



Iniciación al análisis de datos procedentes de técnicas de secuenciación masiva (NGS)

Unidad de Bioinformática (BU-ISCI)I
Unidades Comunes Científico Técnicas – SGAFI-ISCI

17-21 Junio 2019, 7^a Edición
Programa Formación Continua, ISCI

OBJETIVOS DEL CURSO

- ❖ Aproximación a las técnicas de secuenciación masiva (NGS) y a sus aplicaciones
- ❖ Adquirir conocimientos básicos del entorno linux
- ❖ Familiarizarse con los formatos de ficheros generados en el análisis de datos procedentes de la SM
- ❖ Conocer el flujo de análisis de los datos procedentes de la SM



Sesión 1 - Secuenciación Masiva Plataformas de Secuenciación

Isabel Cuesta

Unidad de Bioinformática
Unidades Comunes Científico Técnicas – SGAFI-ISCIII

17-21 Junio 2019, 7^a Edición
Programa Formación Continua, ISCIII

INDICE

- ❖ Unidad de Bioinformática
Servicios ofertados
- ❖ Evolución de la secuenciación
- ❖ Plataformas de secuenciación masiva (NGS)

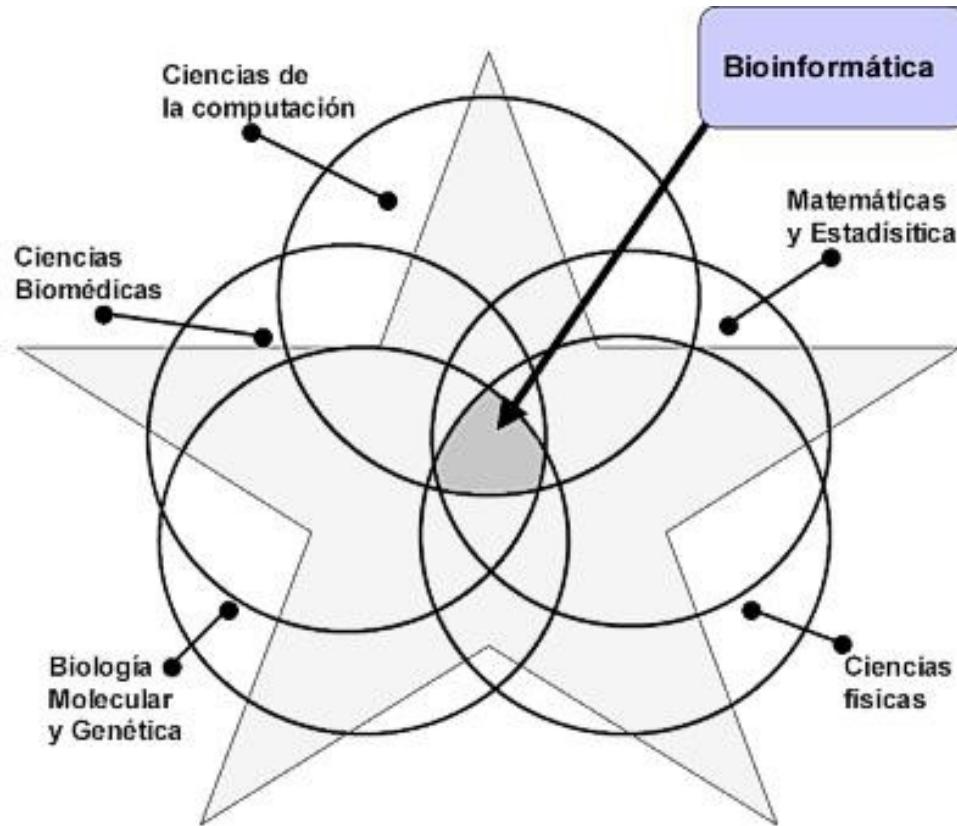
Bioinformatics (*i/baɪənətɪks/*) is the application of statistics and computer science to the field of molecular biology.



(Quantitative+Computable) Molecular Biology

Bioinformatics now entails the creation and advancement of databases, algorithms, computational and statistical techniques and theory to solve formal and practical problems arising from the management and analysis of biological data.

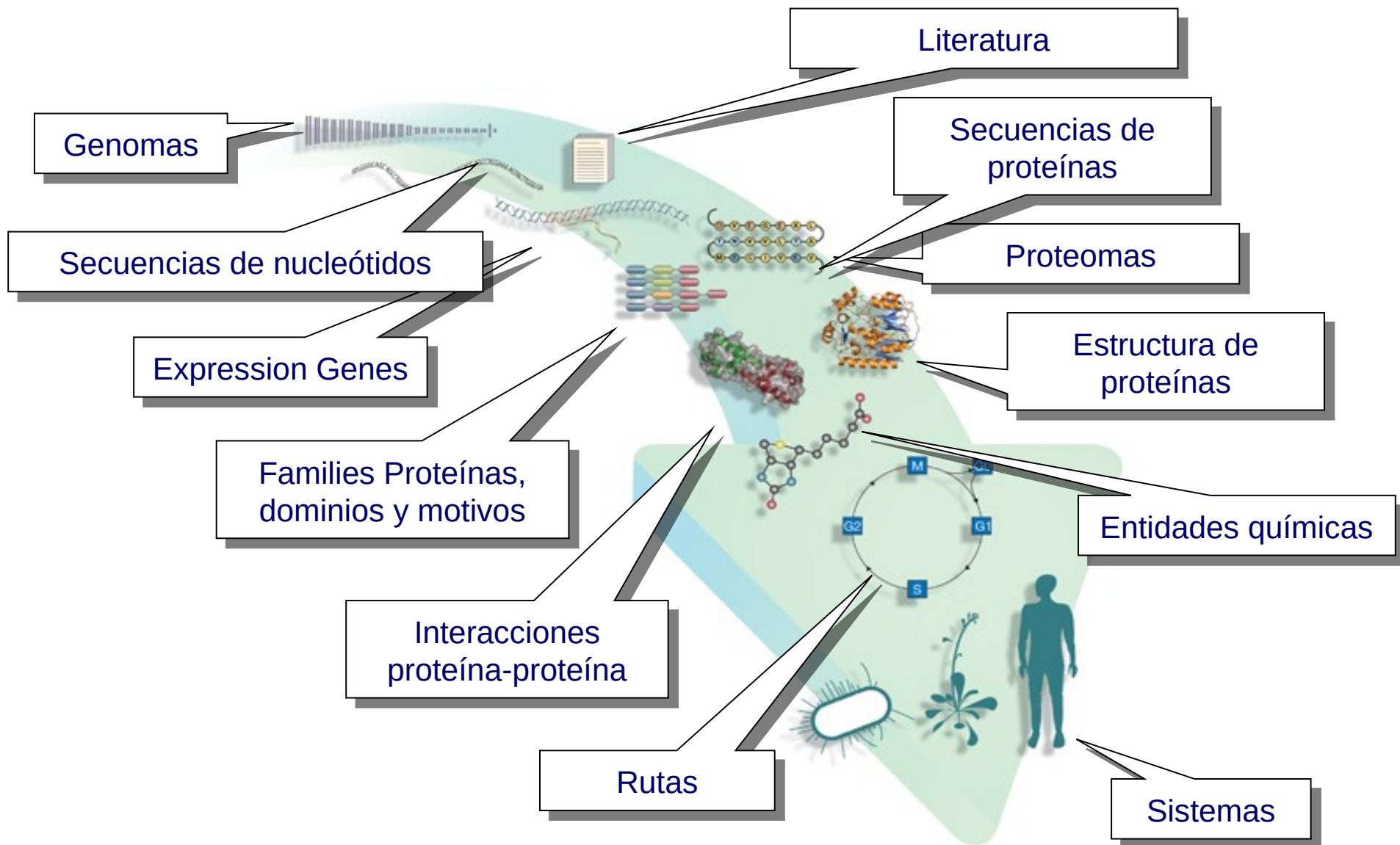
Bioinformática es multidisciplinaria



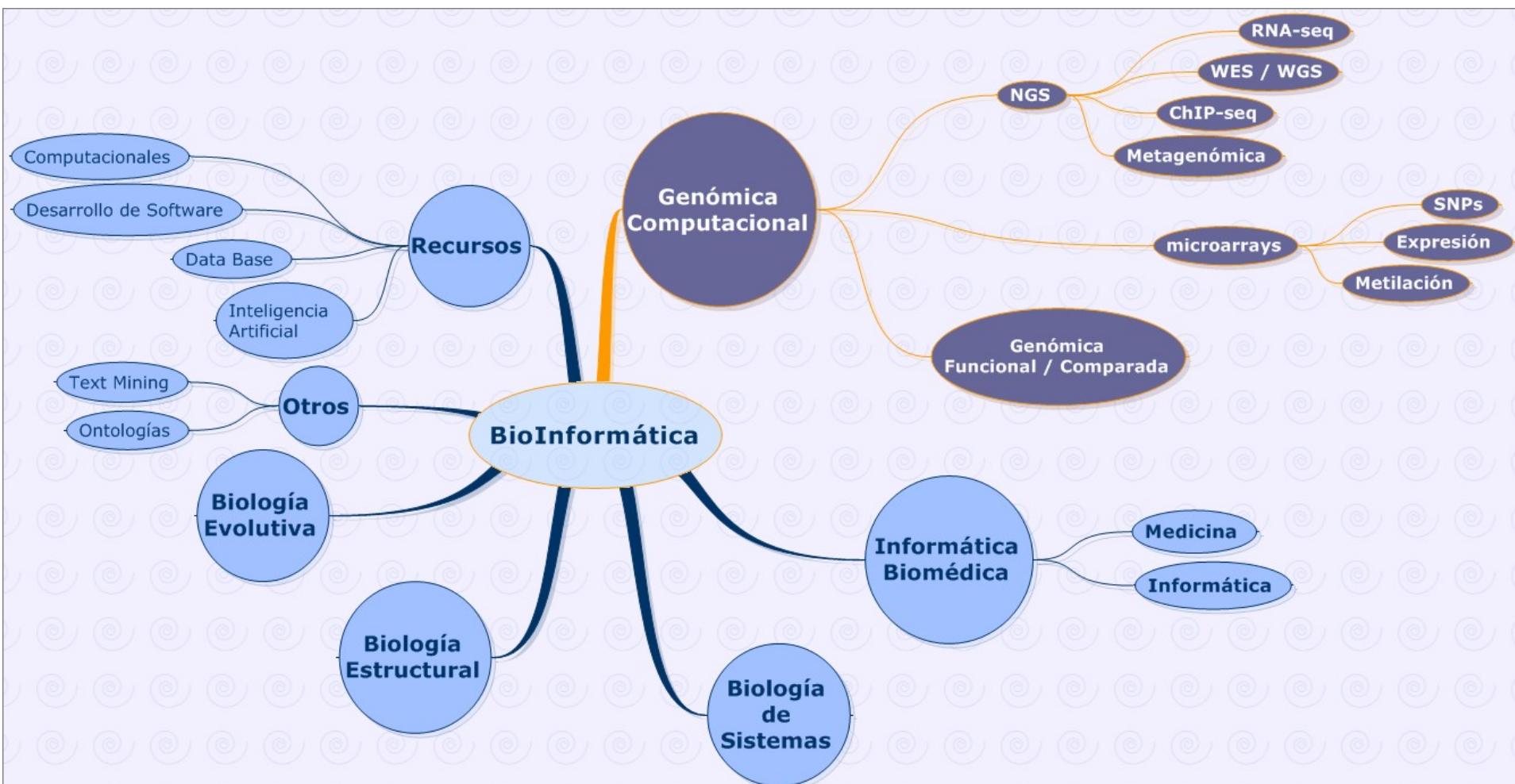
PERSONAL DE LA UNIDAD DE BIOINFORMATICA

	Disciplina	2012	2013	2014	2015	2016	2017	2018	2019	2020
Isabel Cuesta	Dr. Biología Molecular							CIENTIFICO TITULAR OPIS		
Sara Monzón	Biotecnología		CIBERER		PTA MINECO			T.SUPERIOR OPIS		
Bruno Lobo	Administrador Sistemas		PROYECTO		PTA MINECO			SISTEMAS		
Jorge de la Barrera	Informática			PTA FIS						
Miguel Juliá	Matemáticas							PTA MINECO		
José Luis García	Telecomunicaciones							PROYECTO		
Pedro Sola	Biología						U. ANTIBIOTICOS			
????								PROYECTO		

os de datos dan idea de la dimensión de la Bioinformática



ESPECIALIZACIÓN



RECURSOS INFORMÁTICOS

- ❑ Marco de relación con UTIC, establecido en 2013, administración compartida.
- ❑ Workstation (5), 4nucleos, 32Gb Ram, 4TB almacenamiento
- ❑ Servidor de la Unidad, 4-quad, 120Gb Ram,
- ❑ HPC 320 cores, 8TB RAM (16 nodos: Por nodo: 2 procesadores de 10 cores (20 threads) a 2,5 GHZ. 256 GB de RAM. Almacenamiento interno de 500 GB. 2 Interfaces de red a 10 Gbps.)
- ❑ 2 cabinas de almacenamiento, NetApp, 70 TB y 250TB escalable, ubicadas en el nuevo CPD del ISCIII.



