# Microbial Ocean*omics* using High-Throughput DNA Sequencing

Ramiro Logares

Institute of Marine Sciences, CSIC, Barcelona

9th RES Users'Conference – 23 September 2015
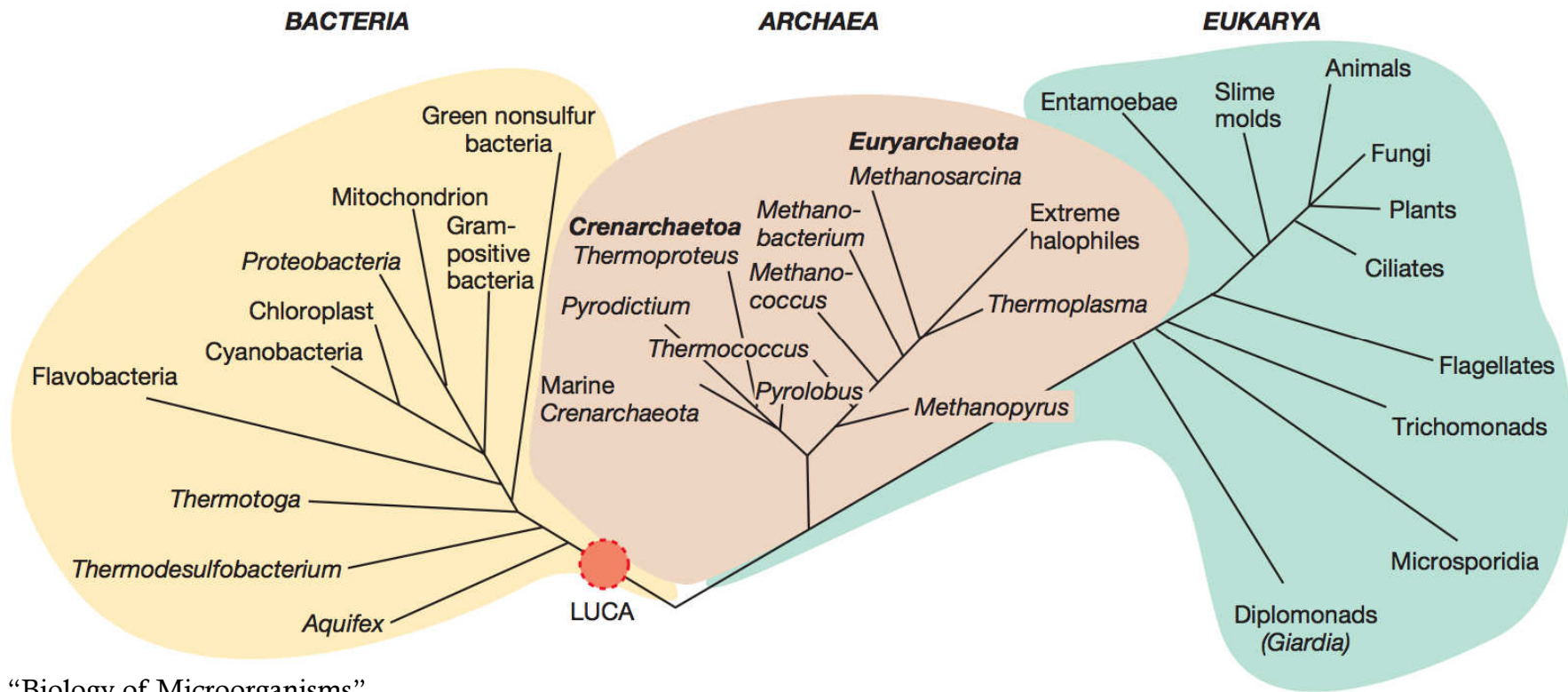
# Importance of microbes in the sunlit ocean

- Phytoplankton: 50% primary production of the Earth (Field et. Al 1998)

