

# BIG4: Biosystematics, informatics and genomics of the big 4 insect groups- training tomorrow's researchers and entrepreneurs

Kick-Off Meeting  
14-18 September 2015  
Copenhagen, Denmark



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 642241

# Next Generation Sequencing (NGS) Revolution

14 of September 2015

**BIG4** Kick-off meeting

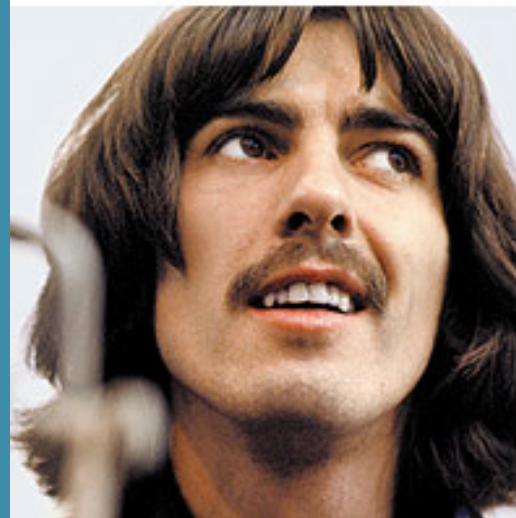
[www.era7bioinformatics.com](http://www.era7bioinformatics.com)



A bit of context ....

What is Era7 ?

(Some people does not understand why we are in a project related with “beatles”)



We think that bioinformatics is not an objective in itself but a tool that has sense to get

**New Biological Knowledge**



The most important thing in Era7 is our multidisciplinary, young and great team:

- Mathematics
- Informatics
- Bioinformatics
- Biochemistry
- Medicine

Covering all the spectrum and understanding different languages and partners what is needed to get

# New Biological Knowledge



Era7 started in 2005

Granada





# From 2010 in Madrid



From 2012 in

Cambridge MA  
USA

Era7 Bioinformatics Inc.



From 2014  
A wet lab in  
Granada



Our mission is to solve ALL the needs a customer could have about the NGS project.



## Innovative Business Model

- Research
- Cloud Computing
- Open Source
- Being focused in NGS

## Research:

- More than 70% of our staff doing R&D
- European, national and regional projects

**ohnosequences!**

era7 bioinformatics R&D group

## Research in Era7 Bioinformatics

### Research Lines:

- Algorithms for assembly
- Bacterial genome annotation
- Cloud Computing Architectures
- Graph Databases for Biological data
- Comparative genomics
- Cancer Genomics
- Genome Plasticity
- Big Data integration and visualisation
- Host Immune System and infection

### Software Projects

- BG7
- MG7
- CG7
- Bio4j
- Biographika
- NEXTMICRO
- Statika
- Nispero

# Cloud Computing

- Scalability (thousands of machines)
- Flexibility adapting computation to project
- Big Data Solutions
- Parallel Computation
- No need of investment:
- Lower costs
- IaaS





## Open Source

- Freedom for our customers
- The possibility of building the best pipelines



## Bio4J Graph oriented database integration (Uniprot GO Refseq)

Bio4j is a bioinformatics graph based DB including most data available in:

- UniProt (SwissProt + Trembl)
- UniRef (50,90,100)
- Gene Ontology (GO)
- RefSeq
- NCBI taxonomy
- Enzyme DB



Big Data:

2.000.000.000 relationships  
400.000.000 nodes  
1.000.000.000 properties

Bio4j project selected this year among 190 projects world wide of all areas for the

## *Google Summer of Code 2014*



# Biographika

An Open Source project about data visualisation

[www.biographika.com](http://www.biographika.com)

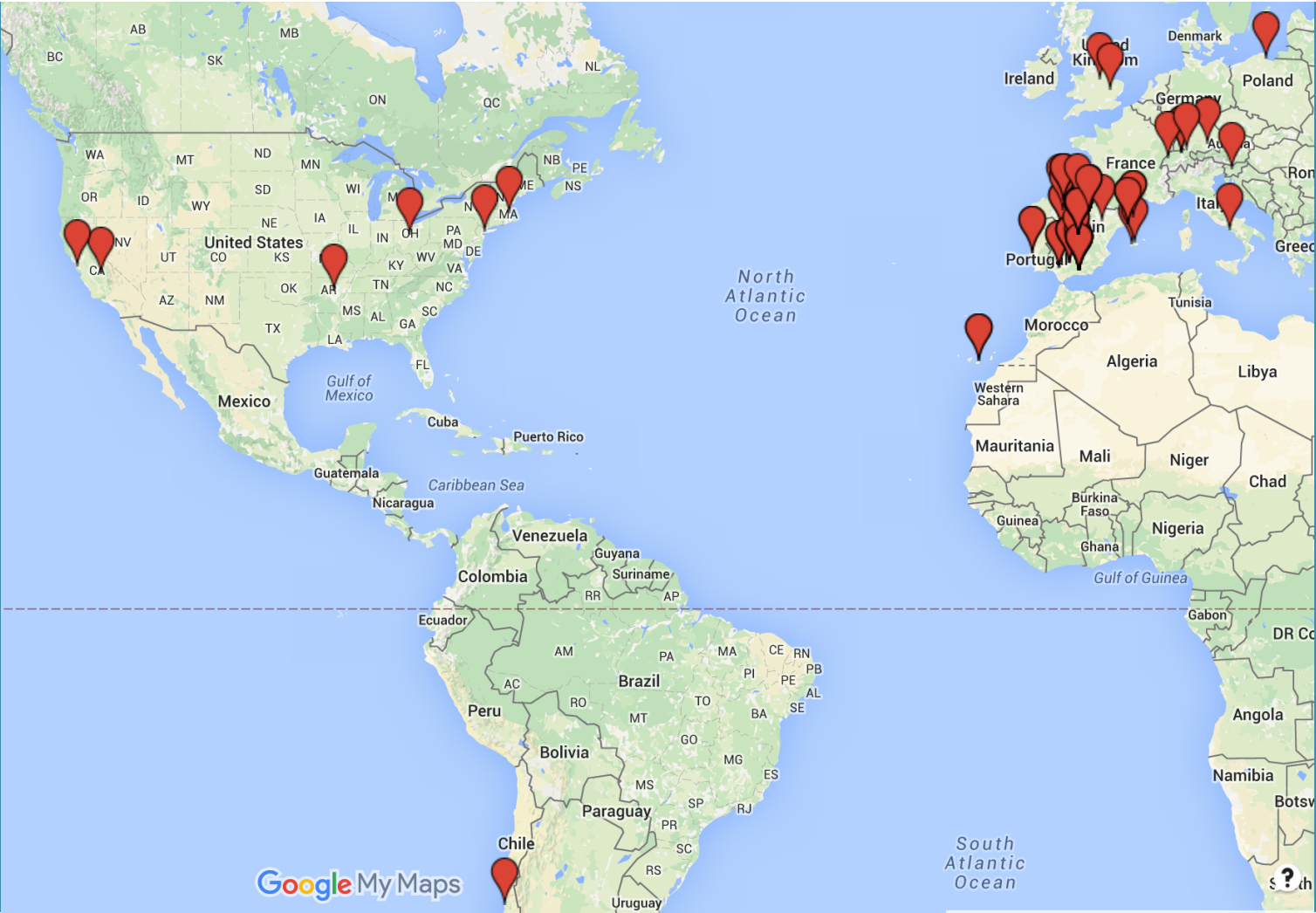
<https://vimeo.com/127019770>



## Our customers are researchers from:

- Research centres
- Universities
- Hospitals
- Biotech Companies
- Pharma Companies
- Agrifood Companies

