



# Introduction to MetaGenomics



Institut de Recherche sur le Cancer et le Vieillessement, Nice  
Institute for Research on Cancer and Aging, Nice  
CNRS UMR 7284 - INSERM U 1081 - UNS



**Olivier Croce** – [croce@unice.fr](mailto:croce@unice.fr)  
Bioinformatics service - IRCAN

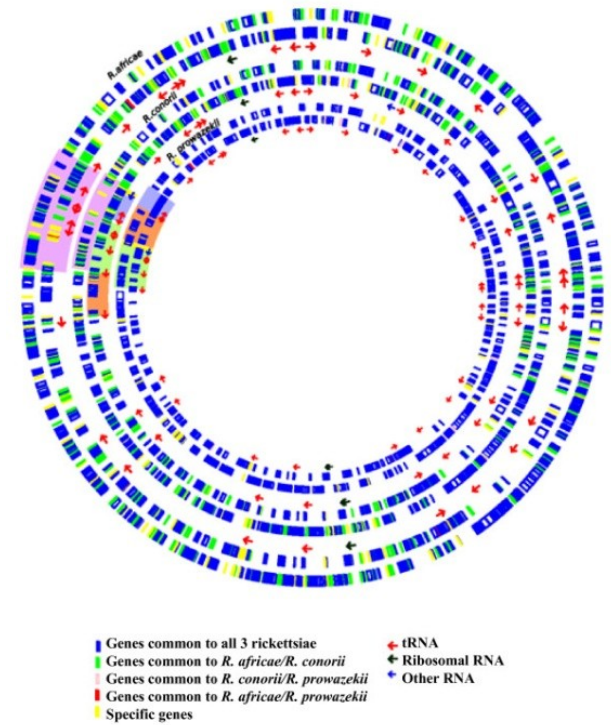
# Definitions

## Genomics:

=> Analysis of the genome from a single organism, phenotype prediction

- Presence/absence of genes
- Chromosomal gene order (synteny)
- Comparison of presence/absence of orthologous genes
- Presence of indels or SNPs in conserved genes
- Repetitive motifs
- ...

*l.g.* for prokaryotes => optimization or design of culture media, resistance to antibiotics, detection of virulence



## MetaGenomics:

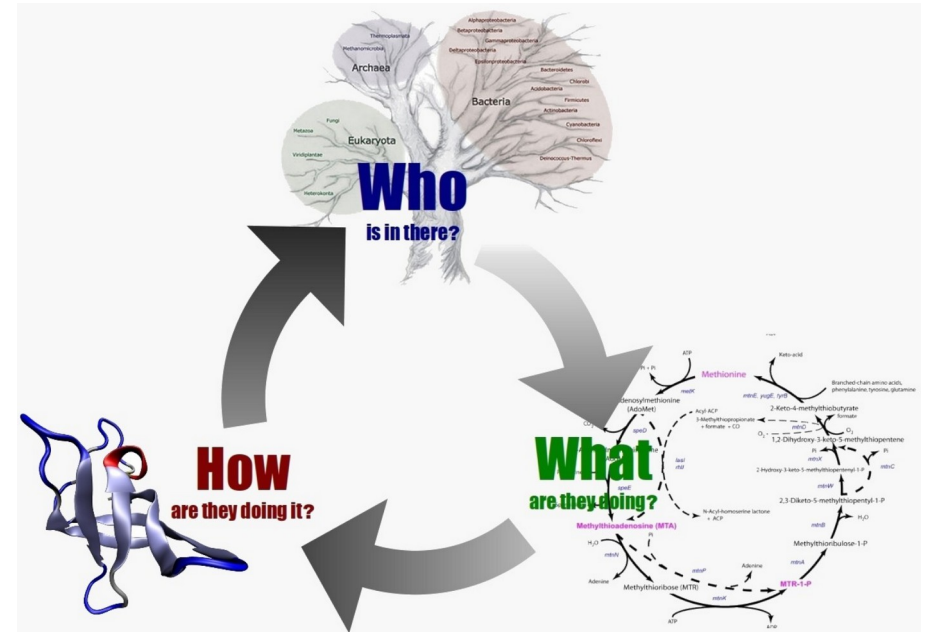
=> Study of the genetic material recovered from an environmental sample. The sample usually contains >1 species !

Metagenome: the collective genome of all the microorganisms in a given environment.