- Microplankton crucial for the marine food chain

- Biogeochemical cycling

- Large phylogenetic and metabolic diversity

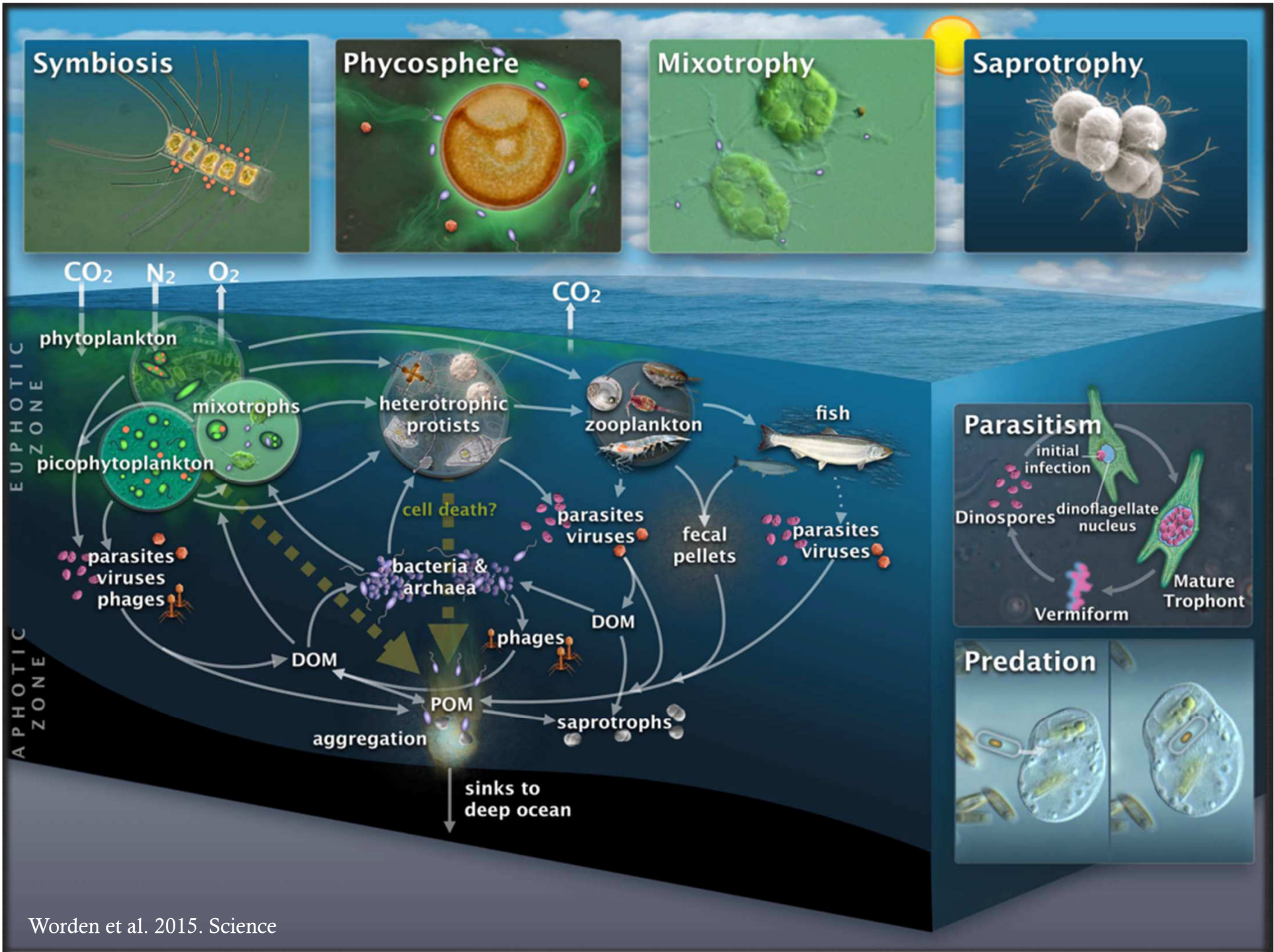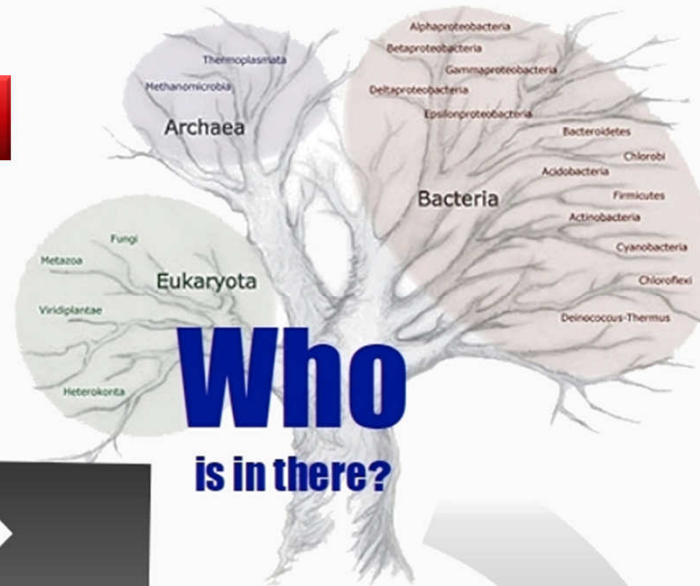*Chlorophyll concentration by SeaWiFS*

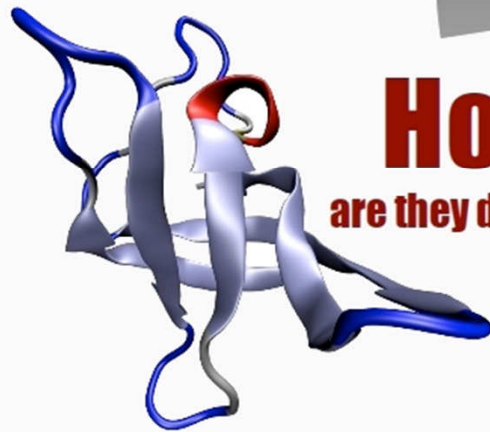*September 1997 – August 2000*

# Microbial phylogenetic diversity
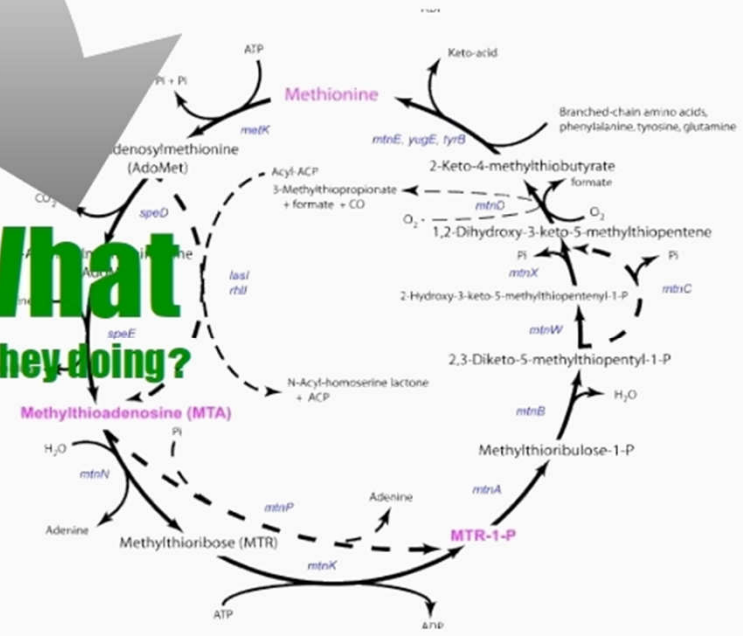


Brock "Biology of Microorganisms"