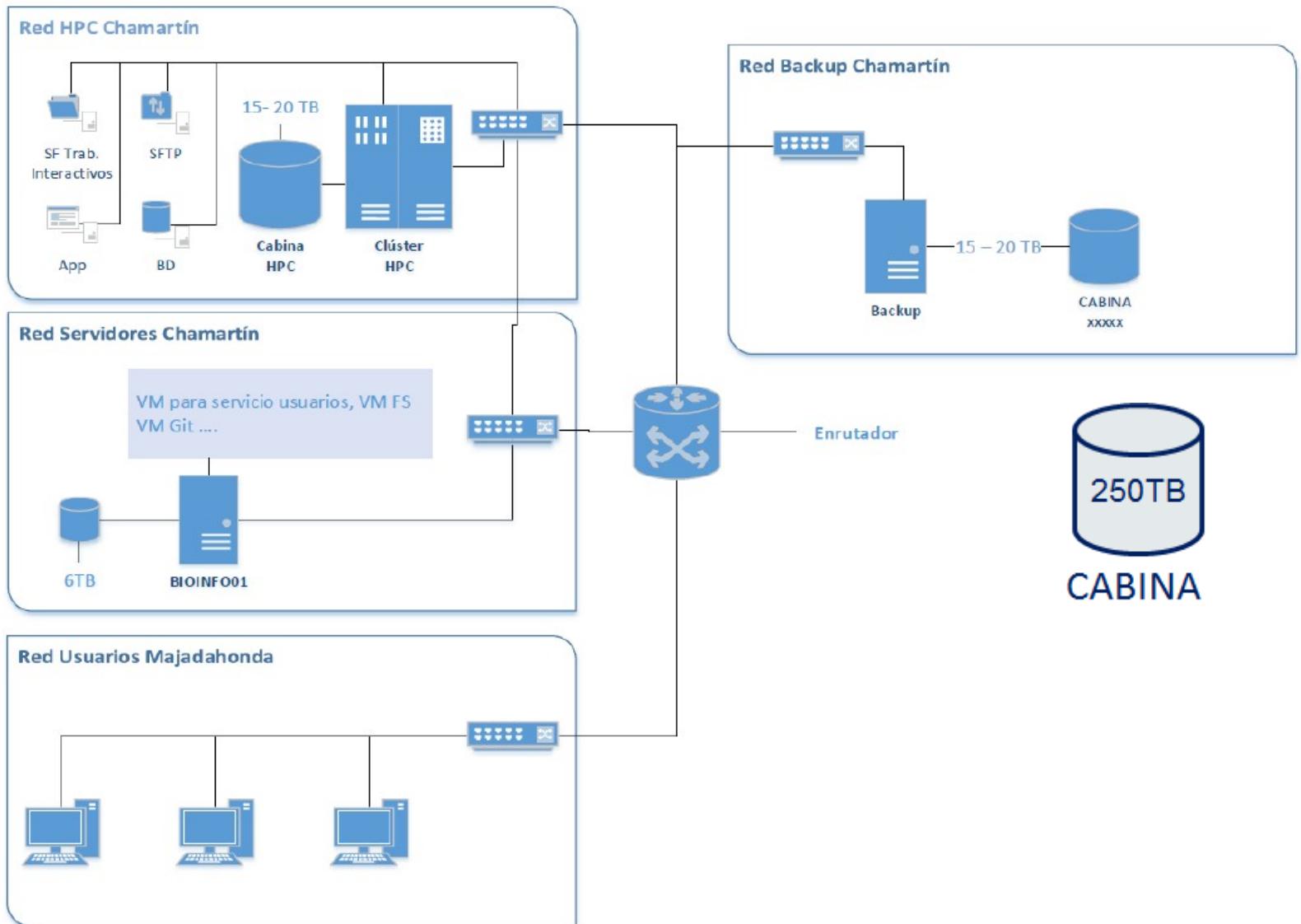
Nuevo CPD

25/06/2018 Pabellón 4 – planta semisótano

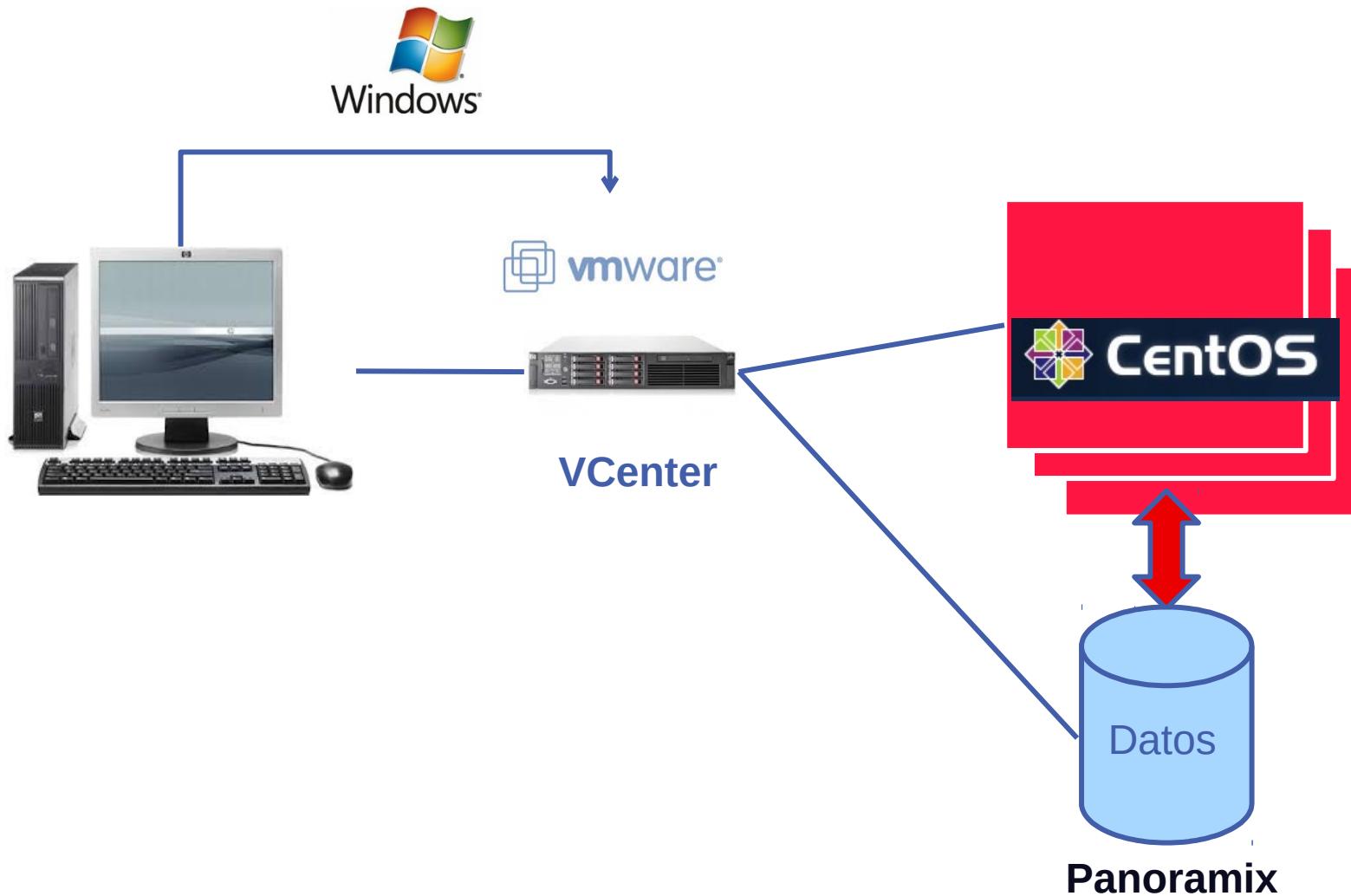


CPD respaldo Majadahonda

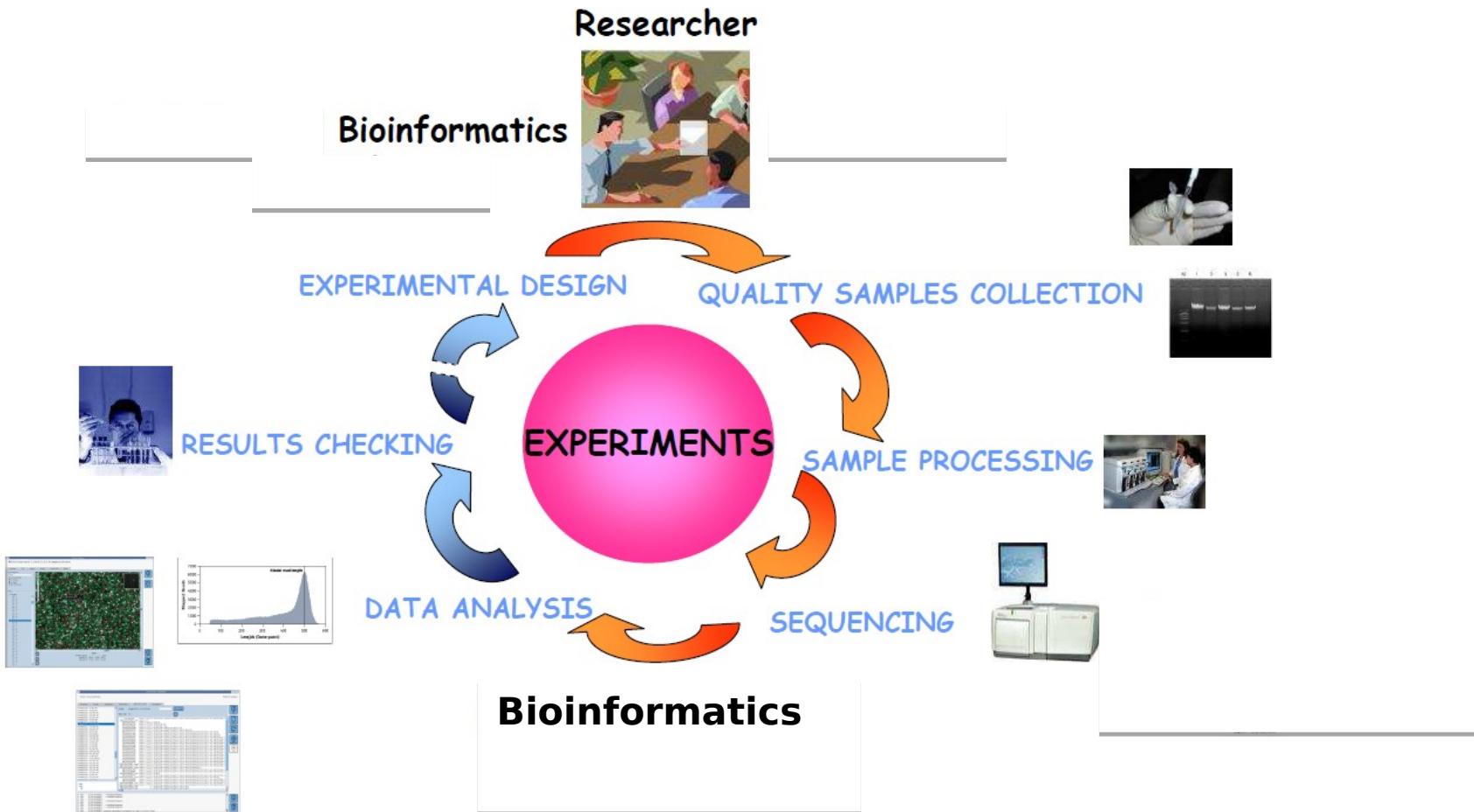
INFRAESTRUCTURA



Recursos Informáticos para el curso



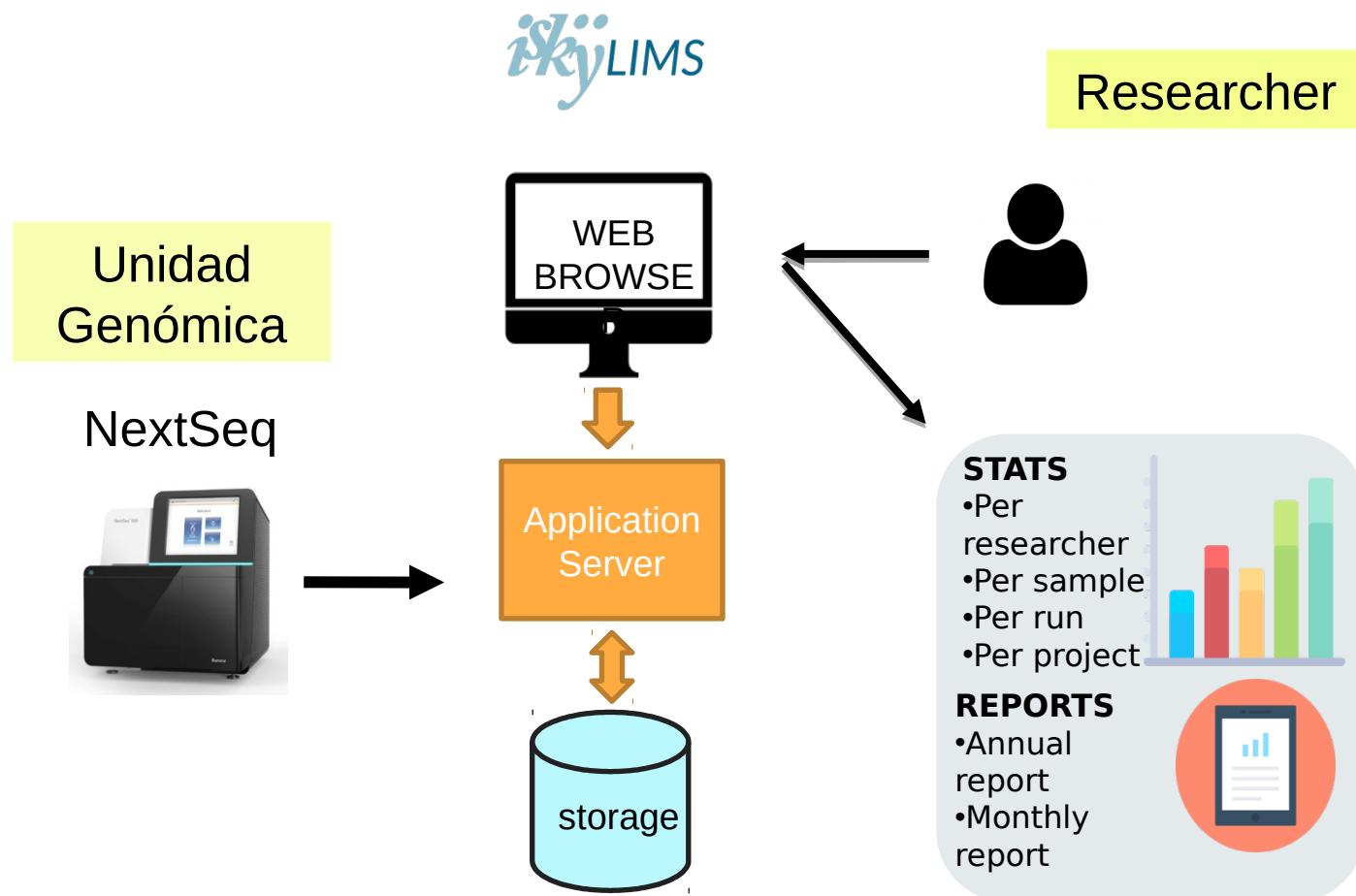
Workflow en NGS



- **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

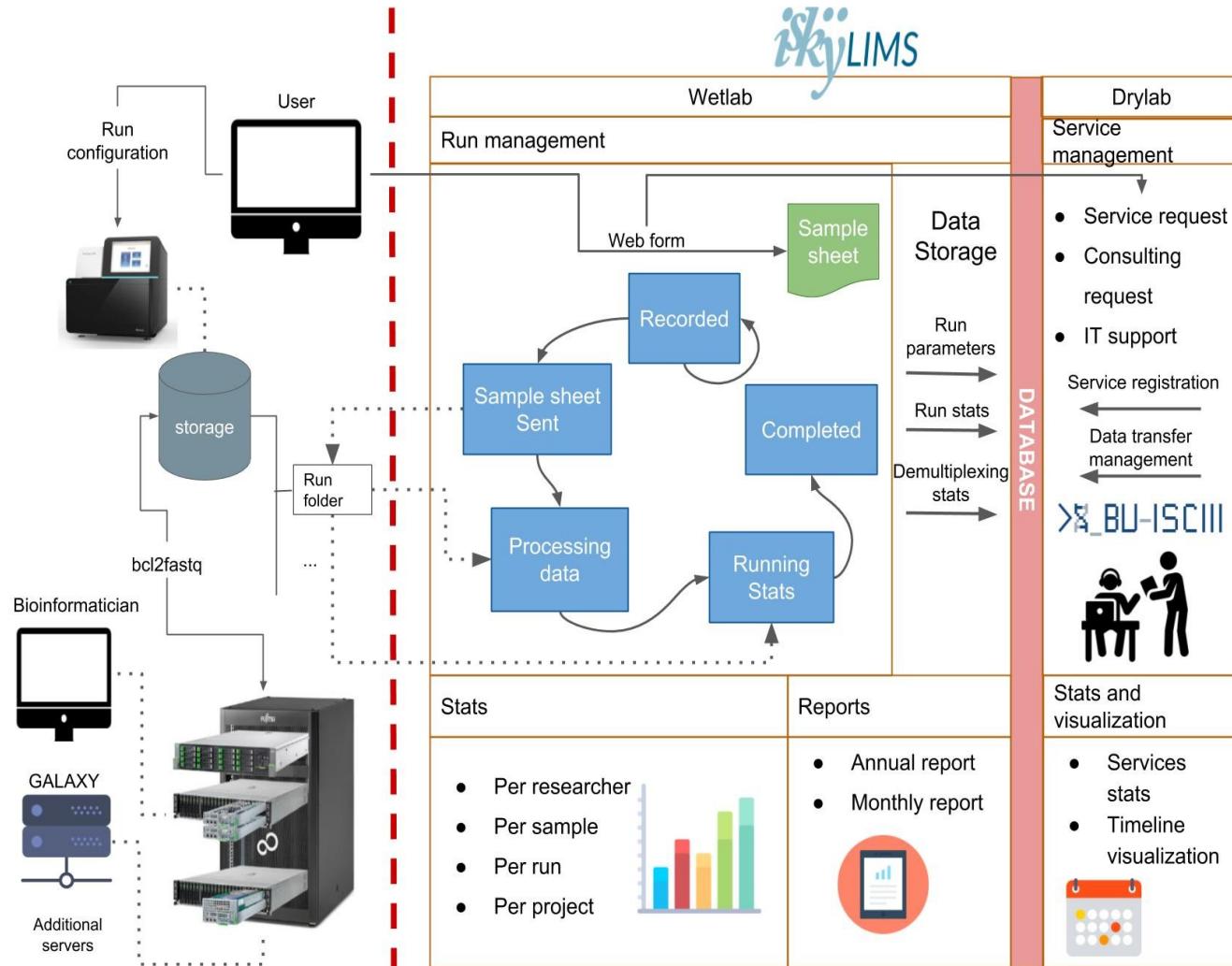
Cartera de Servicios de la Unidad – Análisis de datos Secuenciación Masiva

	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)
	MICROBIAL WGS Amplicon	Report html	<i>De novo</i> / Reference (fasta file)	MLST, Resistance, Virulence	SNPs Phylogenetic analysis	Structural Functional WGSOutbraker Plasmid ID
RNaseq	mRNA	RSQC Report html	<i>De novo</i> (fasta file)	Transcripts coverage / expression	Variants (Vcf file)	mRNA seq
	miRNA				Transcripts annotation	miRNA seq
Metagenomics	16S taxonomic profile	Report html	<i>De novo</i>	Green genes DB		Qiime MeTRS Kaiju
	Shotgun			Genome Ref Seq		PikaVirus

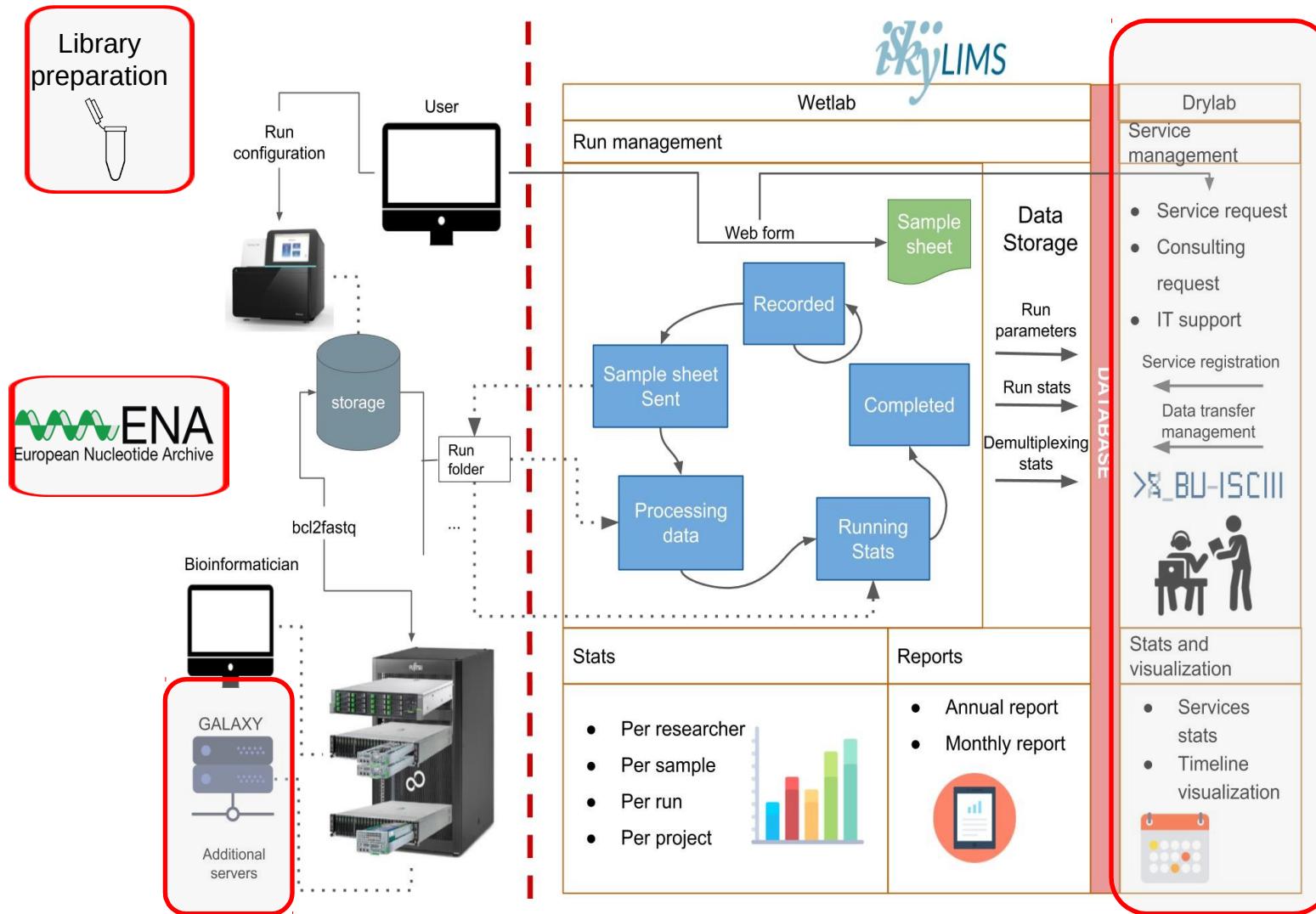


iSkyLIM

S



iSkyLIM



25/06/2018

BU-ISCIII

18

BioInformatics

iSkyLIMS: DryLab

Welcome

This section will allow you to check BU-ISCIII service activity. Available processes are request new services, colaborations, 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.

SERVICES REQUEST

[HOME](#)[SERVICES REQUEST](#) ▾[COUNSELING REQUEST](#)[INFRASTRUCTURE REQUEST](#)[STATISTICS](#)[CONTACT](#)

Sequencing Data

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

SERVICES REQUEST



Service selection

Available Services *

- 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

SERVICES REQUEST



Service Description

Service description file^{*}

No file selected.

Service Notes^{*}

COUNSELING REQUEST



Service selection

AvailableServices *

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

Service Description

Service description file*

No file selected.

Service Notes*

INFRASTRUCTURE REQUEST



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

Service Description

Service description file^{*}

No file selected.

Service Notes^{*}

Galaxy

(| 172.23.2.60)

Galaxy Analyze Data Workflow Shared Data Visualization Help Login or Register

Tools 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

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.

>  

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](#).

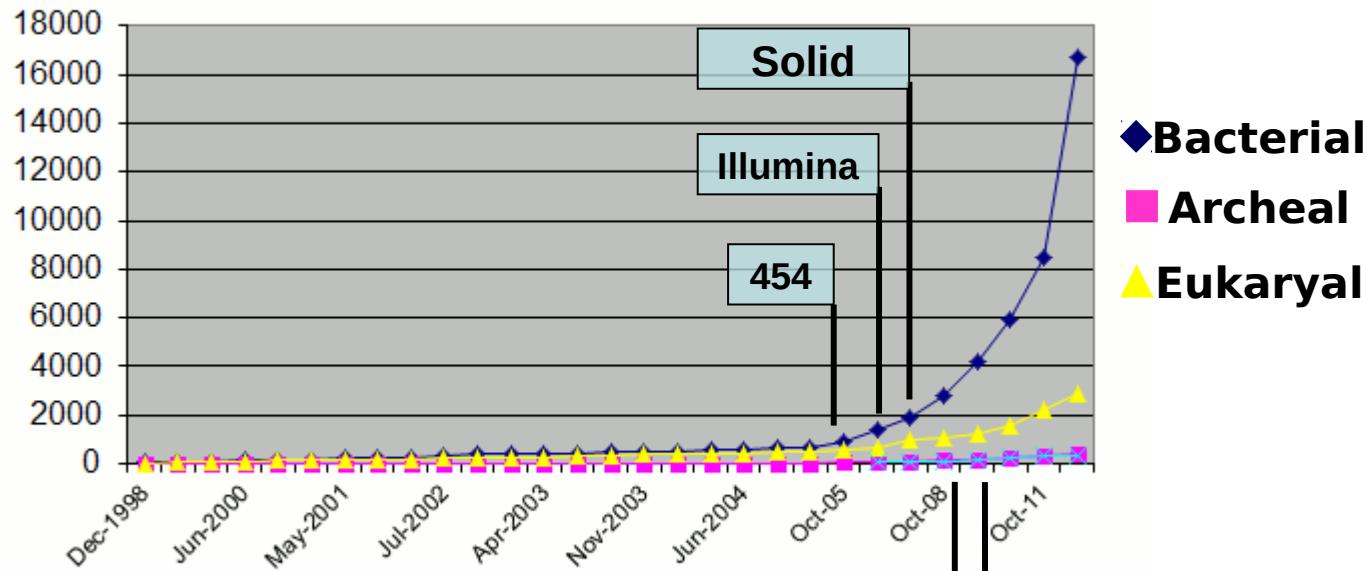
INDICE

- ❖ Unidad de Bioinformática
Servicios ofertados
- ❖ Evolución de la secuenciación
- ❖ Plataformas de secuenciación masiva
(NGS)

Genomics Revolution Era



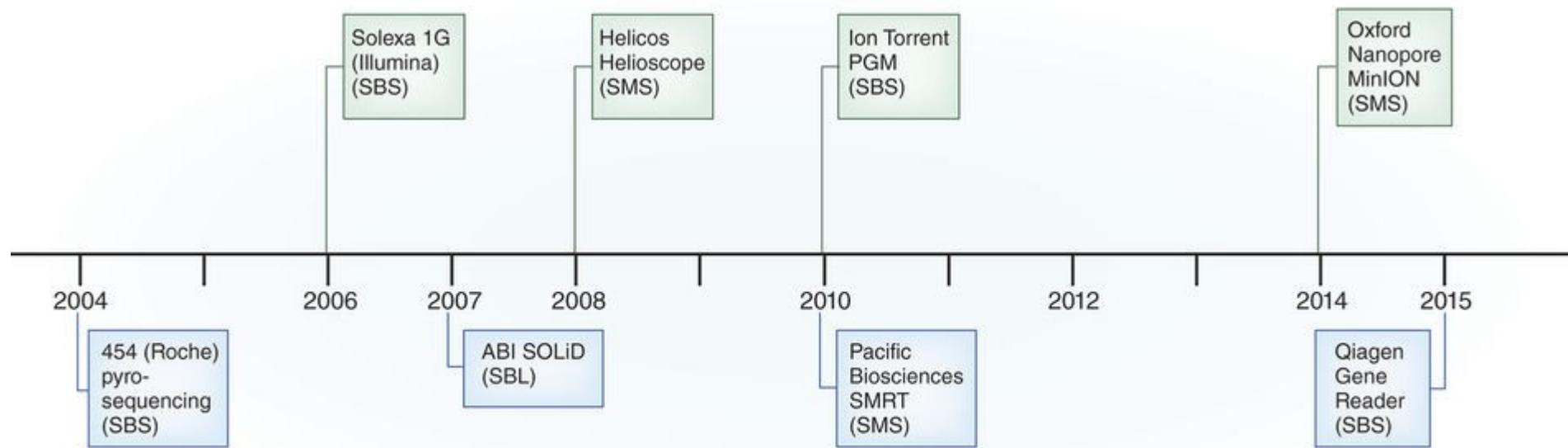
Genome Projects on GOLD according to Phylogenetic Groups ©
October 2012 - 20327 Projects



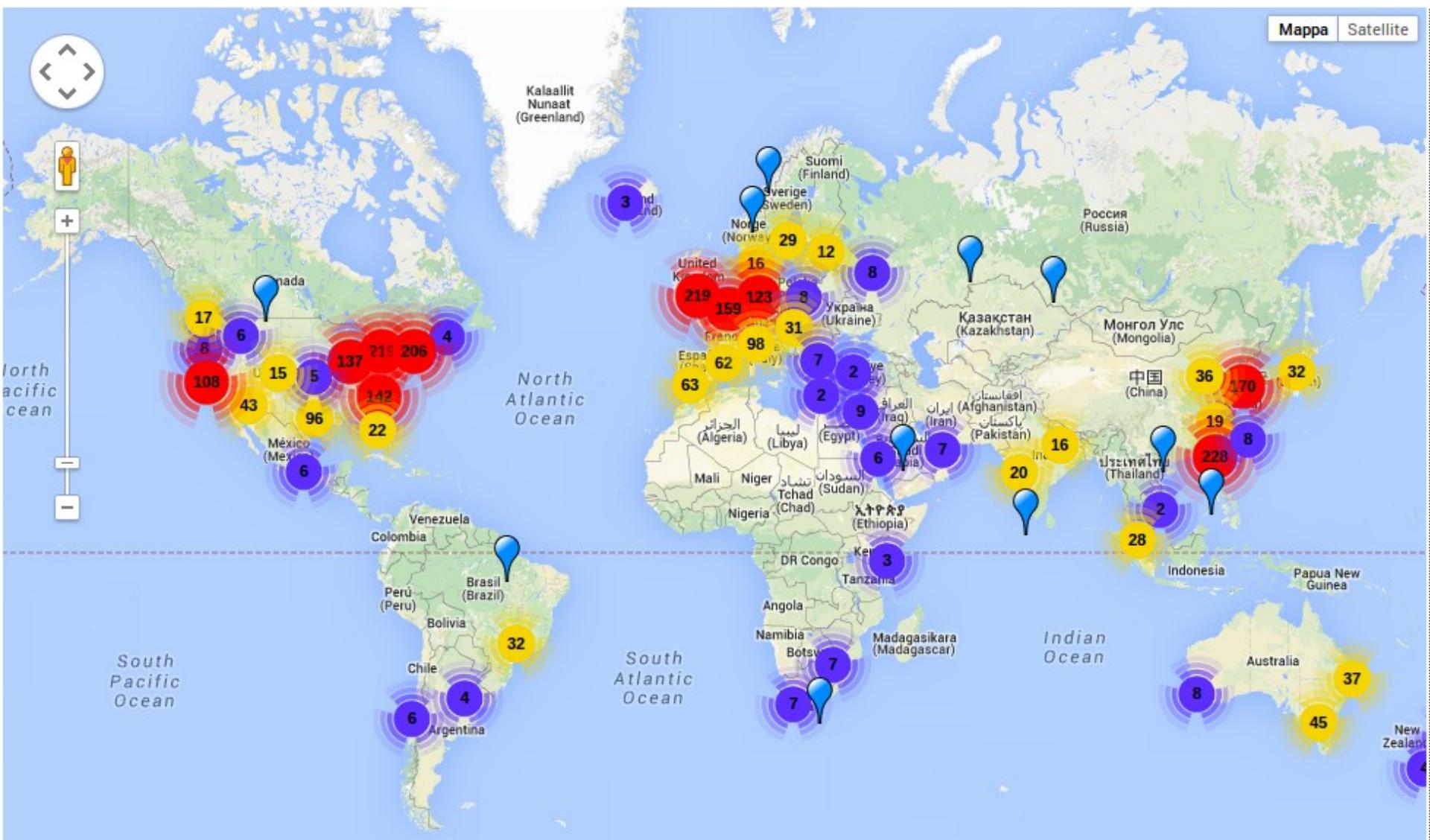
Source: <http://www.genomeonline.org>



NGS PLATFORMS: TIMELINE



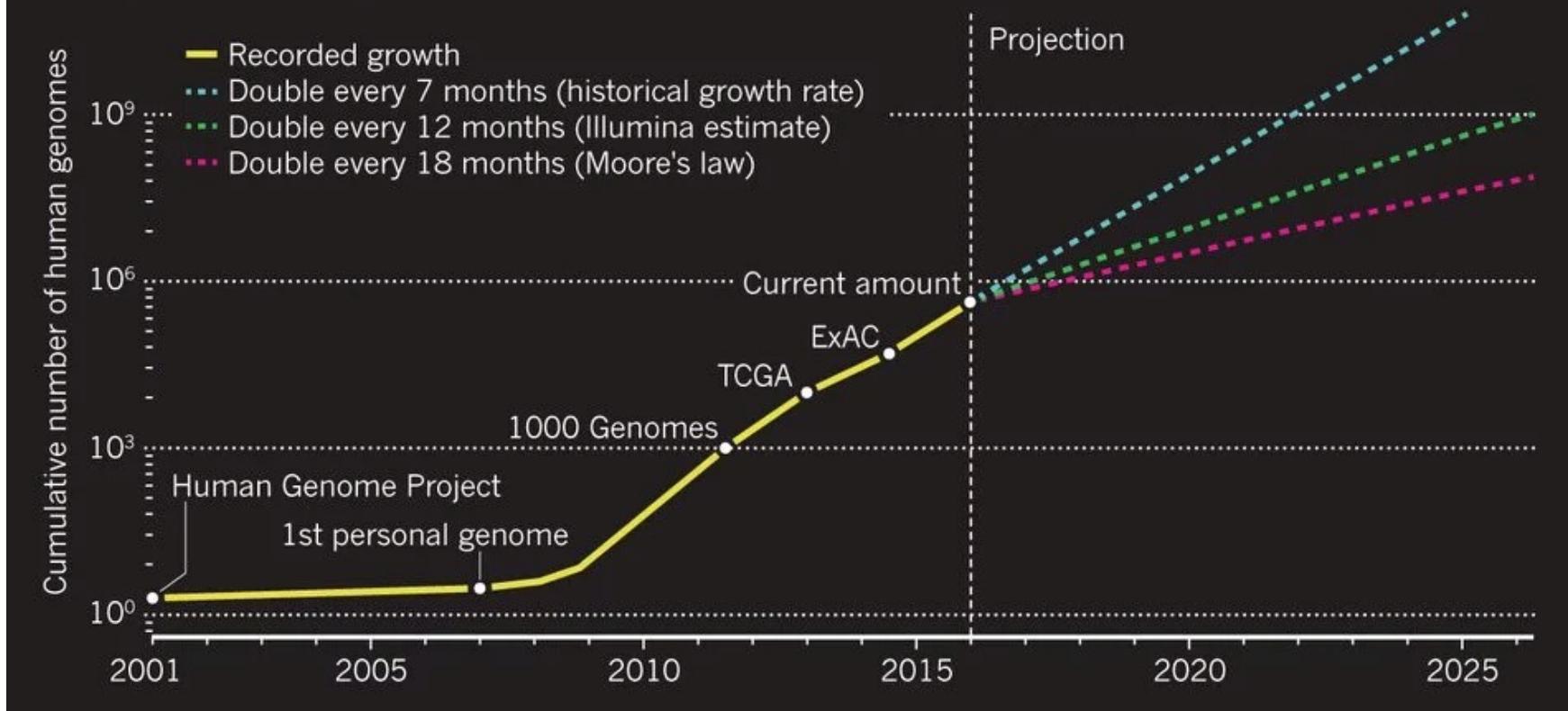
Democratización de la secuenciación. Mapa de los avances



SEQUENCING PROJECTS

DNA SEQUENCING SOARS

Human genomes are being sequenced at an ever-increasing rate. The 1000 Genomes Project has aggregated hundreds of genomes; The Cancer Genome Atlas (TCGA) has gathered several thousand; and the Exome Aggregation Consortium (ExAC) has sequenced more than 60,000 exomes. Dotted lines show three possible future growth curves.