Our Customers



NGS (Next Generation Sequencing)  
is a revolution. It is changing the way we do

- Genomics
- Metagenomics
- Transcriptomics
- Amplicon Sequencing
- .....

# Why a Revolution ?

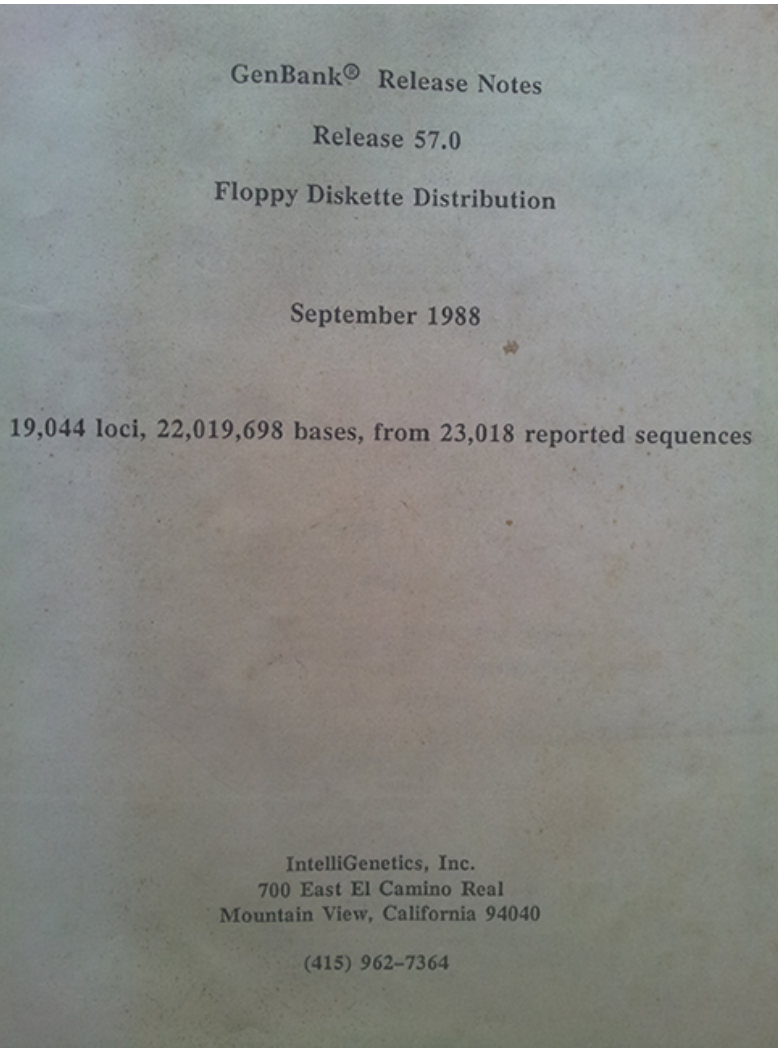
- Quicker and Easier
- Much more throughput
- Much less Price





1988

NCBI did not exist at this moment



1988

19,044	loci
23,018	sequences
22,019,698	bases

1988

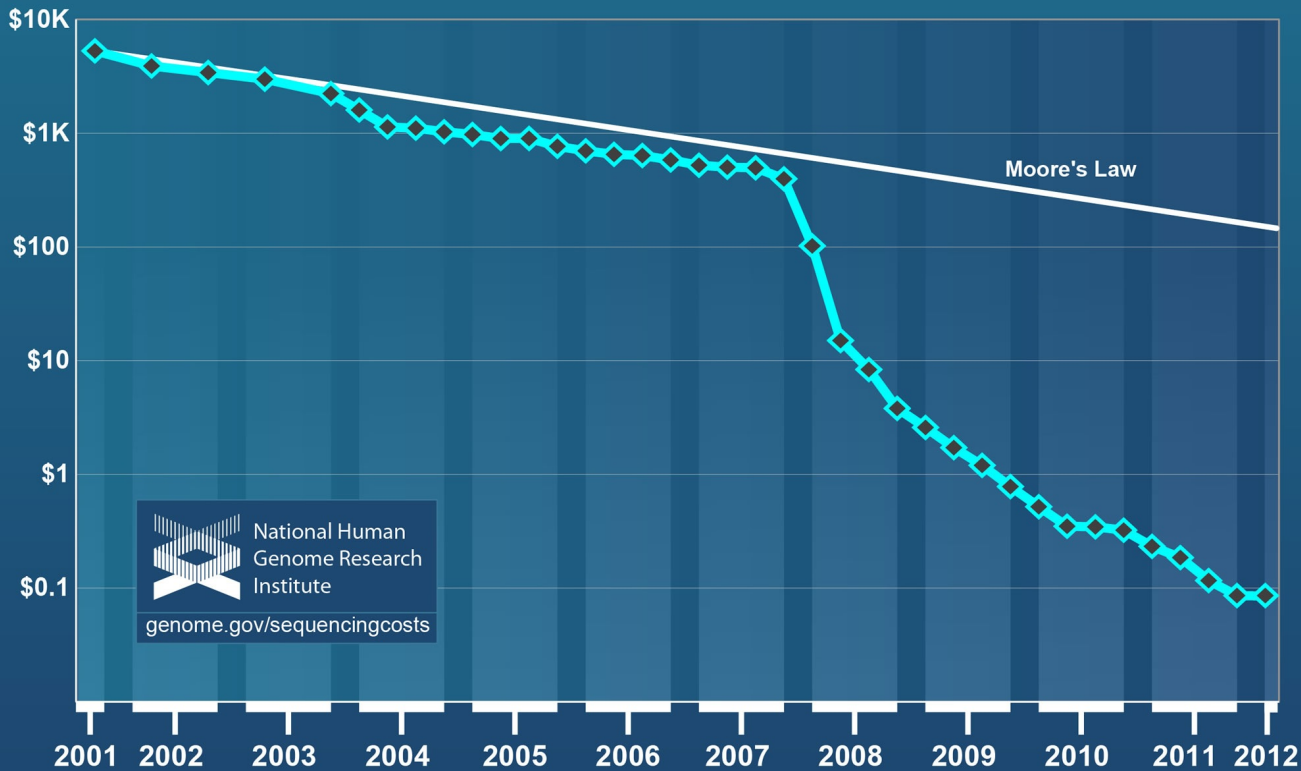


360 Kb

# World Wide Sequencing is today more than 35 Ptb

- Kilo base 1000 nucleotides
- Megabase 1000000 nucleotides
- Gigabase 1000000000 nucleotides
- Terabase 1000000000000 nucleotides
- Petabase 1000000000000000 nucleotides