Protists

Prokaryotes

Viruses

Multicellular

| 0.1 | 1 | 10 | 100 | 1000 | 10,000 | 100,000 | $10^6$ | $10^7$ | $10^8$ |

Body Size [µm]

Source: C. de Vargas

Worden et al. 2015. Science

**Who** is in there?

**What** are they doing?

**How** are they doing it?

# DNA sequencing: an accelerating revolution

| Year | Landmark |
|------|----------|
| 1953 | Discovery of the double helix |
| 1977 | First DNA genome (bacteriophage) |
| 1977 | F. Sanger publishes "chain-terminator" method for DNA sequencing |
| 1987 | First commercial sequencing machine (ABI 370) |
| 1995 | First genome of a free-living organism (bacteria). WGSS initial use |
| 1996 | Nygren & Ronaghi publish "pyrosequencing" |
| 2001 | First draft human genome (3 billion US$) |
| 2004 | 454 pyrosequencing commercialized |
| 2009 | Illumina 50 K US$ per human genome |
| 2010 | Single molecule real time sequencing (SMRT) commercialized |
| 2011 | Human genome for 8000 US$. About 30 K human genomes sequenced |
| 2015 | 1,000 US$ - Human genome (Illumina X10; 18,000 per year-machine) |

# DNA sequencing revolution

| Year | Landmark | |
|------|----------|---|
| 1953 | Discovery of the double helix | |
| | ...cteriophage | |
| | ...ain-terminator" method for DNA sequencing | |
| | ...encing machine (...370) | |
| | ...living organi...). WGSS initial use | |
| | ...olish "pyro...g" | |
| | ...ne (3...) | |
| | ...m | |
| | ...h...ome | |
| | ...ne sequencing (SMRT) commercialized | |
| | ...00 US$. About 30 K human genomes sequenced | |

**62 years**

THE $1,000 GENOME

THE REVOLUTION IN DNA SEQUENCING AND THE NEW ERA OF PERSONALIZED MEDICINE

KEVIN DAVIES

# Sequencing platforms evolution

**1st Generation**

**2nd Generation**

**3rd Generation**

Still used for smaller projects or when high quality is needed

Widely used in most sequencing projects

Not widely used yet, some devices still not in the market

Quality reference

Oxford NANOPORE Technologies

SMRT seq

GridION

MinION

Signal: electricity

Next Generation Genomics: World Map of High-throughput Sequencers

http://omicsmaps.com/hts/centres/imppc/
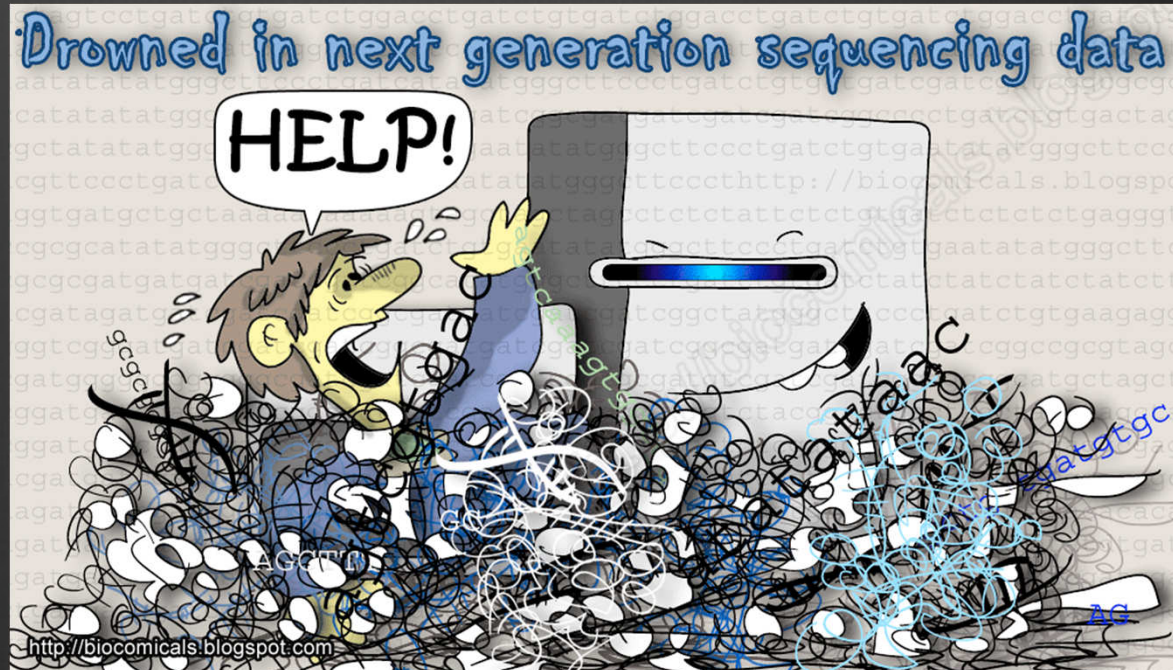
Amount of HTS DNA data produced now in the world:

- ≈ 7,389 functional HTS machines

- ≈35 x 10$^{15}$ bases / year  == 35 PETAbases

- 250,000 human genomes per year

Computing power

# Moore's Law

# Kryder's Law



Computer power tend to double every two years

Storage capacity doubles annually

## DNA Sequencing Is Now Improving Faster Than Moore's Law!

Adrienne Burke, Contributor

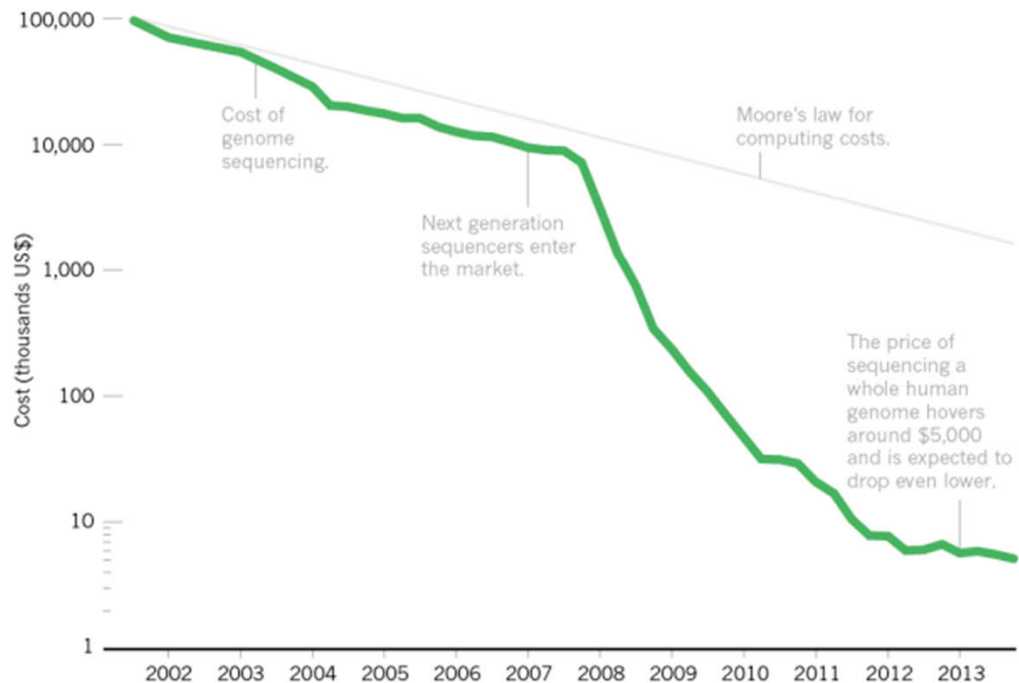+ Comment now

*A "worldwide genomics revolution" is upon us.*

The genomics industry marked a new milestone on Tuesday. As Forbes' Matthew Herper reported in three separate posts and nearly 100 related Tweets, the two leading manufacturers of DNA sequencing instruments announced almost simultaneously at an investors' conference that they would introduce new machines this year capable of sequencing an entire human genome in a single day. Life Technologies said its forthcoming Ion Proton machine, which processes DNA on a semiconductor chip, will do it for a cost of $1,000 per genome.

These advances are not just big news for biotech and medicine, but exciting for all Techonomists. They're proof that the pace of advances in genome sequencing technology has exceeded Moore's Law. The speed of genome sequencing has far better than doubled every two years since 2003, when the

Image by World Economic Forum via Flickr

Forbes Magazine, 2012

## Falling fast

In the first few years after the end of the Human Genome Project, the cost of genome sequencing roughly followed Moore's law, which predicts exponential declines in computing costs. After 2007, sequencing costs dropped precipitously.

Cost of genome sequencing.

Moore's law for computing costs.

Next generation sequencers enter the market.

The price of sequencing a whole human genome hovers around $5,000 and is expected to drop even lower.
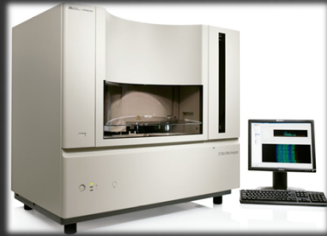
## Worldwide sequencing capacity is growing at about 2-3 times per year

- Only one HiSeq2500 produces about 3 TeraBytes of data per month

- Data processing costs should be considered
  - Electricity
  - Costs for data-admin, and reparation

- Amortization (value decrease) of equipment (3yrs CSIC)

- Data storage cost is not trivial

- What to do with used data? And backups? (maybe cheaper re-sequencing than storing?)

# Data processing and computation

# Minimum needed computer power

| 1st Generation | 2nd Generation | | 3rd Generation |
|---|---|---|---|
| SANGER | 454 Roche | Illumina HiSeq | PacBio |
| Cores= 1-2 | 16-32 | >64 (128) | >16 (32) |
| Mem= 1-4 GB | 32-64 | >64 (128) | >32 (64) |
| Disk= 0.2Tb | >2 Tb | > 10 Tb | >1Tb |

# MareNostrum (Barcelona)

```
Welcome to MareNostrum III

        .  .  -
     .´  ;  . `.
    ; ; | | BSC |
     `. ; . `.´
        `.`..

- All home directories are in GPFS and quotas are enabled
- Applications are located at /apps
- To change password, please login from yo
  to:
    dl01.bsc.es
- Active Archive and transfer management n
    dt01.bsc.es
- For further information read MareNostrum

    http://www.bsc.es/support/MareNostrum

- BSC SUPPORT COMMANDS:

    See 'man bsc' for more information
```

# Question in microbial community ecology

**METAGENETICS**

Community DNA / RNA

Targeted amplification & HTS

Amplicon library

**OR**

Functional Gene

Taxonomic markers 16/18S rDNA

Shotgun HTS (mRNA enrichment)