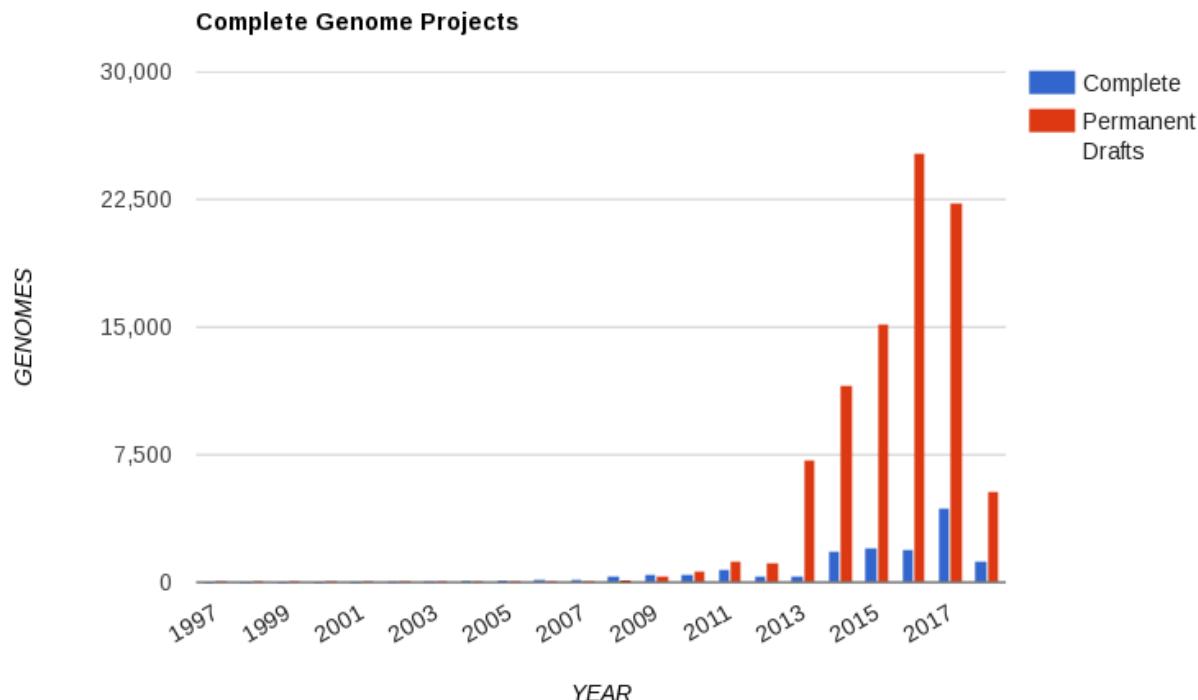
Genomics Revolution Era



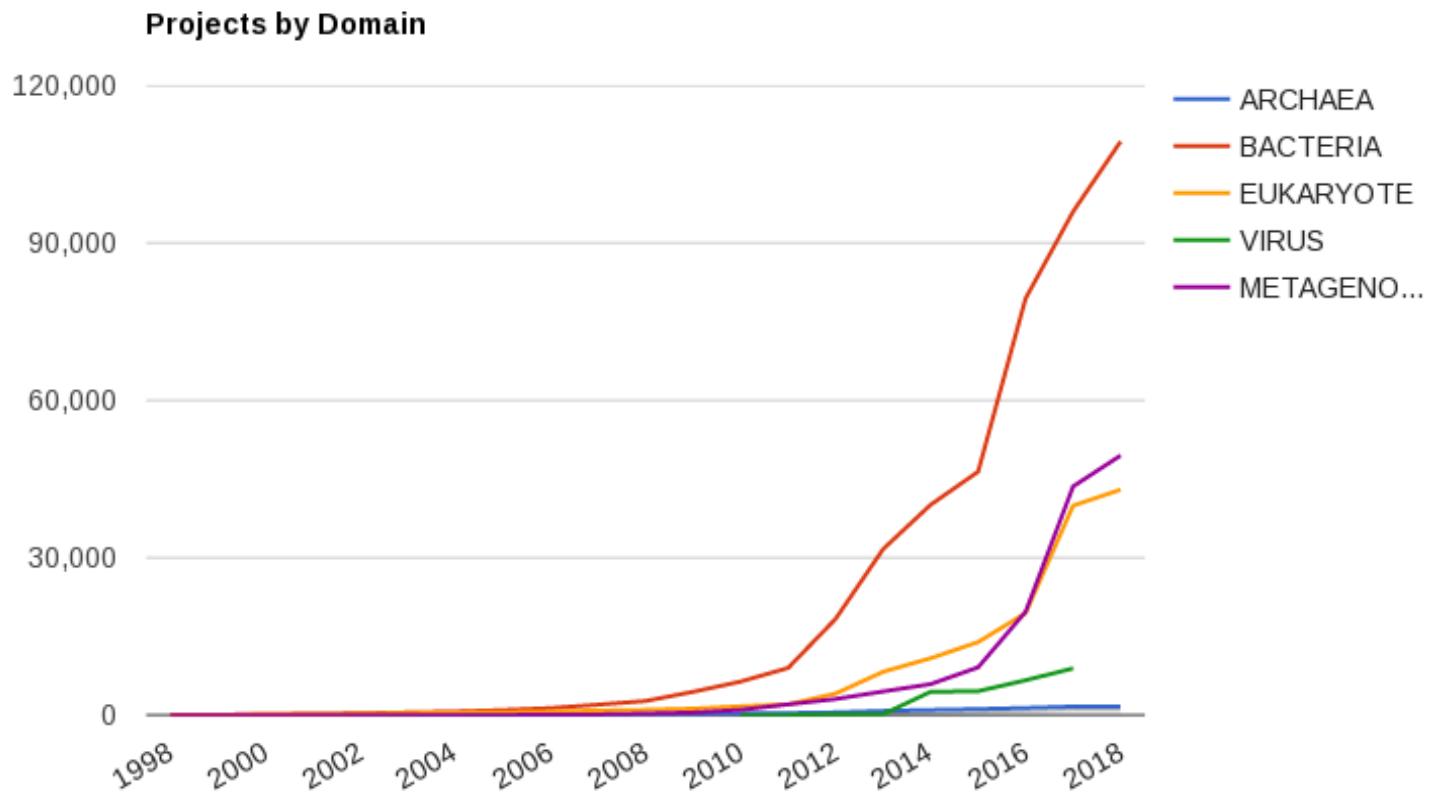
[JGI HOME](#) [LOG IN](#)

Home Search Distribution Graphs Biogeographical Metadata Statistics References Team Help News

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



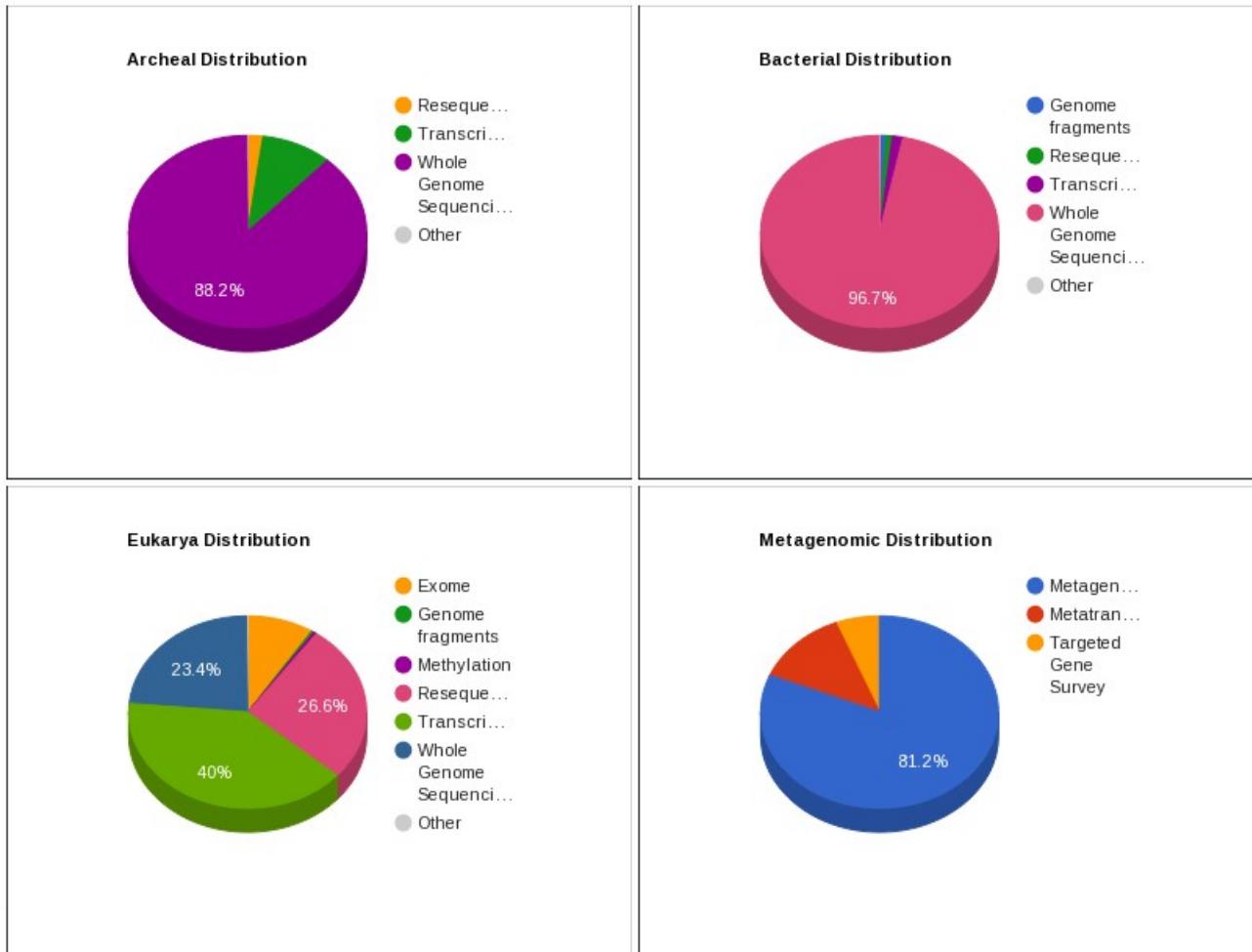
Genomics Revolution Era



Genomics Revolution Era



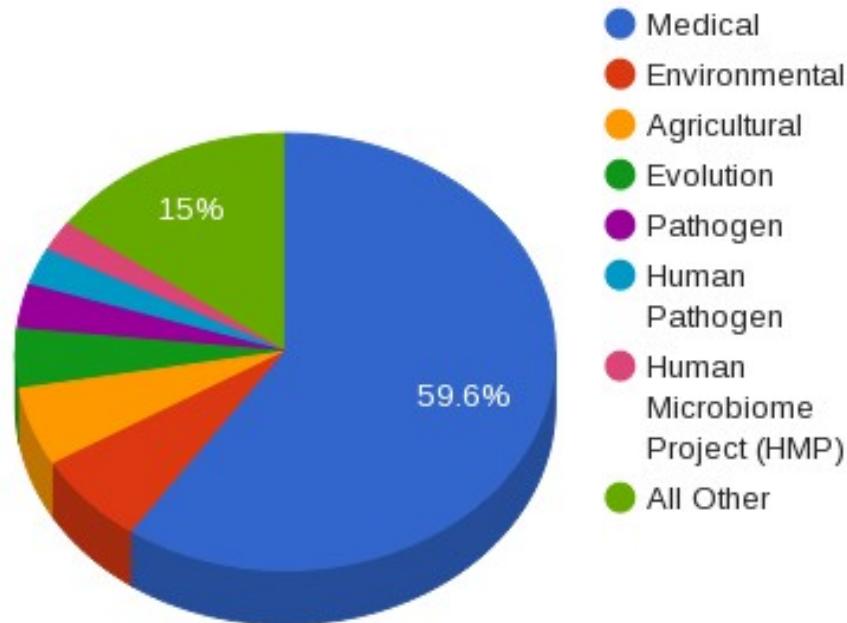
GOLD Project Distributions



Genomics Revolution Era

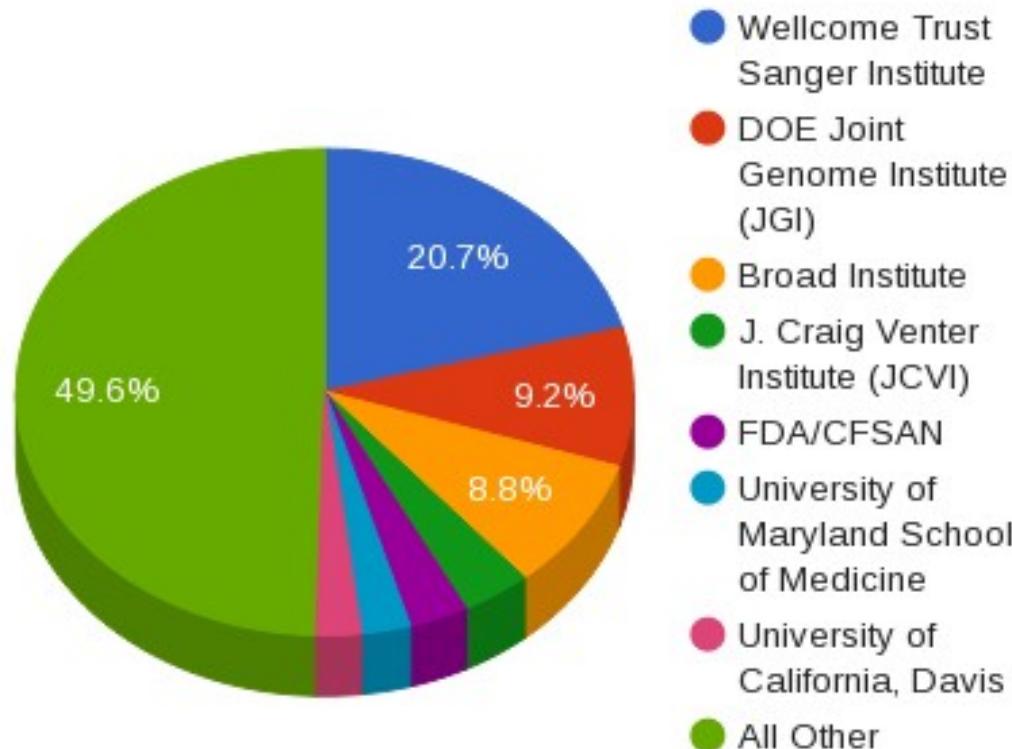


Project Relevance of Bacterial Projects



Genomics Revolution Era

Sequencing Centers for Archaeal and Bacterial Projects



Secuenciadores

Primera generación

- Sanger

Segunda Generación

- 454/Roche
- Solexa/Illumina
- Solid
- Ion Torrent

Tercera Generación

- Pacific Biosciences
- Nanopore

High-Throughput Sequencing Platforms



GS-FLX System



Genome Analyzer IIx



SOLID 3 Plus/4

Sequencing Chemistry	Sequencing by synthesis, pyrosequencing	Sequencing by synthesis with reversible terminators	Sequencing by ligation
Amplification approach	Emulsion PCR	Cluster amplification	Emulsion PCR
DNA support	25-35 µm bead	Flow cell surface	Bead (Solid 3 Plus/4) Flow cell surface (GA5500w)

High-Throughput Sequencing Platforms



Benchtop High-Throughput Sequencing Platform



Roche 454 GS Junior

illumina®



MiSeq



iSeq 100



MiniSeq

BU-ISCIII

life technologies™



Ion Proton™ System



Ion PGM™ System

Illumina Benchtop Sequencers

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



iSeq 100 System



MiniSeq System



MiSeq Series +



NextSeq Series +

Popular Applications & Methods	Key Application	Key Application	Key Application	Key Application
Large Whole-Genome Sequencing (human, plant, animal)				●
Small Whole-Genome Sequencing (microbe, virus)	●	●	●	●
Exome Sequencing				●
Targeted Gene Sequencing (amplicon, gene panel)	●	●	●	●
Whole-Transcriptome Sequencing				●
Gene Expression Profiling with mRNA-Seq				●
Targeted Gene Expression Profiling	●	●	●	
Long-Range Amplicon Sequencing*	●	●	●	
miRNA & Small RNA Analysis	●	●	●	●
DNA-Protein Interaction Analysis			●	●
Methylation Sequencing				●
16S Metagenomic Sequencing		●	●	●
Run Time	9–17.5 hours	4–24 hours	4–55 hours	12–30 hours
Maximum Output	1.2 Gb	7.5 Gb	15 Gb	120 Gb
Maximum Reads Per Run	4 million	25 million	25 million †	400 million
Maximum Read Length	2 × 150 bp	2 × 150 bp	2 × 300 bp	2 × 150 bp

Illumina Production-Scale Sequencers

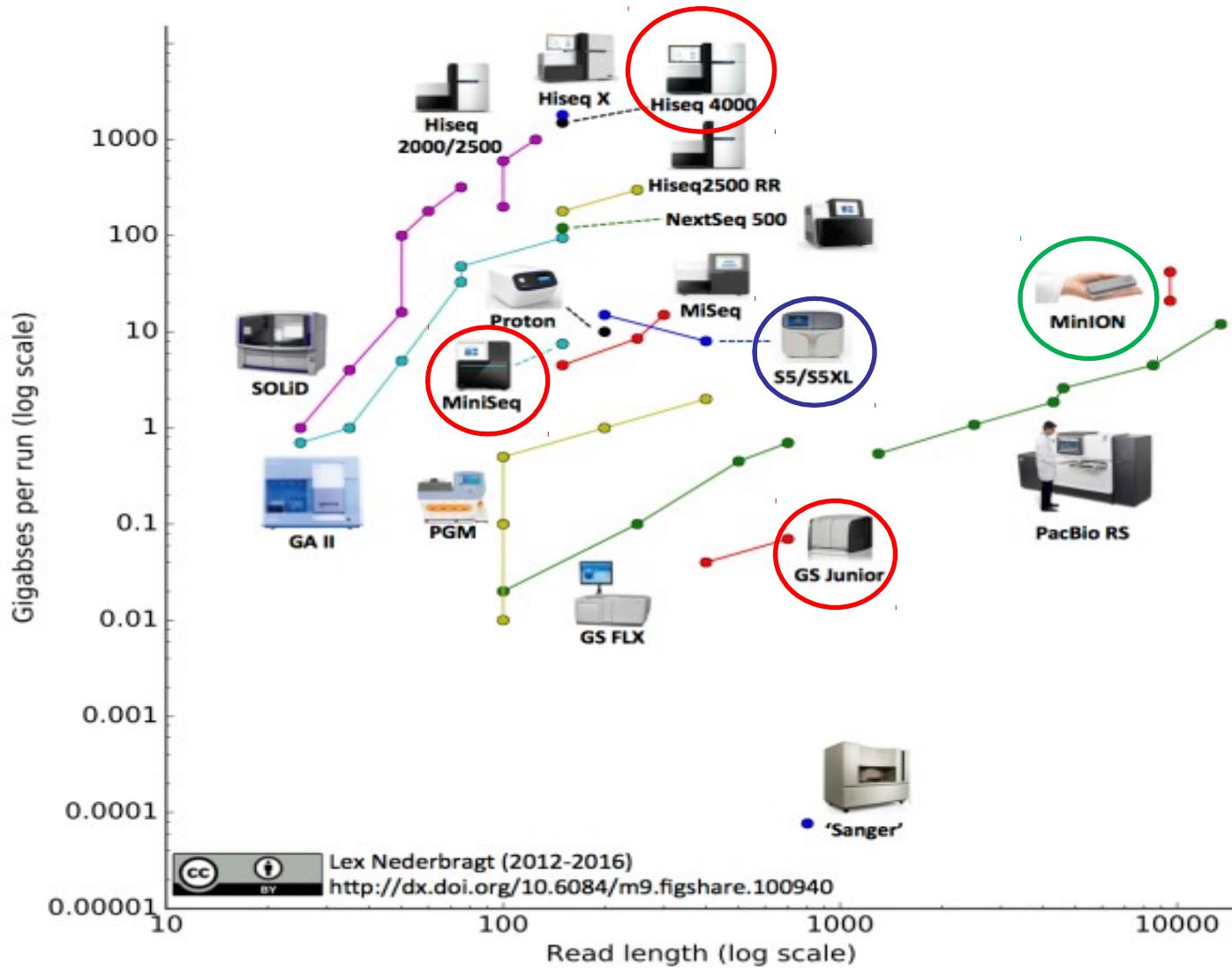
Benchtop Sequencers		Production-Scale Sequencers		
		NextSeq Series 	HiSeq Series 	HiSeq X Series [†] 
Popular Applications & Methods	Key Application 	Key Application 	Key Application 	Key Application 
Large Whole-Genome Sequencing (human, plant, animal)				
Small Whole-Genome Sequencing (microbe, virus)				
Exome Sequencing				
Targeted Gene Sequencing (amplicon, gene panel)				
Whole-Transcriptome Sequencing				
Gene Expression Profiling with mRNA-Seq				
miRNA & Small RNA Analysis				
DNA-Protein Interaction Analysis				
Methylation Sequencing				
Shotgun Metagenomics				
Run Time	12–30 hours	< 1–3.5 days (HiSeq 3000/HiSeq 4000) 7 hours–6 days (HiSeq 2500)	< 3 days	16–36 hours (Dual S2 flow cells) 44 hours (Dual S2 flow cells)
Maximum Output	120 Gb	1500 Gb	1800 Gb	6000 Gb [§]
Maximum Reads Per Run	400 million	5 billion	6 billion	20 billion [¶]
Maximum Read Length	2 × 150 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp

High-Throughput Single Molecule Sequencing Pl

3^a GENERACIÓN

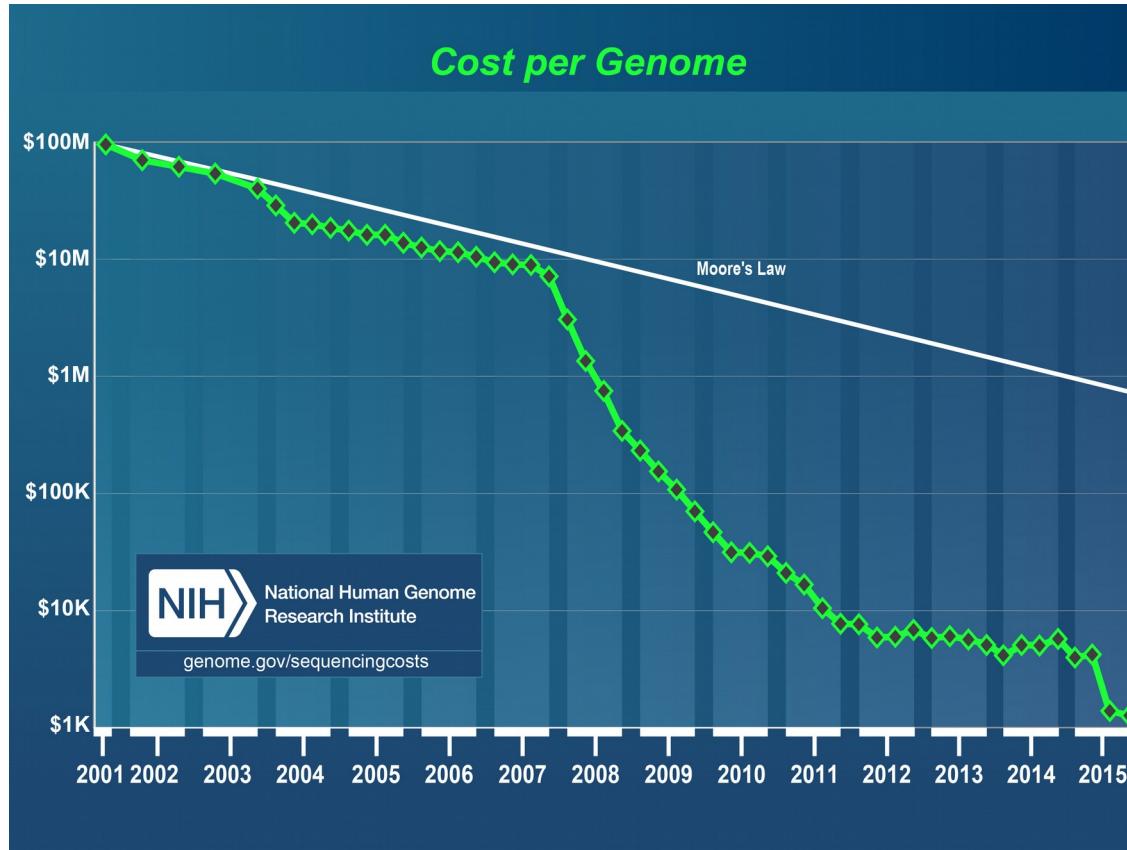


PLATAFORMAS DE SECUENCIACIÓN, 2016 Edición

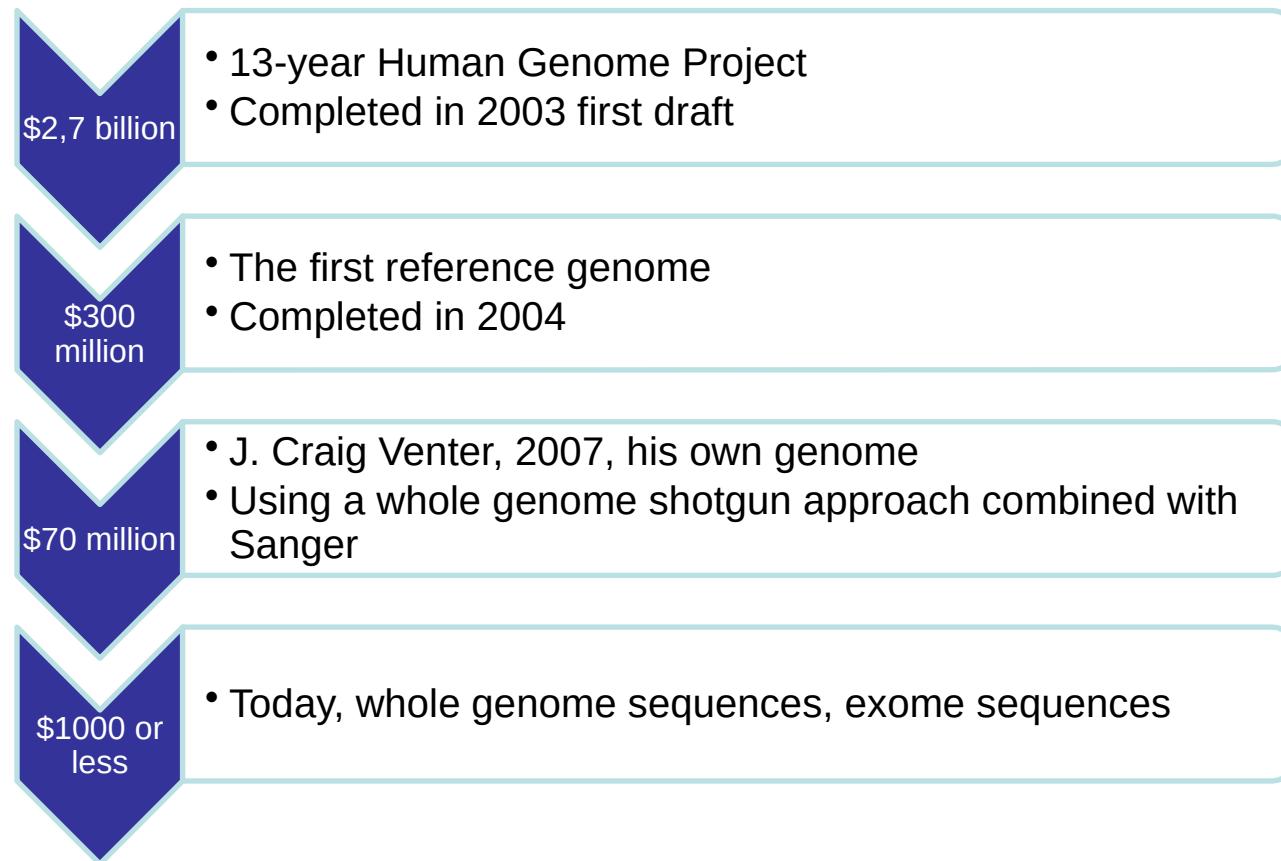


<https://flxlexblog.wordpress.com/>

Coste actual de la secuenciación



Evolución del coste de la secuenciación de u



<https://www.aacc.org/publications/cln/articles/2012/april/sequencing>

Differences Between Sanger and Next Gen samples, tracking, and data relationships

	Sanger	Next Gen
Sequencing Samples	Clones, PCR	DNA libraries
Sample Tracking	Many samples in 96, 384 well plates	Few
Preparation Steps	Few, sequencing reactions clean up	Many, complex procedures
Data collection	Samples in plates 96, 384	Samples on slides 1 - 16+
Data	One read / sample	Thousands and millions reads / sample

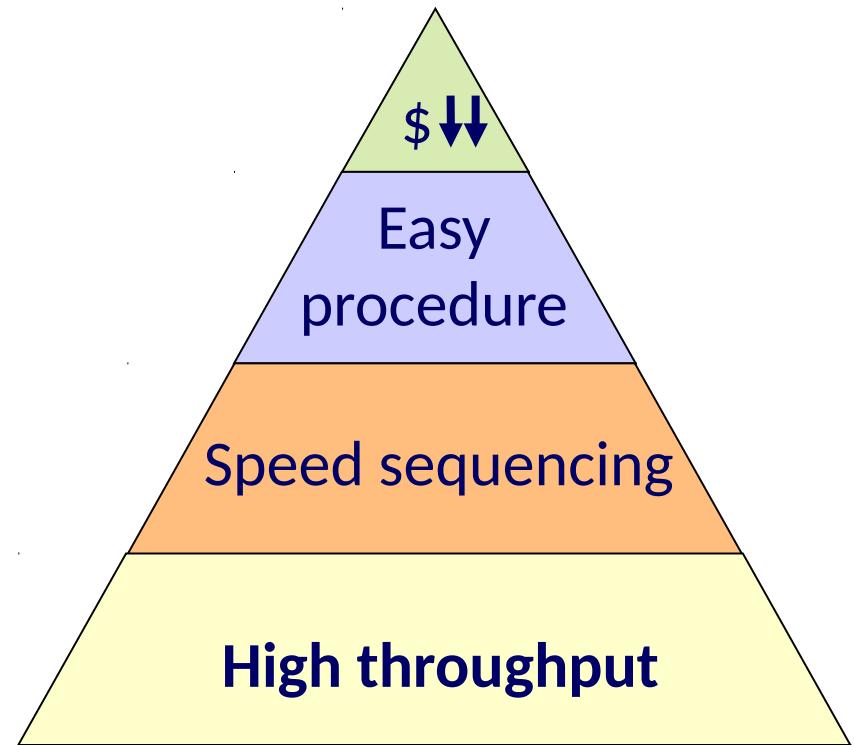
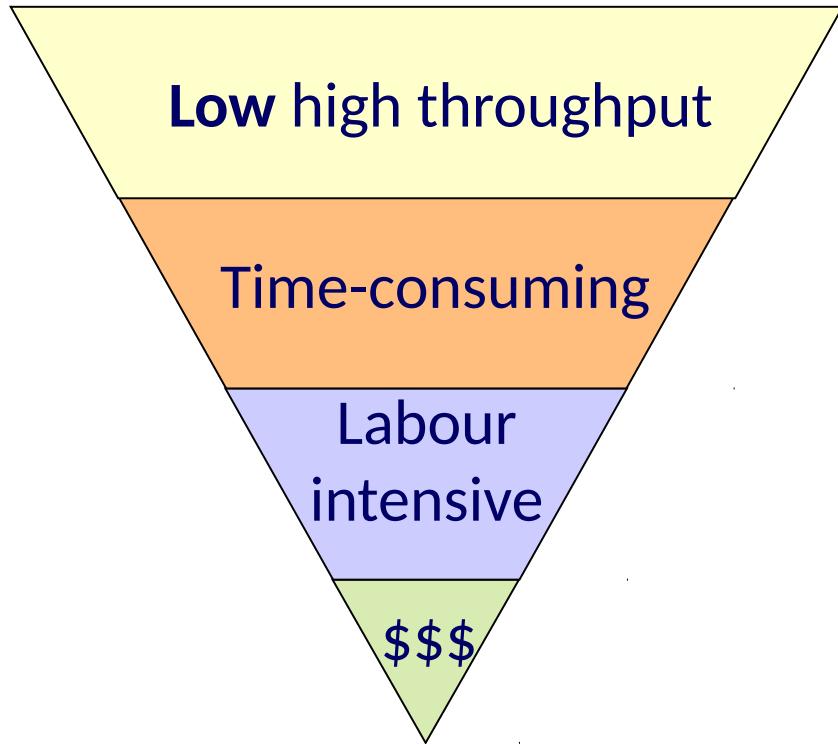
 geospiza™
FROM SAMPLE TO RESULTS™

Copyright 2008

En 2008,
hoy 96+
Kits comerciales

Sanger vs NGS

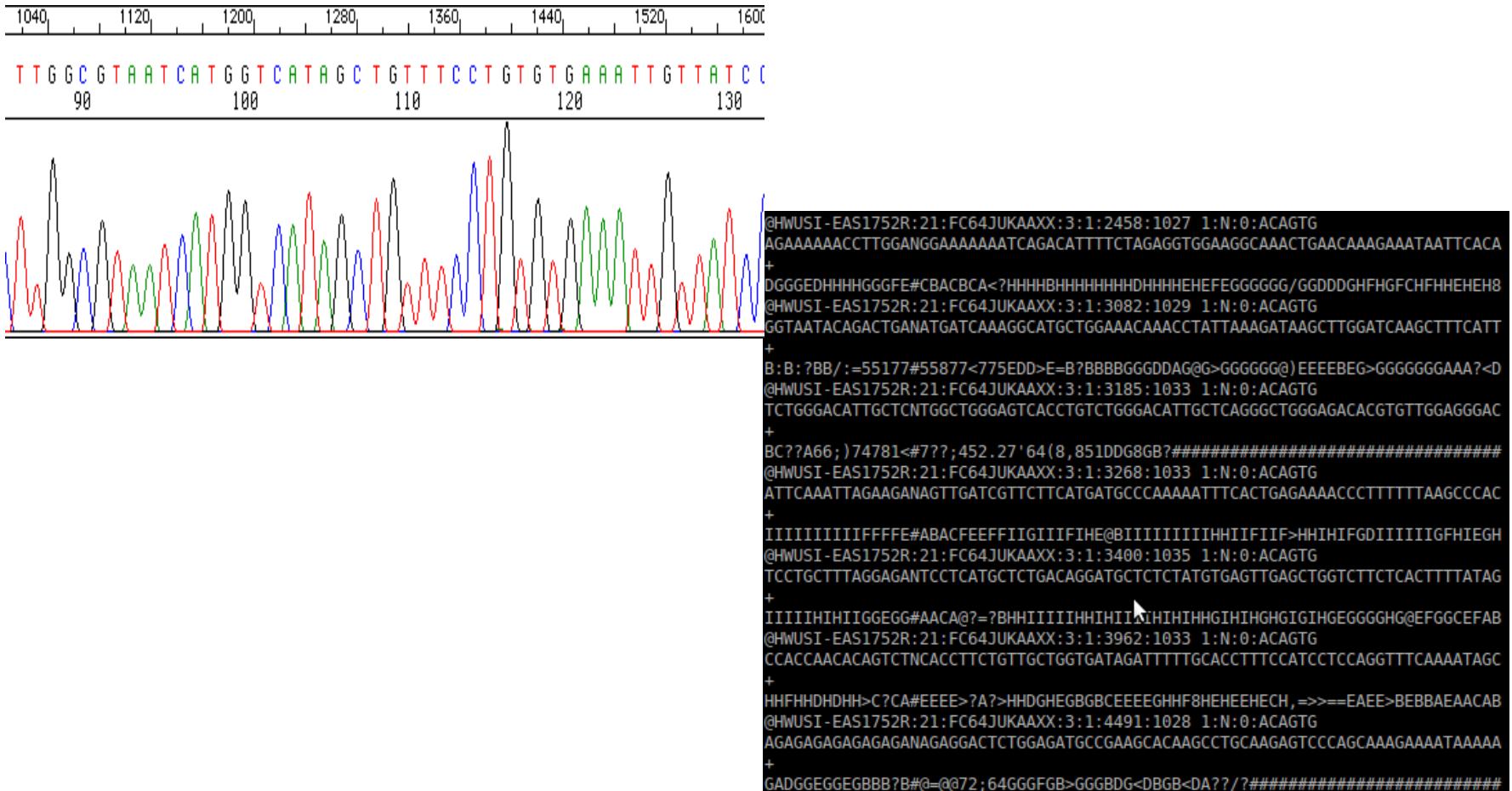
advantages of new technologies



Semiautomatic **Sanger** capillary-based sequencing technology

NGS
Next Generation Sequencing =
Now Generation Sequencing

Nuevo escenario en el análisis de datos



Secuenciadores

Primera generación

- Sanger

Segunda Generación

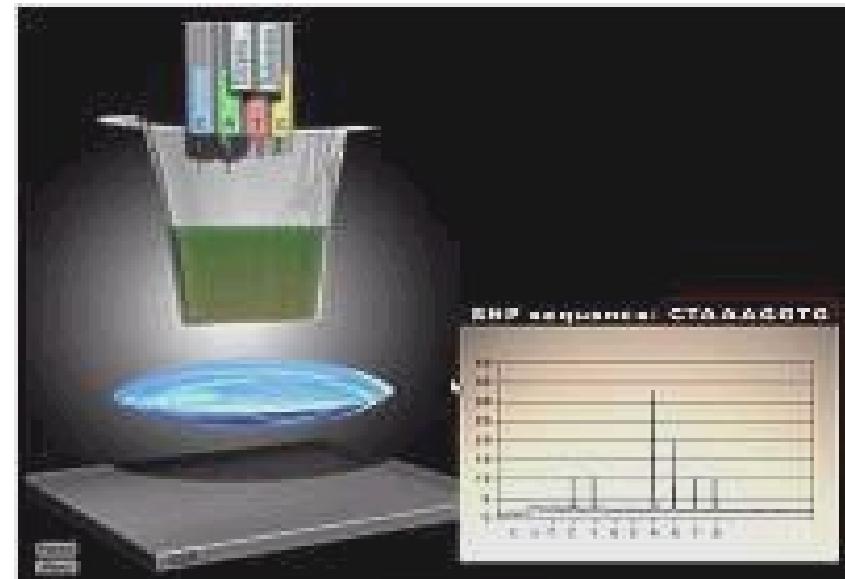
- 454/Roche
- Solexa/Illumina
- Solid
- Ion Torrent

Tercera Generación

- Pacific Biosciences
- Nanopore



PIROSECUENCIACIÓN

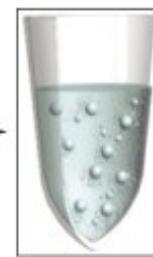
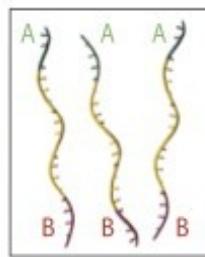
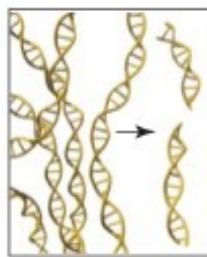




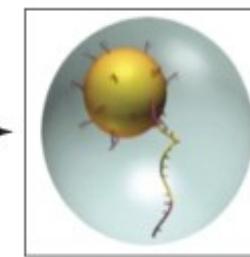
Roche (454) Workflow

Roche (454) GSFLX Workflow:

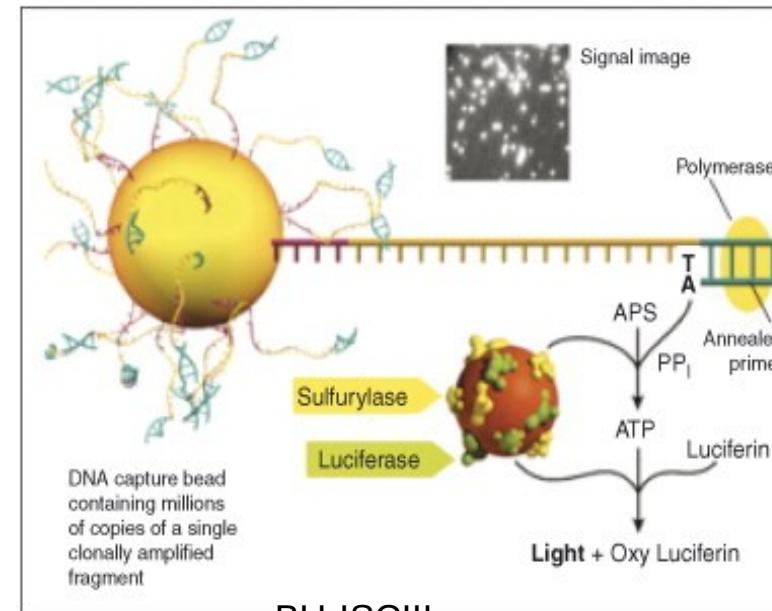
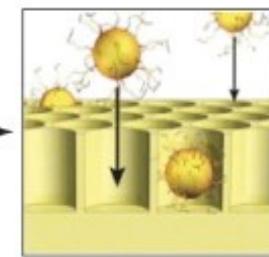
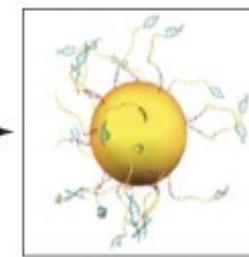
Library construction