### Cost per Raw Megabase of DNA Sequence



A factor of more than 1000X !!

Past 7 years

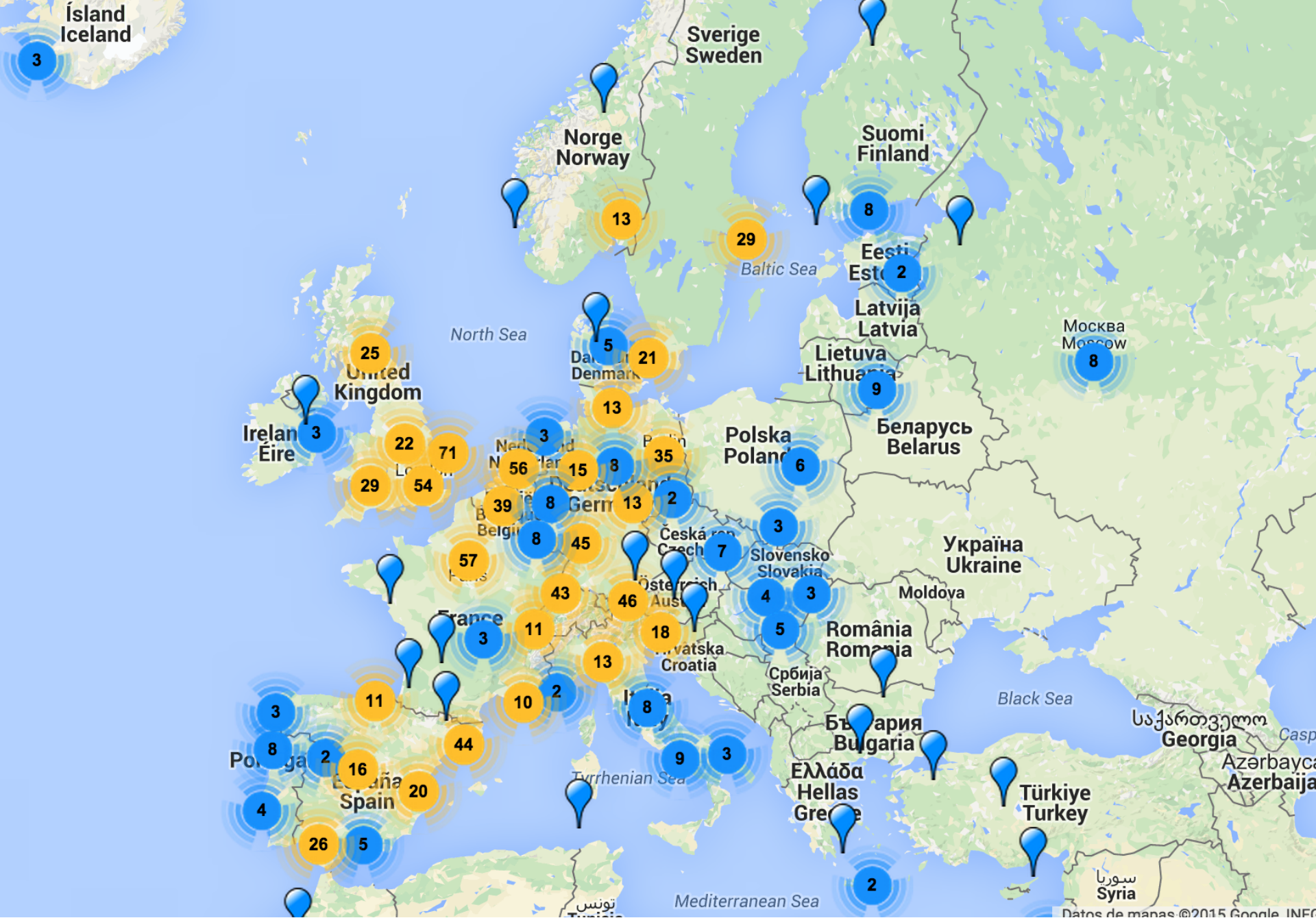


=



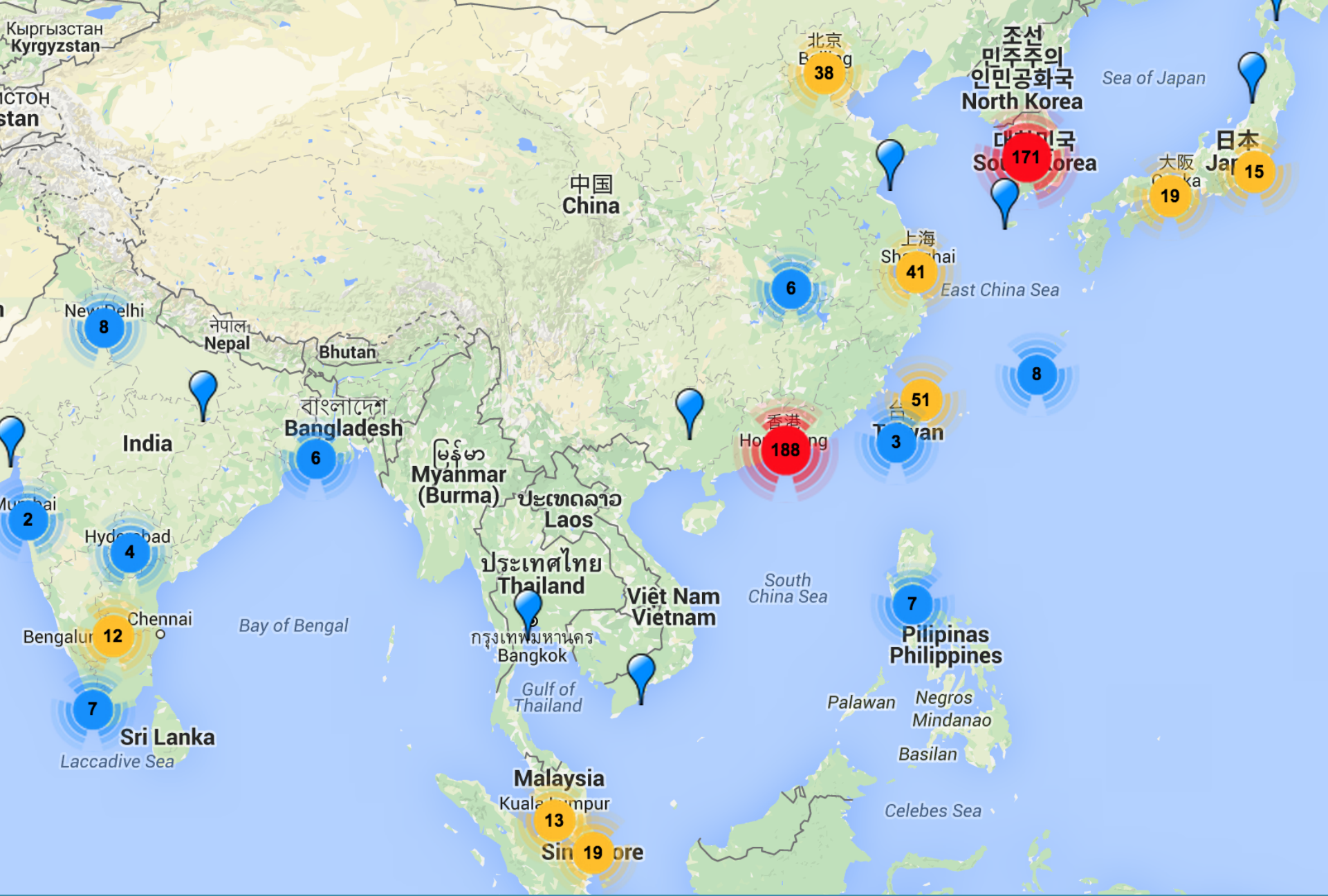
# NGS

- Total amount of sequence data produced **doubling approximately every seven months**
- The OmicsMaps reports that there are more than 2,500 high-throughput instruments located in nearly 1,000 sequencing centers in 55 countries in universities, hospitals, and other research laboratories