# Goals and applications

**Goals:** Who is there ? What are they doing ?  
How are they doing ?  
Characterization of the diversity



**Many applications :**  
fondamental research,  
industrial, clinical  
applications



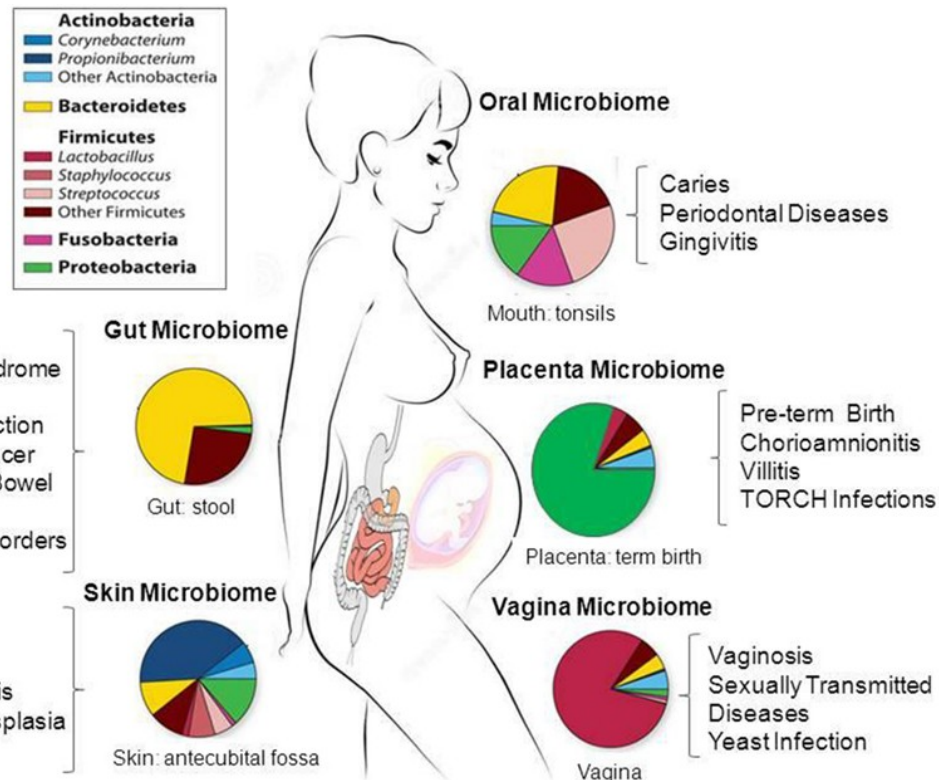
- Growth
- Protein
- Vaccine
- Secondary Metabolite
- Resistance

# Where ?

Almost everywhere !

=> **Microbiome** (gut, skin,...), ocean (TARA project) , soils and deep ground, space ! And more

- More cell in guts than in the rest of the body ! ( $>10^{14}$  )
- Be considered as a real organ



Colonization of the intestine:

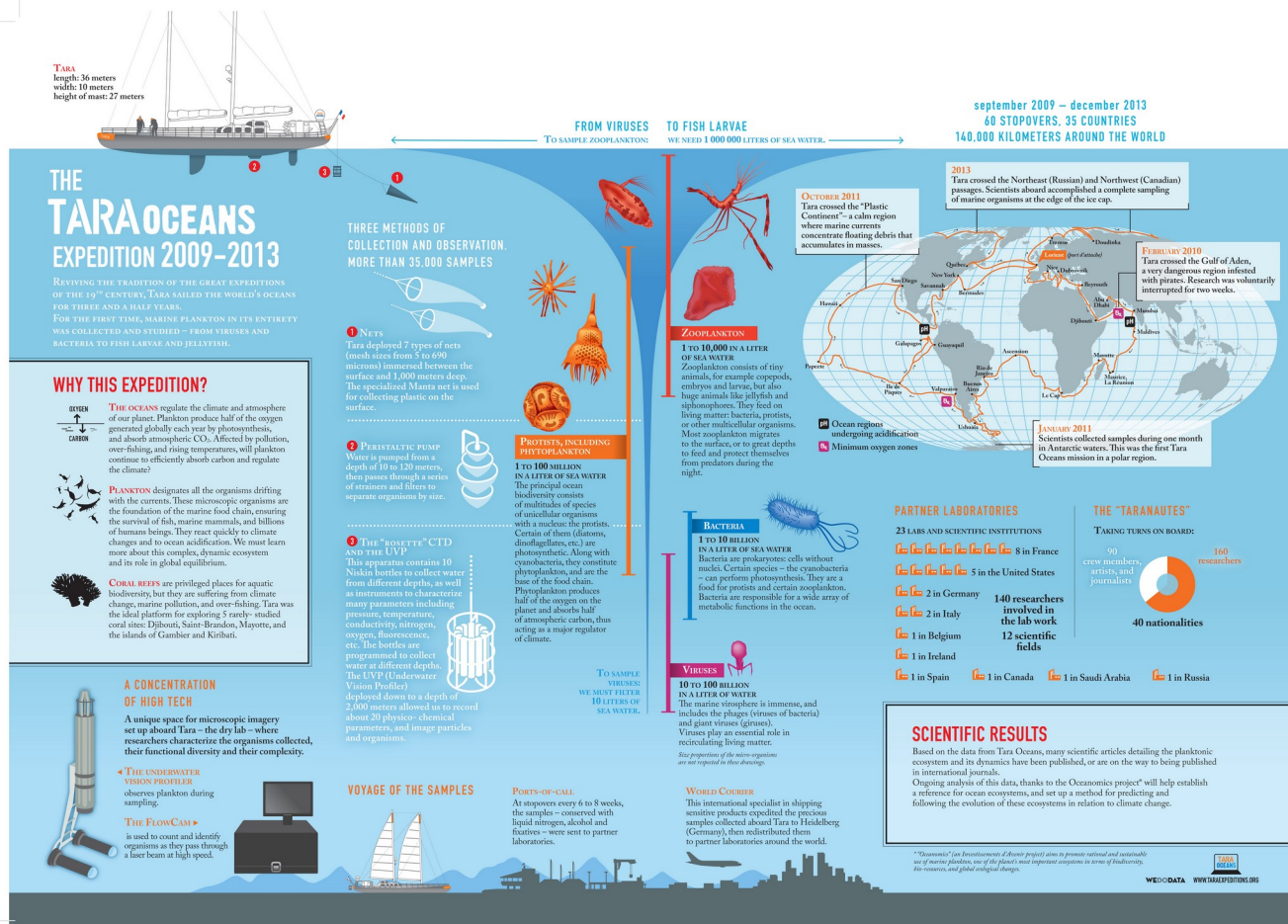
- Initial acquisition during delivery by vagina, stool and skin
- Breastfeeding
- Breast milk: Bifidobacterium (more than 90% of the flora)
- Industrial milk: more Bacteroides and Clostridium
- Stool begins to become hard = adult flora



# Where ?

Almost everywhere !

=> Microbiome (gut, skin,..), **ocean** (TARA project) , soils and deep ground, space ! And more



## Tara Oceans

117 millions of oceanic genes founded !

# Where ?

Almost everywhere !

=> Microbiome (gut, skin,..), ocean (TARA project) , **soils and deep ground**, space ! And more

## High coverage sequencing of DNA from microorganisms living in an oil reservoir 2.5 kilometres subsurface

Hans K. Kotlar,<sup>1\*</sup> Anna Lewin,<sup>2\*</sup> Jostein Johansen,<sup>3</sup> Mimmi Throne-Holst,<sup>4\*</sup> Thomas Haverkamp,<sup>5</sup> Sidsel Markussen,<sup>4</sup> Asgeir Winnberg,<sup>4</sup> Philip Ringrose,<sup>1</sup> Trine Aakvik,<sup>2</sup> Einar Ryeng,<sup>3</sup> Kjetill Jakobsen,<sup>3</sup> Finn Drabløs<sup>3</sup> and Svein Valla<sup>2\*</sup>

<sup>1</sup> Statoil ASA, 7053 Ranheim, Norway.

<sup>2</sup> Department of Biotechnology, Norwegian University of Science and Technology, 7491 Trondheim, Norway.

<sup>3</sup> Department of Cancer Research and Molecular Medicine, Norwegian University of Science and Technology, 7491 Trondheim, Norway.

