| Yields in field are dependent on sample and preparation methods. Users can get outputs in the following ranges per flow cell when utilising the latest chemistries and protocols | 1 - 2 Gb | 10 - 30 - 50 Gb | 10 - 30 - 50 Gb | 10 - 30 - 50 Gb | 100 - 200 - 300 Gb | 100 - 200 - 300 Gb |
| Price per Gb for different flow cell yields (yields vary according to sample and preparation methods) | @1 - 2 Gb \$90 per flow cell: \$90 - 45 | @10 - 30 - 50 Gb \$900 per flow cell: \$90 - 30 - 18 \$790 per flow cell: \$79 - 26 - 16 \$675 per flow cell: \$68 - 23 - 14 \$500 per flow cell: \$50 - 17 - 10 \$475 per flow cell: \$48 - 16 - 9.5 | @10 - 30 - 50 Gb \$900 per flow cell: \$90 - 30 - 18 \$790 per flow cell: \$79 - 26 - 16 \$675 per flow cell: \$68 - 23 - 14 \$500 per flow cell: \$50 - 17 - 10 \$475 per flow cell: \$48 - 16 - 9.5 | @10 - 30 - 50 Gb \$900 per flow cell: \$90 - 30 - 18 \$790 per flow cell: \$79 - 26 - 16 \$675 per flow cell: \$68 - 23 - 14 \$500 per flow cell: \$50 - 17 - 10 \$475 per flow cell: \$48 - 16 - 9.5 | @100 - 200 - 300 Gb \$1,600 per flow cell: \$16 - 8 - 5 \$1,120 per flow cell: \$11 - 6 - 4 \$940 per flow cell: \$9 - 5 - 3.1 \$680 per flow cell: \$7 - 3.4 - 2.3 \$625 per flow cell: \$6 - 3 - 2 | @100 - 200 - 300 Gb \$1,600 per flow cell: \$16 - 8 - 5 \$1,120 per flow cell: \$11 - 6 - 4 \$940 per flow cell: \$9 - 5 - 3.1 \$680 per flow cell: |

Library preparation



Oxford Nanopore has developed VolTRAX – a small device designed to perform library preparation automatically, so that a user can get a biological sample ready for analysis, hands-free. VolTRAX is designed as an alternative to a range of lab equipment, to allow consistent and varied, automated library prep options.

VolTRAX V2 Starter Pack

\$8,000.00

VolTRAX V2 is designed to automate all laboratory processes associated with Nanopore Sequencing from sample extraction to library preparation.

MinIT, Analysis



Eliminating the need for a dedicated laptop
for nanopore sequencing with MinION.
\$2400

MinIT Specifications:

Pre-installed software: Linux OS, MinKNOW, Guppy, EPI2ME

Bluetooth and Wi-Fi enabled; you can control your experiments using a laptop, tablet or smartphone

fastq or fast5 files are written to Onboard storage: 512 GB SSD

Processing: GPU accelerators (ARM processor 6 cores, 256 Core GPU), 8 GB RAM.

Small footprint, 290g

1 x USB 2.0 port, 1 x USB 3.0 port and 1 x Ethernet port (1 Gbit capacity)

MinIT has now been replaced by the MinION Mk1C, which combines the real-time, portable sequencing of MinION, with powerful integrated compute, a high-resolution touchscreen, and full connectivity.

SmidgION, Mobile analysis



Oxford Nanopore has now started developing an even smaller device, SmidgION.

potential applications may include remote monitoring of pathogens in a breakout or infectious disease; the on-site analysis of environmental samples such as water/metagenomics samples, real time species ID for analysis of food, timber, wildlife or even unknown samples; field-based analysis of agricultural environments, and much more.

PacBio sequencing and its applications

Rhoads & Au, Gen Prot Bioinf 2015

Performance comparison of sequencing platforms of various generations

| Method | Generation | Read length (bp) | Single pass error rate (%) | No. of reads per run | Time per run | Cost per million bases (USD) | Refs. |
|-----------------------------------|------------|--|----------------------------|-------------------------------|--------------|------------------------------|------------|
| Sanger ABI 3730×1 | 1st | 600–1000 | 0.001 | 96 | 0.5–3 h | 500 | [14,18–21] |
| Ion Torrent | 2nd | 200 | 1 | 8.2×10^7 | 2–4 h | 0.1 | [15,25] |
| 454 (Roche) GS FLX+ | 2nd | 700 | 1 | 1×10^6 | 23 h | 8.57 | [14,17,27] |
| Illumina HiSeq 2500 (High Output) | 2nd | 2×125 | 0.1 | 8×10^9 (paired) | 7–60 h | 0.03 | [9,16,26] |
| Illumina HiSeq 2500 (Rapid Run) | 2nd | 2×250 | 0.1 | 1.2×10^9 (paired) | 1–6 days | 0.04 | [9,16,26] |
| SOLiD 5500×1 | 2nd | 2×60 | 5 | 8×10^8 | 6 days | 0.11 | [14,24] |
| PacBio RS II: P6-C4 | 3rd | $1.0\text{--}1.5 \times 10^4$ on average | 13 | $3.5\text{--}7.5 \times 10^4$ | 0.5–4 h | 0.40–0.80 | [5,12,15] |
| Oxford Nanopore MinION | 3rd | $2\text{--}5 \times 10^3$ on average | 38 | $1.1\text{--}4.7 \times 10^4$ | 50 h | 6.44–17.90 | [22,23] |

Characteristics, strengths and weaknesses of commonly used sequencing platforms

Table 2

Characteristics, strengths and weaknesses of commonly used sequencing platforms

| Platform \ Instrument | Throughput range (Gb) ^a | Read length (bp) | Strength | Weakness |
|----------------------------|------------------------------------|------------------|--|---------------------------------------|
| <i>Sanger sequencing</i> | | | | |
| ABI 3500/3730 | 0.0003 | Up to 1 kb | Read accuracy and length | Cost and throughput |
| <i>Illumina</i> | | | | |
| MiniSeq | 1.7–7.5 | 1×75 to ×150 | Low initial investment | Run and read length |
| MiSeq | 0.3–15 | 1×36 to 2×300 | Read length, scalability | Run length |
| NextSeq | 10–120 | 1×75 to 2×150 | Throughput | Run and read length |
| HiSeq (2500) | 10–1000 | ×50 to ×250 | Read accuracy, throughput, | High initial investment, run |
| NovaSeq 5000/6000 | 2000–6000 | 2×50 to ×150 | Read accuracy, throughput | High initial investment, run |
| <i>Ion Torrent</i> | | | | |
| PGM | 0.08–2 | Up to 400 | Read length, speed | Throughput, homopolymers ^c |
| S5 | 0.6–15 | Up to 400 | Read length, speed, | Homopolymers ^c |
| Proton | 10–15 | Up to 200 | Speed, throughput | Homopolymers ^c |
| <i>Pacific BioSciences</i> | | | | |
| PacBio RSII | 0.5–1 ^b | Up to 60 kb | Read length, speed (Average 10 kb, N50 20 kb) | High error rate and initial |
| Sequel | 5–10 ^b | Up to 60 kb | Read length, speed (Average 10 kb, N50 20 kb) | High error rate |
| <i>Oxford Nanopore</i> | | | | |
| MINION | 0.1–1 | Up to 100 kb | Read length, portability | High error rate, run length, |

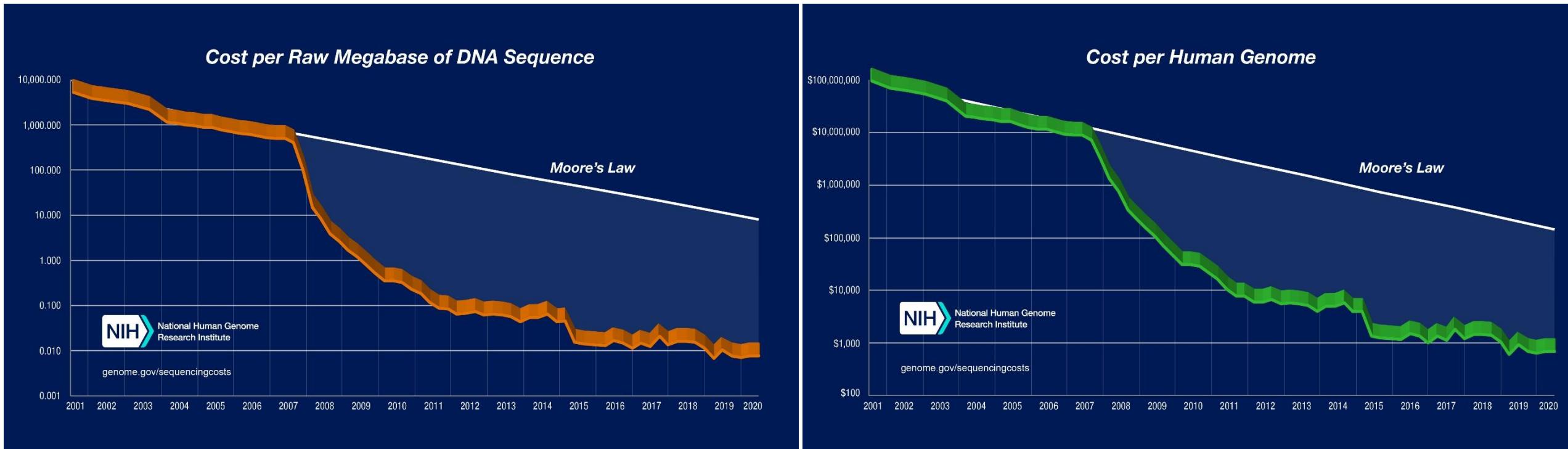
^a The throughput ranges are determined by available kits and run modes on a per run basis. As an example of a 15-GB throughput, thirty-five 5-MB genomes can be sequenced to a minimum coverage of 40× on the Illumina MiSeq using the v3 600 cycle chemistry.

^b Per one single-molecule real-time cell.

^c Results in increased error rate (increased proportion of reads containing errors among all reads) which in turn results in false-positive variant calling.