# NGS

~32,000 microbial genomes, ~5,000 plant and animal genomes, and ~250,000 individual human genomes that have been sequenced or are in progress thus far

- 100.000 genomes UK project
- 1.000.000 genomes Obama Precision Medicine Initiative
- 1.000.000 plant and animals genomes. (BGI)

(The proof of the anthropocentric vision)

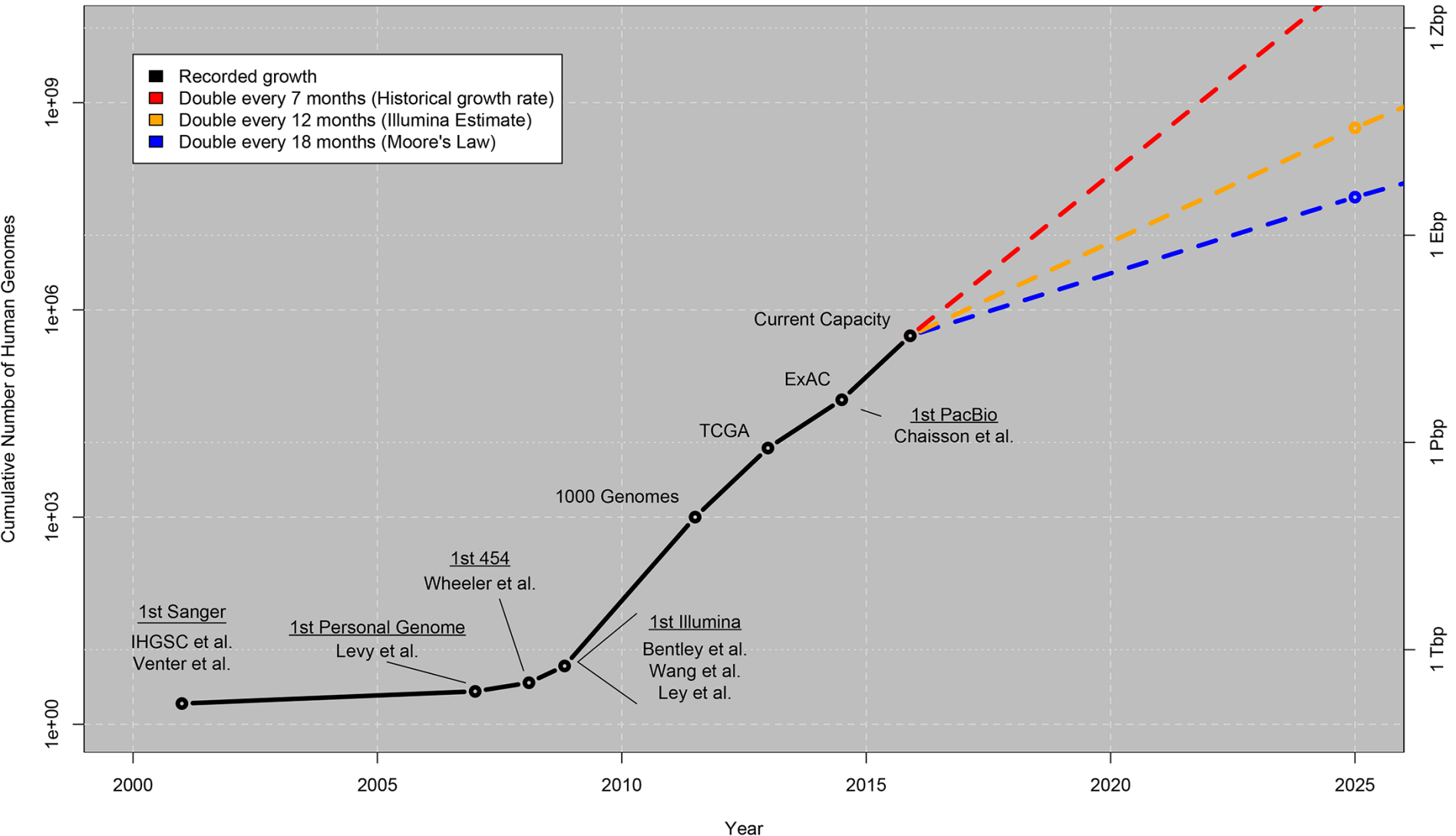
Plos Biology 10.1371/journal.pbio.1002195