Emulsion PCR



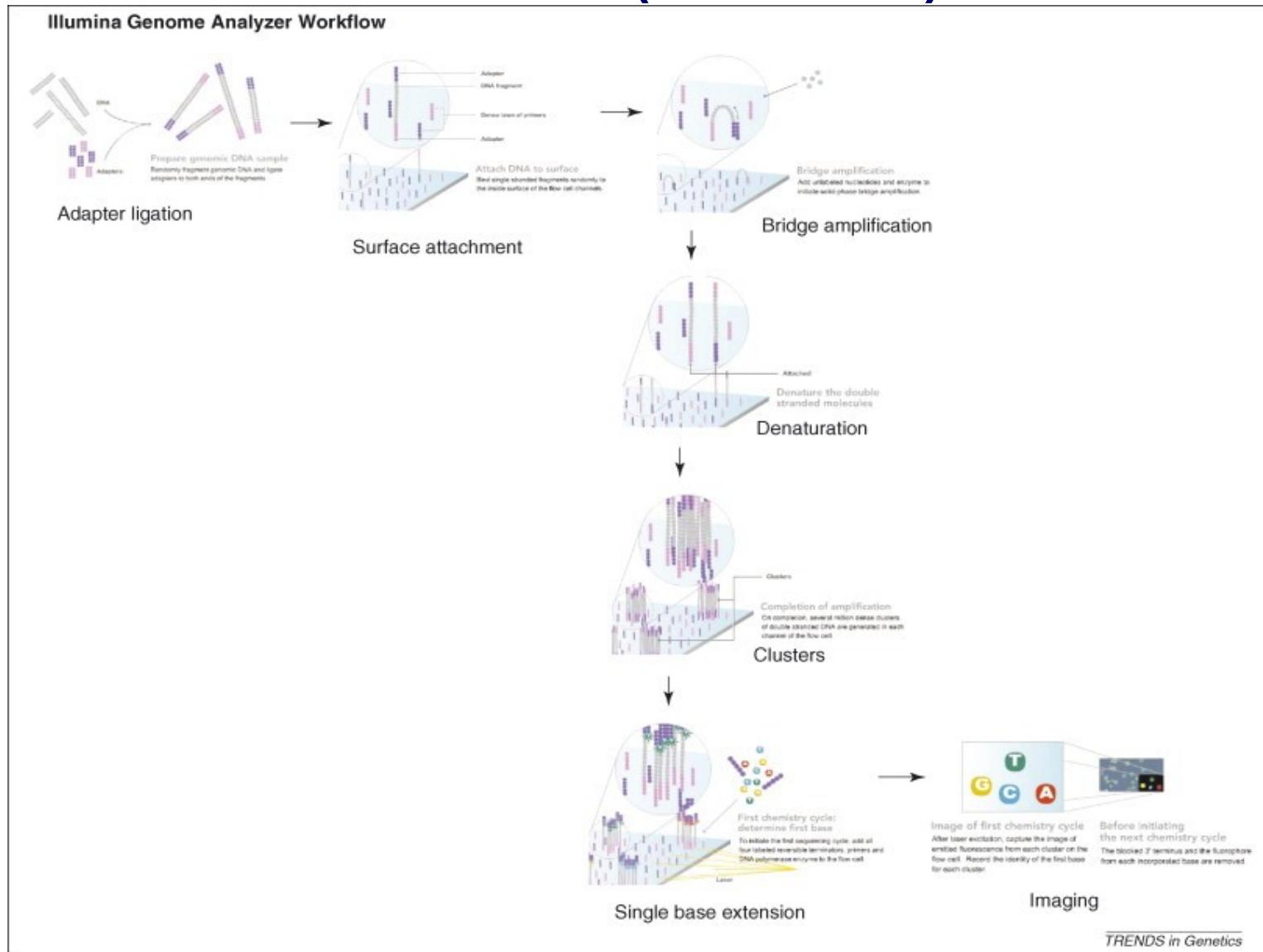
PTP loading



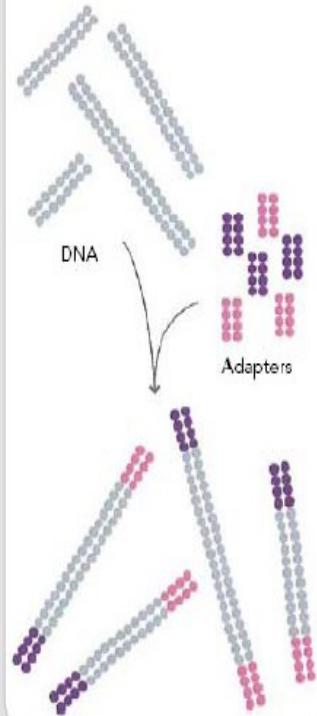
BU-ISCHII

Pyrosequencing reaction

Illumina (Solexa) Workflow

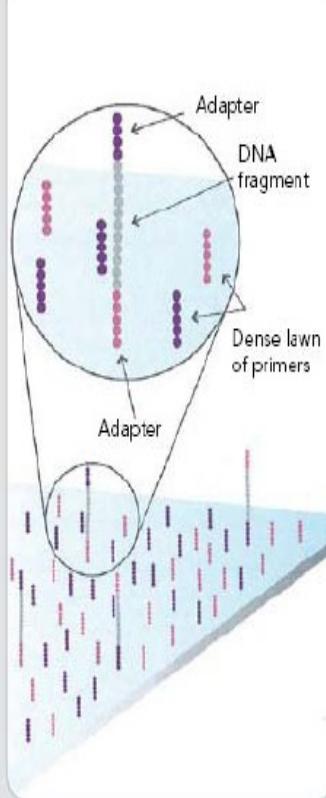


1. PREPARE GENOMIC DNA SAMPLE



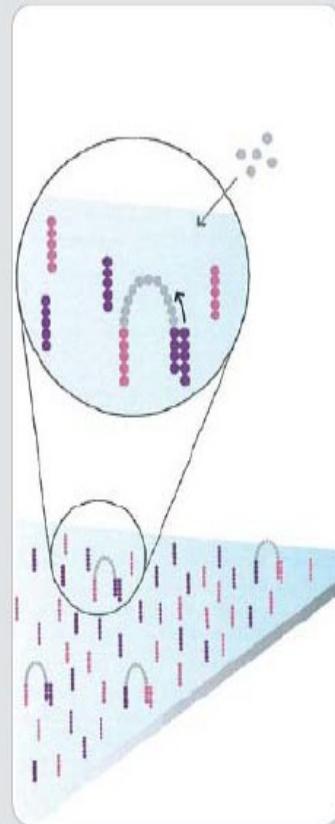
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



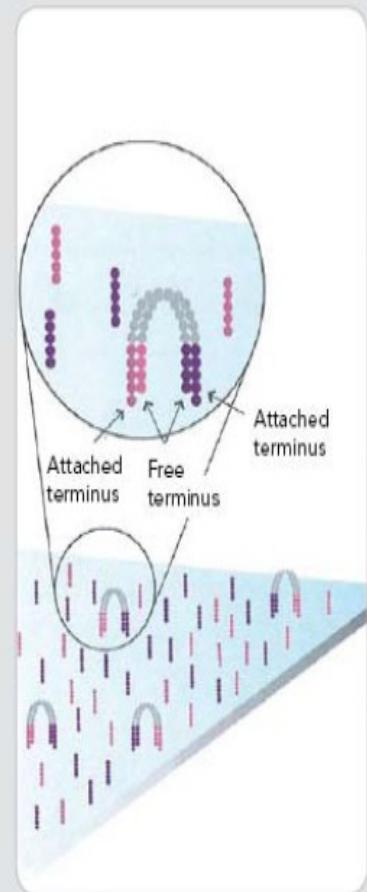
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

3. BRIDGE AMPLIFICATION



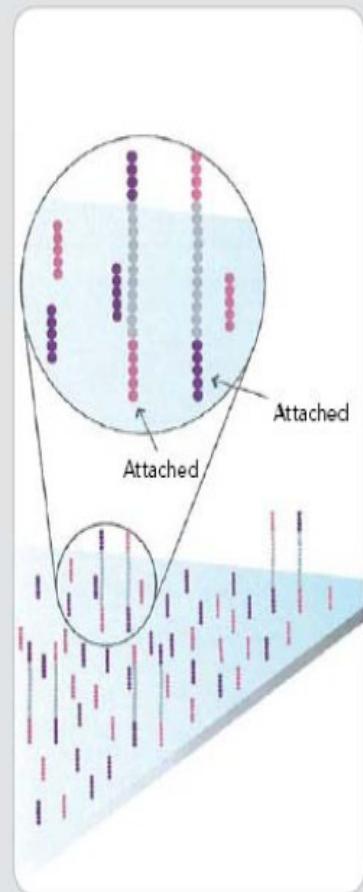
amplificación en fase sólida

4. FRAGMENTS BECOME DOUBLE-STRANDED



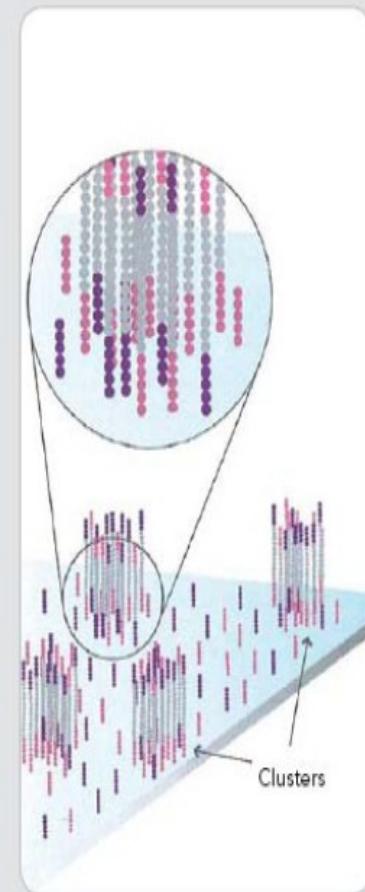
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



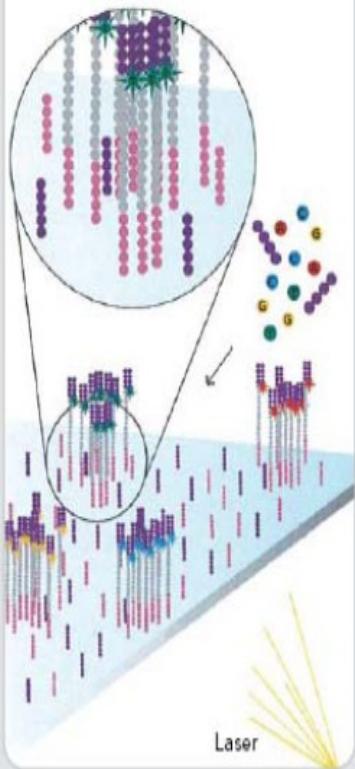
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



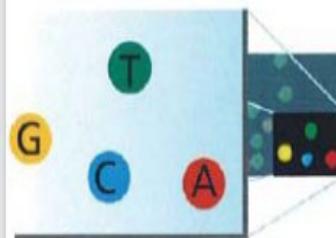
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

7. DETERMINE FIRST BASE



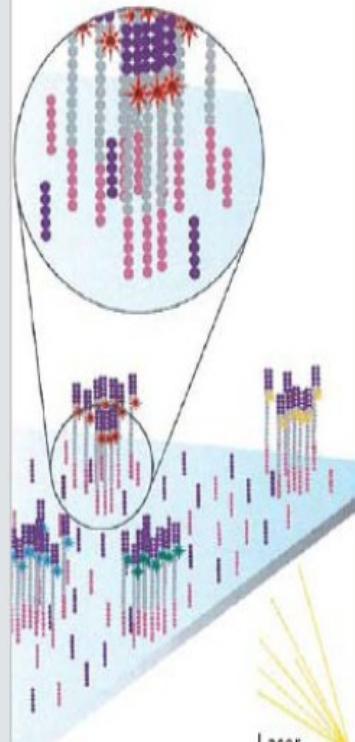
The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.

8. IMAGE FIRST BASE



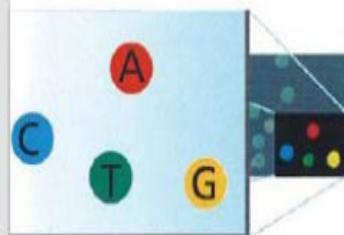
After laser excitation, the emitted fluorescence from each cluster is captured and the first base is identified.

9. DETERMINE SECOND BASE



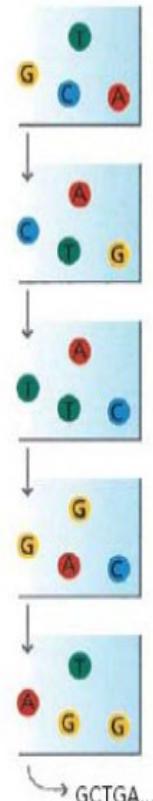
The next cycle repeats the incorporation of four labeled reversible terminators, primers, and DNA polymerase.

10. IMAGE SECOND CHEMISTRY CYCLE



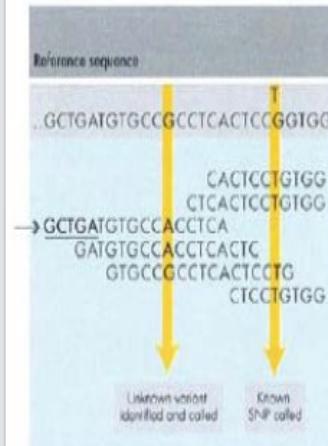
After laser excitation, the image is captured as before, and the identity of the second base is recorded.

11. SEQUENCING OVER MULTIPLE CHEMISTRY CYCLES



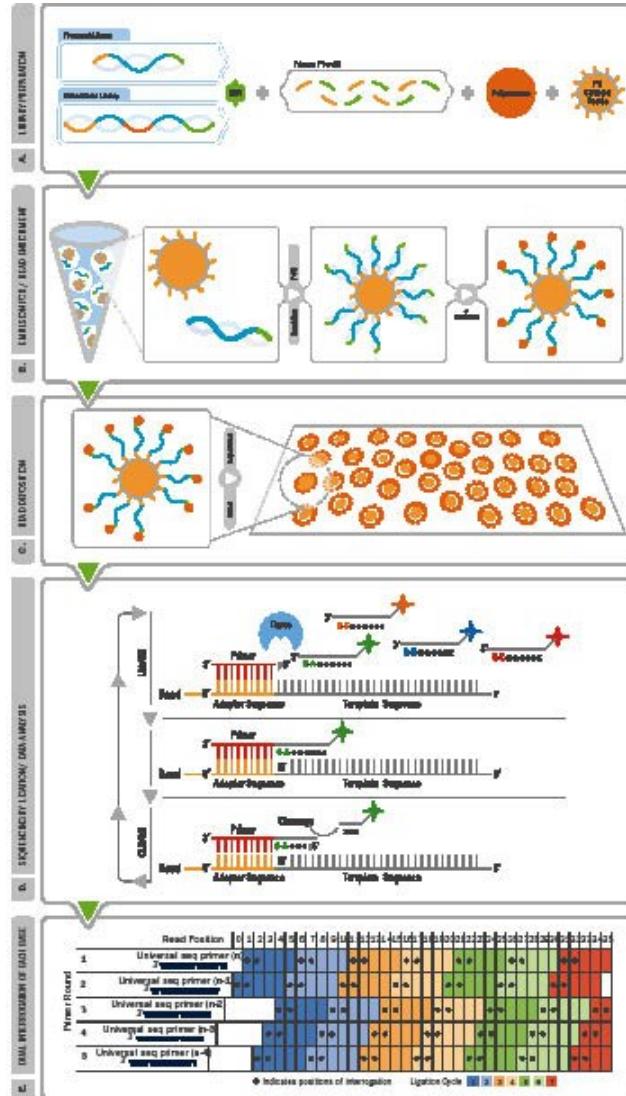
The sequencing cycles are repeated to determine the sequence of bases in a fragment, one base at a time.

12. ALIGN DATA



The data are aligned and compared to a reference, and sequencing differences are identified.

Life Technologies / Solid



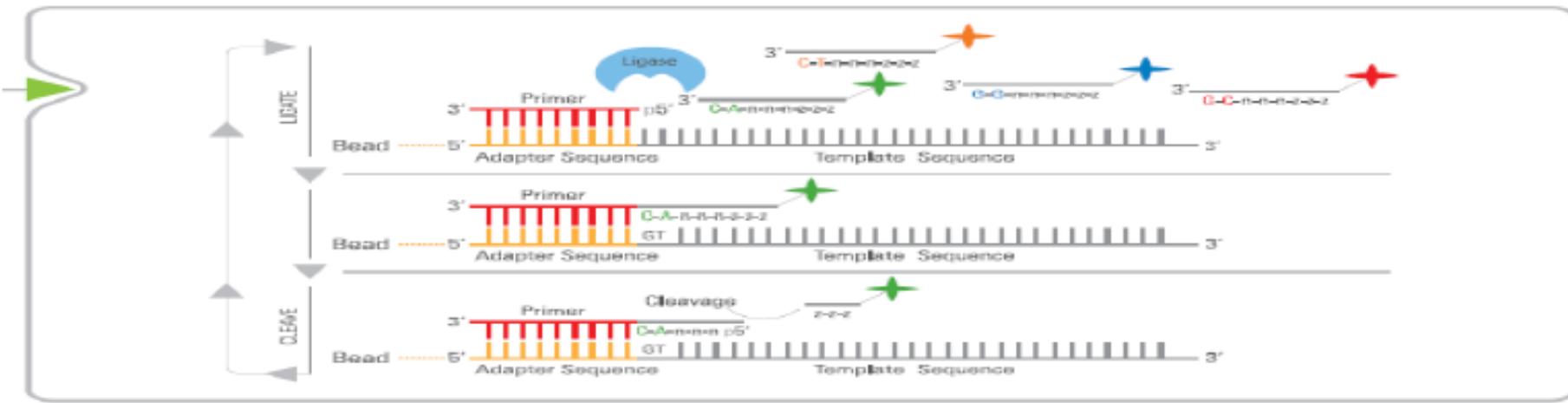
Library Construction

Emulsion PCR

The beads are deposited onto a slide

Life Technologies / Solid

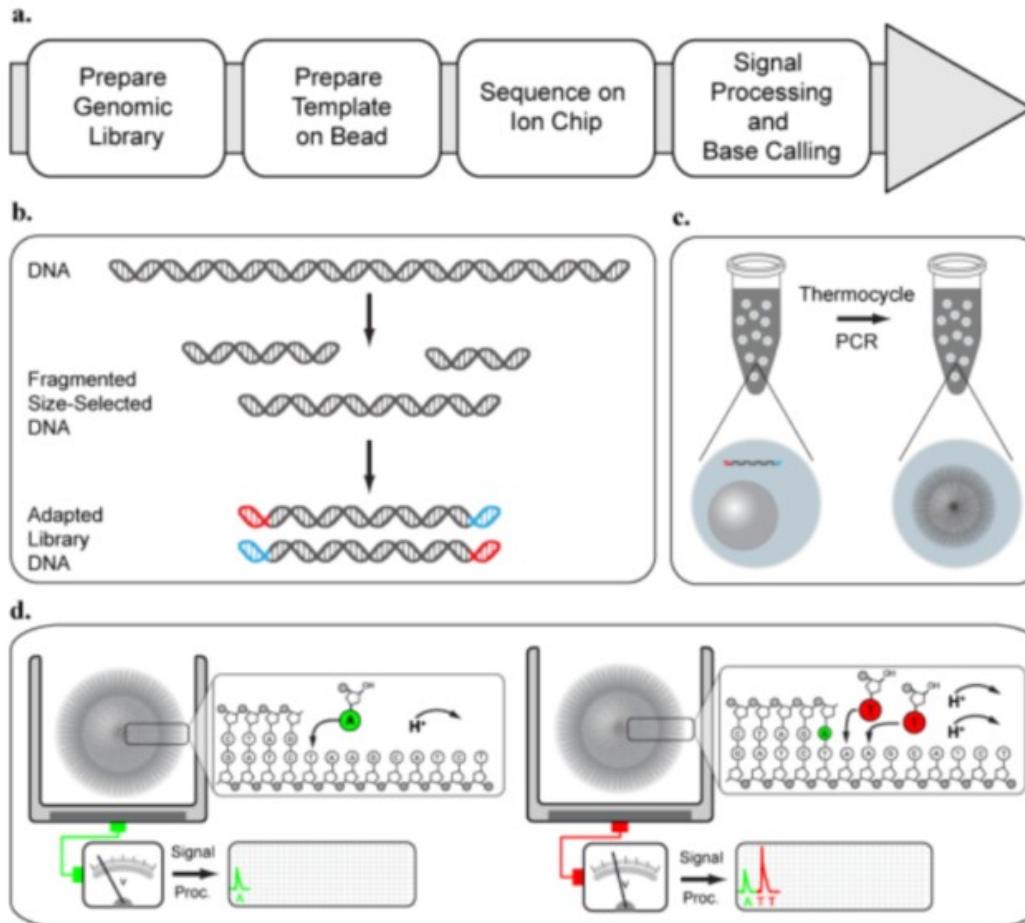
Sequencing by Ligation



Primers hybridize to the P1 adapter sequence within the library template. A set of four fluorescently labeled di-base probes compete for ligation to the sequencing primer. Specificity of the di-base probe is achieved by interrogation every 1st and 2nd base in each ligation reaction. Multiple cycles of ligation, detection and cleavage are performed.

Ion Torrent PGM

Personal Genome Machine



Secuenciadores

Primera generación

- Sanger

Segunda Generación

- 454/Roche
- Solexa/Illumina
- Solid
- Ion Torrent

Tercera Generación

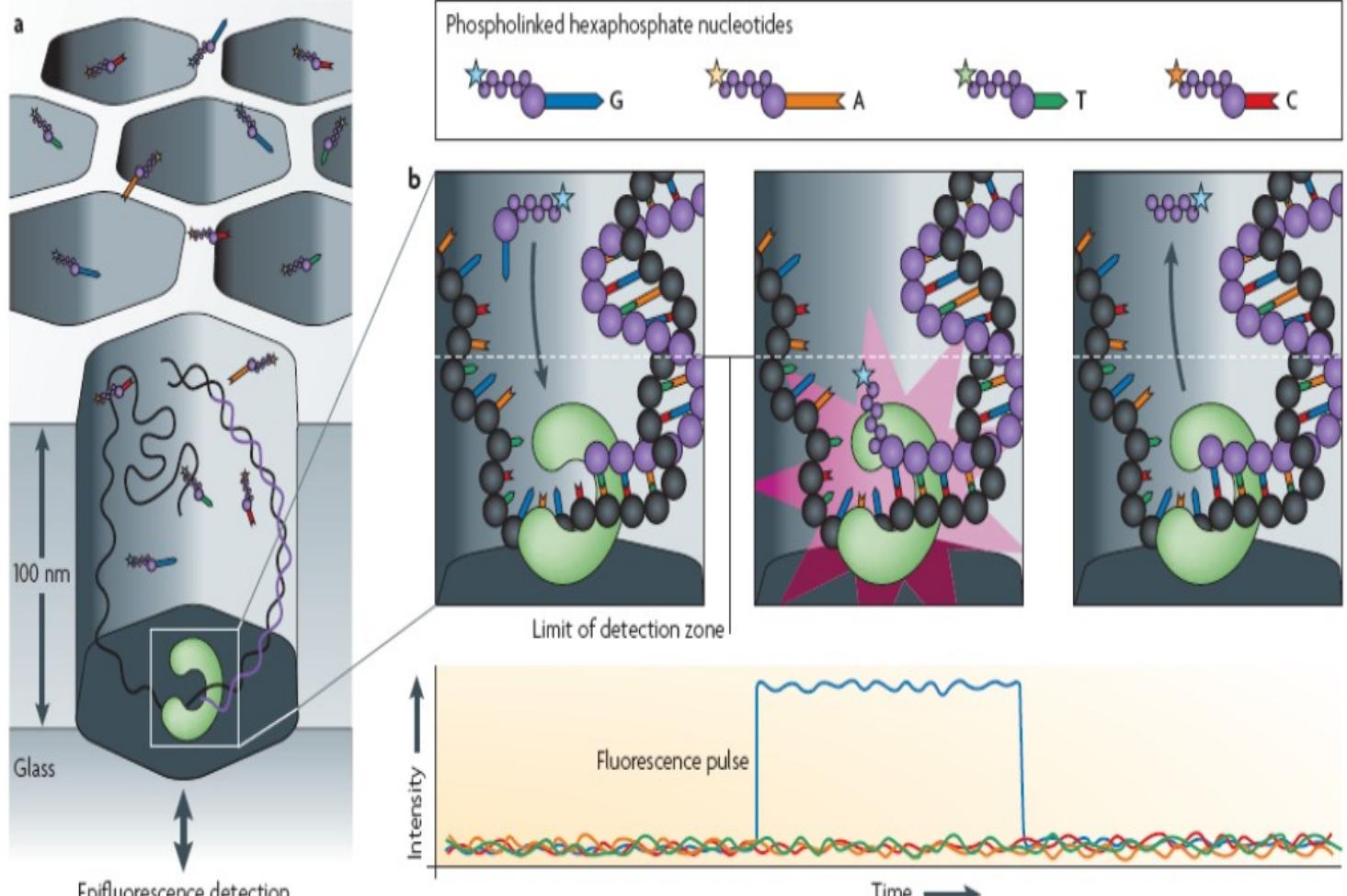
- Pacific Biosciences
- Nanopore

3^a GENERACIÓN: LECTURAS MAS LARGAS

- PacBio, PACIFIC BIOSCIENCE
- Moleculo, ILLUMINA
- MinION, GridION, OXFORD NANOPORE

ULTRASECUENCIACIÓN A TIEMPO

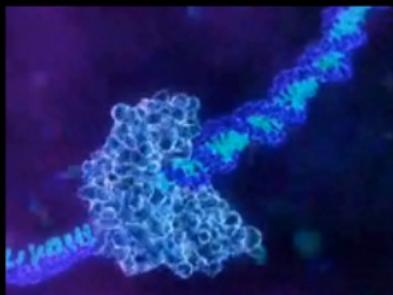
Pacific Biosciences — Real-time sequencing



ULTRASECUENCIACIÓN A TIEMPO

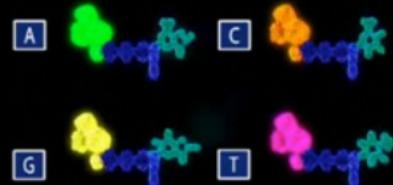
Next-gen sequencing: Pacific Biosciences