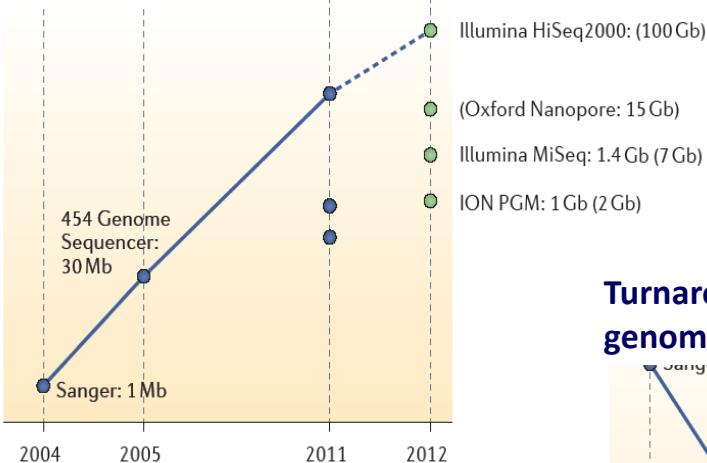
Besser et al., Clin Micr Infect, 2018

Coste actual de la secuenciación

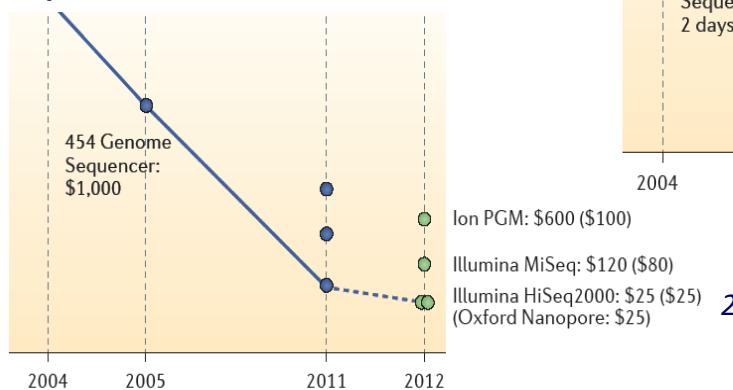


Sequencing platforms in Microbiology

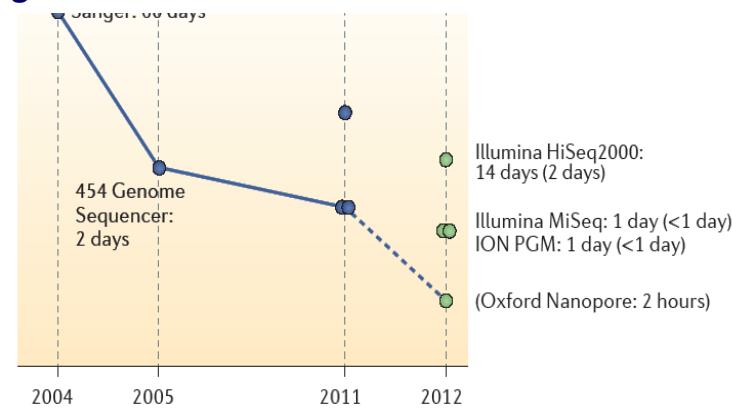
Raw daily output



Cost per Mb assembled sequence



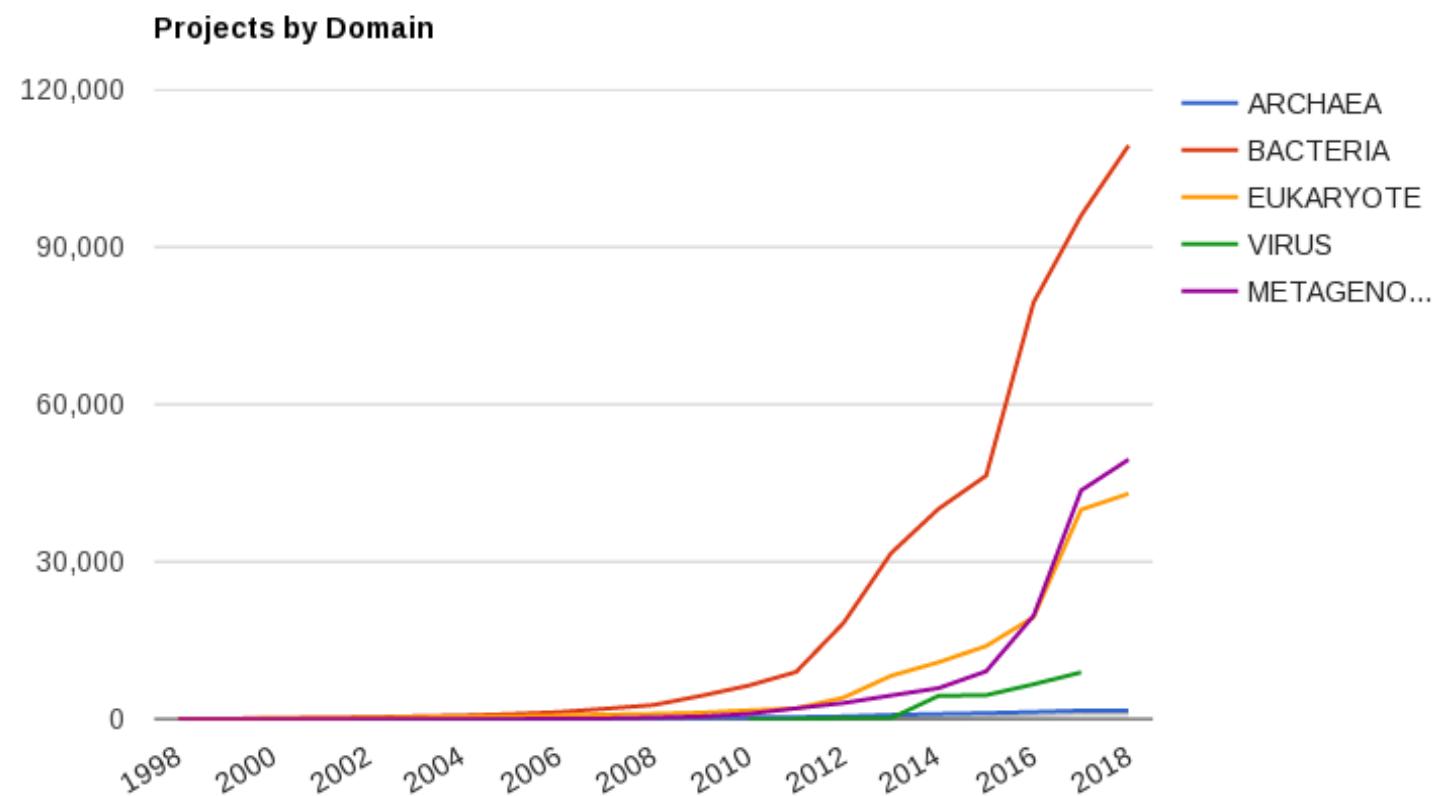
Turnaround time: bacterial genome



Sequencing projects

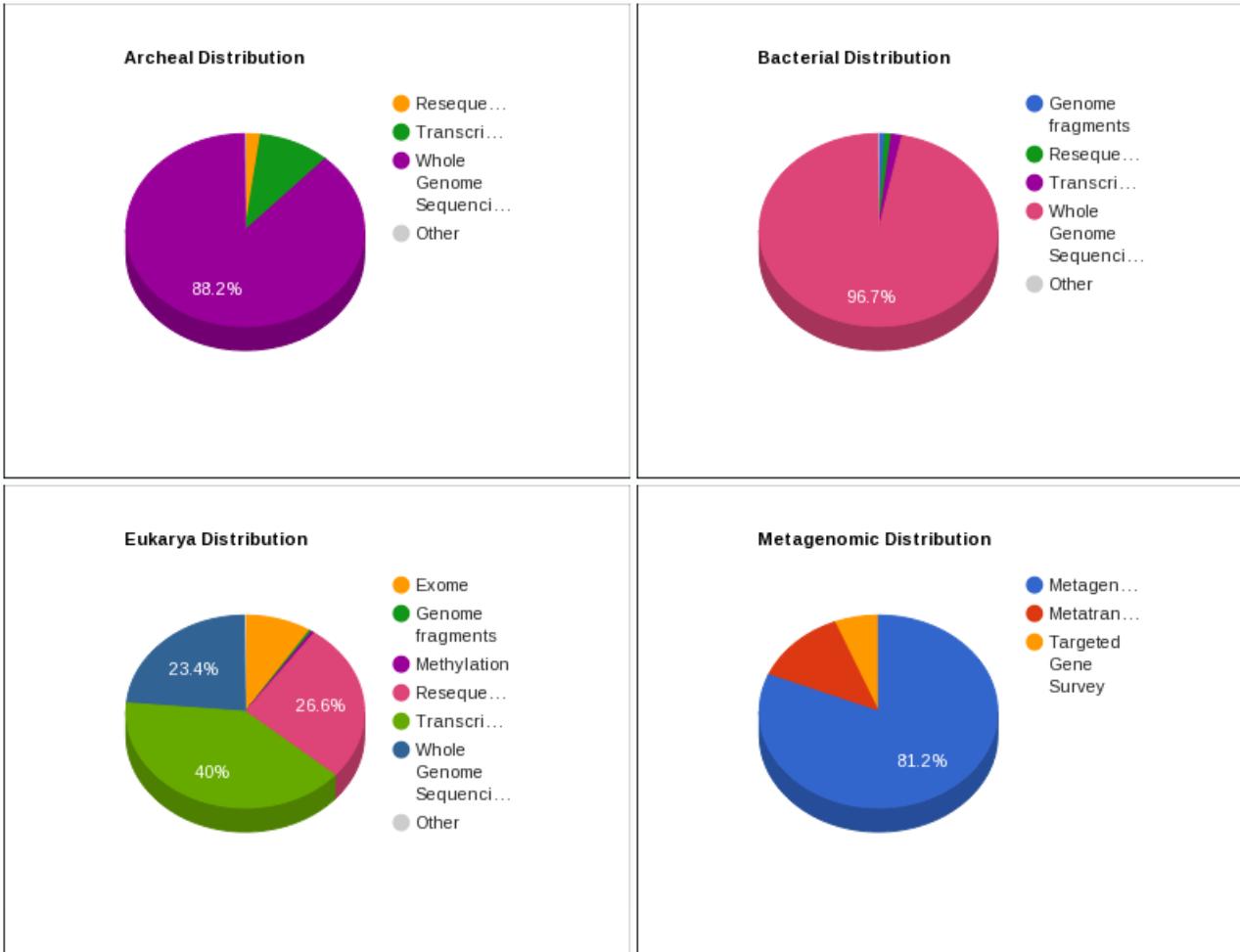
<https://gold.jgi.doe.gov/>

GOLD, Genome Online DataBase

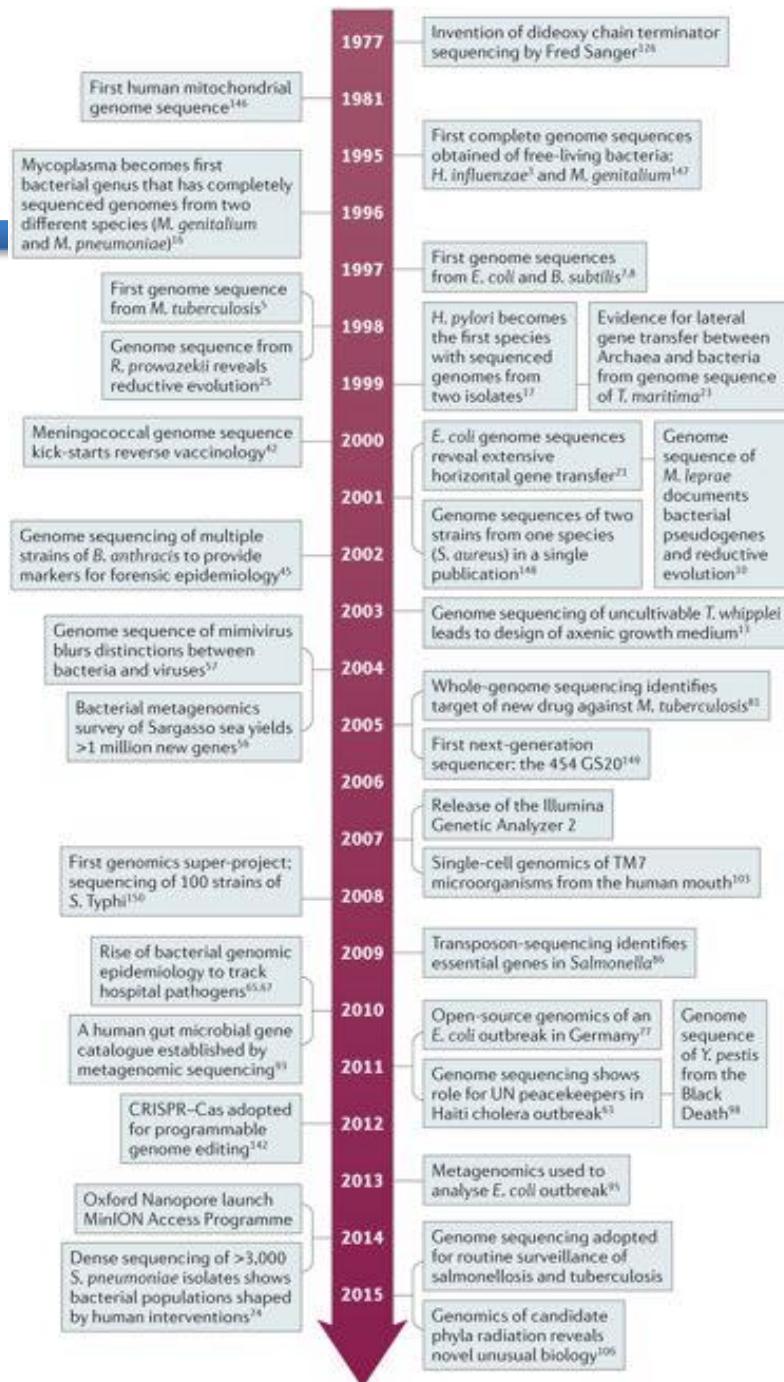


Organisms distribution of Sequencing projects

GOLD Project Distributions



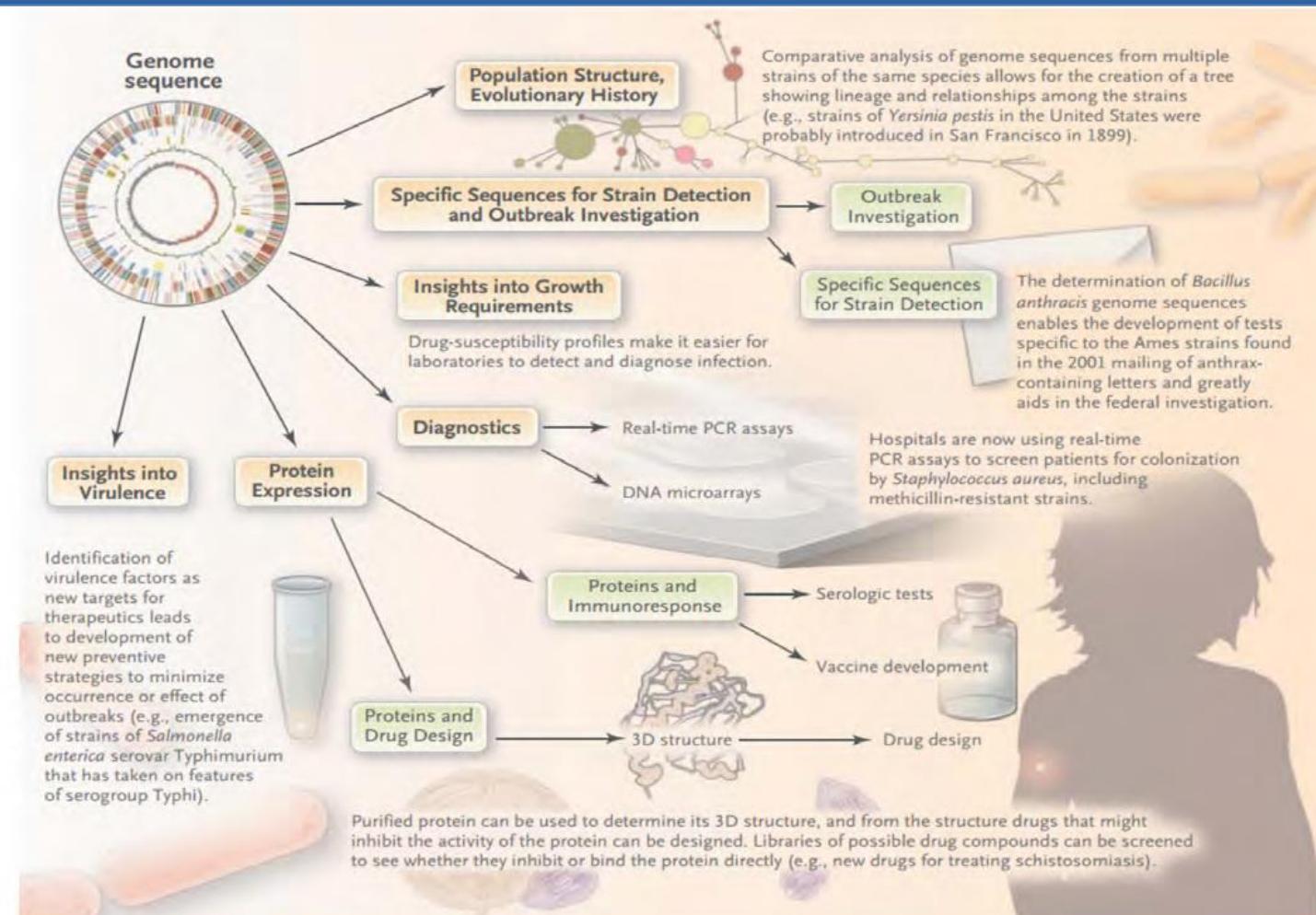
Brief history of the major events that have shaped the sequencing and analysis of bacterial genomes in the past two decades



de genomas bacterianos: herramientas y aplicaciones

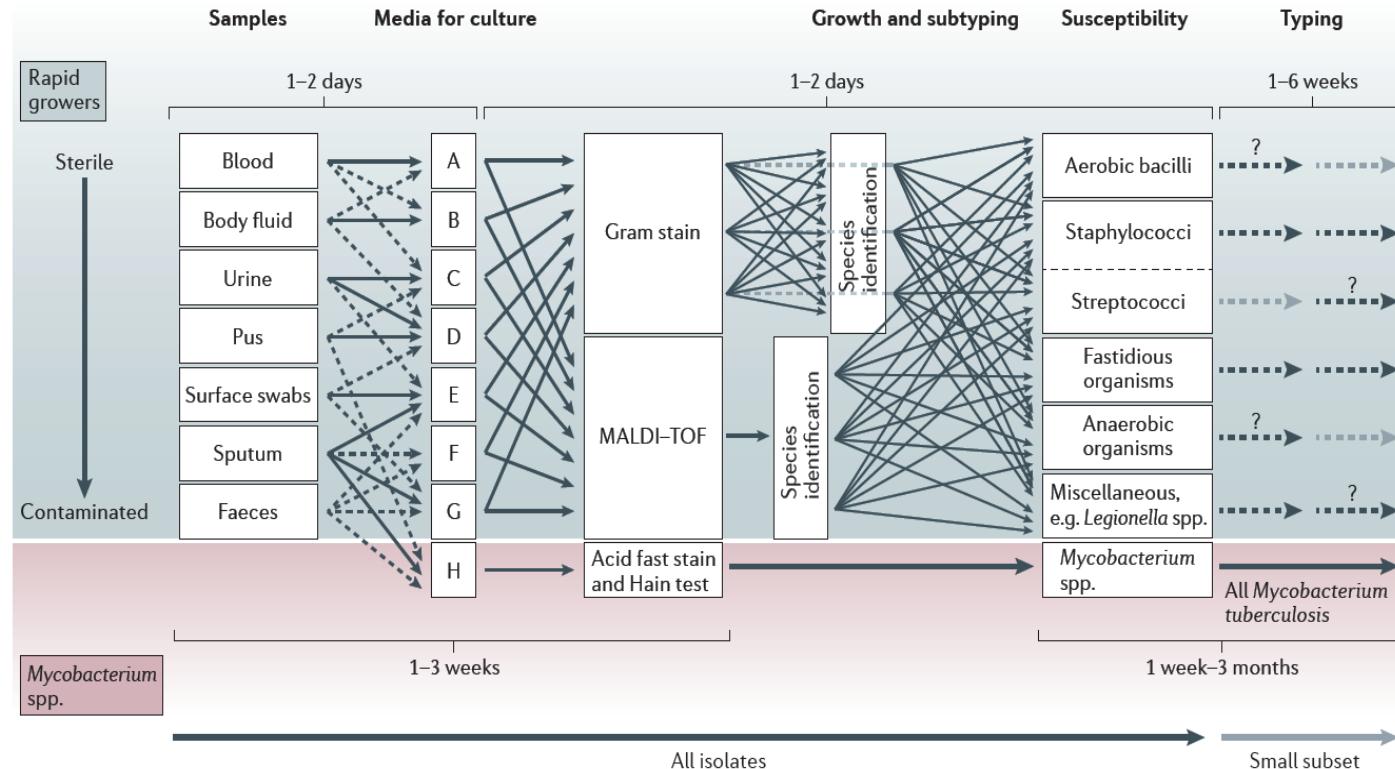
Use of microbial genomics for tool development

Report from The American Academy of Microbiology, 2015



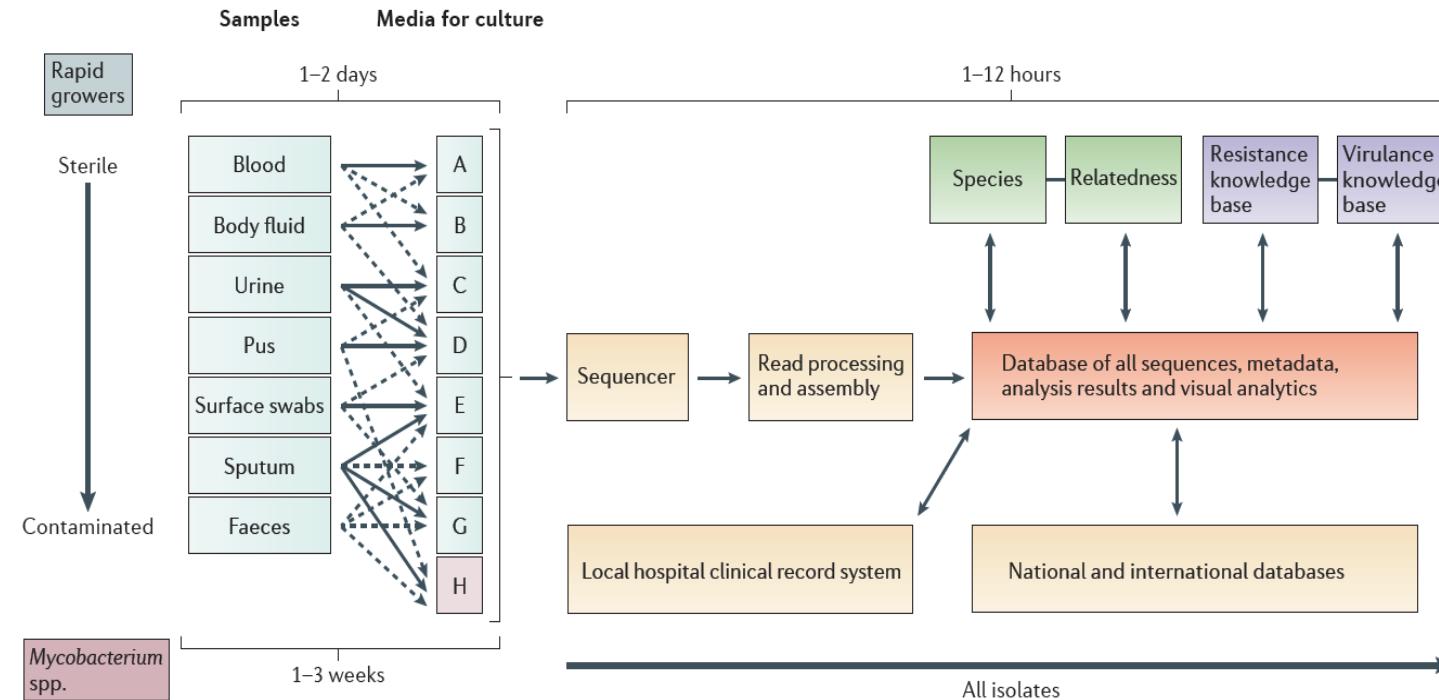
Workflow for processing samples for bacterial pathogens

Didelot et al., Nature Genet Review 2012, 13:601-612



Ongoing developments in DNA-sequencing technologies are likely to affect the diagnosis and monitoring of all pathogens, including viruses, bacteria, fungi and parasites.