# NGS

We therefore estimate between 100 million and as many as 2 billion human genomes could be sequenced by 2025 far exceeding the growth for the three other Big Data domains. Indeed, this number could grow even larger, especially since new *single-cell genome sequencing* technologies are starting to reveal previously unimagined levels of variation, especially in cancers, necessitating sequencing the genomes of thousands of separate cells in a single tumor

Plos Biology 10.1371/journal.pbio.1002195

**Growth of DNA Sequencing**



<b>Data Phase</b>	<b><u>Astronomy</u></b>	<b><u>Twitter</u></b>	<b><u>YouTube</u></b>	<b><u>Genomics</u></b>
<b>Acquisition</b>	25 zetta-bytes/year	0.5–15 billion tweets/year	500–900 million hours/year	1 zetta-bases/year
<b>Storage</b>	1 EB/year	1–17 PB/year	1–2 EB/year	2–40 EB/year
<b>Analysis</b>	In situ data reduction	Topic and sentiment mining	Limited requirements	Heterogeneous data and analysis
	Real-time processing	Metadata analysis		Variant calling, ~2 trillion central processing unit (CPU) hours
	Massive volumes			All-pairs genome alignments, ~10,000 trillion CPU hours
<b>Distribution</b>	Dedicated lines from antennae to server (600 TB/s)	Small units of distribution	Major component of modern user's bandwidth (10 MB/s)	Many small (10 MB/s) and fewer massive (10 TB/s) data movement

doi:10.1371/journal.pbio.1002195.t001

# NGS Technologies:

- illumina
- Ion Torrent
- PacBio
- Oxford Nanopore
- Complete Genomics
- (SOLiD)
- (454)

# illumina is the market leader

More than 90% of sequenced nucleotides  
has been sequenced with illumina  
technology



**Sequencing**



**MiSeq**

Focused Power



**MiSeqDx**

Focused Power



**MiSeq FGx**

Focused Power



**NextSeq 500/550**

Flexible Power



**HiSeq 2500**

Production Power



**HiSeq 3000**

Production Power



**HiSeq 4000**

Production Power

**Sequencing Library Prep**



**HiSeq X Five**

Population Power



**HiSeq X Ten**

Population Power



**NeoPrep**

Powerfully Simple

**Arrays**



**HiScan**

Flagship Scanner



**iScan**

Cutting-edge Scanner



**MiSeq**  
**Focused power.** Speed and simplicity for targeted and small genome sequencing.



**NextSeq 500**  
**Flexible power.** Speed and simplicity for everyday genomics.



**HiSeq 2500**  
**Production power.** Power and efficiency for large-scale genomics.



**HiSeq X\***  
**Population power.** \$1,000 human genome and extreme throughput for population-scale sequencing.

Key applications	Small genome, amplicon, and targeted gene panel sequencing.	Everyday genome, exome, transcriptome sequencing, and more.		Production-scale genome, exome, transcriptome sequencing, and more.		Population-scale human whole-genome sequencing.
Run mode	N/A	Mid-Output	High-Output	Rapid Run	High-Output	N/A
Flow cells processed per run	1	1	1	1 or 2	1 or 2	1 or 2
Output range	0.3-15 Gb	20-39 Gb	30-120 Gb	10-180 Gb	50-1000 Gb	1.6-1.8 Tb
Run time	5-65 hours	15-26 hours	12-30 hours	7-40 hours	< 1 day - 6 days	< 3 days
Reads per flow cell†	25 Million‡	130 Million	400 Million	300 Million	2 Billion	3 Billion
Maximum read length	2 × 300 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp	2 × 125 bp	2 × 150 bp

# Ion Torrent

It has some applications

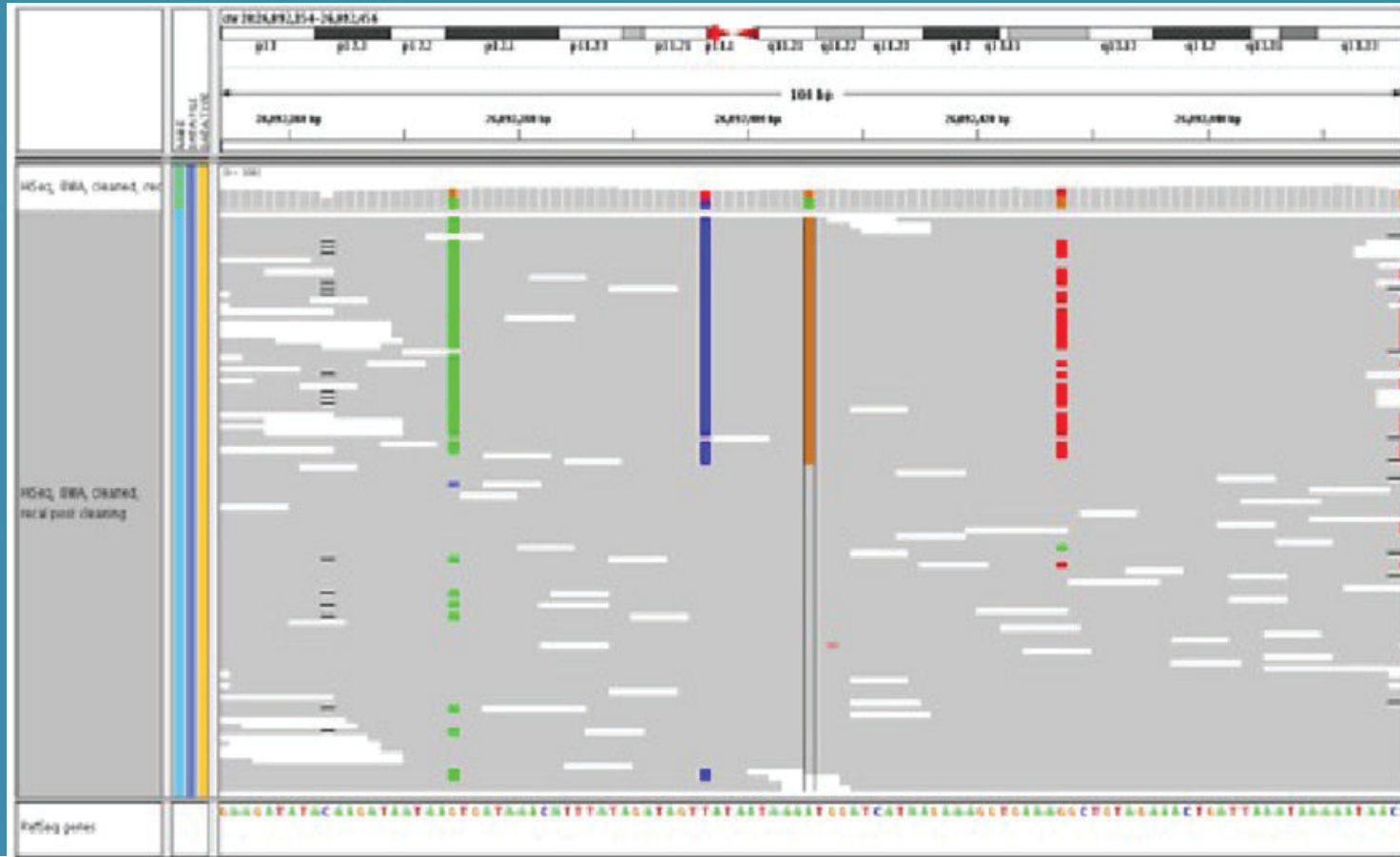
Ion PGM™ Chip	Run time		Output	
	200 bp read	400 bp read	200 bp read	400 bp read
Ion 314™ Chip v2	2.3 hr	3.7 hr	30–50 Mb	60–100 Mb
Ion 316™ Chip v2	3.0 hr	4.9 hr	300–500 Mb	600 Mb–1 Gb
Ion 318™ Chip v2	4.4 hr	7.3 hr	600 Mb–1 Gb	1.2–2 Gb
Ion Proton™ Chip	Run time		Output	
	200 bp read		200 bp read	
Ion PI™ Chip	2–4 hr		Up to 10 Gb	