Metagenome / Metatranscriptome

Extraction

**METAGENOMICS / METATRANSCRIPTOMICS**

Metabolic potential or actual transcription of the community

# Alternatively: Single Cell Genomics

# Some results using MareNostrum

# Malaspina 2010 expedition
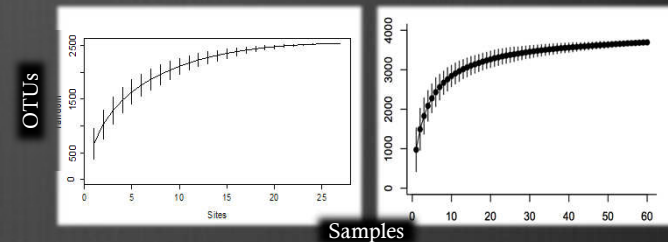
Assembly per sample ➔ eukaryotic few contigs

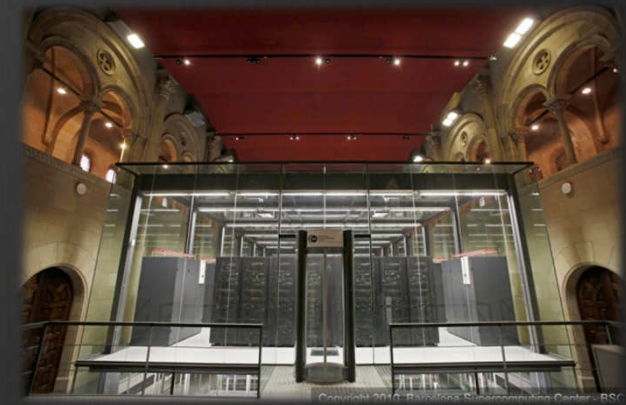…but, we knew the same OTUs were found in several samples

Co-Assembling all samples together (58 metagenomes) should generate longer contigs

Co-assembly of all samples generated longer contigs

MareNostrum Supercomputer
2,048 processors with
Ray assembler

| | One Sample | Sum of All Samples | Assembly All Samples |
|---|---|---|---|
| Contigs > 2Kb | 1,055 | 102,705 | 152,175 |
| Mean Coverage (> 2Kb) | 24.87 | - | 139.8 |
| Contigs > 10Kb | 21 | 5,823 | 23,086 |
| Largest Contig (bp) | 40,779 | 207,037 | 925,604 |
| % assembled reads (> 2Kb) | ~5% | - | ~40% |
| Largest Scaffold (bp) | 40,779 | - | 1,275,015 |

# Metagenomes

*1,500- 4,000m*

Co-assembly of 1,500-4,000m samples:

- 58 | 4,000m (5Gb each)

- 29 | 1,500-4,000m (20-40 Gb each)

- Ray assembly with 2,048 threads @MN (18hs)

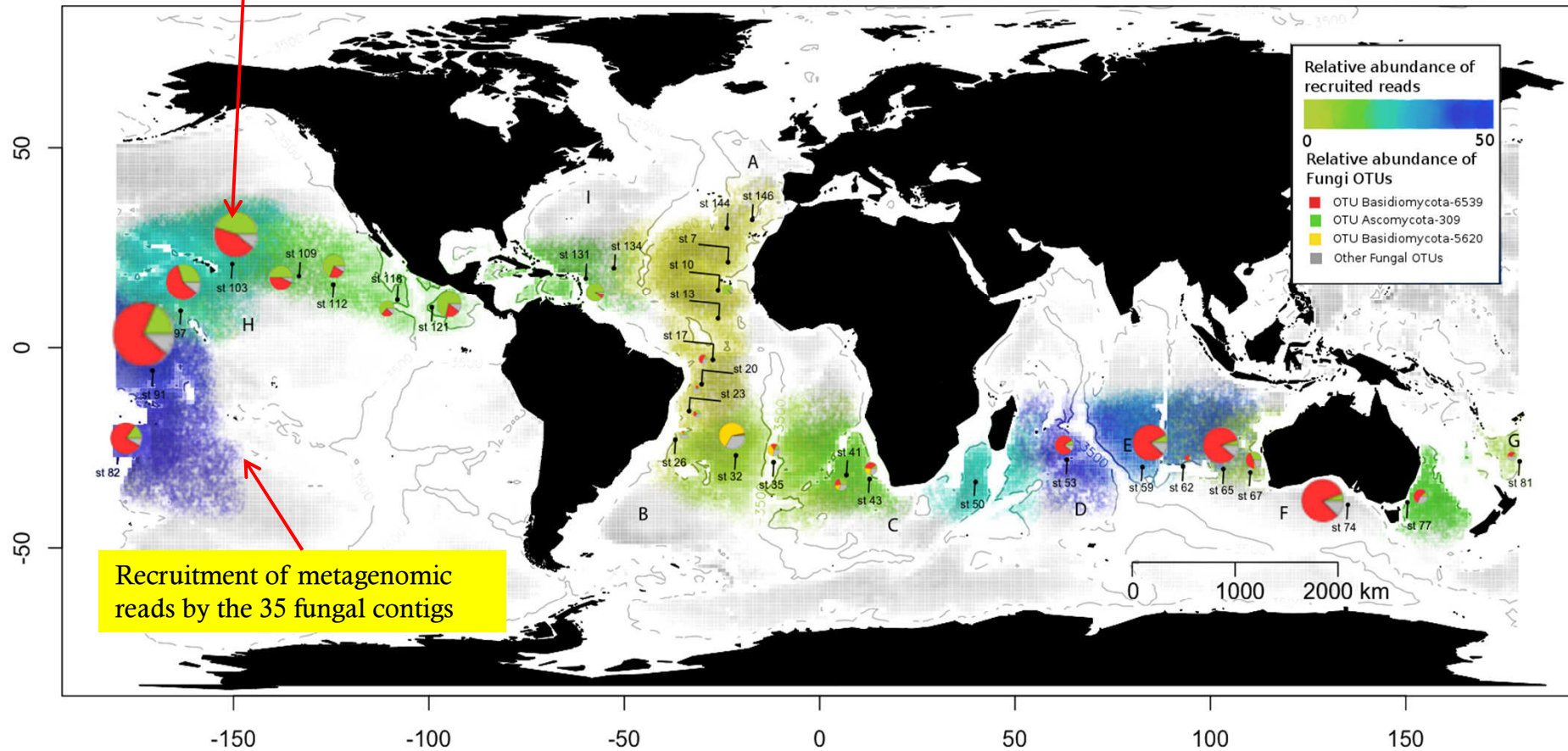|  | 4,000m | 1,500 - 4,000m |
|---|---|---|
| # Contigs > 2kb | 152,175 | 339,898 |
| Largest contig (bp) | 925,604 | 1,119,237 |