The diagnostic and clinical applications of bacterial WGS

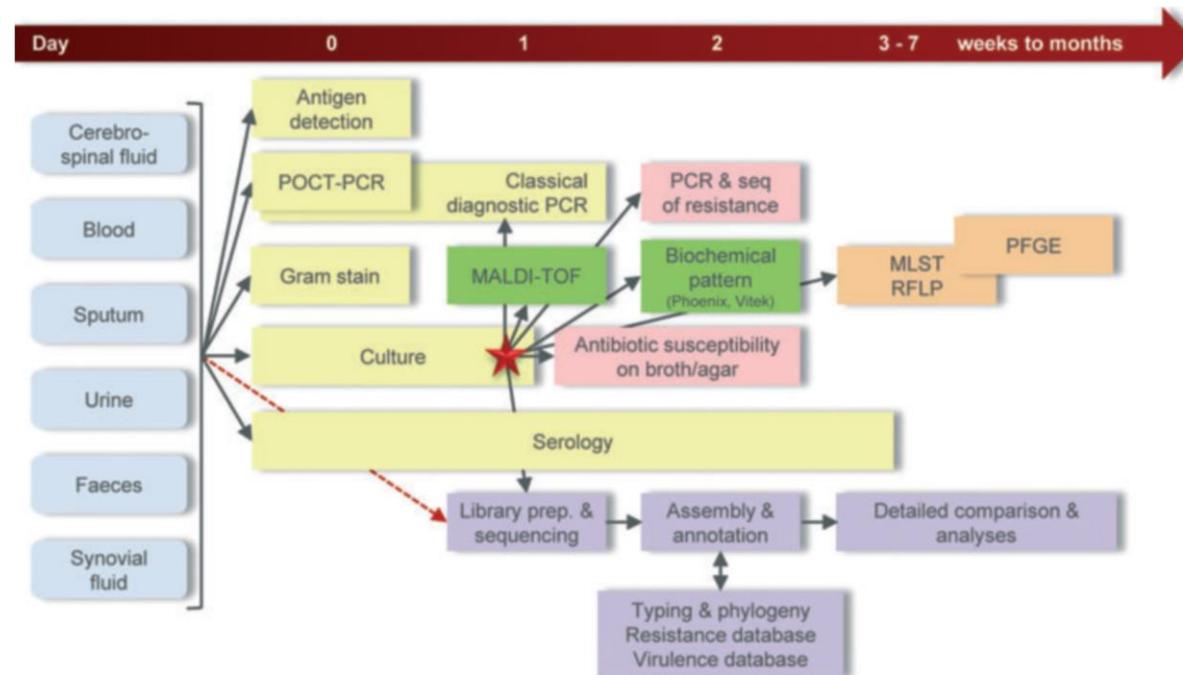


Didelot et al., Nature Genet Review 2012, 13:601-612

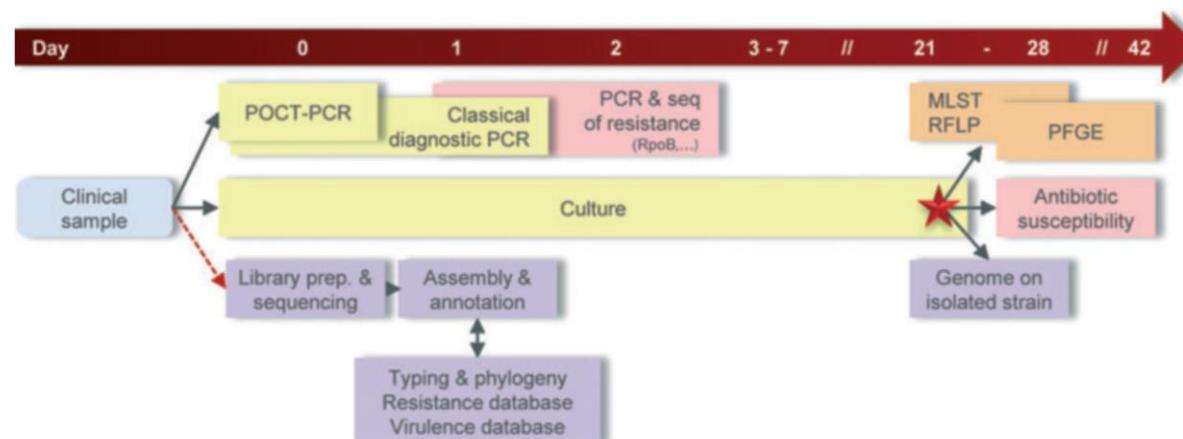
Schematic representation of the timeline for the processing of clinical samples with classical pathogens

Bertelli and Greub, Clin Microb and Infect, 2013

(a) classical
pathogens



(b) slow-growing
bacteria such as
Mycobacterium
tuberculosis



Foodborne outbreak identification “Crisis del pepino”

2011

Mayo

- 24 Primera muerte en Alemania
- 26 Alemania acusa a los pepinos españoles
- 30 Prohibición de importaciones de verduras de España y Alemania
- 31 Laboratorios alemanes desmienten oficialmente que los pepinos españoles sean el foco de infección

Junio

- 10 Resolución de la crisis

**Secuenciación
Genoma**

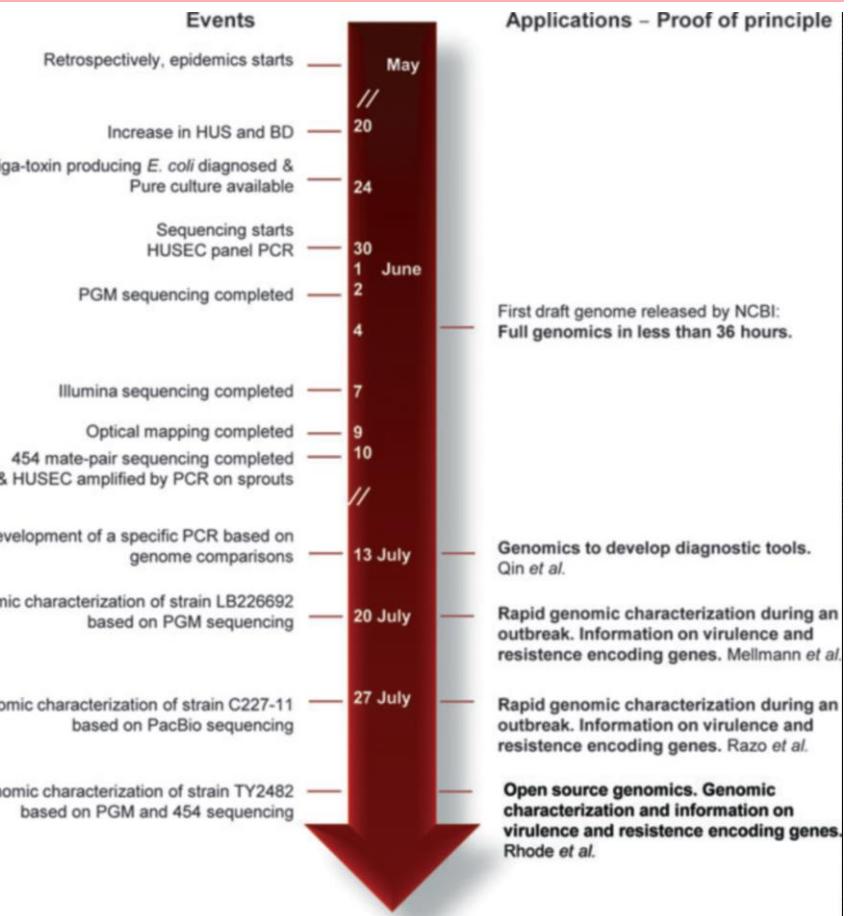


华大基因
BGI

Universitätsklinikum
Hamburg-Eppendorf

Secuenciación de genomas bacterianos: herramientas y aplicaciones

The Escherichia coli O104:H4 epidemics: event timeline and major outputs



Foodborne outbreak identification “Crisis del pepino”

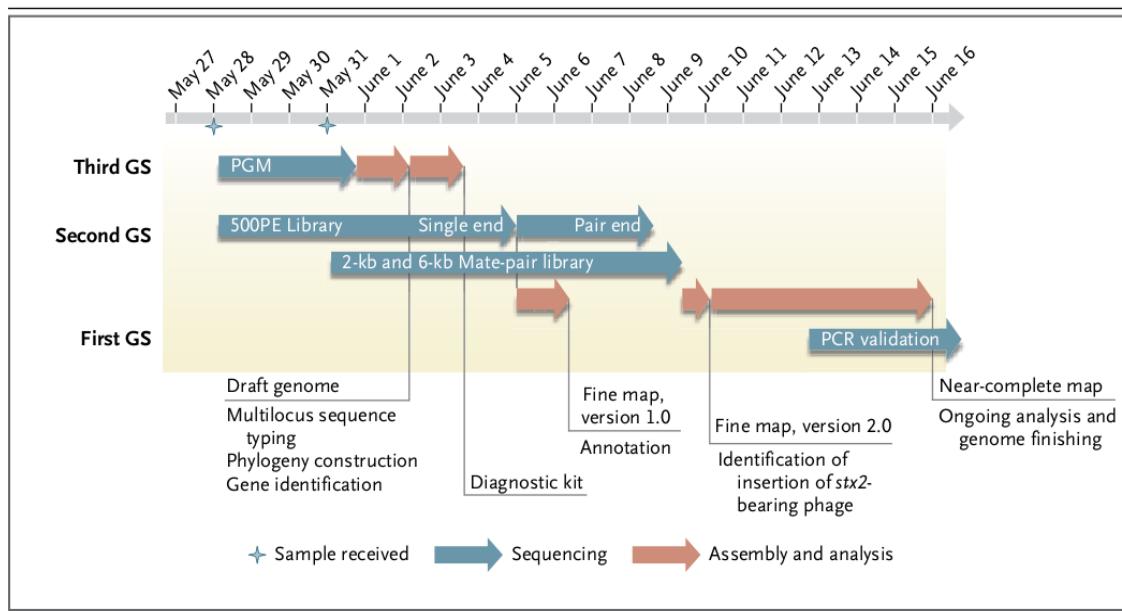


Figure 1. Timeline of the Open-Source Genomics Program.

After receiving the first batch of DNA samples on May 28, 2011, sequencing runs with the use of the Ion Torrent Personal Genome Machine (PGM) and Illumina (small-insert library) were initiated simultaneously. On May 31, the second batch of DNA was received and used for Illumina large-insert sequencing. An assembly of the Ion Torrent reads was released on June 2, which enabled subsequent analyses (multilocus sequence typing, phylogenetic analysis, and genome comparisons). Errors in the Ion Torrent data were corrected with the use of later Illumina data, and a high-quality draft genome sequence was created. GS denotes generation of sequencing technology. The symbols at May 28 and May 31 in the timeline indicate the arrival of DNA samples.

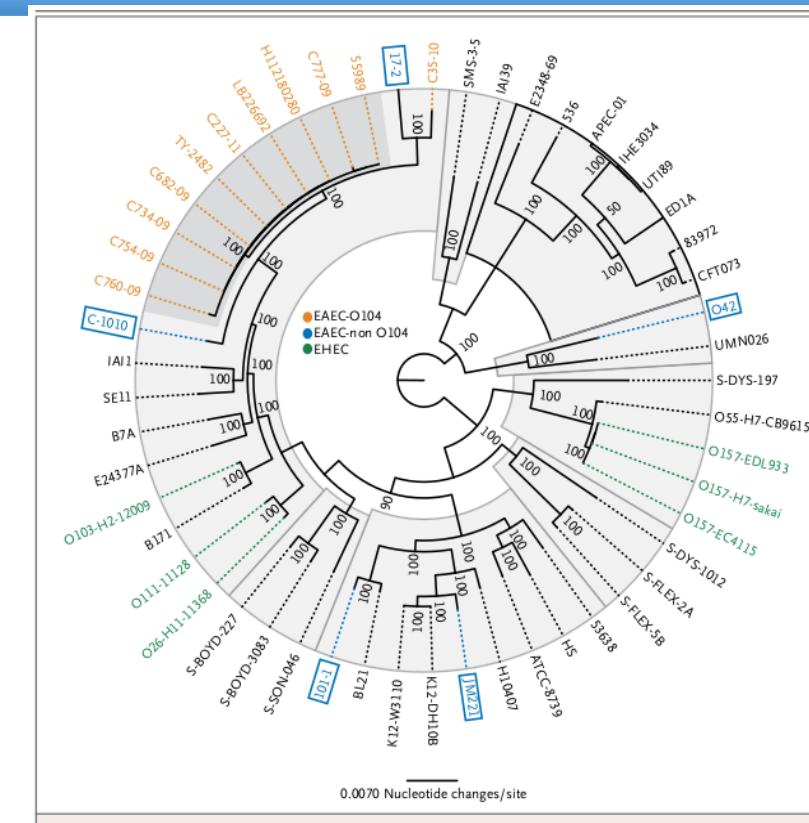


Figure 2. Phylogenetic Comparisons of 53 *Escherichia coli* and *Shigella* Isolates.

Genomic sequences were compared with the use of 100 bootstrap calculations, as described by Sahl et al.³⁵ The species-based phylogeny was inferred with the use of 2.56 Mbp of the conserved core genome. The O104:H4 isolates are shown in orange, the reference enteroaggregative *E. coli* (EAEC) isolates in blue, and the enterohemorrhagic *E. coli* isolates in green. (The classification of the other strains is shown in Fig. 4 and Table 4 in the Supplementary Appendix.) The O104:H4 isolates cluster into a single clade (dark gray); in contrast, the reference EAEC isolates are extremely divergent and are represented throughout the phylogeny.

Andalusian Listeria Outbreak

Actualización de información sobre el brote de intoxicación alimentaria causado por *Listeria monocytogenes*.