# Pacific Biosciences

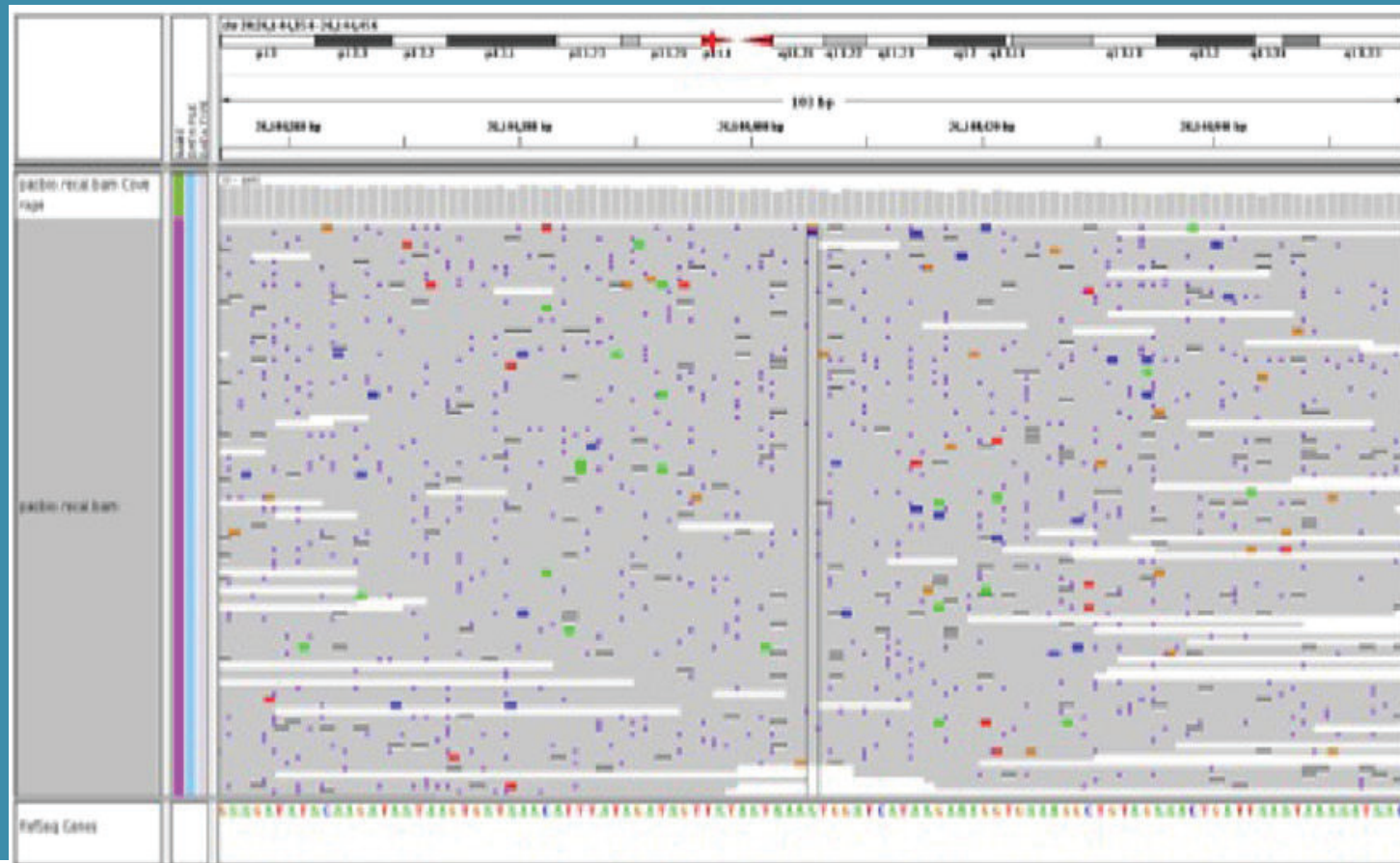
The advantage of long reads



# SYSTEMATIC ERROR in illumina



# RANDOM ERROR in PacBio





# Oxford Nanopore MinION

The advantage of not being really a machine

Oxford Nanopore  
Long reads,  
still a lot of  
errors.