1



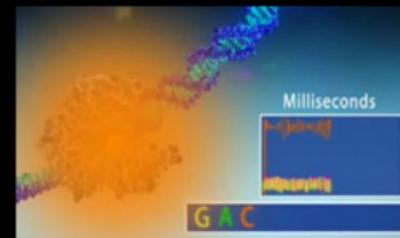
DNA polymerase wrapped around DNA chain

2

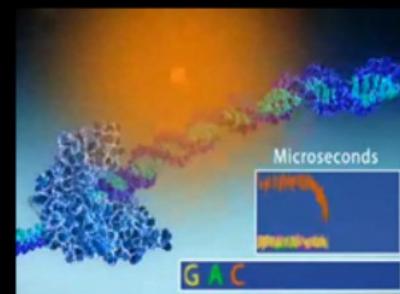


Phospholinked nucleotides

3a



3b

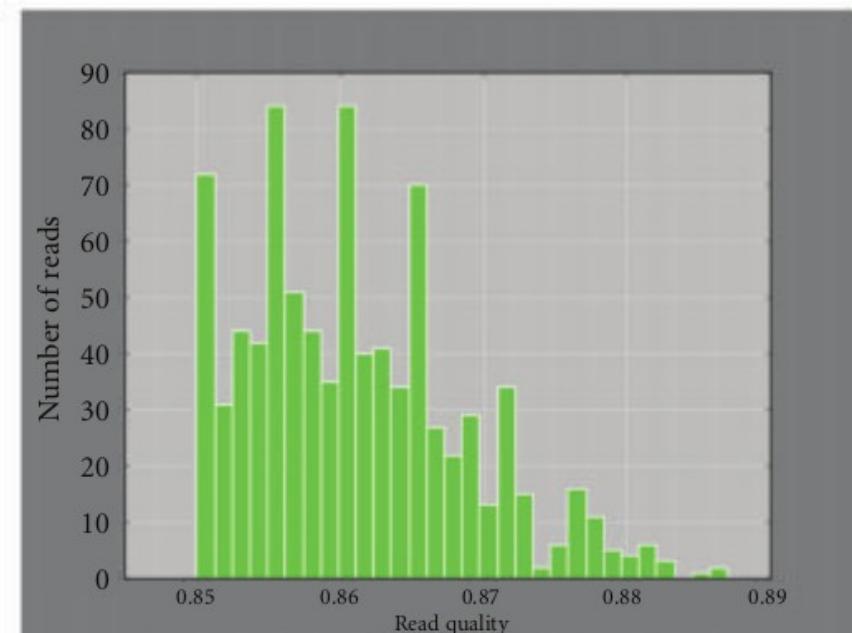
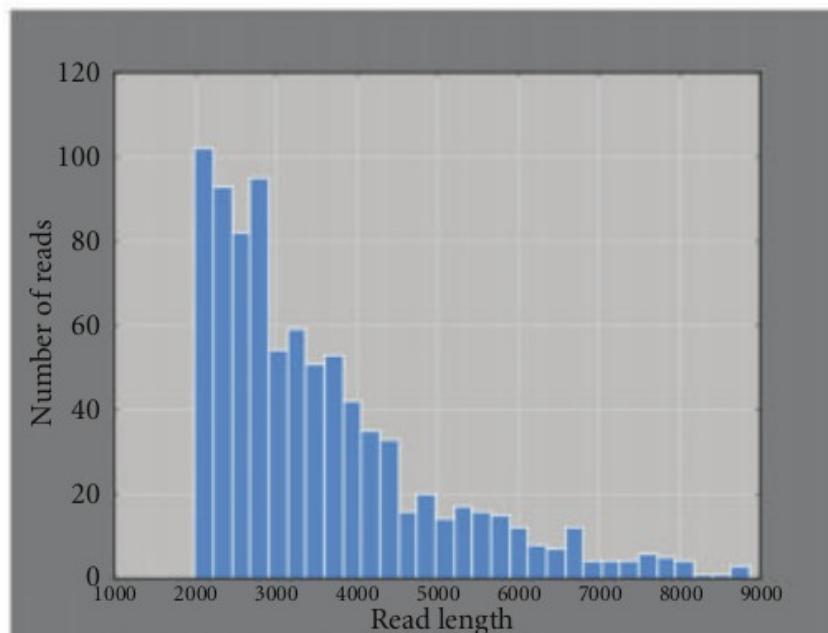


Phospholinked nucleotide binds, fluoresces and detaches as nucleotide base is read

Secuenciación de un fosmido DNA usando el secuenciador PacBioscience

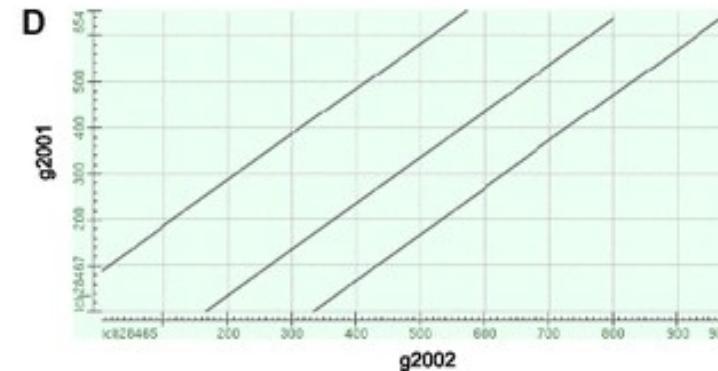
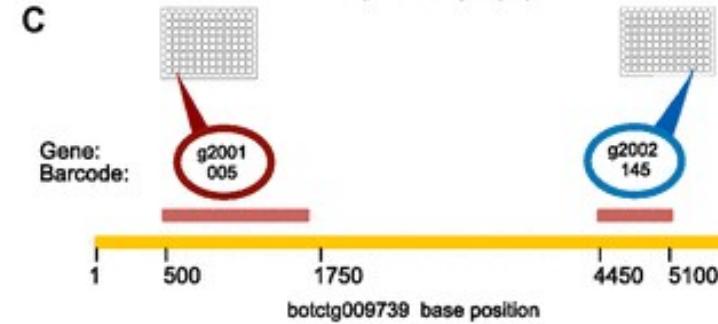
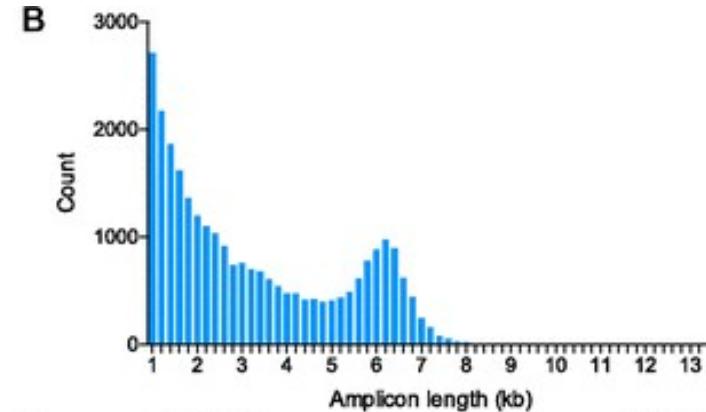
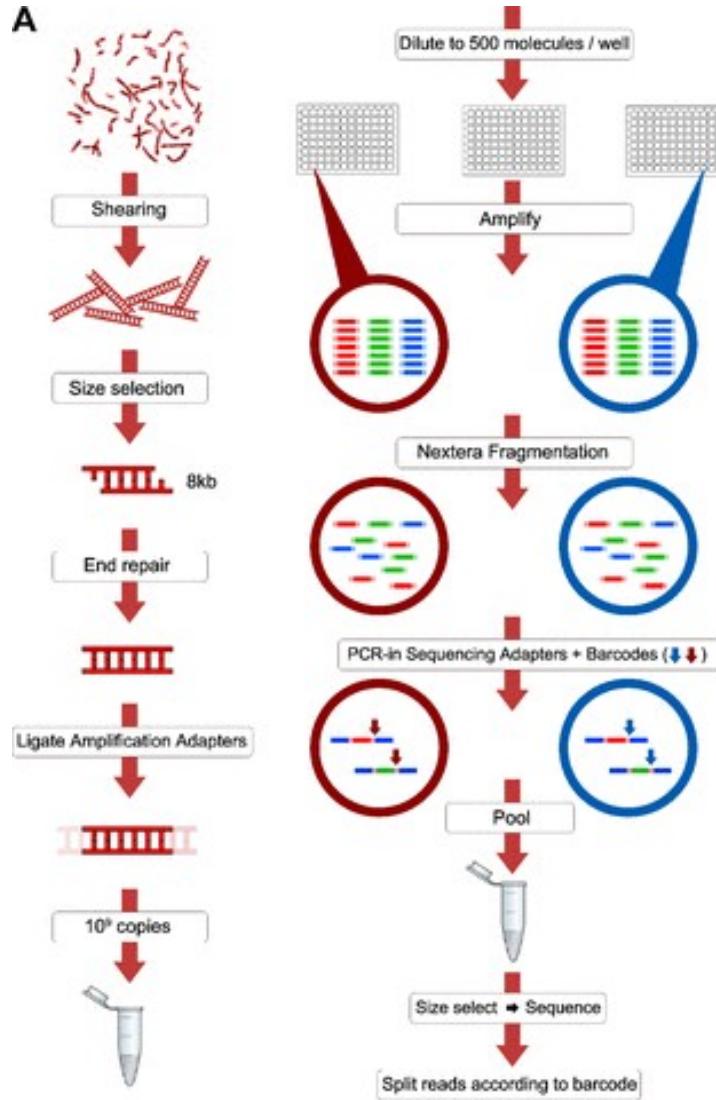
	Prefilter	Post-QC filter*
Number of bases	84,110,272 bp	22,373,400 bp
Number of reads	46,861	6,754
Mean read length	513 bp	2,566 bp
Mean read score	0.144	0.819

* MinRL = 50, MinRS = 0.75



Liu et al., J. Biomedicine and Biotechnology, 2012

MOLECULO, ILLUMINA

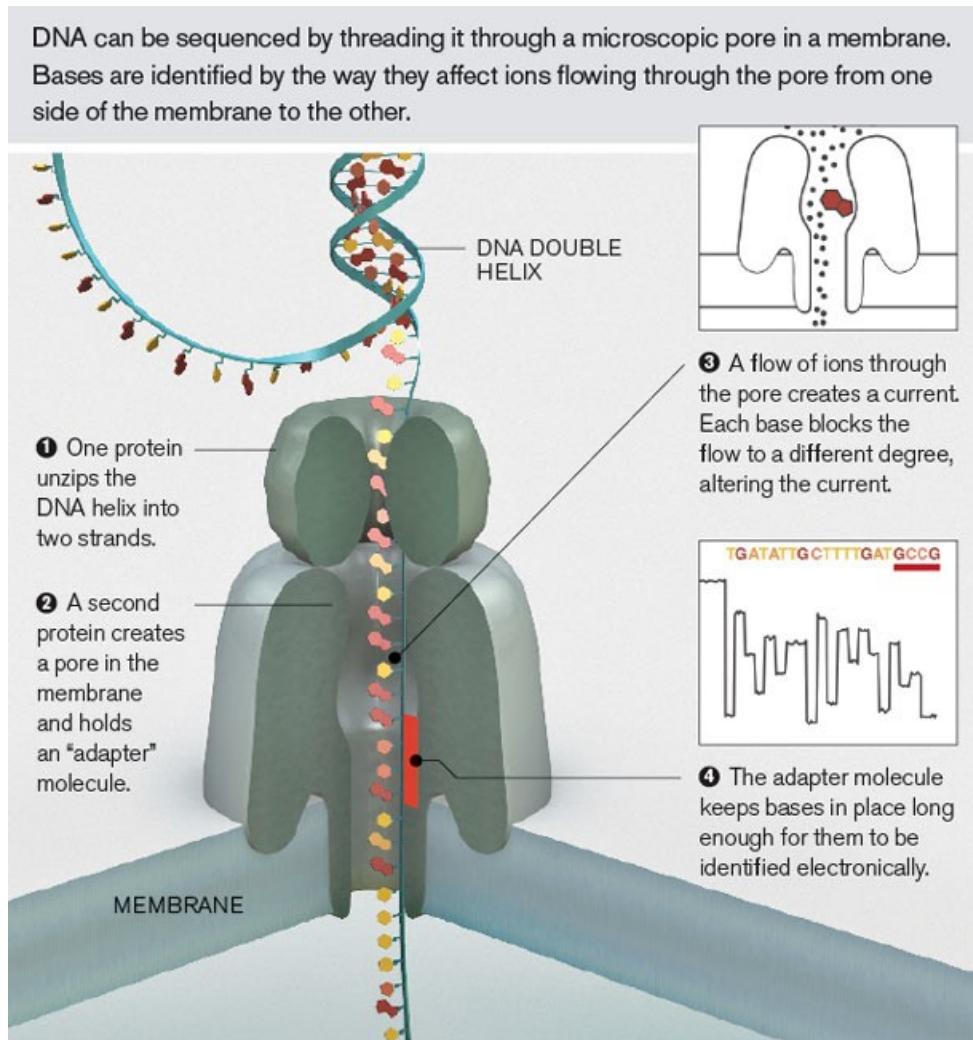


Moleculo, acquired by Illumina in late 2012, developed an innovative technology for generating long reads that combines a new library prep method and genome analysis tools

BU-ISCM

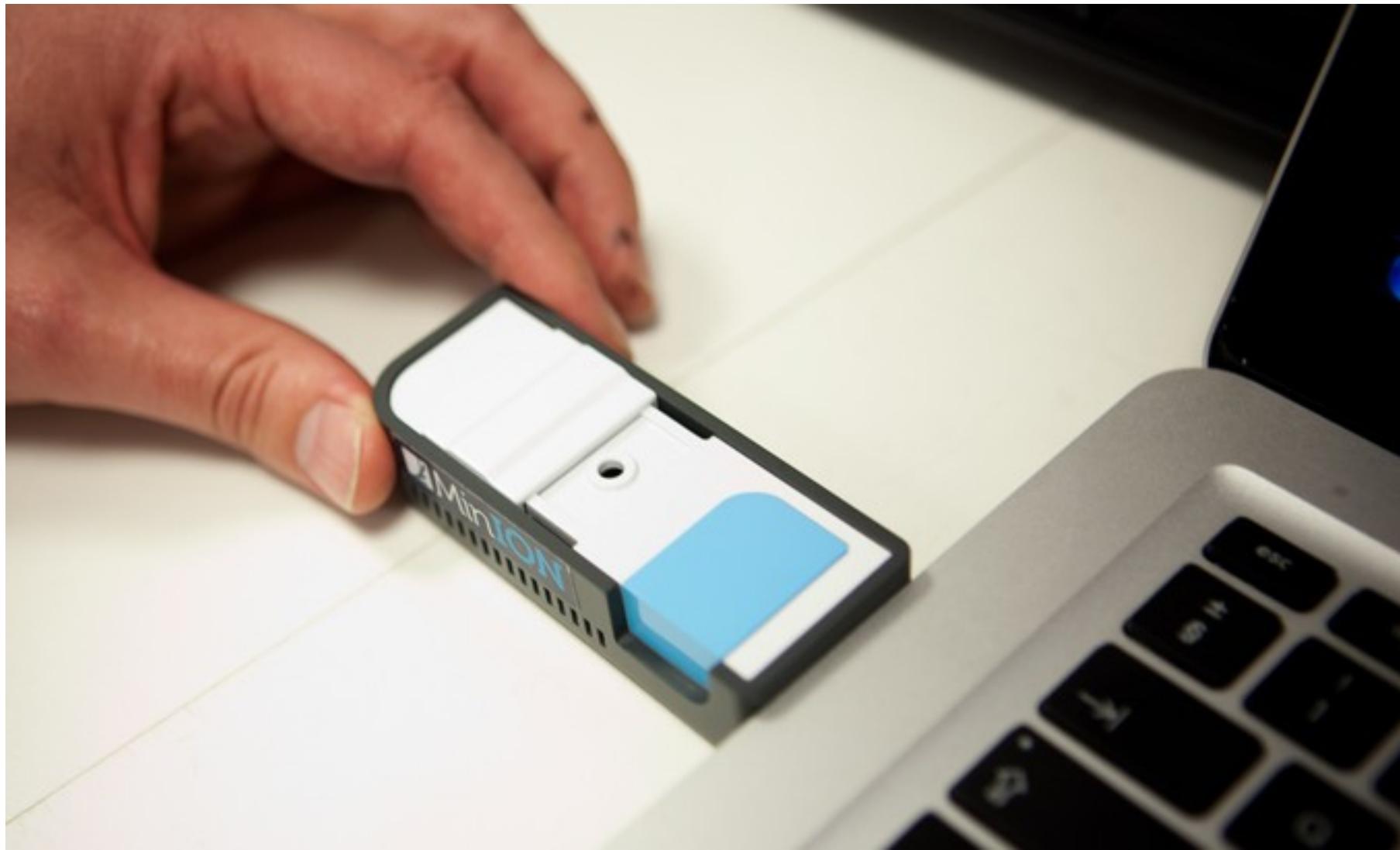
Voskoboynik et al., elife 2013;2:e00569

MinIon, OXFORD NANOPORE



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

MinIon, OXFORD NANOPORE



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

OXFORD NANOPORE



MinION

- Pocket-sized
- Up to 512 nanopore channels
- Simple 10-minute sample prep available
- Many publications illustrate broad usage
- Commercially available

[Learn more](#)

[Buy MinION](#)



GridION

- Multiple sequencing devices, one compute module
- Use up to five MinION Flow Cells at a time
- Benchtop processor capable of handling high data volumes in real time
- Rapid, real-time applications such as Read Until ...

[Learn more about GridION](#)



PromethION

- Benchtop system
- Up to 48 flow cells of 3,000 nanopore channels each
- Simple 10-minute sample prep available
- Flow cells may be run individually or concurrently

[Learn more](#)

[Register for early access](#)

Community

[The MinION Access Programme](#)

Specifications

MinION Flow cell pricing

MinION progress so far

MAP FAQs

MAP registration form

Joining MAP – what you need to know

Publications from the MAP

The MAP community login

PEAP (PromethION Early Access Programme)

The MinION Access Programme

A Community to support your use of MinION



To start to use a MinION, you must register to [join the MinION Access Programme](#). The MinION Access Programme (MAP) started as a broad, community-focused access project which started in spring 2014. MinION is now commercially available and the Community provides support and discussions around development of new tools, applications and methods.

Listening to this community helps Oxford Nanopore provide continuous improvements to our products and support. MinION is currently enabled for nucleic acid analysis, and in time is expected to encompass the direct analysis of proteins.

To join the MAP you will need to pay the \$1,000 access fee and agree to the Terms and Conditions. Subsequently your MinION starter pack will be shipped to you and you will be able to purchase additional supplies immediately.

As part of the MAP you will also receive periodic free supplies of consumable items. Purchased items are warranted and take delivery priority over purchased goods.

[Register an account](#) now to join experimental teams around the world or read the [latest publications](#) that have come out of the MAP community.

The MinION Access Programme (MAP) is a community-focused technology access programme that started in April 2014 and includes participants in more than 30 countries. Launched as an 'early-access' community the MAP seeks to enable a broad range of people to explore how the MinION may be useful to them in nucleic acid analysis, to contribute to developments in analytical tools and applications and to share their experiences and collaborate. The MAP is now open to any registrant. See tools and publications at <http://publications.nanoporetech.com>

How to start using MinION

1 What you need to know

Useful information including payment details, laboratory and IT requirements

[Before you register](#) >

2 Register an account

We will confirm your account within 3 working days. You will then be able to pay the access fee of \$1,000 and accept the terms & conditions. You will then have full access to the MAP community and will receive a MinION as part of the starter pack.

The Nanopore analysis channel on F1000Research will be a central, open platform on which scientists can publish and discuss new applications and analyse workflows for nanopore sequencing data. People can access and contribute data easily, so that the wider life-science community can realise the full potential of this new technology quickly."



NANOPORE
APPLICATIONS

Sequencing, analysing and counting short tandem repeats on the MinION™

MinION sequencing of STRs with accurate analysis and counting, without assembly or complex statistics, enables a range of biological applications

Content: minionnanopore.com. More information at www.nanoporetech.com and publications nanoporetech.com

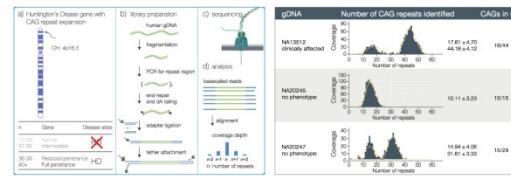


Fig. 1: Expansion of Huntington's Disease gene

Clinical relevance of sequencing STRs

Huntington's Disease (HD) is a triplet expansion disease. The 'D' gene (HTT), on chromosome 4p16.3, contains a series of consecutive CAG nucleotides (Fig. 1a). A polyC tract >30 causes the formation of a mutant protein that accumulates in the nucleus, leading to progressive cognitive decline, motor impairment, and the number of repeats correlates to severity and age of onset. We PCR-amplified the locus, prepared libraries (Fig. 1b) and sequenced (Fig. 1c). The coverage depth and number of repeats with different numbers of repeats (Fig. 1d). The repeat count corresponds to the peak coverage.

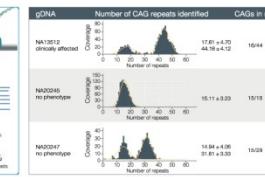
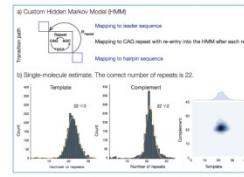


Fig. 2: Proof-of-concept for HD using cell-line gDNA

Counting STRs in Huntington's Disease cell-line samples

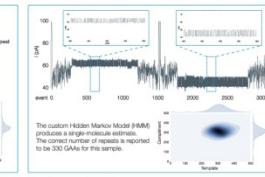
STRs can be very long, over 8k, and some repeats are impossible to sequence due to heterozygosity or other factors. We sequenced and analysed three Huntington's Disease (gDNA) samples, acquired from the NIGMS human genetic cell repository (Fig. 2). The samples were homozygous (Na13212), heterozygous (Na32045), and homozygous (Na20247). The heterozygous and homozygous samples generate double and single peaks of coverage respectively, as expected. The positions of the coverage peaks are in good agreement with the reference value also given in the 3rd column.



Repeat estimation in 'squiggle-space' using a custom HMM

Instead of performing the repeat-length analysis on binned data, we can apply a custom Hidden Markov Model (HMM) (Fig. 3a) to the pre-binned 'squiggle' data. The HMM predicts the number of repeats based on counting the repeats in each individual read. Complex repeats and arbitrary structure variations are addressable by varying the structure of the HMM.

For example, we show a repeat estimate for the DMPK gene, implicated in myotonic dystrophy type 1 (DM1). DM1 is caused by an expansion of a GAA triplet repeat, which is usually symptomless, and longer repeats are associated with early onset, and greater severity.



Counting long repetitive regions using raw data and custom HMM

Friedreich's ataxia (FRDA) is an inherited neurodegenerative disease caused by a GAA triplet expansion in the first intron of the frataxin gene. This expansion causes gene silencing by inducing a heterochromatin structure in the DNA at this point. Patients typically have >700 copies of the triplet repeat. We sequenced a sample containing 330 copies of the repeat (Fig. 4). Using the raw data and the custom HMM, we were able to sequence and analyse a sample containing 330 copies of the repeat (Fig. 4). The distribution of repeat counts in our analysis is likely to have arisen from sample heterogeneity and PCR leakage.

News Press releases

Overview | Leadership | Funding | Background | Collaborations | Jobs | People & groups | **News** | Events | Visit us | Contact

About us > News > Press releases > Mini DNA sequencer tests true

Mini DNA sequencer tests true



- o Public access to Oxford Nanopore's MinION™ miniature sensing device enabled an international consortium to evaluate the technology and provide a standard protocol for its use;
- o Preliminary analysis of data generated in five very different laboratories indicates the performance and accuracy of the device is consistently good;
- o Data are freely available for re-analysis and innovation in the Nanopore analysis channel on [F1000Research](#).