Publica: Agencia Española Seguridad alimentaria y Nutrición
Fecha: 29 agosto 2019
Sección: Seguridad Alimentaria

Jueves 29 de agosto de 2019, 12.00 horas

ACTUALIZACIÓN EN RELACIÓN CON LA DISTRIBUCIÓN DE PRODUCTOS RELACIONADOS CON LA ALERTA.

La Agencia Española de Seguridad Alimentaria y Nutrición (AESAN) recomienda a las personas que tengan en su domicilio algún producto de la marca "La Mechá" se abstengan de consumirlo. Si se dispone del producto se debe devolver al punto de compra y, de no ser posible, desecharlo.

Brote de listeriosis: sube el número de afectados y se apunta a la falta de higiene en la carne como causa

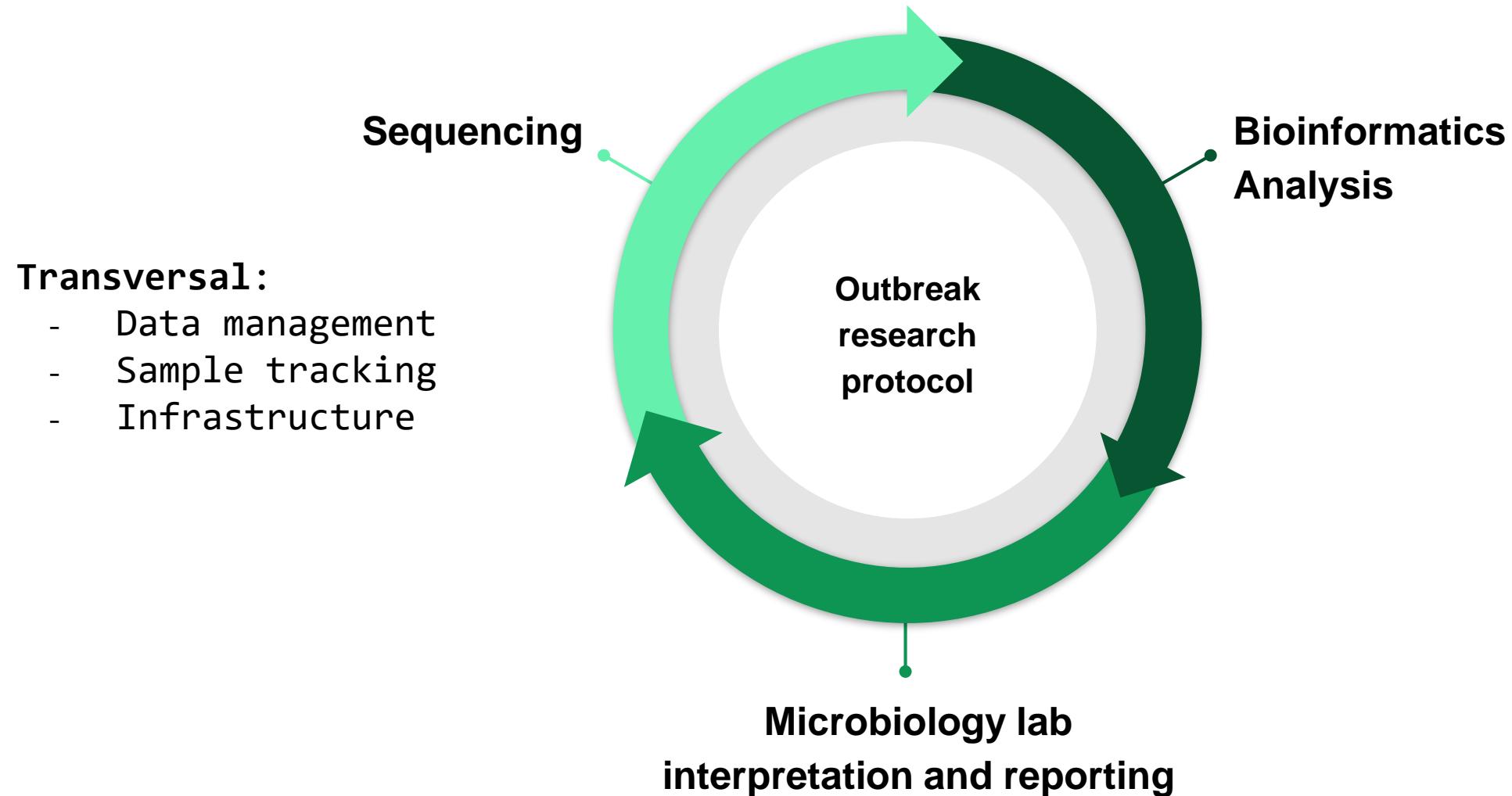
EFE 25.08.2019

- Tres nuevos casos, en Sevilla y Cádiz, dejan el número de personas afectadas en Andalucía en 192.
- [La carne con listeria de la marca blanca se vendió en los municipios de Sevilla.](#)
- La empresa que vendió la marca blanca de Magrudis dice que cumple los protocolos.



- Meat “La Mechá”. Margulis S.L.
- 250 cases related.
- Meat “"La Montanera del Sur". INCARYBE S.L”, suspicion. (Cádiz)
- Meat “Sabores de Paterna” (Málaga)

Andalusian Listeria Outbreak

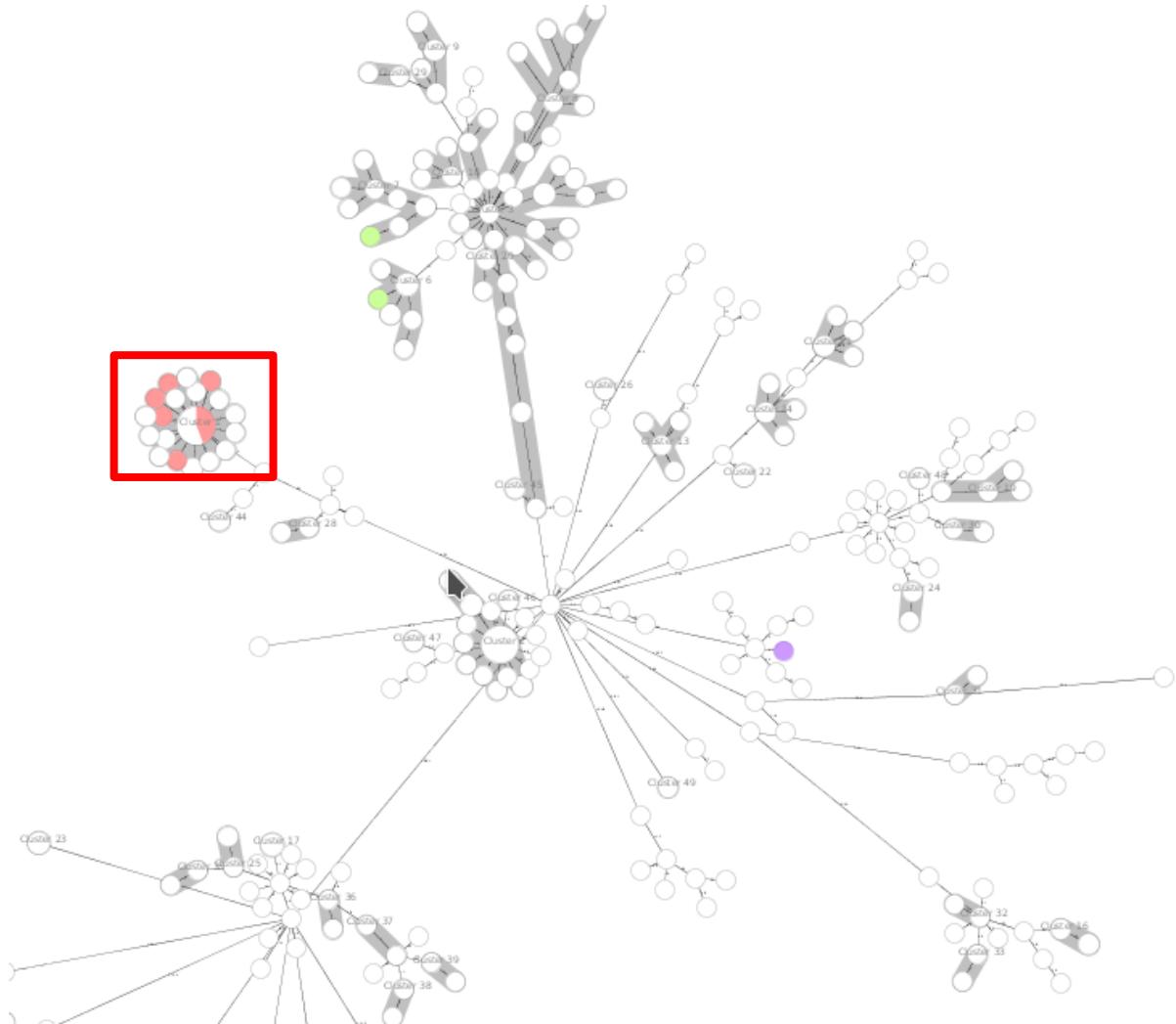


Andalusian Listeria Outbreak

- 625 listeria samples already sequenced
 - 258 suspected to be related to the outbreak (mid august to mid september)

Results:

- 233 related to the outbreak, confirmed to be caused by the meat “La Mechá”
 - 25 sporadic cases not related to the outbreak.



PREPARACIÓN LIBRERÍA, estrategias

SECUENCIACIÓN GENOMA, EXOMA, TRANSCRIPTOMA

1. Sin amplificación
2. Amplificación con PCR
3. Sondas captura

- Tamaño de fragmento
- Longitud de la lectura
- Single o Paired-end
- Número de bases por muestra
- Profundidad de cobertura x