# Oxford Nanopore Ebola Sequencing

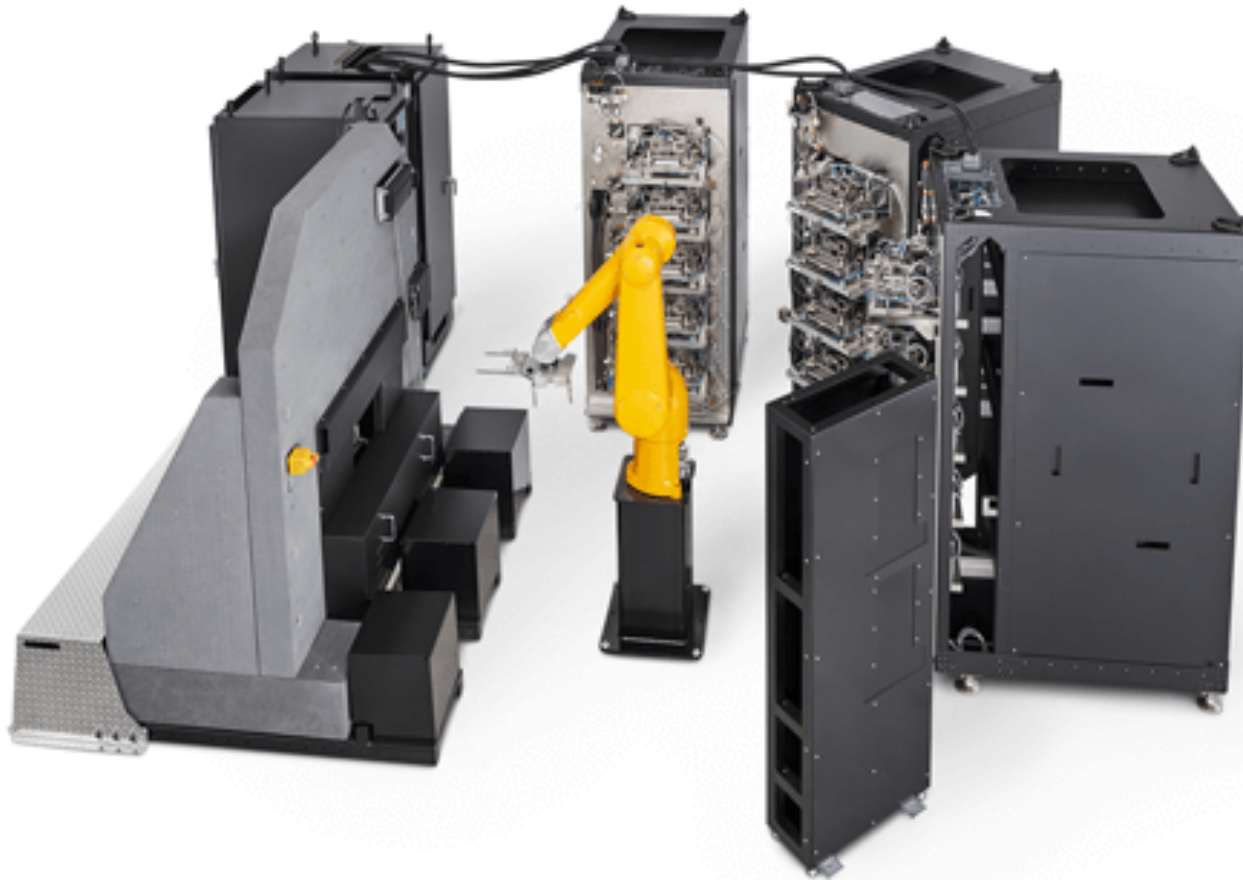
NATURE NEWS

05 May 2015



# Complete Genomics (BGI)

Now, you can buy machines....



Revolocity: From Blod to VCF Files BGI NEWS 05 June 2015

## PacBio :

- Really LONG READS (8-10 kb)
- Nucleotide Modifications (no PCR)
- Large amount of DNA

# PacBio for De novo Bacterial Genomics:

## The Best !!!

- “Finished” Genome
- Chromosome and plasmids assembled in one experiment !!

## PacBio :

- To improve Complex Genome Assembly
- Impressive information about Isoforms
- Large amplicons (haplotype info)



# With such amount of sequencing data

- Storage
- Analysis
- Interpretation

There is a growing need of

# Bioinformatics

## Some concepts:

- De novo Genome
- Resequencing
- Shotgun
- amplicon sequencing
- RNA-seq
- Exome
- Enrichment
- Targeted sequencing
- Metagenomics
- Metatranscriptomics
- Libraries
- Run, Lane, multiplexing...

# Some bioinformatics concepts

- Assembly
- Mapping
- Annotation
- .....

# Some kind of applications

# Bacterial Genomics:

- Design
- Sequencing
- Assembly
- Annotation
- Submitting
- Comparative Genomics (SNPs and Whole Genome)
- Interactive Data Visualisation
- Databases and Cloud Solutions

# Bacterial Genomics:

- illumina
- PacBio
- Strategies in a set of isolates

# Comparative genomics

- Genome alignment
- Ortholog tables
- Strain exclusive genes and genomic regions
- SNPs
- Innovative interactive visualisation



# Metagenomics:

- 16S
- Shotgun
- (Consortium)

# Metagenomics:

- 16S
  - Sequencing.
  - Analysis:
    - alfa – diversity
    - beta – diversity

# Bacterial Metagenomics:

- 16S

Primers

Target regions

Read length....

illumina, MiSeq (2x300) and now HiSeq (2x250) (454)

# Metagenomics:

- Shotgun
  - Sequencing
  - Analysis:
    - Functional Profile
    - Specific Functional Profiles:

# Metagenomics:

- Consortium,
  - A few species
  - Usually  $< 10$
  - Potential assembly
  - Useful:
    - Infections
    - Industrial Biotechnology

# Application of Metagenomics:

## HUMAN MICROBIOME AND HEALTH

### Microbiome influence in health:

- Microbiome of mother and child
- Microbiome of centenary individuals
- Microbiome of lean individuals

### Microbiome and diseases:

- Microbiome and Cancer
  - In cancer prevention
  - In cancer pathogeny
  - Metagenomics as biomarker
- Microbiome and Obesity
  - Gut Microbiome as biomarker of chronic inflammatory diseases

# Application of Metagenomics:

## Microbiome and diseases:

- Oral microbiome and cardiovascular diseases
- Nervous System and gut microbiome
- Aging and gut microbiome
- Microbiome and allergy
- Microbiome and Inflammatory Bowel Diseases
- Microbiome and autoimmune diseases

## Microbiome and infections

- Microbiome and antibiotic treatments
- Study of Commensal/Pathogen microbial organisms
- Biofilms
- Virome study
- Mobilome study
- Resistome study

# Application of Metagenomics:

## Microbiome and infections

De novo Diagnosis of nonculturable microbial pathogens

Microbial fungi

Viruses

Bacteria

## METAGENOMICS IN VETERINARY

Metagenomics in Veterinary has equivalent applications as in Human health

In many cases human diseases have a disease model in animals

Vector-borne pathogens studies



# Application of Metagenomics:

## ENVIRONMENTAL SCIENCES

Diversity in earth microenvironments to analyze:

- Climatic changes

- Pollution induced changes

- Seasonal changes

- Plague defence

- Water analysis

- Waste water treatment plants

# Application of Metagenomics:

## METAGENOMICS IN AGRIFOOD SECTOR

Conservation of food (contamination, toxins, ..)