# Metagenomes: deep ocean fungi
*4,000m samples*

MALASPINA 2010

FRA of 2.6Mb of what seems to be a widespread fungus in the deep ocean

Pyrotag data

Recruitment of metagenomic reads by the 35 fungal contigs

Relative abundance of recruited reads

0          50

Relative abundance of Fungi OTUs
- OTU Basidiomycota-6539
- OTU Ascomycota-309
- OTU Basidiomycota-5620
- Other Fungal OTUs

# Co-assembly stats

| | % Genome recovery (CEGMA) | Assembly size (Mb; contigs >1,000bp) | Contigs (> 1,000bp) | Max. contig | N50 (>1,000bp) |
|---|---|---|---|---|---|
| **Co-Assembly 14 SAGs MAST-4 clade A (SPAdes)** | **89.1** | **47.5** | **14,564** | **57,905** | **4,563** |
| Co-Assembly 14 SAGs MAST-4 clade A (MegaHit) | 80.6 | 42.5 | 15,158 | 51,080 | 3,475 |
| *MAST-4 clade A single SAG assembly (mean \| SD)* | *20.6 \| 10.2* | *9.1 \| 4.5* | *1,694 \| 763* | *72,570 \| 20,347* | *11,041 \| 3,121* |
| **Co-Assembly 9 SAGs MAST-4 clade E (SPAdes)** | **68.5** | **32.3** | **5,677** | **104,912** | **9,991** |
| *MAST-4 clade E single SAG assembly (mean \| SD)* | *14.3 \| 5.5* | *6.2 \| 2.4* | *1,098 \| 350* | *63,915 \|18,608* | *10,567 \|1,920* |

# Continuing analyses with the co-assembly

1) Gene prediction [Augustus]

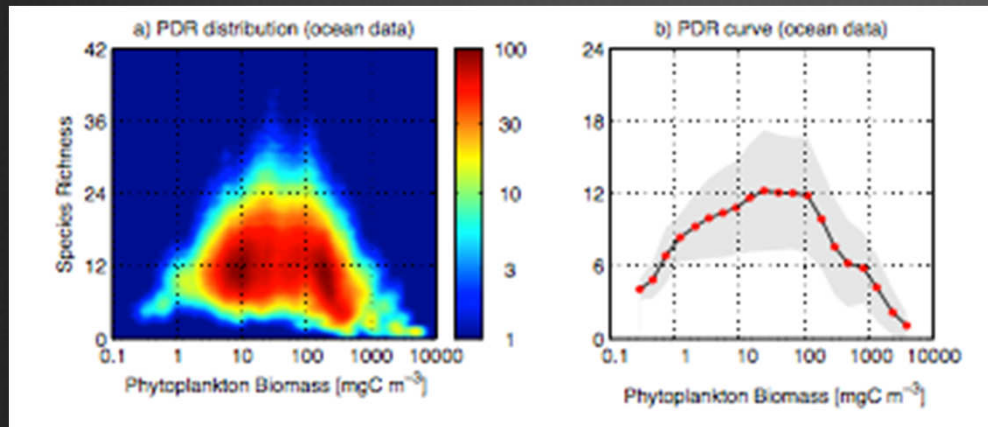2) Annotation (KEGG, KOG, Pfam, eggNOG, OMRGC, MMETSP)

General metabolic pathways
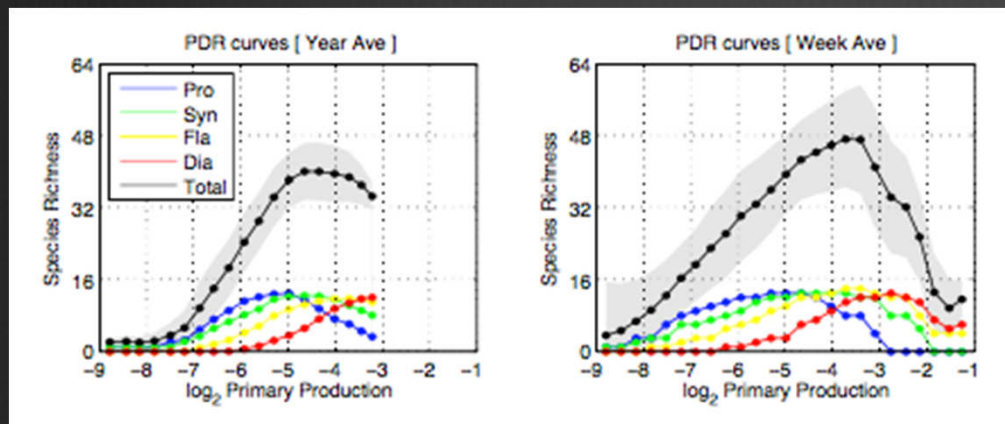
In Red: mapped MAST4 pathways/functions according to KO

# Large phylogenetic analyses



G. Salazar

⊛ Deep ocean Bacteria [Malaspina]

⊛ 3,500 sequences (16S rRNA)

⊛ 60 samples

⊛ About 1,000 processes

⊛ Stramenopiles (protists)

⊛ 3,835 sequences (>1,100bp) [18S ]

⊛ About 1,000 processes

# Simulations

## Primary production vs. richness

**Real data**



**MIT ecological selection model**

1. Global marine ecosystem model
2. Several plankton and nutrient types
   - 64 phytoplankton species (small, large)
   - 2 zooplankton generic (small, large)
   - 4 nutrients (N,P,Si,Fe)
3. Four phytoplankton functional groups with trade-offs
   - slow growth niche specialists (*Prochlorococcus*, *Synechococcus*)
   - fast growth niche opportunists (flagellates, diatoms)
4. Self-assembly of the phytoplankton community
   - ecological selection by resource competition
   - survival of the most adapted to the environment