SECUENCIACIÓN GENOMAS

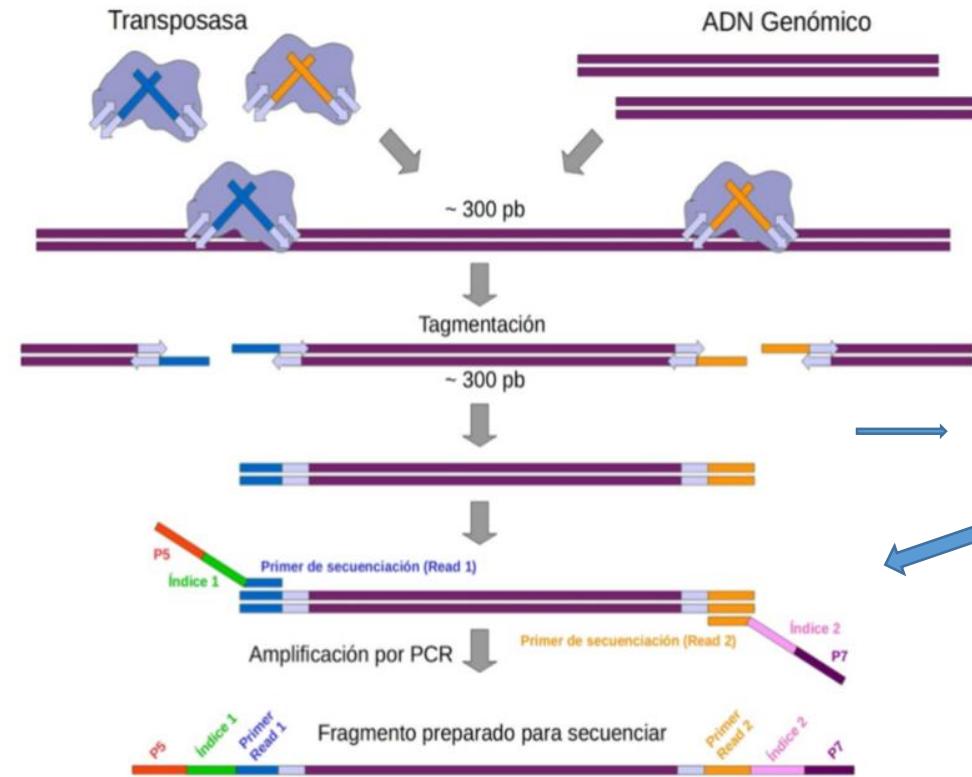
1. Metagenómica

IDENTIFICACIÓN MICROORGANISMOS

1. Metataxonomía

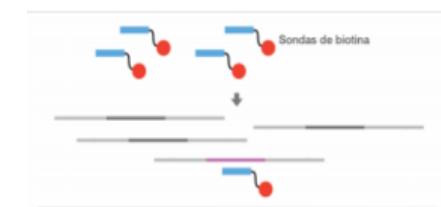
PREPARACIÓN LIBRERÍA

ENZIMÁTICA FÍSICA



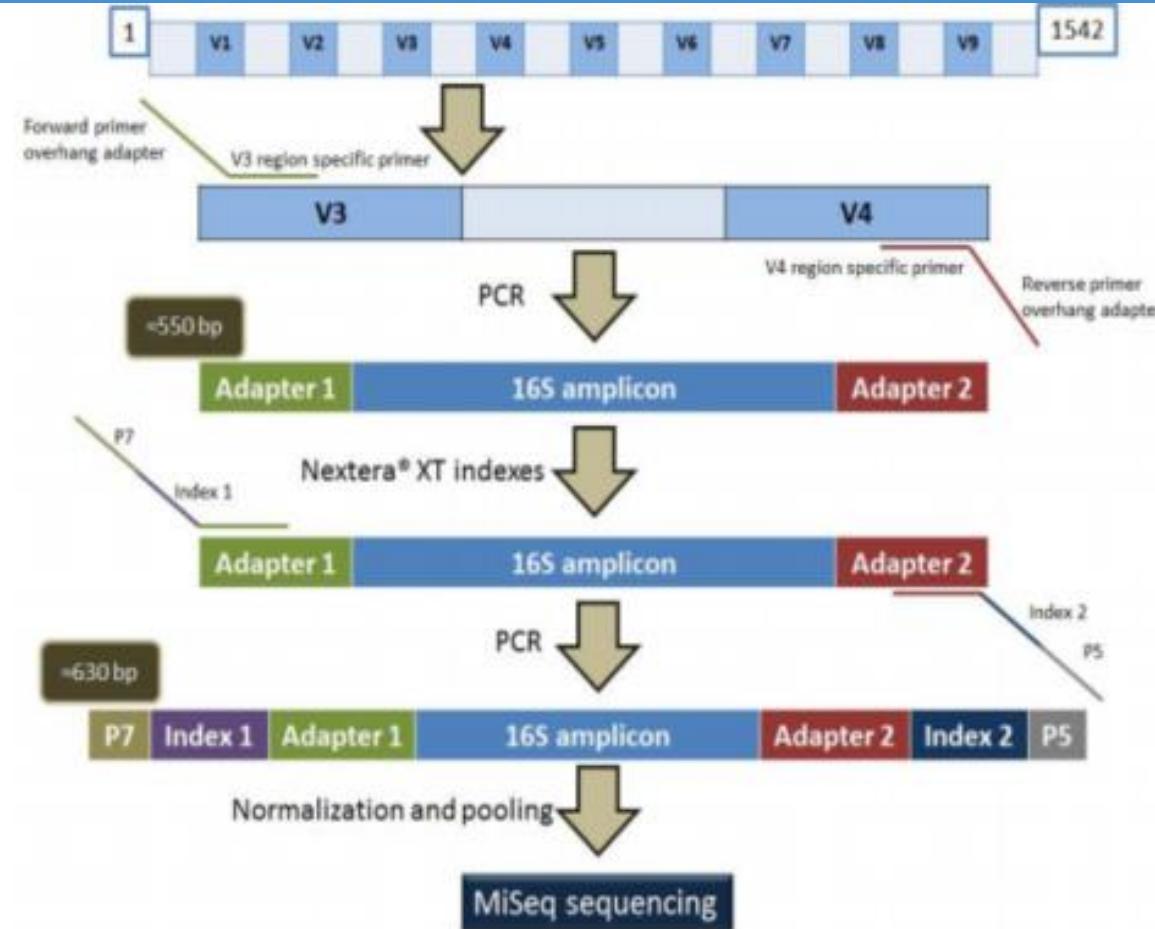
RNA → cDNA

ENRIQUECIMIENTO:
PCR
CAPTURA SONDAS



Guia Práctica Genómica https://www.uv.es/varnau/GM_Cap%C3%ADtulo_2.pdf

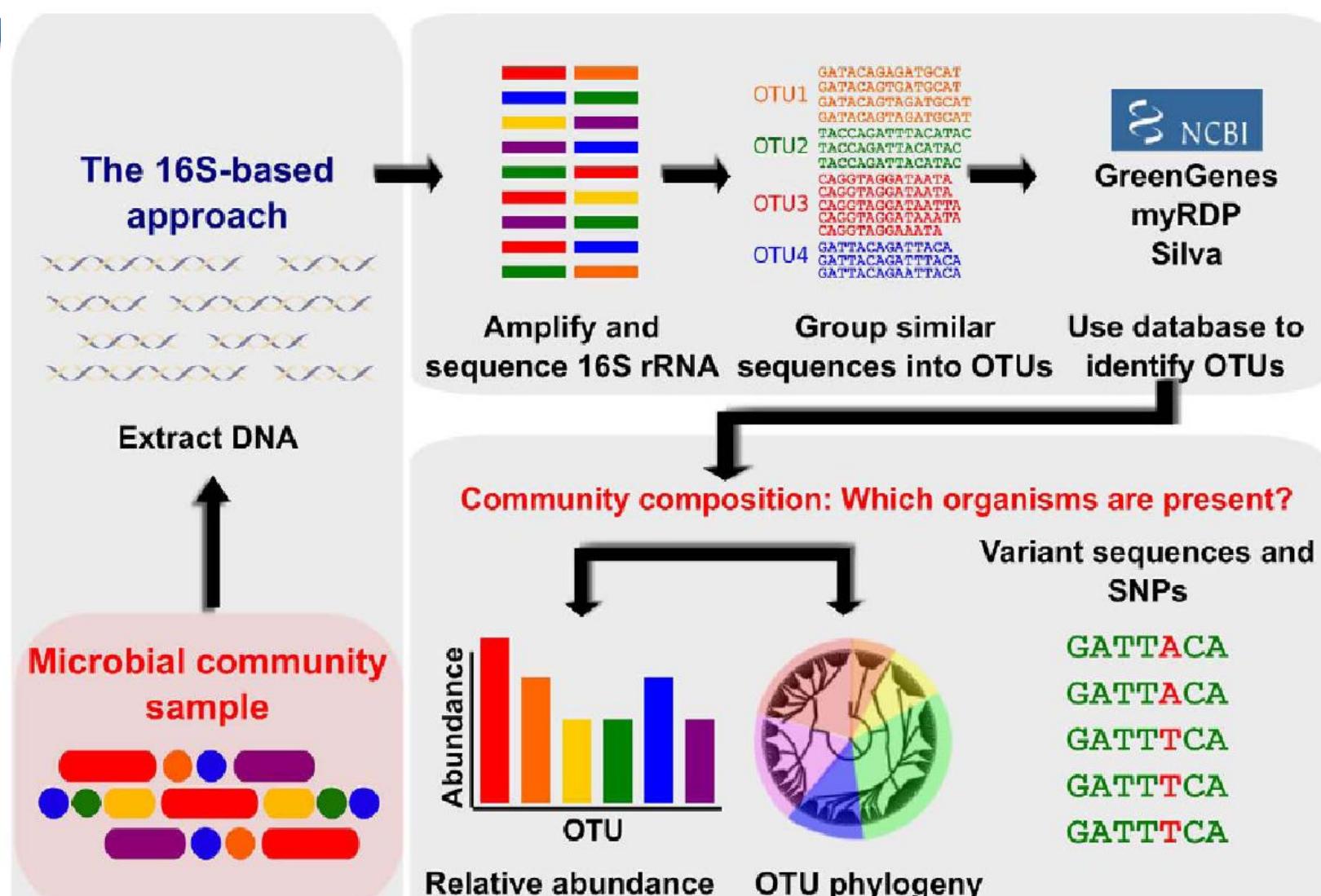
PREPARACIÓN LIBRERÍA, rRNA 16S, caracterización microbiota



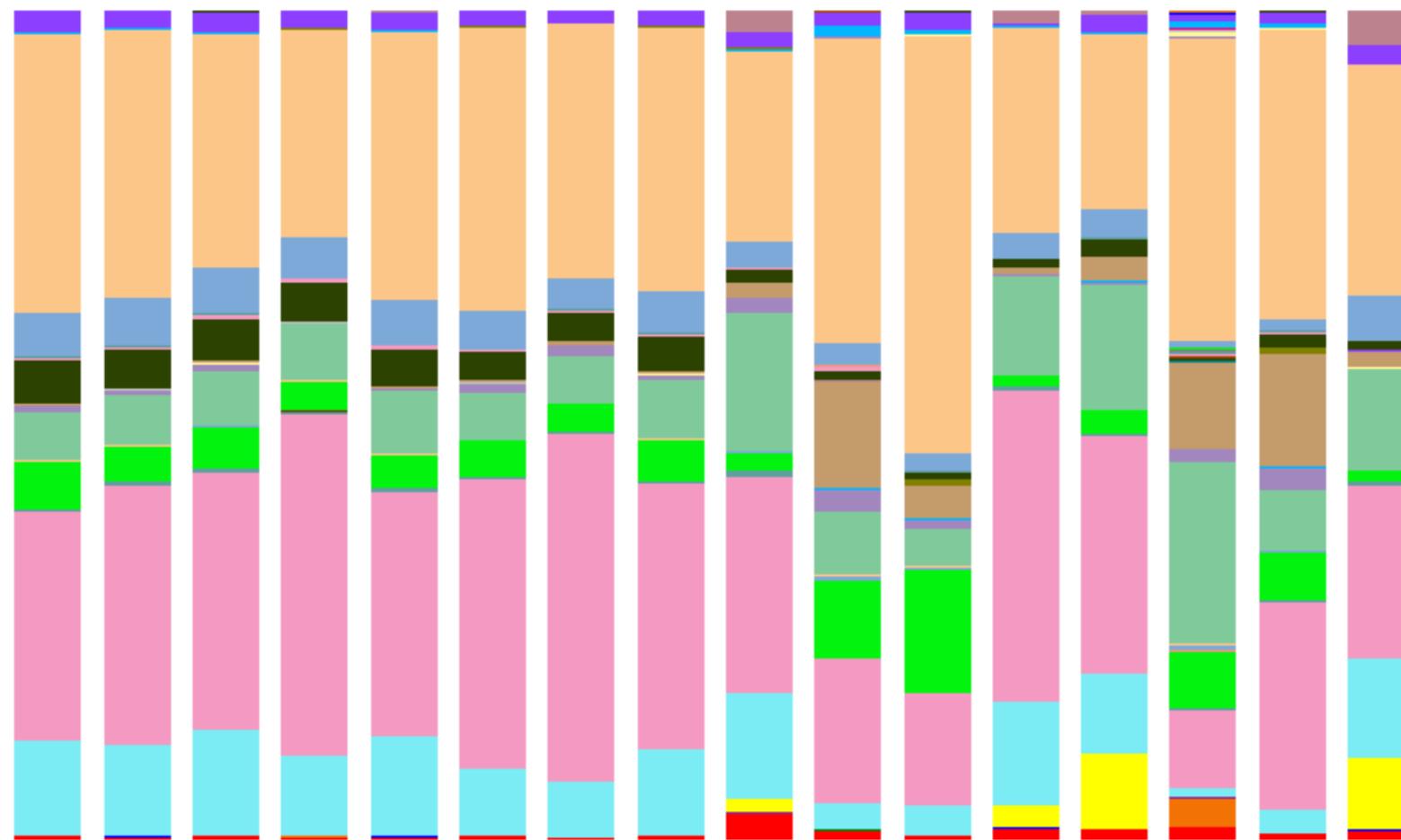
Metataxonomics vs Metagenomics (16S vs Shotgun)

| | Metagenetics | Metagenomics |
|------------------------------------|----------------|--------------|
| Amplified sequence | Marker regions | Whole genome |
| Computing time | Usually short | Usually long |
| Taxonomic composition | Yes | Yes |
| New pathogen detection | No | Yes |
| Genome coverage information | No | Yes |

Metataxonomics



Taxonomy summary (i.e. phylum level)

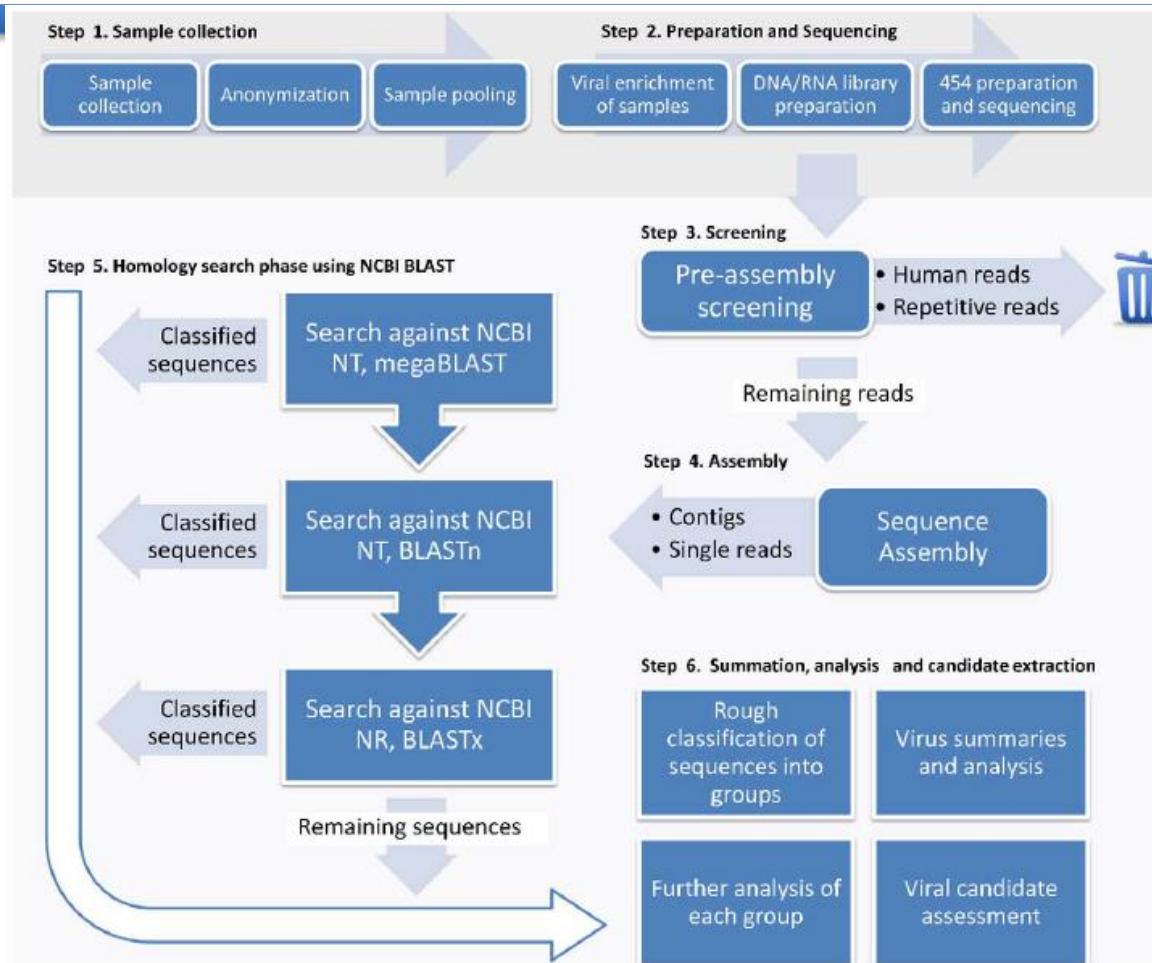


Metataxonomics

Problemas:

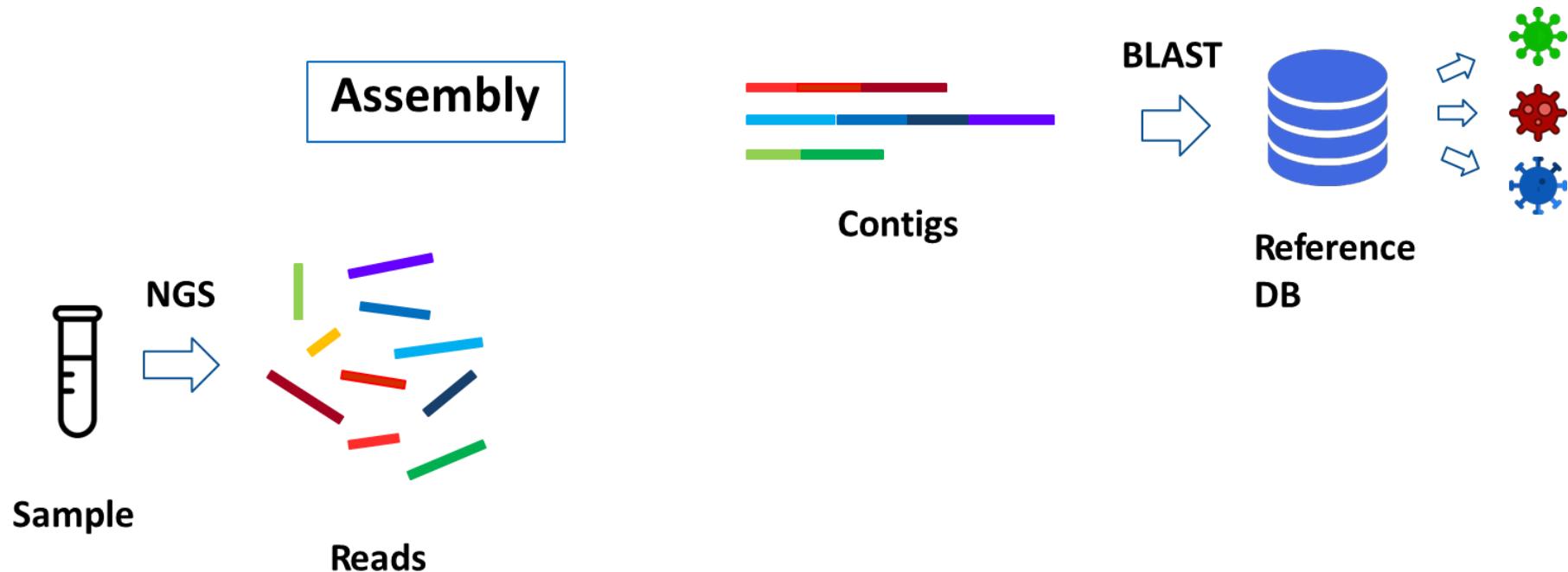
- Raros en el genoma (< 0.1%)
- Los trozos similares dificultan el ensamblado correcto de lecturas pequeñas
- No todos los rRNA se amplifican en la misma medida con los *primers* universales
- Especies con diversas copias de sus genes rRNA
- **No se conoce un umbral fijo de similitud que separe especies**
- **Tendencia a producirse quimeras en la PCR**