15 October 2015 – The MinION, a handheld DNA-sequencing device developed by Oxford Nanopore, has been tested and evaluated by an independent, international consortium coordinated by EMBL's European Bioinformatics Institute (EMBL-EBI). The innovative device opens up new possibilities for using sequencing technology in the field, for example in tracking disease outbreaks, testing packaged food or the trafficking of protected species.

<http://www.ebi.ac.uk/about/news/press-releases/mini-dna-sequencer-tests-true>

SmidgION, OXFORD NANOPORE



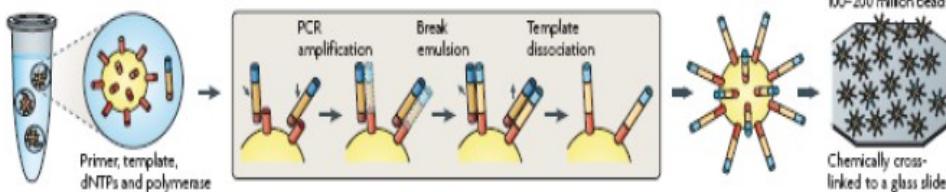
One of the biggest highlights from the talk is SmidgION, the smallest sequencing device that can be plugged into a smartphone for sequencing and analysis. SmidgION will have 256 pores, can sequence for four hours with a smartphone, and produce 230 Mb/hour. SmidgION is expected to be available in late 2017.

<http://nextgenseek.com/2016/05/oxford-nanopore-introduces-smidgion-at-london-calling-2016/>

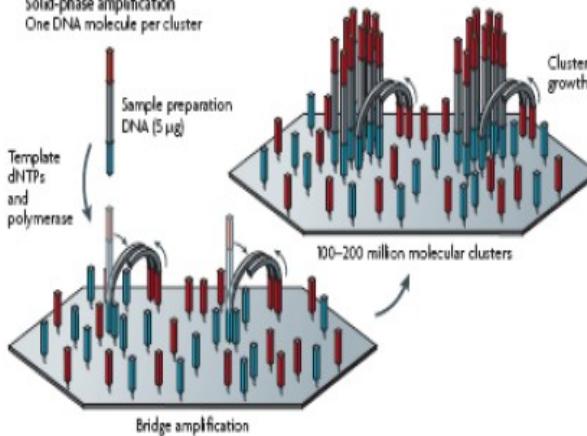
ESTRATEGIAS DE INMOVILIZACIÓN

a Roche/454, Life/PG, Polonator
Emulsion PCR

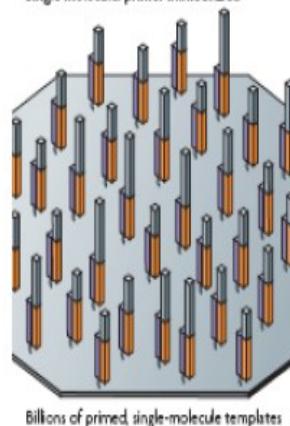
One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion



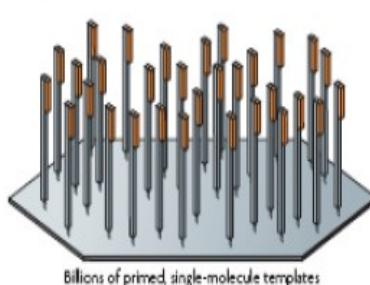
b Illumina/Solexa
Solid-phase amplification
One DNA molecule per cluster



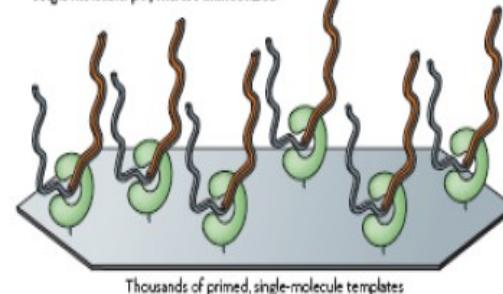
c Helicos BioSciences: one-pass sequencing
Single molecule: primer immobilized



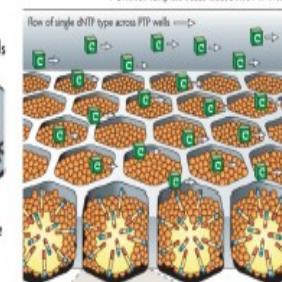
d Helicos BioSciences: two-pass sequencing
Single molecule: template immobilized



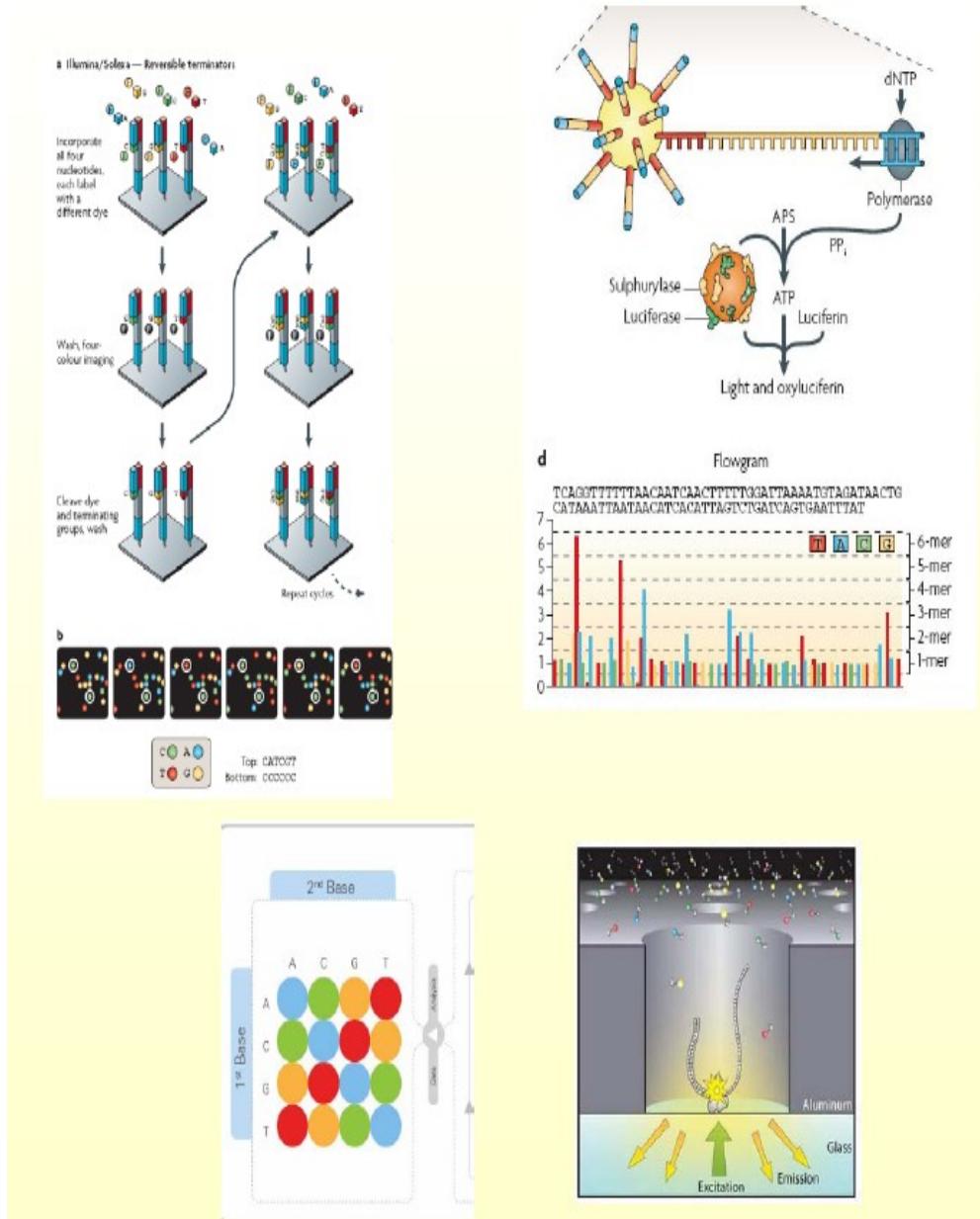
e Pacific Biosciences, Life/Vision, LI-COR Biosciences
Single molecule: polymerase immobilized



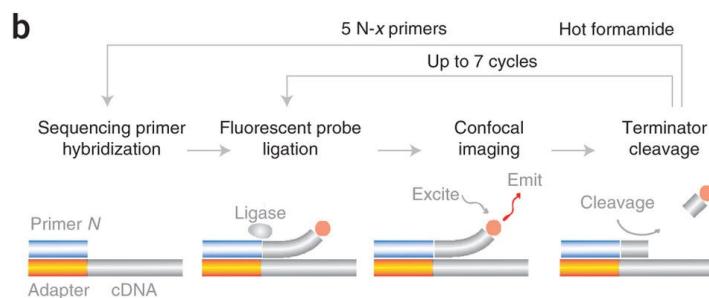
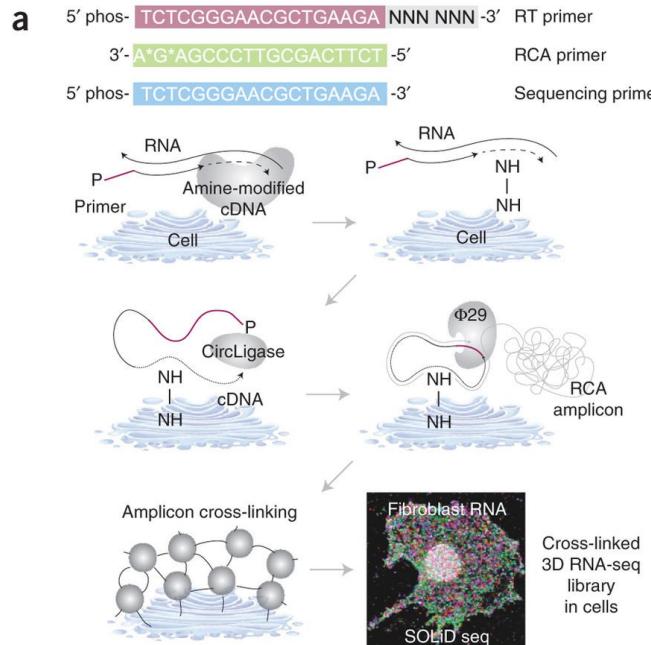
c Roche/454 — Pyrosequencing
1–2 million template beads loaded into PTP wells



ESTRATEGIAS DE LECTURA



Secuenciación in situ!



Fluorescent in situ sequencing (FISSEQ) of RNA for gene expression profiling in intact cells and tissues

Lee et al., Nat protocols 10 442-458 (2015)

Comparative of NGS platforms

Next-Gen Sequencer	Machine Cost	Cost per run	Minimum Throughput	Sequencing Run Time	Cost Per Mb
Illumina MiSeq	\$125,000	\$750	1500 Mb (2 x 150 Bases)	27 Hours	\$0.5
454 GS Junior	\$108,000	\$1,100	35 Mb (400 Bases)	8 Hours	\$31
Ion Torrent PGM - 314 Chip	\$80,490	\$225	10Mb (100 Bases)	3 Hours	\$22.5
Ion Torrent PGM - 316 Chip	\$80,490	\$425	100Mb	3 Hours	\$4.25
Ion Torrent PGM - 318 Chip	\$80,490	\$625	1000Mb	3 Hours	\$0.63

<http://nextgenseek.com/2012/08/comparing-price-and-tech-specs-of-illumina-miseq-ion-torrent-pgm-454-gs-junior-and-pacbio-rs/>

Comparative of NGS platforms

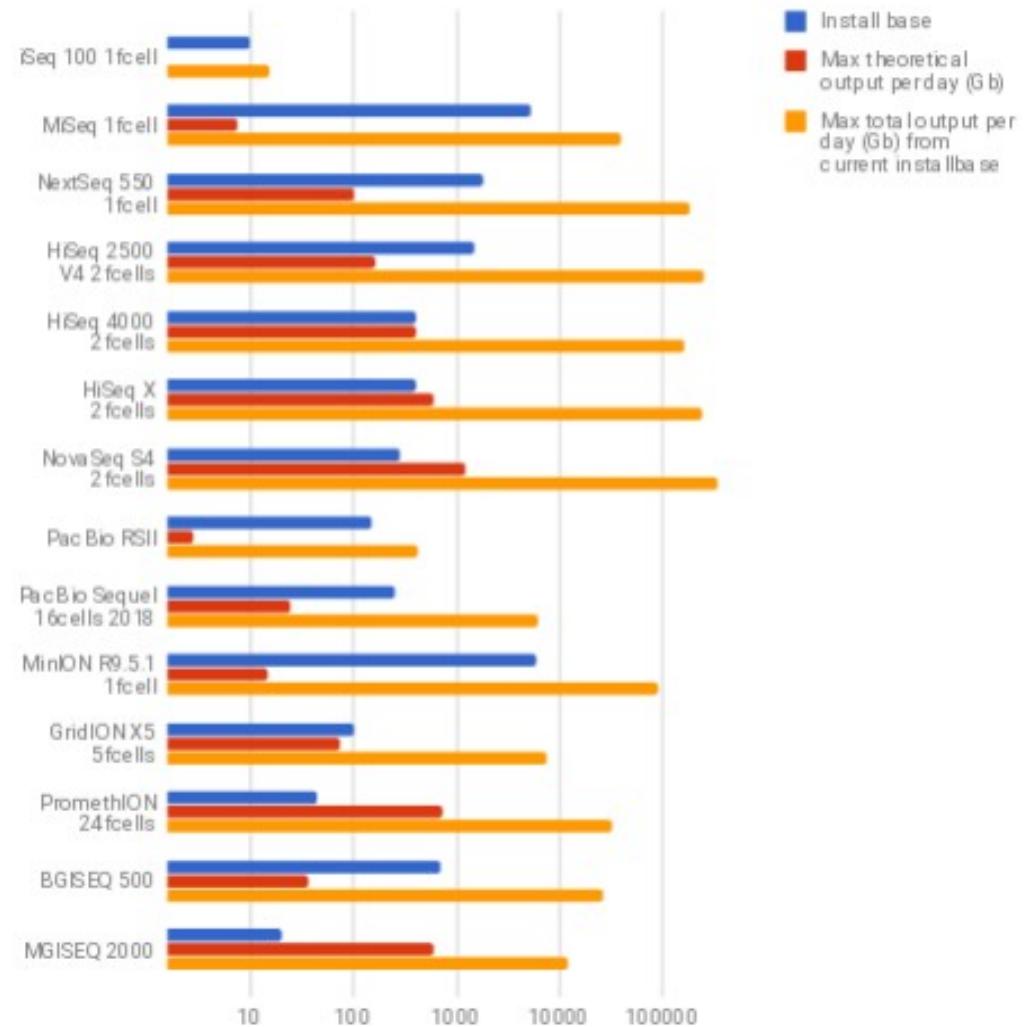
	Machine Cost	Throughput Per Run	Sequencing Cost per Gb	Sequencing Run Time	Observed Error rate
Illumina MiSeq	\$128,000	1.5-2 Gb	\$502	27 Hours	0.8%
PacBio RS	\$695,000	100 Mb	\$2000	2 Hours	12.86%
Ion Torrent -318 Chip	\$80,000	1Gb	\$1000	2 Hours	1.71%
Illumina GAIIx	\$256,000	30 Gb	\$148	10 Days	0.76%
Illumina HiSeq 2000	\$654,000	600 Gb	\$41	11 Days	0.26%

<http://nextgenseek.com/2012/08/comparing-price-and-tech-specs-of-illumina-miseq-ion-torrent-pgm-454-gs-junior-and-pacbio-rs/>

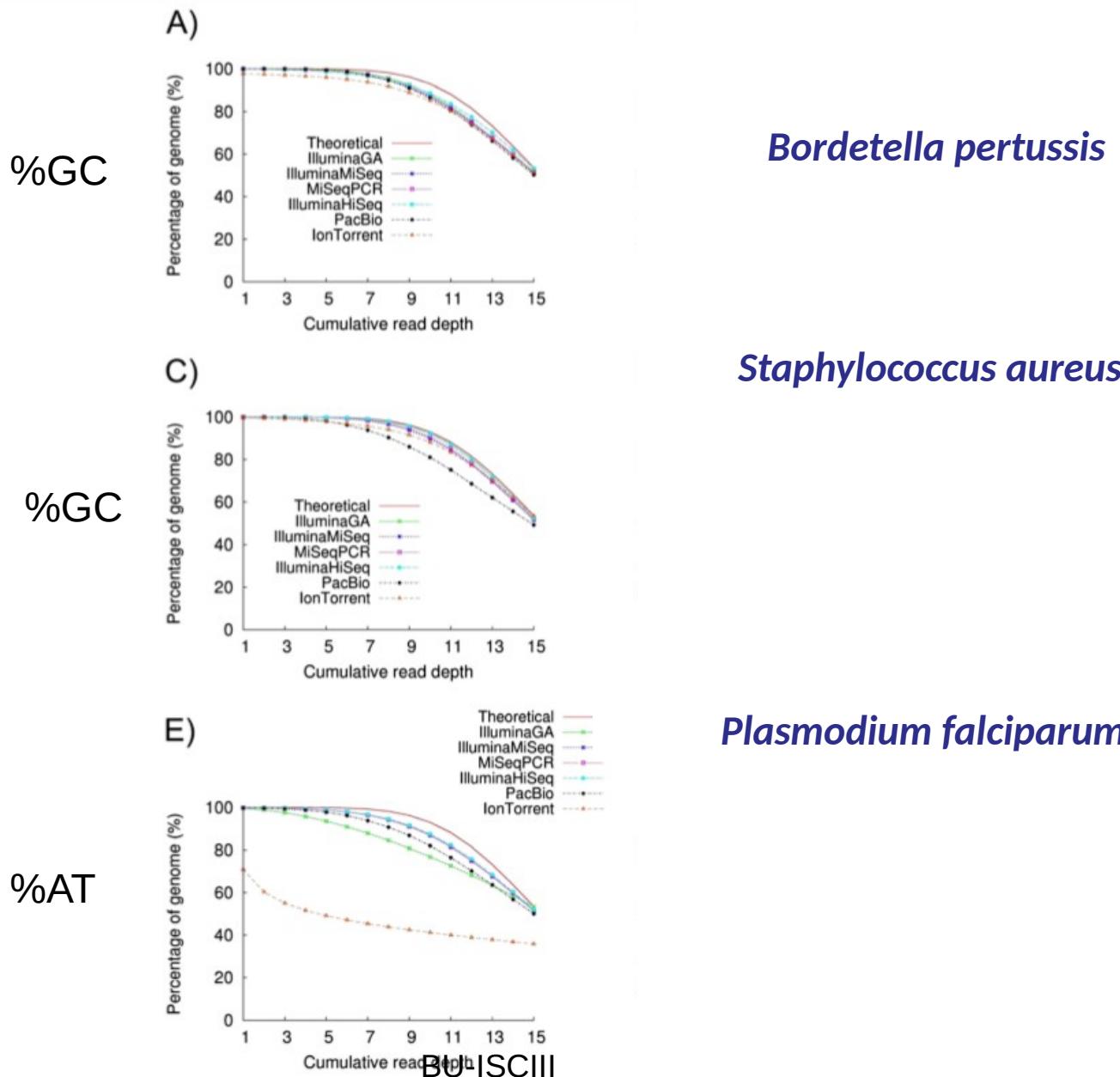
Comparative of NGS platforms

[https://docs.google.com/spreadsheets/d/1GMMfhyLK0-q8Xklo3YxIWaZA5vVMuhU1kg41g4xLkXc/edit?
hl=en_GB&hl=en_GB#gid=1569422585](https://docs.google.com/spreadsheets/d/1GMMfhyLK0-q8Xklo3YxIWaZA5vVMuhU1kg41g4xLkXc/edit?hl=en_GB&hl=en_GB#gid=1569422585)

Next Generation Sequencing Install Base
<http://tinyurl.com/ngsspecs>



Uniformidad de cobertura a lo largo del genoma



Variación en la cobertura dependiendo del Sistema de Secuenciación

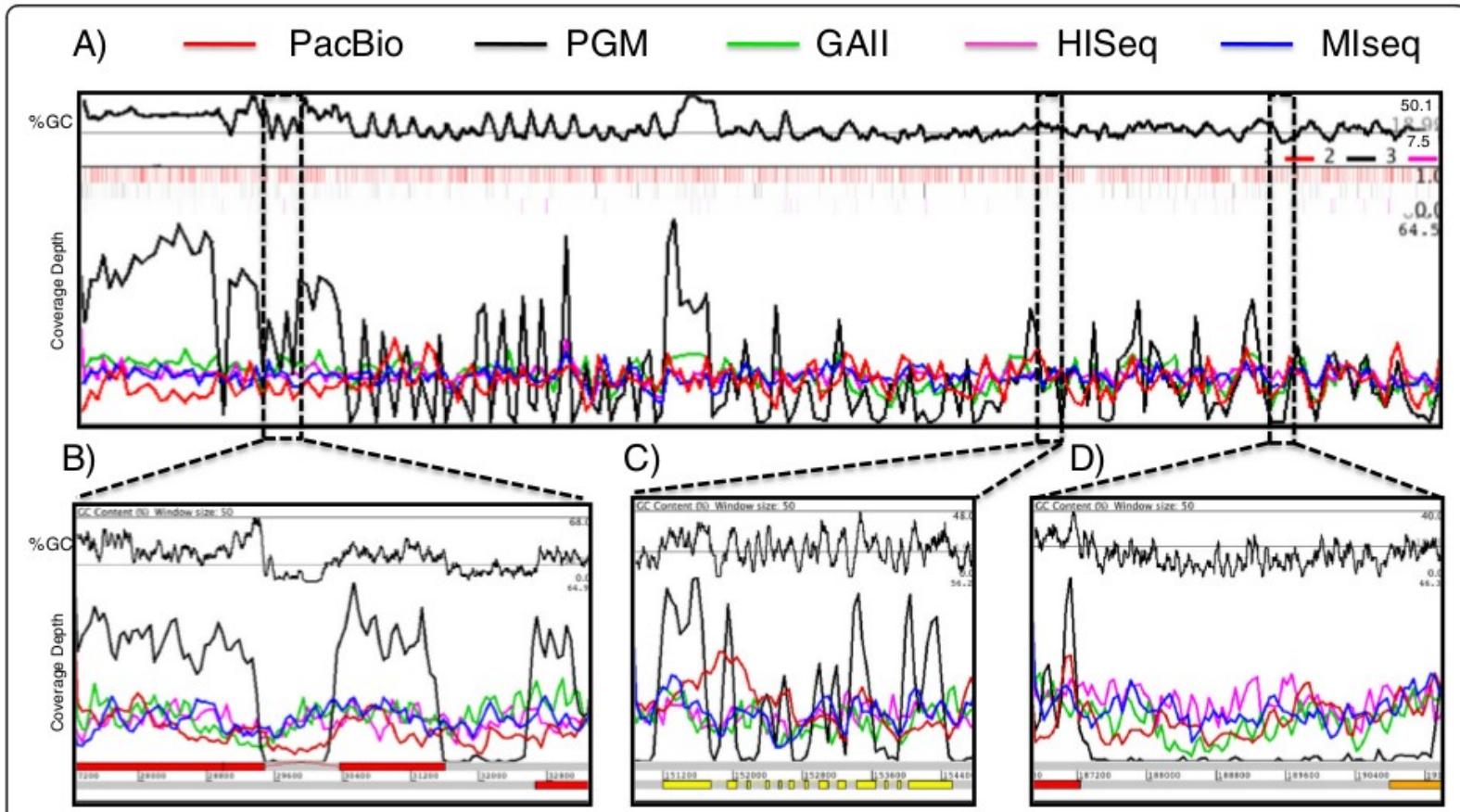
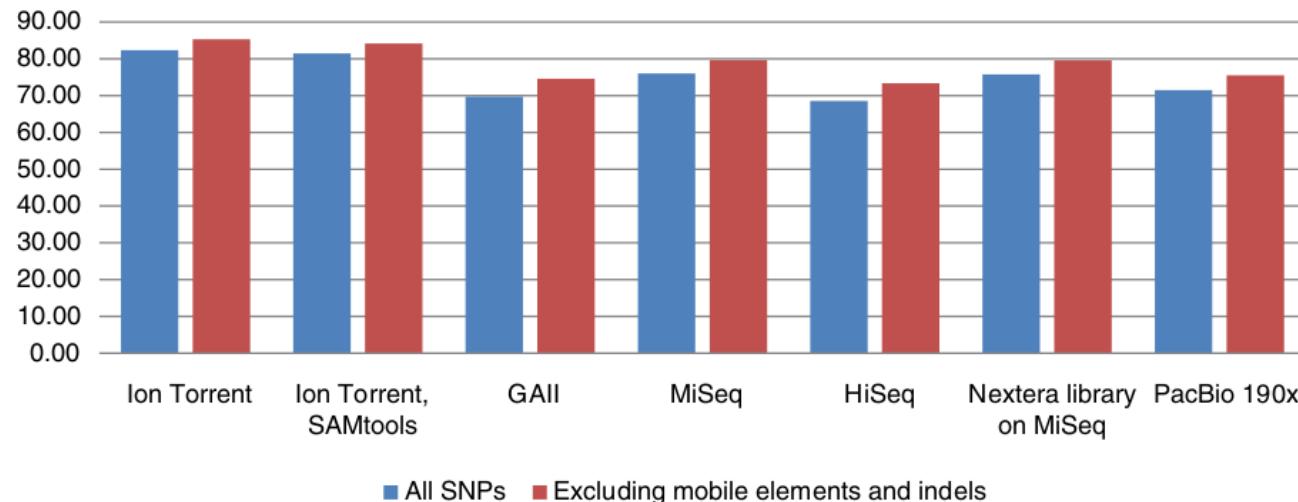


Figure 2 Artemis genome browser [8] screenshots illustrating the variation in sequence coverage of a selected region of *P. falciparum* chromosome 11, with 15x depth of randomly normalized sequence from the platforms tested. In each window, the top graph shows the percentage GC content at each position, with the numbers on the right denoting the minimum, average and maximum values. The middle graph in each window is a coverage plot for the dataset from each instrument; the colour code is shown above graph a). Each of the middle graphs shows the depth of reads mapped at each position, and below that in B-D are the coordinates of the selected region in the genome with gene models on the (+) strand above and (-) strand below. **A)** View of the first 200 kb of chromosome 11. Graphs are smoothed with window size of 1000. A heatmap of the errors, normalized by the amount of mapping reads is included just below the GC content graph (PacBio top line, PGM middle and MiSeq bottom). **B)** Coverage over region of extreme GC content, ranging from 70% to 0%. **C)** Coverage over the gene PF3D7_1103500. **D)** Example of intergenic region between genes PF3D7_1104200 and PF3D7_1104300. The window size of B, C and D is 50 bp.