**Simulation**



Vallina et al., 2014. Nat. Comm

# Summary of results with RES support since 2011

- 11 published papers

- 2 in revision


- Contributing mostly
  - Metagenomics
  - Genomics
  - Phylogenetics
  - Modelling

# Conclusions:
# Microbial ecology

- Massive amounts of DNA data need powerful computers as well as programs that can deal with them

- Future developments require further integration with high-performance computers and quantitative methods

- Analysis of large datasets will likely unveil patterns of genomic functioning as well as interactions between marine microbes

microbial Malaspina@*ICM*

MALASPINA 2010

MASSIMO PERNICE

GUILLEM SALAZAR

SILVIA G. ACINAS

RAMON MASSANA

JULIA PERERA BEL

PABLO SANCHEZ

PEP GASOL

CARLES PEDRÓS-ALIÓ

JOSE A. MARTIN CANO

FRANCISCO CORNEJO

MARTA SEBASTIAN

CATERINA RODRÍGUEZ

**Single Molecule Real Time (SMRT) sequencing**

- Average read length: 4,200-8,500 bp (longest read 30Kbp)
  - P4/C2: shorter reads, higher accuracy
  - P5/C3: longer reads, lower accuracy
  - 200-300 Mbp from each SMRT cell for 15-20kb insert size libraries
  - 100-150 Mbp for >20Kbp libraries
  - No multiplexing in genomics libraries (multiplexing in amplicons)
  - Library preparation 400-1200€
  - About 350 € per SMRT cell
  - Signal: colors

# What if you are not interested in the whole community but in one species?

**If you have a clonal culture:**

Genomics and/or RNA-Seq

**If you don't have a clonal culture:**

Single-Cell genomics

illumina®

**MiSeq**
- $25 \times 10^6$ reads
- 15Gb/run
- 2x300bp

**NextSeq 500**
- $400 \times 10^6$ reads
- 120Gb/run
- 2x150bp

**HiSeq 2500**
- $4 \times 10^9$ reads
- 1000 Gb/run
- 2x125bp

Signal: colors
$1Gb = 1 \times 10^9$ bases

# Oceanic microbial community

- Includes all species occurring at a particular site and their abundances

# Multiple metabolisms



Brock "Biology of Microorganisms"
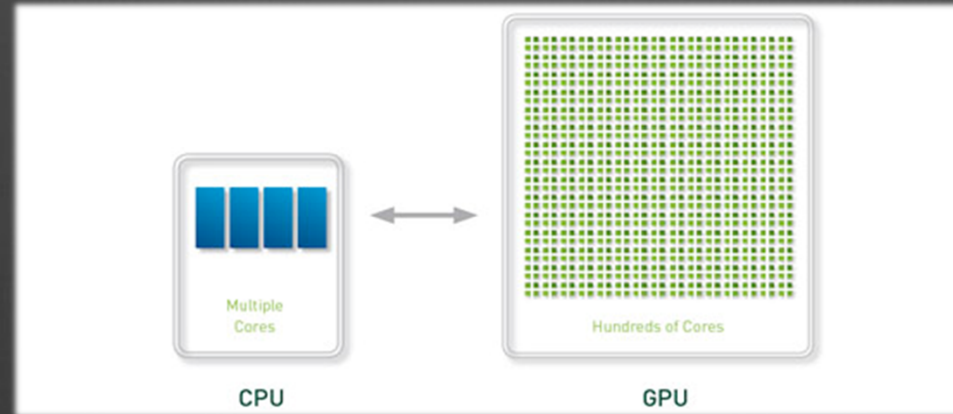
# Comparison of platforms

| | Run time | Mreads/run | Read length | Mb/run | €/Mb |
|---|---|---|---|---|---|
| Sanger (3730xl) | 2h | 0.000096 | 400-900 | 0.06 | 1500 |
| 454 FLX Titanium | 10h | 1 | 400 | 400 | 15 |
| 454 FLX+ | 18-20 h | 1 | 700 | 900 | 9 |
| Ion Torrent | 2h | 80 | 400 | 32,000 | 1 |
| PacBio | 0.5-2h | 0.005 | 4-8 K | 300 (SMRT) | 0.33-1 |
| Illumina MiSeq | 55h | 25 | 2x300 | 15,000 | 0.1 |
| Illumina GAIIx | 14 days | 320 | 2x150 | 96,000 | 0.12 |
| Illumina HiSeq2500 | 1-11 days | 4000 | 2x125 | 1,000,000 | 0.05 |

$1Gb=1x10^9$ bases

Glenn 2011 updated in
http://www.molecularecologist.com/next-gen-table-2-2014/

Endocytosis

Clathrin-dependent endocytosis

Clathrin-independent endocytosis

In Red: mapped
MAST4
pathways/functions
according to KO

04144 7/23/14
(c) Kanehisa Laboratories

# GPU (graphics processing unit) computing



Serial part of an application runs on a CPU and the computationally-intensive part runs on a GPU



GPU Pipeline for HTS sequencing

Centro de Investigación Príncipe Felipe

http://docs.bioinfo.cipf.es/projects/ngs-gpu-pipeline/wiki

# Cloud Computing...

- Purchase needed computer power
- Scalable (few to thousands of processors)
- No maintenance costs