Metagenómica, pipeline de análisis

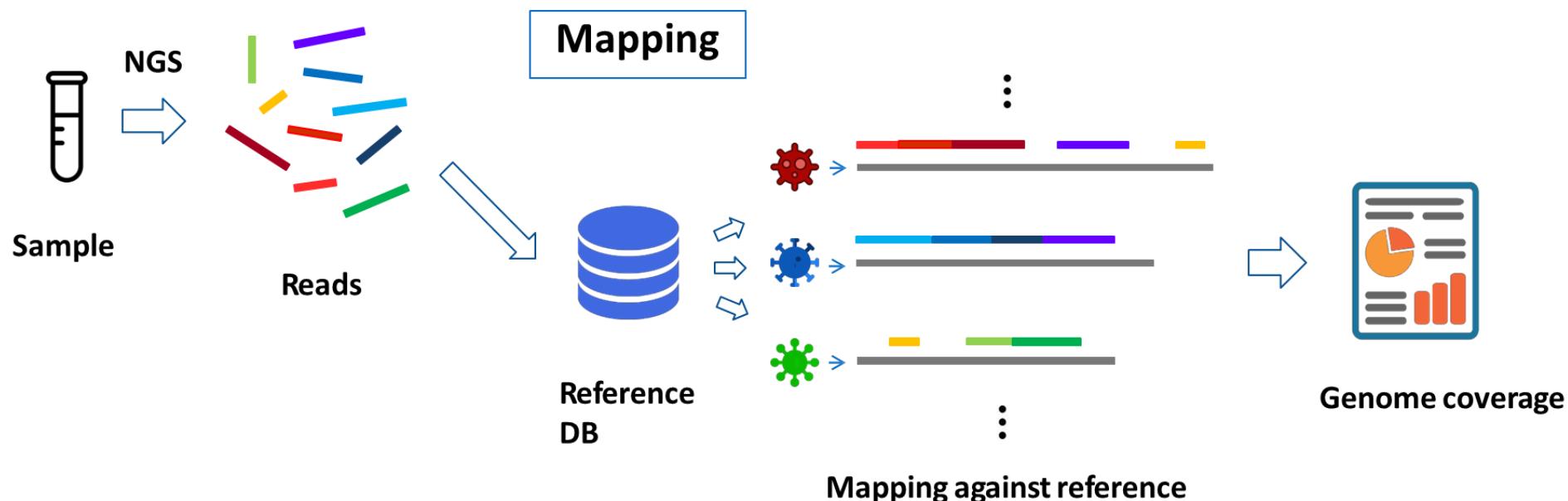


Lysholm et al., Plos One 2012:7,2, e30875

Metagenomic analysis approaches



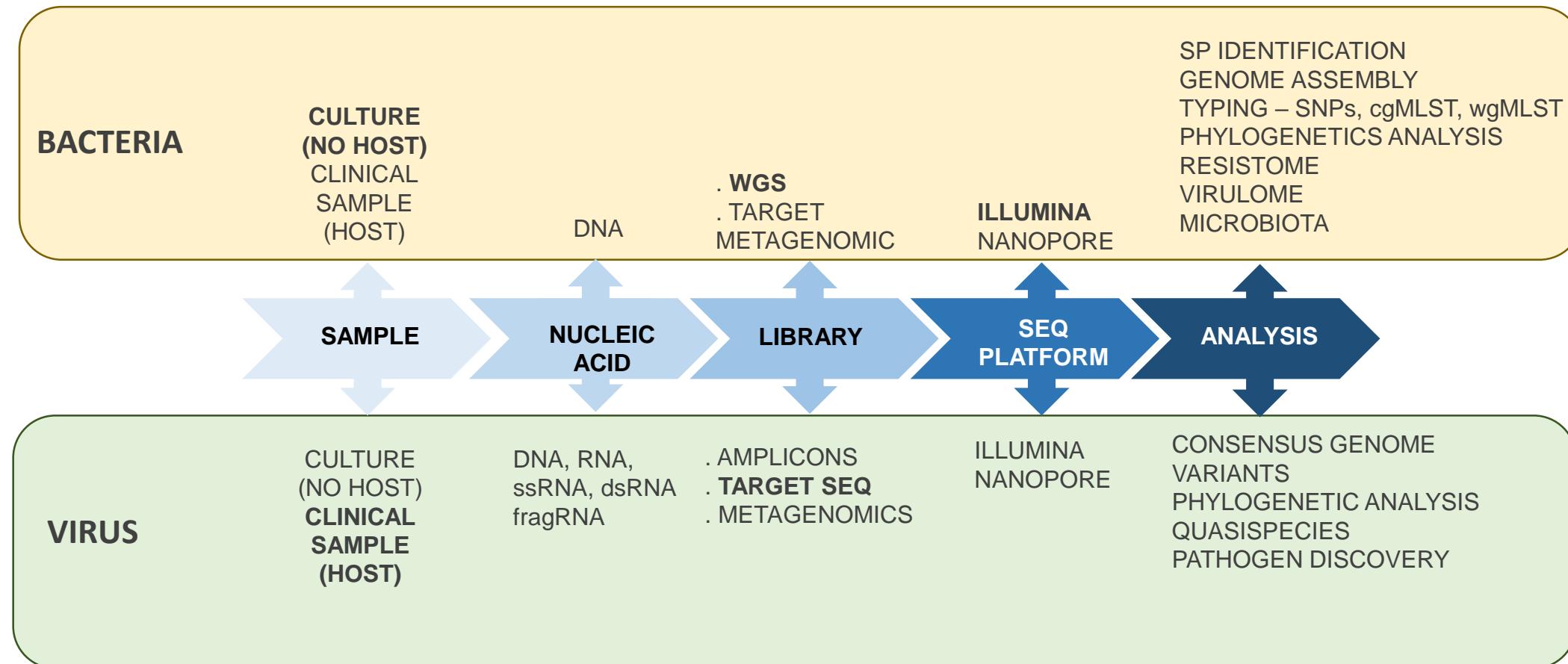
Metagenomic analysis approaches



Metataxonomics vs Metagenomics (16S vs Shotgun)

| Software | Organism | Genetic portion used | | Binning algorithm used | | | Genome coverage | Novel pathogen discovery |
|-------------|--------------------------------|----------------------|--------------|------------------------|---------|----------|-----------------|--------------------------|
| | | Genetic markers | Whole Genome | Clustering | Mapping | Assembly | | |
| Mothur | Bacteria | X | | X | | | No | No |
| QIIME | Bacteria | X | | X | | X | No | No |
| MEGAN | Bacteria | | X | | | X | No | No |
| Platypus | Bacteria | | X | | X | | No | No |
| SURPI | Virus | | X | | | X | No | Yes |
| Virus-TAP | Virus | | X | | | X | No | Yes |
| VIP | Virus | | X | | X | | No | Yes |
| Pathosphere | Virus, Bacteria, Eukarya | | X | | | X | No | Yes |

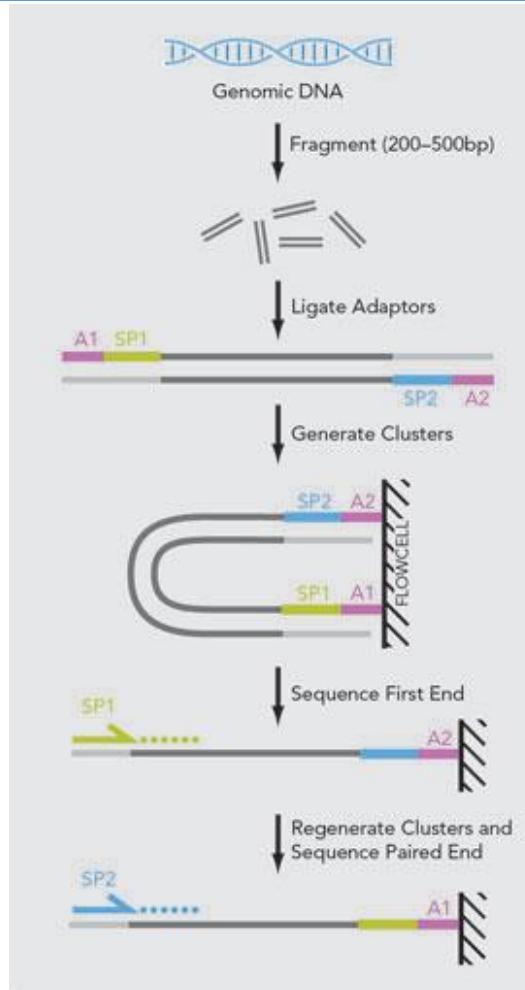
Bacterial and viral Genome Sequencing



Bioinformatics analysis in microbial genomics

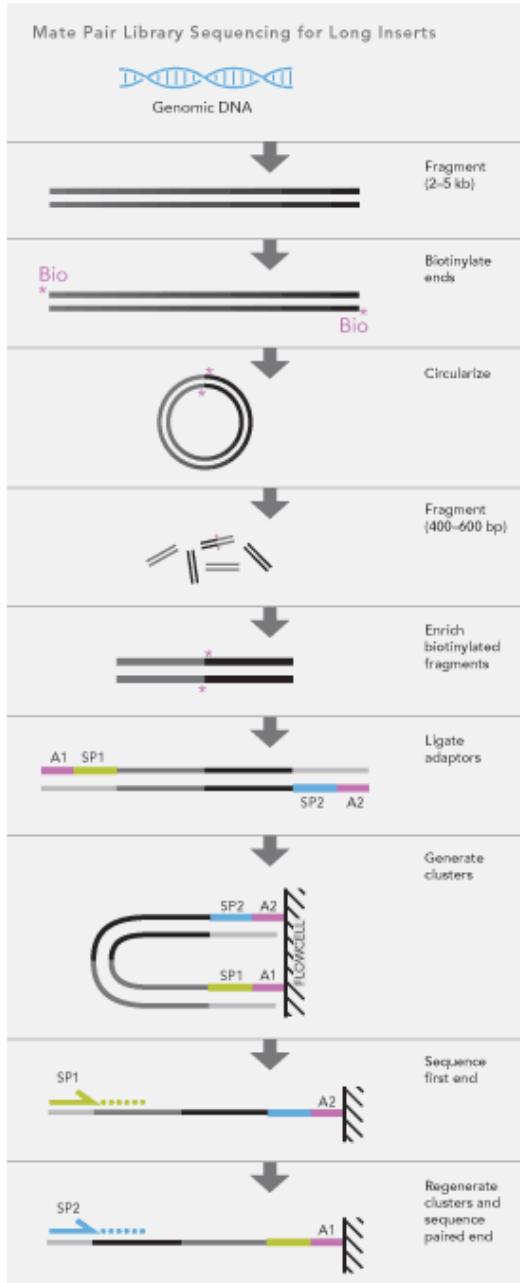
- SPECIE IDENTIFICATION
 - WGS - Kmers analysis
 - TARGET METAGENOMIC, rRNA - MICROBIOTA
- ASSEMBLY GENOME
 - de NOVO or REFERENCE -BASED
 - cgMLST, wgMLST - MINIMUM SPANING TREE
 - METAGENOMIC - HOMOLOGY -BASED
- VARIANT CALLING
 - REFERENCE GENOME SELECTION
 - HAPLOID GENOME
 - LOW FREQUENCY VARIANT - QUASISPECIES
 - SNPs MATRIX - PHYLOGENETIC ANALYSIS
- STRUCTURAL AND FUNCTIONAL ANNOTATION
 - RESISTOME, VIRULOME, SEQUENCE-TYPE

Que es Pair-end?



Secuenciación de un fragmento (bp)

**Modificación de single-read DNA,
Leyendo por ambos extremos, forward y reverse**



Mate Pair library preparation is designed to generate short fragments that consist of two segments that originally had a separation of several kilobases in the genome. Fragments of sample genomic DNA are end-biotinylated to tag the eventual mate pair segments. Self-circularization and refragmentation of these large fragments generates a population of small fragments, some of which contain both mate pair segments with no intervening sequence. These Mate Pair fragments are enriched using their biotin tag. Mate Pairs are sequenced using a similar two-adapter strategy as described for paired-end sequencing.

Que es Mate-pair?

Secuenciación de dos fragmentos separados kb.

Util:
Secuenciación de un Genoma de novo
Finalizar un genoma
Detección de variantes estructurales

Sequencing terms

Depth of coverage

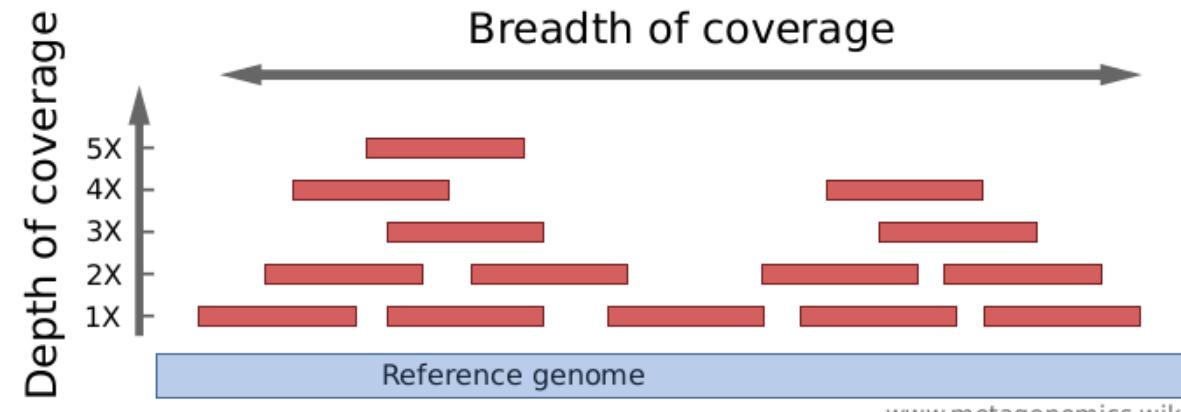
How strong is a genome "covered" by sequenced fragments (short reads)?

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).

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.



Calculo de cobertura: número de lecturas

Estimating Sequencing Runs

Coverage Equation

The Lander/Waterman equation is a method for computing coverage¹.

The general equation is:

$$C = LN / G$$

- C stands for coverage
- G is the haploid genome length
- L is the read length
- N is the number of reads

So, if we take one lane of single read human sequence with v3 chemistry, we get

$$C = (100 \text{ bp}) * (189 \times 10^6) / (3 \times 10^9 \text{ bp}) = 6.3$$

This tells us that each base in the genome will be sequenced between six and seven times on average.

Sequencing Coverage Calculator

Support Center:
Sequencing Coverage Calculator

Application or product: Whole-Genome Sequencing

Coverage: 100 x

Duplicates: 2 %

Genome or region size (in million bases): 3300 Mb

Total read length (e.g. 200 for 2x100): 600 cycles

Benchtop Sequencers Production-Scale Sequencers

iSeq NextSeq 500/550

MiSeq NovaSeq 6000

MiSeq / MiSeq Dx in RUO mode HiSeq 3000/4000

NextSeq 500/550 HiSeq 1500/2500 Rapid Run

Exceeds maximum read length?

Number of units per sample (flow cell or lane)

Samples per unit (flow cell or lane)

Comments

Products

Support Center:
Sequencing Coverage Calculator

Thank you for using the Illumina coverage estimator.

The results were calculated based on: **coverage needed**. Explain the estimations

| Run type | MiSeq | MiSeq | MiSeq | MiSeq |
|--|---|---|--|-------|
| v3 Reagents | v2 Reagents | v2 Nano Reagents | v2 Micro Reagents | |
| 25,000,000 per flow cell | 15,000,000 per flow cell | 1,000,000 per flow cell | 4,000,000 per flow cell | |
| 15,000,000 per flow cell | 9,000,000 per flow cell | 600,000,000 per flow cell | 2,400,000,000 per flow cell | |
| Does not exceed maximum (2x300) | Read length exceeds maximum of 2x250 | Read length exceeds maximum of 2x250 | Read length exceeds maximum of 2x150 | |
| Number of units per sample (flow cell or lane) | 22,449 flow cells | 37,415 flow cells | 561,224 flow cells | |
| Samples per unit (flow cell or lane) | -0/flow cell | -0/flow cell | -0/flow cell | |
| Comments | Upgraded software: MCS v2.3 or later; MiSeq | Upgraded hardware or from September 2012 and later: MCS v2.0 or later; MiSeq Reagent Kit v2 (150/600) | Upgraded hardware or from September 2012 and later: MCS v2.0 or later; MiSeq Reagent Nano Kit v2 (300/500) | |
| Products | MiSeq Reagent Kit v3 | MiSeq Reagent Kits v2 | MiSeq Reagent Kits v2 | |

Get the results in a comma-separated values (CSV) report.

https://emea.support.illumina.com/downloads/sequencing_coverage_calculator.html

The idea that via a **One Health** approach infectious diseases can be better controlled and prevented



Global Microbial Identifier



Spanish National Microbiology Center (CNM)



Mission: Provide support to the National Health System and the different Spanish Regions in the diagnosis and control of infectious diseases. In order to fulfill this mission it acts as Reference center offering a series of scientific activities:

- Diagnosis
- **Surveillance** →
- Infectious diseases research
- Training

Outbreak research:
Molecular source
detection

ECDC roadmap and international commitment



EUROPEAN CENTRE FOR
DISEASE PREVENTION
AND CONTROL

**ECDC roadmap for integration
of molecular and genomic
typing into European-level
surveillance and epidemic
preparedness**

Version 2.1, 2016–2019

www.ecdc.europa.eu

**ECDC strategic framework for
the integration of molecular and
genomic typing into European
surveillance and multi-country
outbreak investigations**