A)

Percentage of correctly called true SNPs



B)

Number of incorrect SNP calls

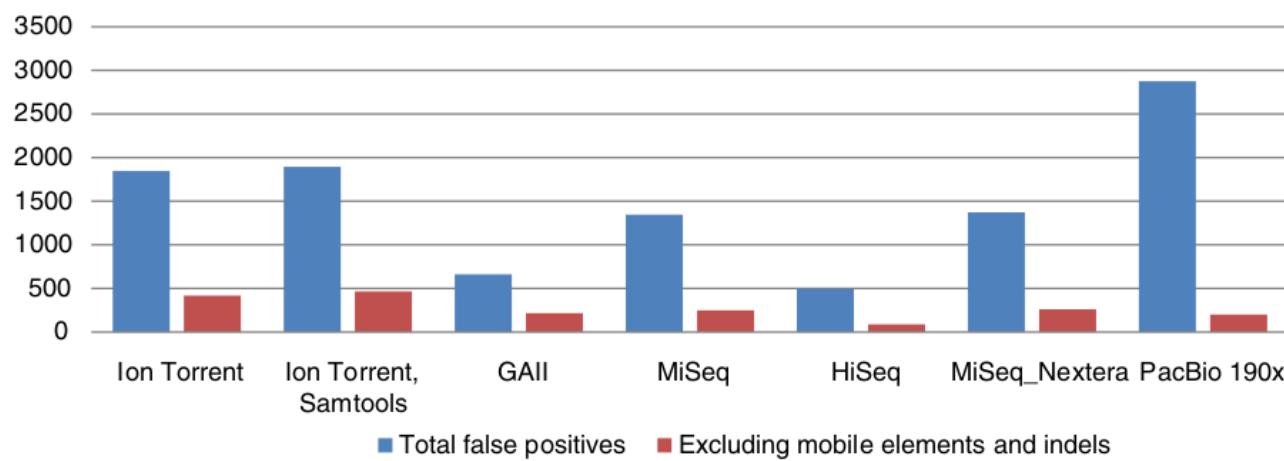


Figure 5 Accuracy of SNP detection from the *S. aureus* datasets generated from each platform, compared against the reference genome of its close relative *S. aureus* USA300_FPR3757. Both the Torrent server variant calling pipeline and SAMtools were used for Ion Torrent data; SAMtools was used for Illumina data and SMRT portal pipeline for PacBio data. **A)** The percentage of SNPs detected using each platform overall (blue bar), and outside of repeats, indels and mobile genetic elements (red bar). **B)** The number of incorrect SNP calls for each platform overall (blue bar), and outside of repeats, indels and mobile genetic elements (red bar).

Errores específicos de plataforma

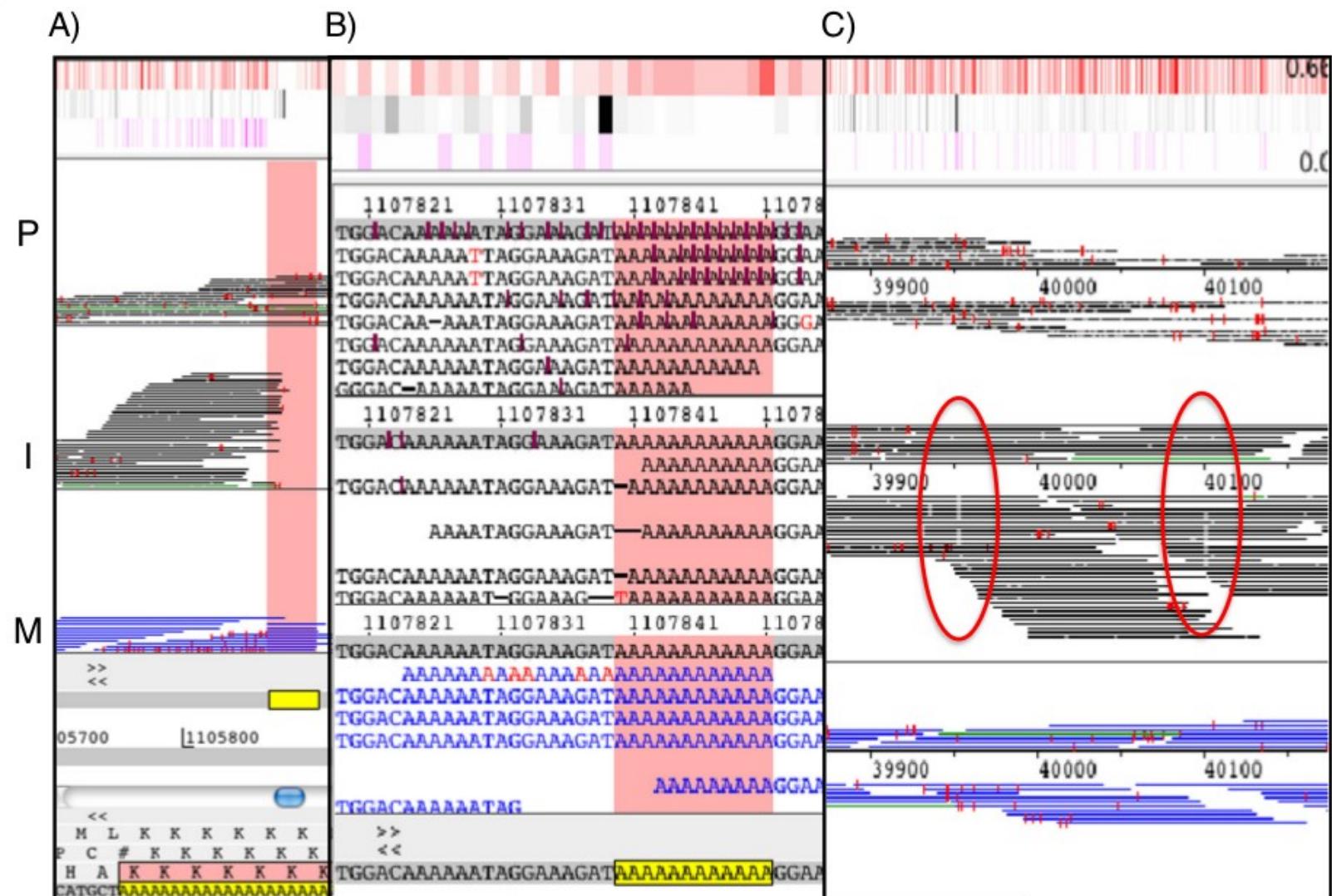


Figure 4 Illustration of platform-specific errors. The panels show Artemis BAM views with reads (horizontal bars) mapping to defined regions of chromosome 11 of *P. falciparum* from PacBio (P; top), Ion Torrent (I; middle) and MiSeq (M; bottom). Red vertical dashes are 1 base differences to the reference and white points are indels. **A**) Illustration of errors in Illumina data after a long homopolymer tract. Ion torrent data has a drop of coverage and multiple indels are visible in PacBio data. **B**) Example of errors associated with short homopolymer tracts. Multiple insertions are visible in the PacBio Data, deletions are observed in the PGM data and the MiSeq sequences read generally correct through the homopolymer tract. **C**) Example of strand specific deletions (red circles) observed in Ion Torrent data.

Conclusiones

A tale of three next generation sequencing platforms: comparison of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers
Quail et al., BMC Genomics 2012, 13:341

- Ion Torrent no es recomendable para secuenciar genomas con bajo contenido en GC
- Pac Bio:
 - Requiere elevada cantidad de DNA (no amplificación)
 - No recomendable para aplicaciones de counting (RNAseq, Chipseq, exoma)
 - No valido para identificar SNPs
 - Necesario adaptar software para aprovechar la longitud de lectura en el assembly.
- Illumina, preparación de librería adecuada a amplicones para generar correctamente los clusters
- Aplicaciones dependientes de plataforma

Conclusiones

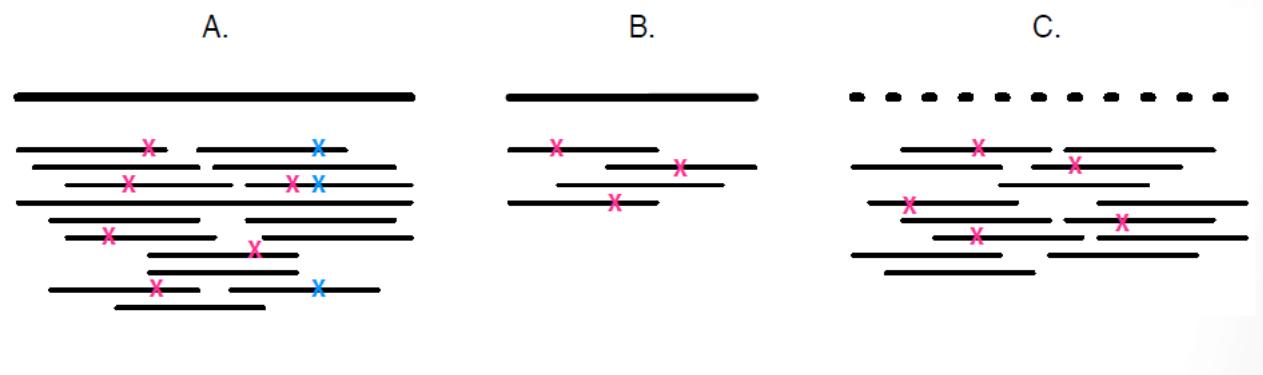
A tale of three next generation sequencing platforms: comparison of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers
Quail et al., BMC Genomics 2012, 13:341

- Conclusiones obtenidas con reactivos y plataformas a 2011
- Hay errores intrínsecos a la plataforma
- Otros errores se solucionaran con la actualización de las plataformas

Algunos conceptos en secuenciación

Básicamente tres problemas

Resecuenciación, Conteo y ensamblado



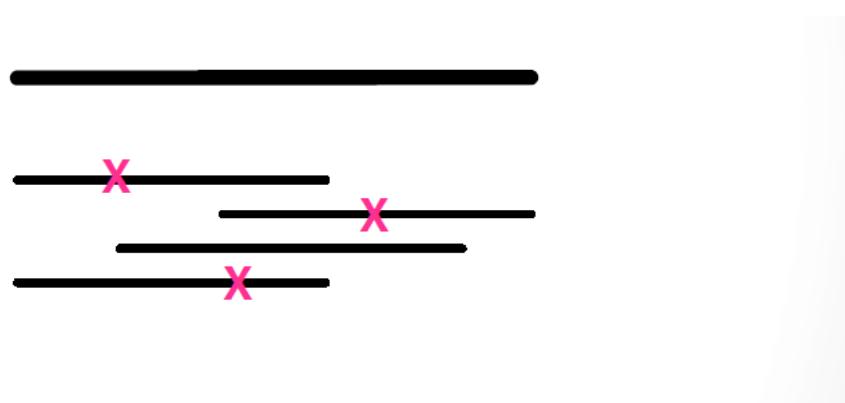
Resecuenciación

Conocemos el genoma, genoma de referencia, y queremos identificar variaciones (azul), en un background de errores (rosa)



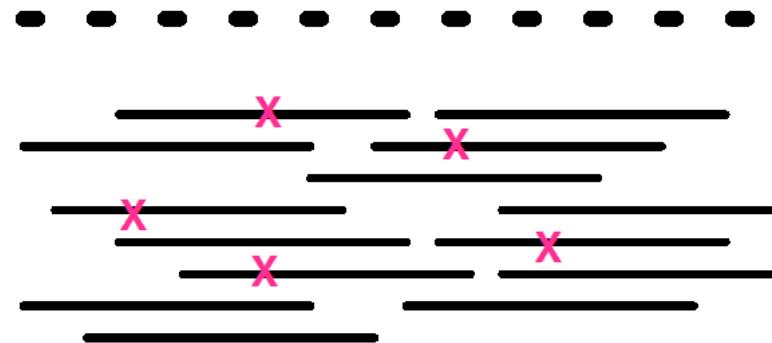
Conteo

Número de lecturas de un gen (amplicón) o mRNA (RNAseq). Equivalente a expresión en Microarrays.



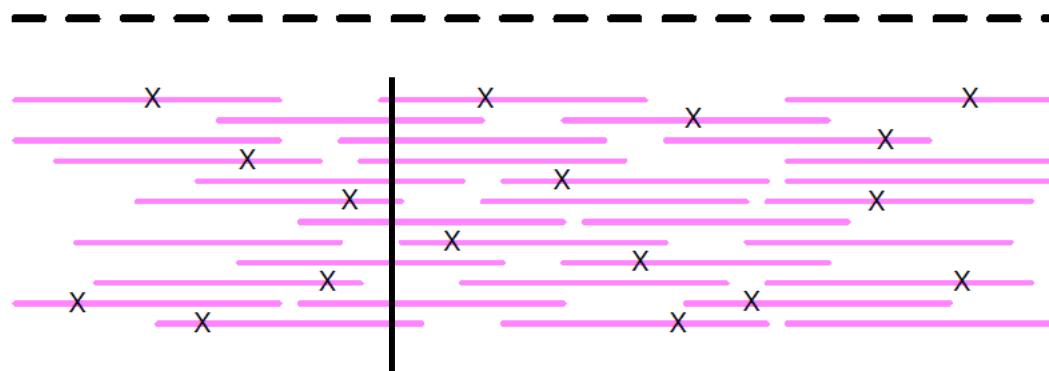
Ensamblado

No hay genoma de referencia y lo construimos de novo



Cobertura (depth of coverage)

Número medio de lecturas por base. i.e. 10x



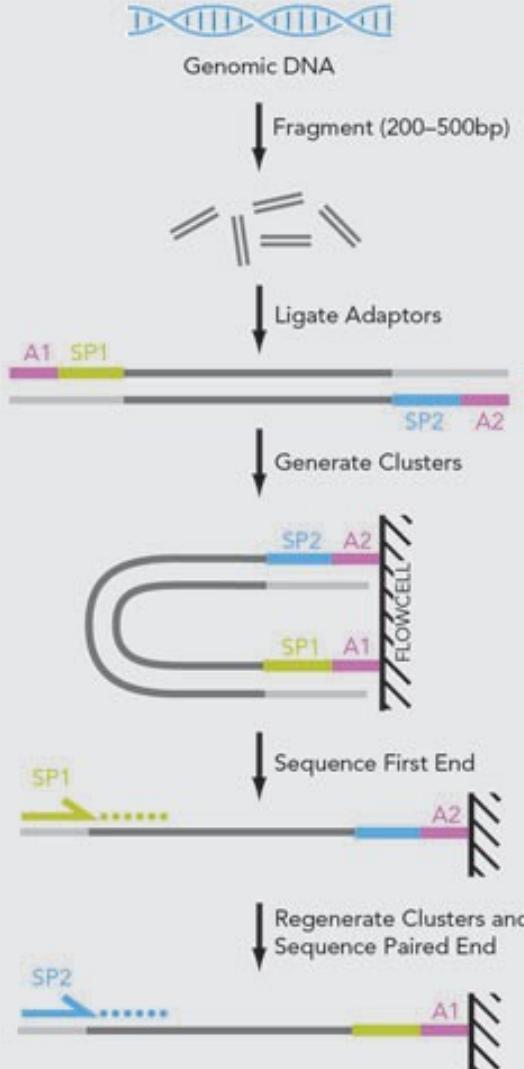
Sequencing coverage

as the average number of reads that align known reference bases

Number of reads x read length / target size

assuming that reads are randomly distributed across the genome.

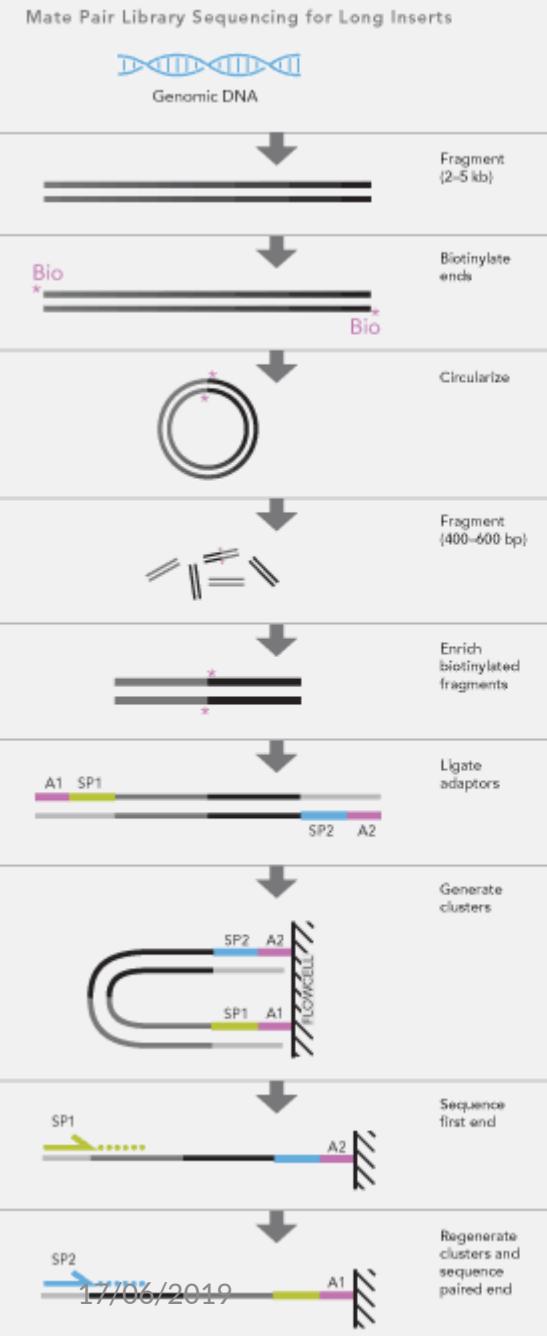
Que es Pair-end?



Secuenciación de un fragmento (bp)

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

Que es Mate-pair?



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.

Secuenciación de dos fragmentos separados kb.

Util:

**Secuenciación de un Genoma de novo
Finalizar un genoma
Detección de variantes estructurales**

Coverage and Read Depth Recommendations by Sequencing Application

Table 1: Coverage and Read Recommendations by Application

Category	Detection or Application	Recommended Coverage (x) or Reads (millions)	References
Whole genome sequencing	Homozygous SNVs	15x	Bentley et al., 2008
	Heterozygous SNVs	33x	Bentley et al., 2008
	INDELs	60x	Feng et al., 2014
	Genotype calls	35x	Ajay et al., 2011
Whole exome sequencing	CNV	1-8x	Xie et al., 2009; Medvedev et al., 2010
	Homozygous SNVs	100x (3x local depth)	Clark et al., 2011; Meynert et al., 2013
	Heterozygous SNVs	100x (13x local depth)	Clark et al., 2011; Meynert et al., 2013
Transcriptome Sequencing	INDELs	not recommended	Feng et al., 2014
	Differential expression profiling	10-25M	Liu Y. et al., 2014; ENCODE 2011 RNA-Seq
	Alternative splicing	50-100M	Liu Y. et al., 2013; ENCODE 2011 RNA-Seq
	Allele specific expression	50-100M	Liu Y. et al., 2013; ENCODE 2011 RNA-Seq
	De novo assembly	>100M	Liu Y. et al., 2013; ENCODE 2011 RNA-Seq

<https://genohub.com/recommended-sequencing-coverage-by-application/>

Coverage and Read Depth Recommendations by Sequencing Application

DNA Methylation Sequencing	CAP-Seq	>20M	Long, H.K. et al., 2013
	MeDIP-Seq	60M	Taiwo, O. et al., 2012
	RRBS (Reduced Representation Bisulfite Sequencing)	10X	ENCODE 2011 Genome
	Bisulfite-Seq	5-15X; 30X	Ziller, M.J et al., 2015; Epigenomics Road Map
RNA-Target-Based Sequencing	CLIP-Seq	10-40M	Cho J. et al., 2012; Eom T. et al., 2013; Sugimoto Y. et al., 2012
	iCLIP	5-15M	Sugimoto Y. et al., 2012; Rogelj B. et al., 2012
	PAR-CLIP	5-15M	Rogelj B. et al., 2012
	RIP-Seq	5-20M	Lu Z. et al., 2014
Small RNA (microRNA) Sequencing	Differential Expression	~1-2M	Metpally RPR et al., 2013; Campbell et al., 2015
	Discovery	~5-8M	Metpally RPR et al., 2013; Campbell et al., 2015

<https://genohub.com/recommended-sequencing-coverage-by-application/>

Resumen

Cobertura es importante para la llamada a variantes,
RNAseq. **Plataforma con mayor rendimiento Illumina**

La longitud de las lecturas es importante para el ensamblado
PacBio y Moleculo mayor longitud de lecturas (corrección de
errores con Illumina)

**Gracias por la
atención
Preguntas ???**



Isabel Cuesta

Unidad de Bioinformática – Unidades Científico Técnicas - ISCIII

isabel.cuesta@isciii.es