- GALAXY (g2.bx.psu.edu)
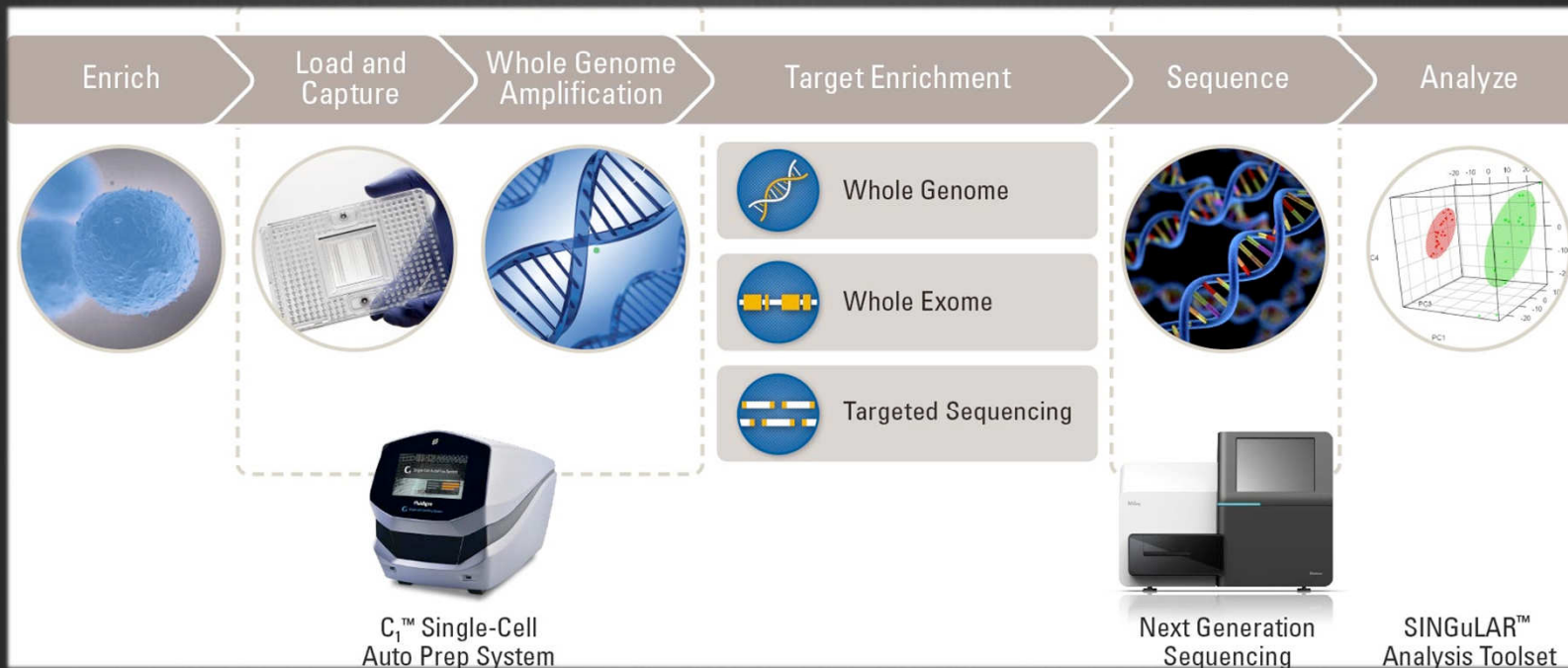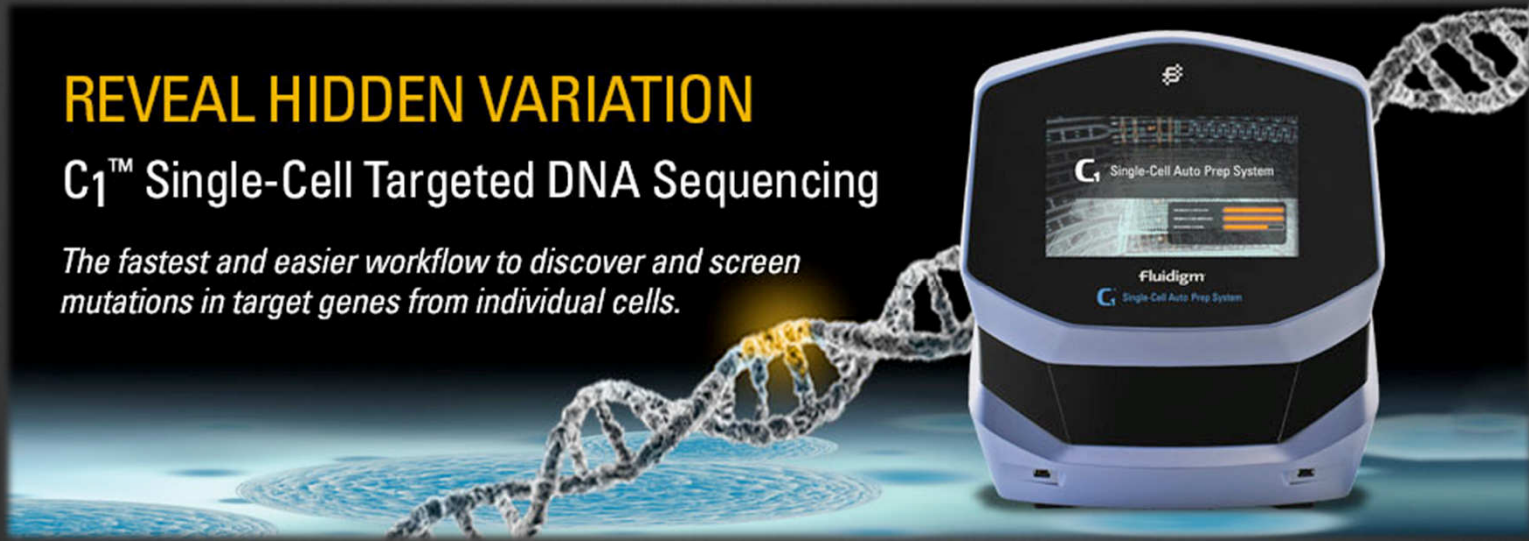- N3phele: HTS analyses at the
    EC2 with QIIME

- **Molecular biology + computers + stats**

- **The next 20 years of genome research**

- **M. Schatz (2015)**

- http://biorxiv.org/content/early/2015/06/02/020289

# General SAG construction strategy