2019–2021

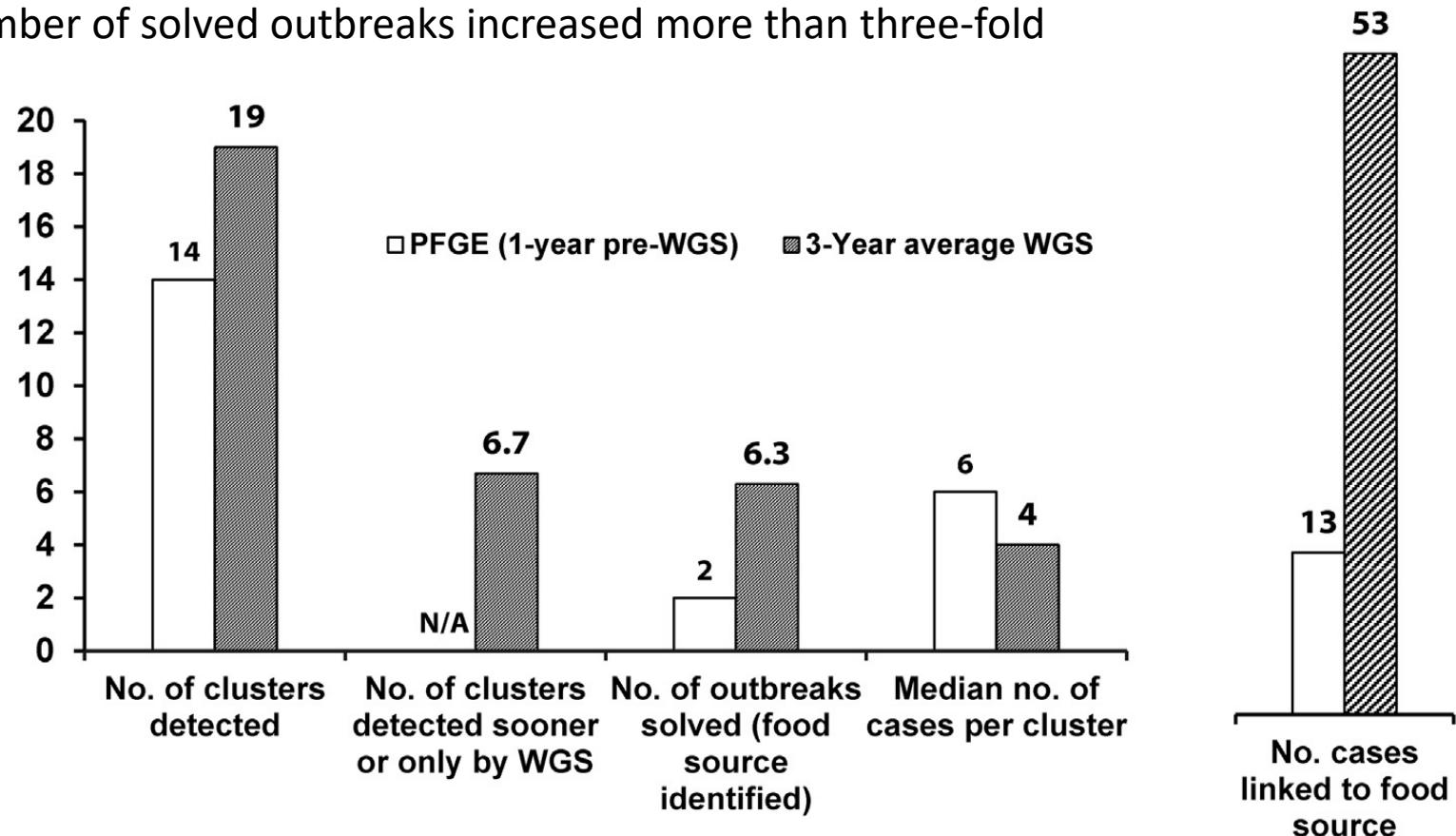
www.ecdc.europa.eu

- **Operationalisation of EU-wide WGS-based surveillance systems in the near term:** start implementation of WGS-based surveillance for *Listeria monocytogenes*, *Neisseria meningitidis*, Carbapenemase-producing *Enterobacteriaceae* and antibiotic-resistant *Neisseria gonorrhoeae*; 2018

Early data from surveillance of listeriosis in the USA

Besser et al., Clin Micr Infect, 2018

The number of outbreaks detected increased 36% after implementation of real-time WGS based surveillance, and likewise the number of solved outbreaks increased more than three-fold



WGS provides higher resolution and accuracy than classical molecular typing methods, such as PFGE or MLVA, contributing to a better understanding of infectious disease and drug resistance transmission patterns and thereby improving the effectiveness of interventions for their control.

ECDC technical report: Monitoring the use of wgs in infectious disease surveillance in Europe 2015-2017

Figure 1. National public health reference laboratories use of WGS-based typing for national surveillance of at least one human pathogen

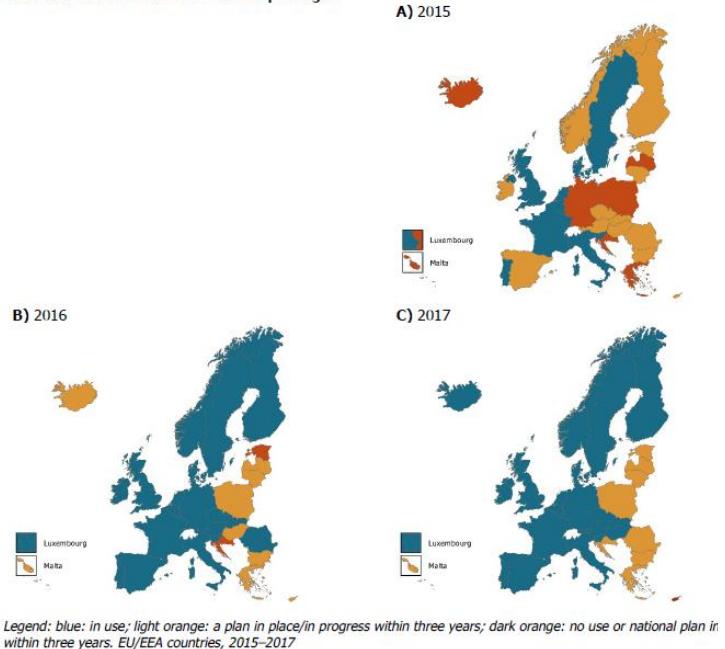


Figure 3. Number of EU/EEA countries using WGS-based typing as first or second-line method for routine surveillance and outbreak investigations in National Public Health Reference Laboratories by disease group and pathogen, 2017

Foodborne pathogens
Listeria monocytogenes
Salmonella enterica
Shiga toxin-producing E. coli (STEC)
Antimicrobial-resistant pathogens
Carbapenemase-producing *Enterobacteriaceae* (CPE)
Antibiotic resistant *Neisseria gonorrhoeae*
Multidrug-resistant *Mycobacterium tuberculosis*
Vaccine-preventable pathogens
Invasive Neisseria meningitidis
Human Influenza virus

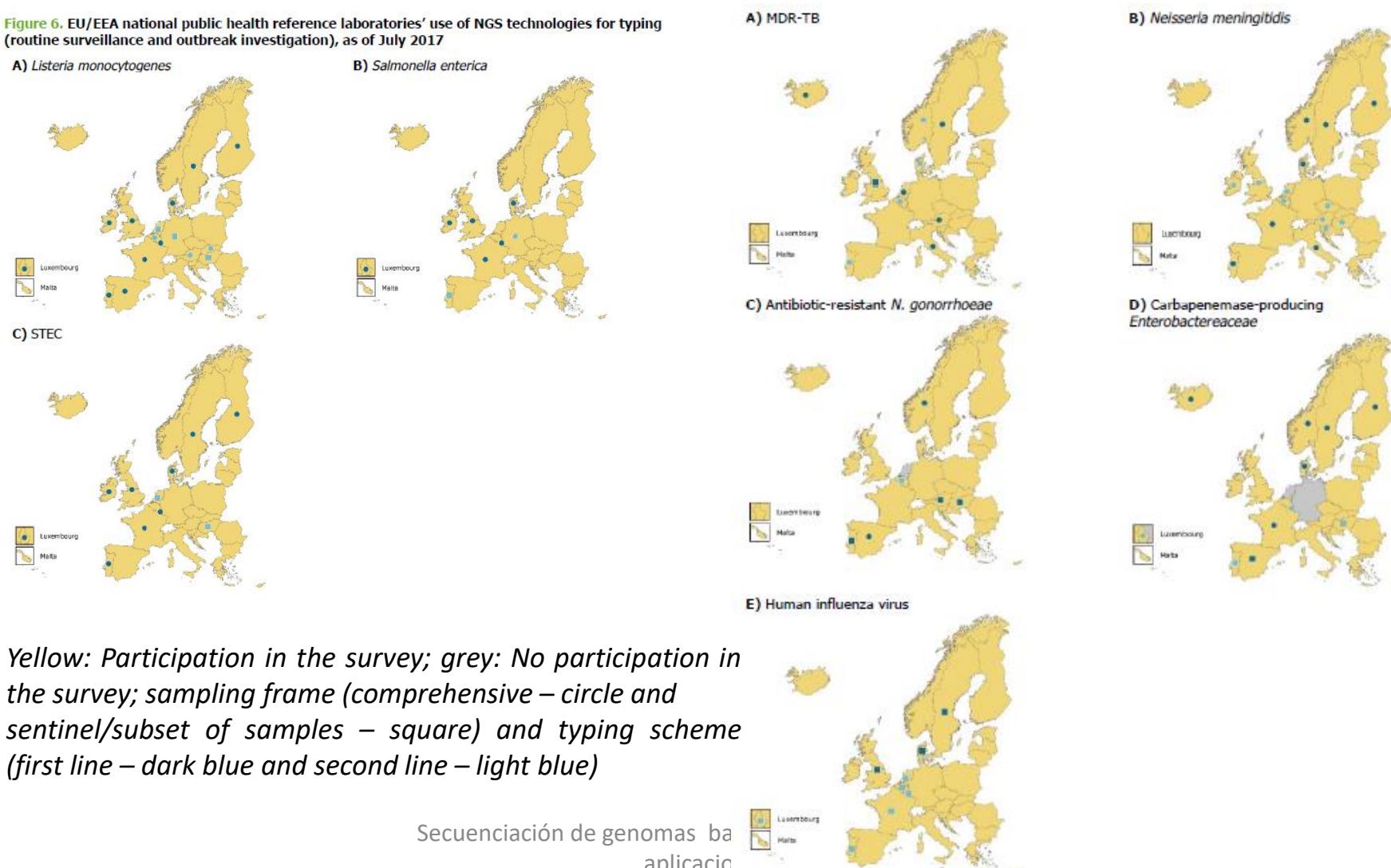
| Number of countries by use of WGS-typing in 2017 or being planned by 2019, per pathogen | | | | | | |
|---|----|----|---|---|---|--|
| 4 | 8 | 3 | 6 | 1 | 9 | |
| 4 | 10 | 9 | 2 | 5 | | |
| 7 | 8 | 5 | 2 | 8 | | |
| 2 | 4 | 12 | 2 | 3 | 7 | |
| 1 | 10 | 7 | 6 | 1 | 5 | |
| 10 | 7 | 3 | 4 | 6 | | |
| 4 | 9 | 2 | 8 | 7 | | |
| 6 | 14 | 2 | 5 | 3 | | |

Legend

- No information
- No, it is not used for public health operations or planned by 2019
- No, it is not used for public health operations or planned by 2019
- Only for outbreak investigations
- Routine surveillance and outbreak investigations - Second-line typing
- Routine surveillance and outbreak investigations - First-line typing

ECDC technical report: Monitoring the use of wgs in infectious disease surveillance in Europe 2015-2017

Figure 6. EU/EEA national public health reference laboratories' use of NGS technologies for typing (routine surveillance and outbreak investigation), as of July 2017



ECDC technical report: Monitoring the use of wgs in infectious disease surveillance in Europe 2015-2017

Table 1. Number of EU/EEA countries with one or more national public health reference laboratories having access to next-generation sequencing (NGS) technologies for routine public health operations, by technology and instrument used, 2017

| NGS technology | Instrument | Foodborne pathogens | | | Antimicrobial-resistant pathogens | | Vaccine-preventable diseases | |
|------------------------------|----------------|-------------------------|--------------------|------|-----------------------------------|---------------------------|------------------------------|------------------------|
| | | <i>L. monocytogenes</i> | <i>S. enterica</i> | STEC | CPE | AR- <i>N. gonorrhoeae</i> | MDR-TB | <i>N. meningitidis</i> |
| Illumina | HiSeq series | 3 | 3 | 3 | 2 | 1 | 2 | 3 |
| | HiSeq X series | | | | | | 1 | |
| | MiniSeq | 1 | 2 | 3 | 2 | | 4 | 1 |
| | MiSeq series | 12 | 10 | 7 | 7 | 10 | 7 | 13 |
| | NextSeq | 2 | 2 | 2 | 3 | 1 | 3 | 7 |
| Ion Torrent | S4 | | | | | | | 1 |
| | S5 | 1 | 1 | 1 | 1 | | 1 | |
| | S5 XL | | | | | | | 1 |
| | PGM | | 1 | 1 | | 1 | | 1 |
| | Proton | | | | | | 1 | 1 |
| Oxford Nanopore Technologies | MinION | | 1 | 1 | 2 | | 1 | 1 |
| Pacific Biosciences-PacBio | PacBio RS II | | 1 | | 2 | | | |
| Other not specified | - | | 3 | 2 | 2 | 1 | 1 | 1 |

Table 2. Bioinformatics tools used by the National Public Health Reference Laboratories using WGS-based typing for surveillance and outbreak investigations of foodborne pathogens, July 2017*

| Tools used for sequence analysis | Number of EU/EEA countries | | |
|------------------------------------|-----------------------------------|-----------------------------|---------------|
| | <i>L. monocytogenes</i> (n=14) | <i>S. enterica</i> (n=7) | STEC (n=9) |
| Commercial software | 9 | 4 | 4 |
| Open source software | 4 | 3 | 5 |
| In-house suite of customised tools | 4 | 2 | 2 |

* Not mutually exclusive

ECDC technical report: Monitoring the use of wgs in infectious disease surveillance in Europe 2015-2017

Table 3. Number of EU/EEA countries using WGS-based typing for surveillance and outbreak investigations in the national public health reference laboratories and respective typing scheme, sampling frame, bioinformatics analysis, and raw data storage practice by pathogen, 2017

| | 2017 | Foodborne pathogens | | | Antimicrobial resistant pathogens | | | Vaccine preventable pathogens | |
|--|---|-------------------------|--------------------|-------|-----------------------------------|---------------------------|--------|-------------------------------|------------------------|
| | | <i>L. monocytogenes</i> | <i>S. enterica</i> | STE C | CP E | AR- <i>N. gonorrhoeae</i> | MDR TB | Human influenza virus | <i>N. meningitidis</i> |
| Number of countries using WGS for routine surveillance and outbreak investigations | | 14 | 7 | 9 | 10 | 6 | 10 | 8 | 15 |
| Typing scheme | First-line WGS | 9 | 5 | 8 | 7 | 5 | 6 | 3 | 7 |
| | Second-line WGS | 5 | 2 | 1 | 3 | 1 | 4 | 5 | 8 |
| Sampling frame | Continuous comprehensive | 12 | 6 | 8 | 7 | 2 | 9 | - | 15 |
| | Sentinel/ subset of case samples | 2 | 1 | 1 | 3 | 4 | 1 | 8 | - |
| Bioinformatic analysis * | cMLST | 12 | 6 | 5 | 6 | 4 | 5 | - | 12 |
| | SNP | 7 | 5 | 5 | 5 | 2 | 7 | - | 5 |
| | Resistome prediction | 4 | 5 | 7 | 8 | 4 | 6 | - | 3 |
| | wgMLST | 5 | 3 | 3 | 2 | 2 | 2 | - | 2 |
| | Virulome/ mobilome prediction | 4 | 2 | 9 | 5 | 1 | - | - | 1 |
| | MLST prediction | 12 | 6 | 8 | 3 | - | - | - | 2 |
| | Serogroup prediction | 7 | 6 | 9 | 1 | - | - | - | 2 |
| | NG-MAST | - | - | - | - | 3 | - | - | - |
| | Speciation | - | 1 | 1 | - | - | 3 | - | 1 |
| | Hemagglutinin and neuraminidase sequence prediction | - | - | - | - | - | - | 4 | - |
| | Phylogenetic relationship | - | 1 | 1 | 1 | - | - | 7 | 1 |
| | Identification of specific point mutations | - | 1 | 1 | - | - | - | 6 | 1 |
| | rMLST | - | - | - | - | - | - | - | 5 |
| | MLST+porA VR1 and VR2+fetA | - | - | - | - | - | - | - | 12 |
| | Vaccine antigen prediction | - | - | - | - | - | - | - | 9 |
| | Other not specified | - | - | - | 1 | - | 3 | 3 | 1 |
| Raw sequence data storage * | Dedicated closed database(s) | 13 | 5 | 7 | 10 | 6 | 10 | 6 | 12 |
| | Publicly available database(s) | 1 | 2 | 2 | - | 1 | 1 | 2 | 3 |

* Not mutually exclusive

Conclusions

This emerging mainstream practice should enable **pan-European WGS-derived data exchange in the medium-term**, subject to harmonisation of sequence analysis pipelines for output compatibility, agreement on international WGS derived type nomenclature and development of **secure and efficient international data sharing** and management platforms.