Fermentations:

- Beer

- Wine

- Bread

- Cheeses and other fermented dairy products (consortiums involved in artisanal products)

Probiotics

Crops:

- Study of Rhizosphere communities:

  - Mycorrhizas to improve crops

  - pathogens forming consortiums or complex biofilms

- Bioremediation

- Microbial studies of compost

- Crop microbial plagues research

# Application of Metagenomics: BIOTECHNOLOGY AND DRUG DISCOVERY

Specific habitats Diversity:

- Searching for specific activities for bioenergy industry

- Digestions in industrial reactors

Research in Marine Samples

- Searching for new active compounds:

  - Cancer research: New anticancer drugs

  - New antibiotics

  - New drugs

- Marine Diversity:

  - Secondary metabolism cluster detection from samples of marine banks

Metagenomes of extreme microenvironments

- New enzymes (especially from extremophile bacteria) for advanced biotechnological uses

- Microorganisms with special resistance to salt, temperature, toxics, pollutants,..

- Astrobiology

# Bacterial RNA – seq

- Sequencing
- Analysis

## “Dual” RNA – seq:

- One experiment providing RNA – seq info:
  - Bacteria (Micro-organism)
  - Host
- Analysis:
  - Time Points
  - Deciphering networks

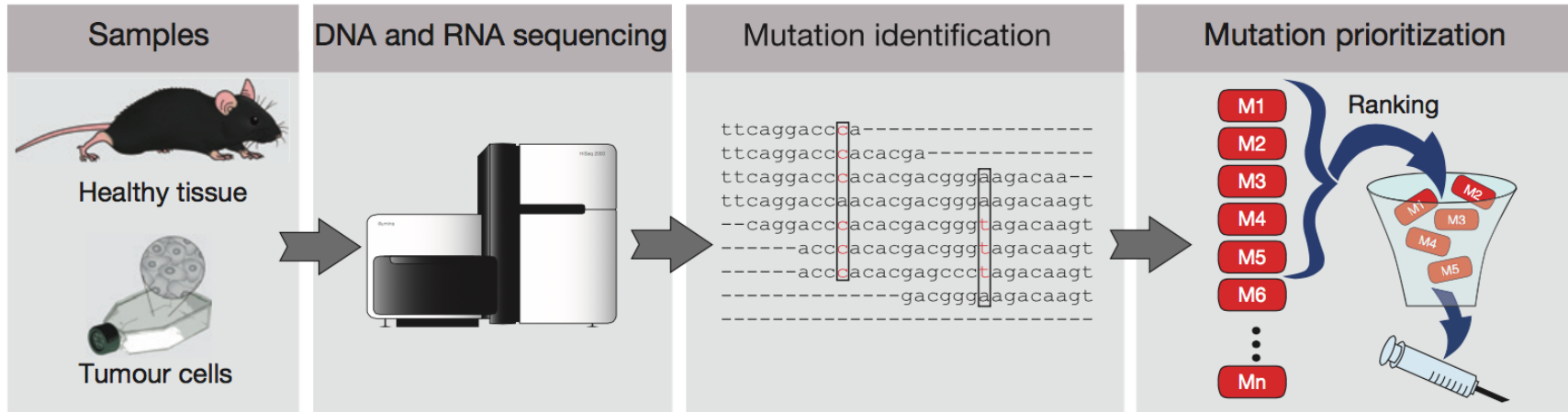
- Bacterial Genomics in medicine:
  - Infectious Diseases
  - Autoimmune Diseases
  - MICROBIOME

- Human (and mouse) Genome
- Human exome
- RNA-seq

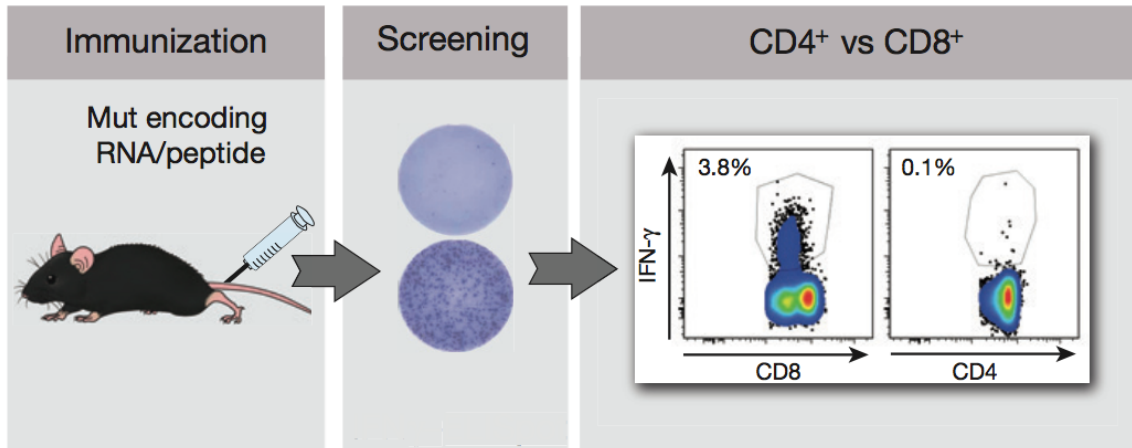
# Cancer Immunotherapy based on neoantigens



### Mutation discovery and prioritization



### Immunogenicity testing



# Cancer Immunotherapy based on neoantigens. NGS and Bioinformatics.

- RNA-seq analysis
- Bioinformatics selection
- Future new criteria (not only immunogenicity)
- The “real” Personalised Medicine

# RNA – seq

- De novo
- With model organism

# RNA – seq

- mRNA
- miRNA
- lncRNA

# Viral Genomics and Metagenomics

## Viral Amplicons (HIV etc.)

# Different kinds of amplicons

- HLA
- Different Large genes

Remember the possible advantages of PacBio Technology

# Immunogenomics

- T cell Repertoires
- B cell Repertoires
- HLA

# Projects

- de novo Genomics
- Genomics with reference (resequencing)
- Exome
- Metagenomics
- Comparative Genomics
- Transcriptomics (mRNA, miRNA, lncRNA)
- Dual RNA-seq (Pathogen and Host)
- SNPs (and whole genome comparison)
- ChiP-Seq
- Cancer Genomics
- Big Data analysis :  
**public data sets**
- Evolutionary studies
- Antibiotic resistance
- Bacteriocins
- Metabolic Engineering
- Biofuels
- Agrifood



## Next future in Hospitals:

- Microbiology:
  - Metagenomics and diagnostics
    - 16S for bacteria
    - Shotgun agnostic
  - Genomics (epidemiology, low cost soon)
- Oncology
- Immunology
- Many: Microbiome(s)

RNA-seq blood analytics?

# An example of use of NGS

Big4 project

Thank you for your  
attention !

Eduardo Pareja

[epareja@era7.com](mailto:epareja@era7.com)

[info@era7.com](mailto:info@era7.com)

[www.era7bioinformatics.com](http://www.era7bioinformatics.com)