Current bottlenecks mainly relate to development of expertise in **epidemiological-WGS data integrative analysis** and access to user-friendly international nomenclature

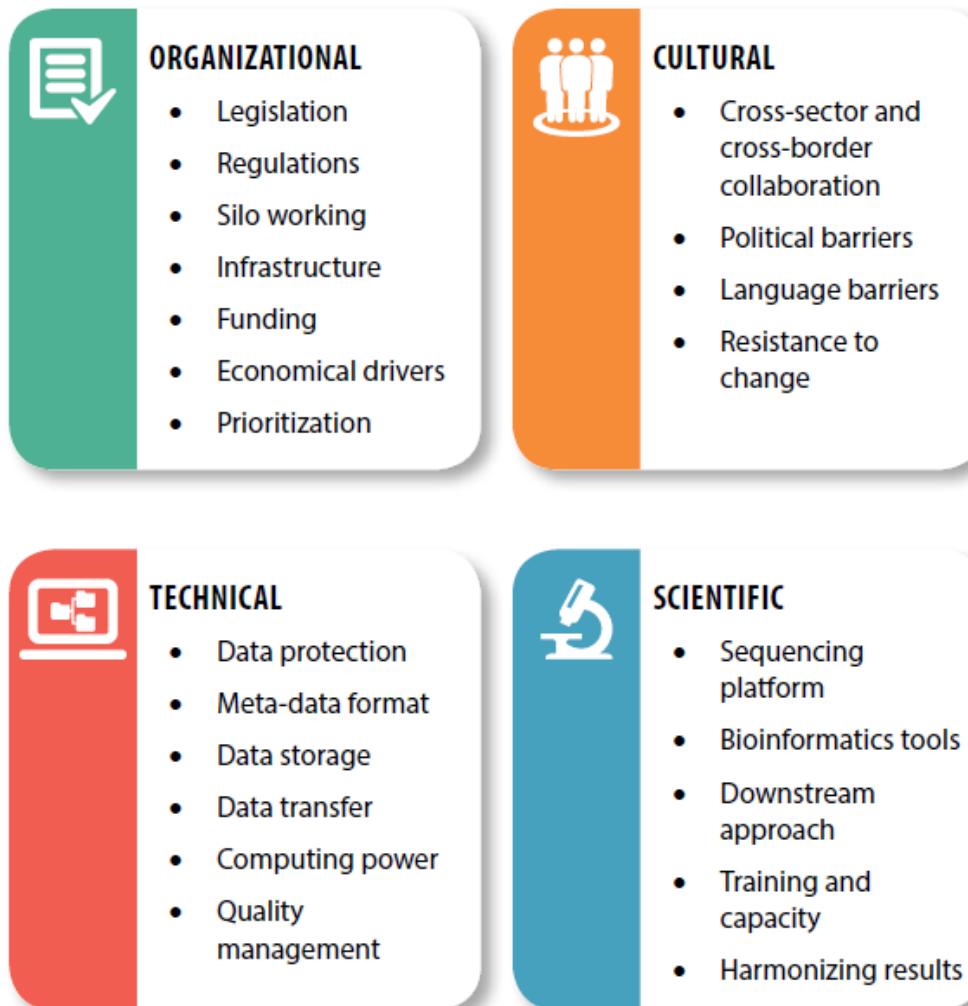
Skills needed to translate WGS data into public health action



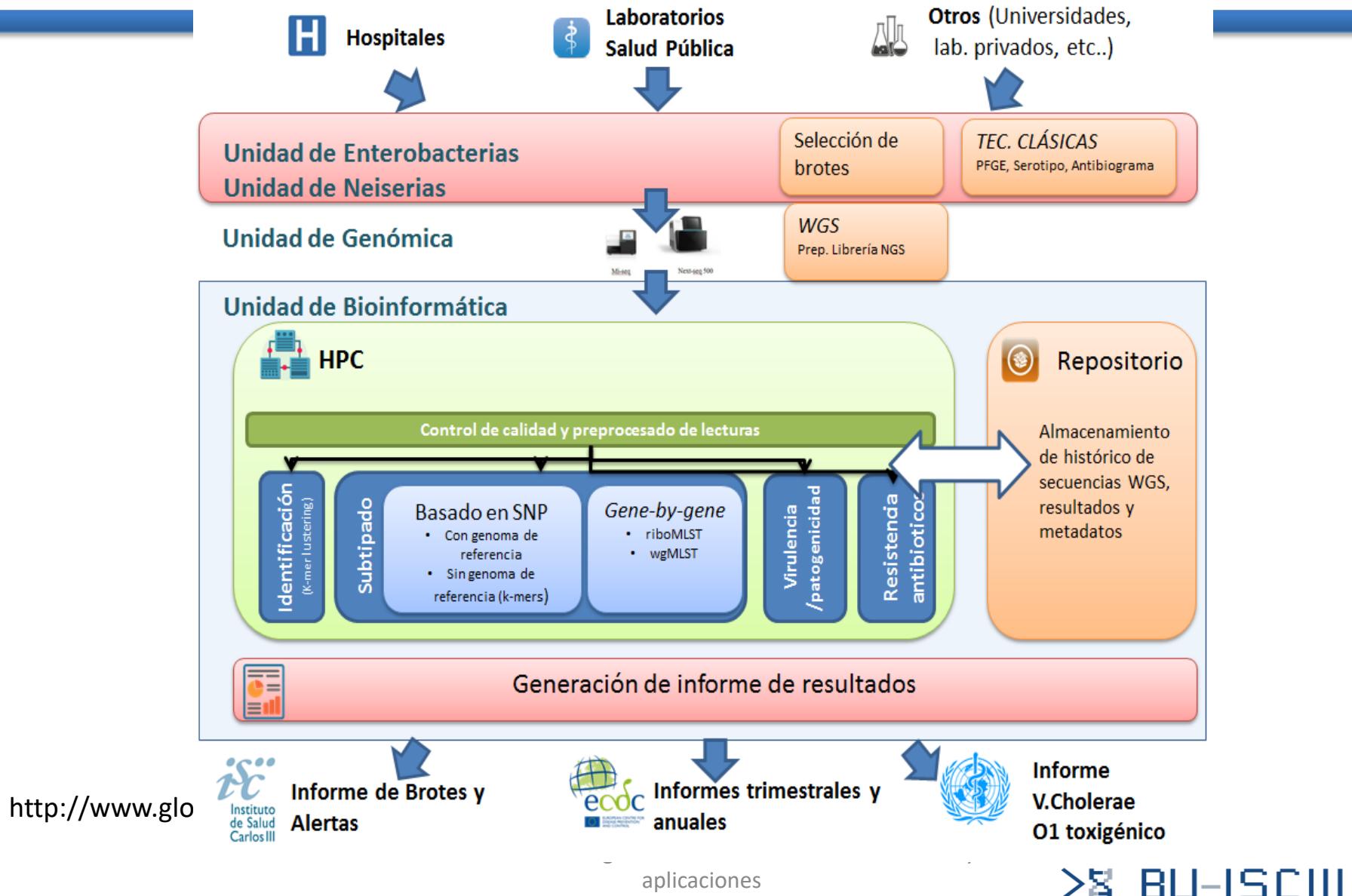
| Bioinformatician | Epidemiologist | Microbiologist |
|---|---------------------------------------|--|
| Algorithms for genome mapping, assembly and comparisons | Epidemiology of communicable diseases | Microbiological diagnostics |
| Inferences from genomic data | Statistical analysis | Subtyping of pathogens |
| Genomic data handling and processing | Case-control studies | Pathogen genomics and evolution |
| Genome data visualization and integration | Health data linkage | Access to culture collections with epidemiological context |
| | Risk assessment and communication | |

FIGURE 2.1

Challenges of coordinating WGS for integrated food chain surveillance

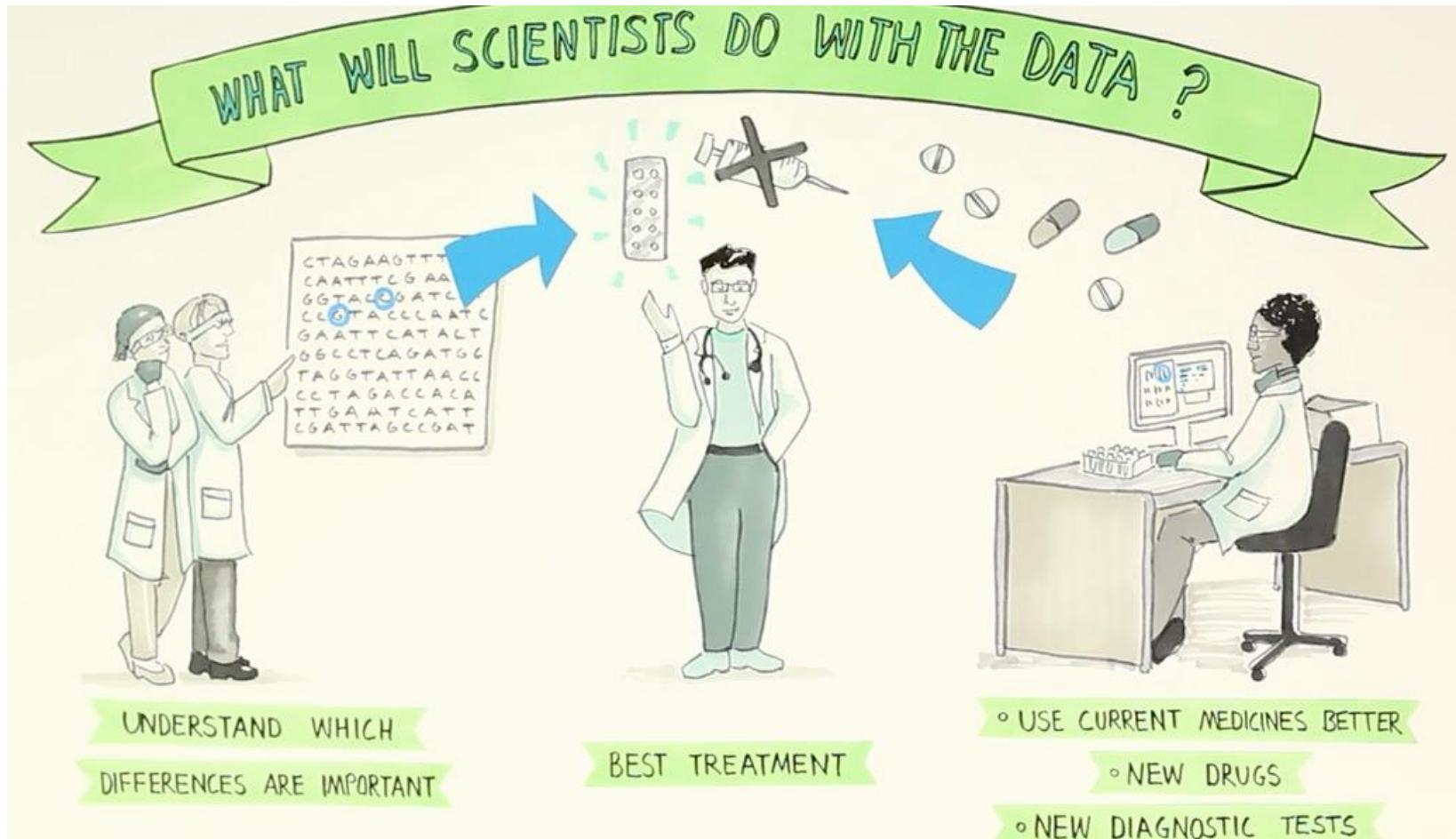


ISCIII Project Ongoing, 2017-2019

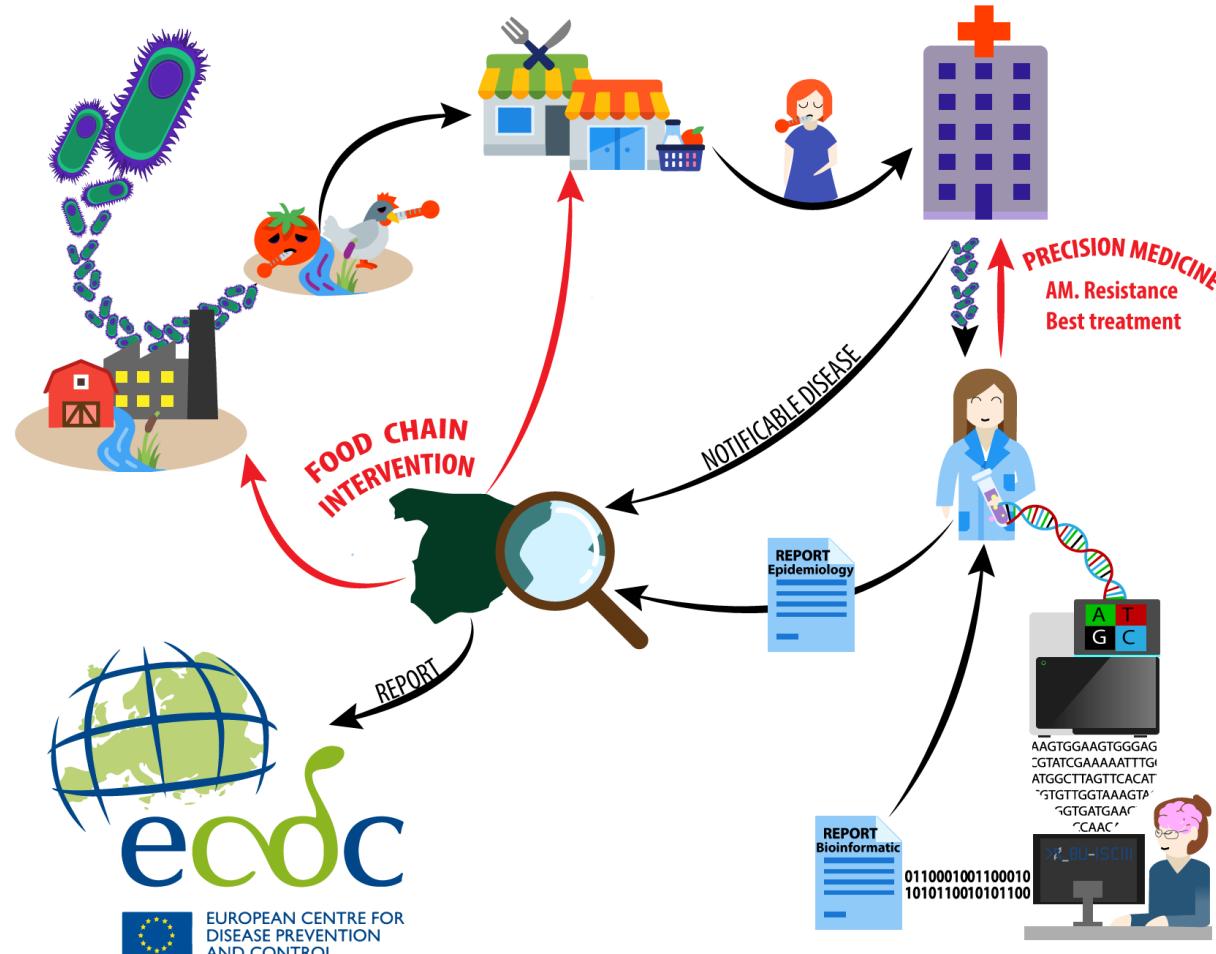


Precision Medicine

<https://labiotech.eu/features/genome-sequencing-review-projects/>



WGS for infectious diseases



Thanks for your attention!

Questions???