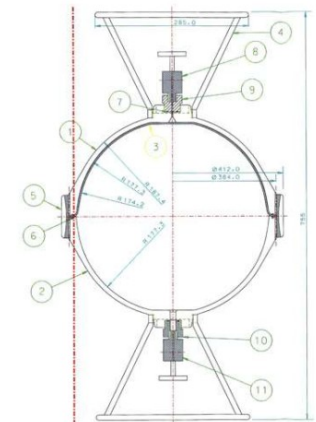
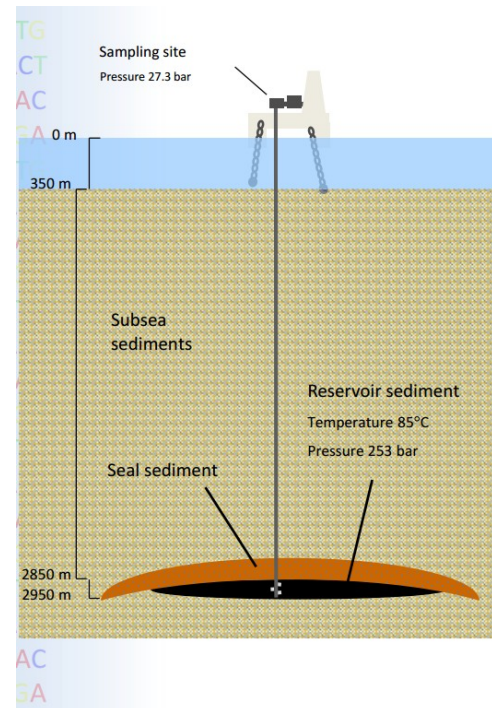
<sup>4</sup> SINTEF Materials and Chemistry, Department of Biotechnology, 7465 Trondheim, Norway.

<sup>5</sup> CEES and MERG, Department of Biology, University of Oslo, 0316 Oslo, Norway.

erant than for a corresponding *E. coli* enzyme, consistent with the conditions in the oil reservoir.

### Introduction

The diversity of environments on the earth is enormous, ranging from, e.g. extremely cold to hot, dry, or acidic conditions, and studies of microbes inhabiting such extreme environments are interesting from a basic biological point of view, for applied biotechnology and to evaluate the outer boundaries for the existence of life. Oil reservoirs located deep into the earth crust attribute a combination of high pressure, temperature, salinity as well as physical barriers to life on the surface. Petroleum oil is generated from organic materials deposited millions of years ago, buried under layers of sediments of gradually increasing



Xpand Pressure Flask

Sampling an oil field

# Where ?

Almost everywhere !

=> Microbiome (gut, skin,..), ocean (TARA project) , soils and deep ground, **space** ! And more

## Survival of *B.subtilis* spores :

Unprotected : several seconds

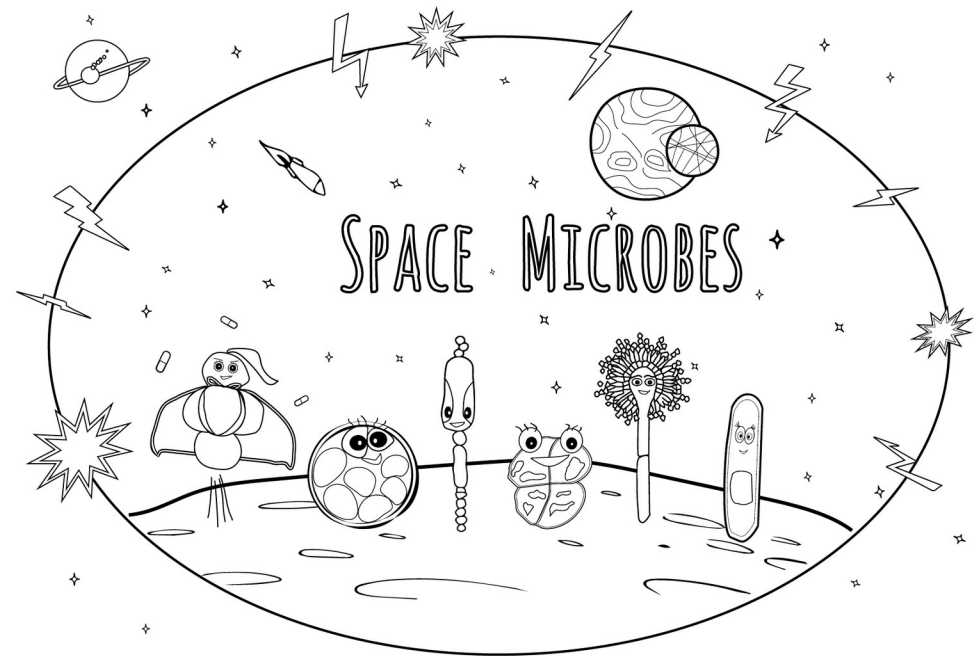
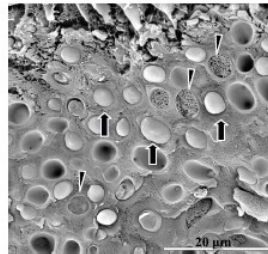
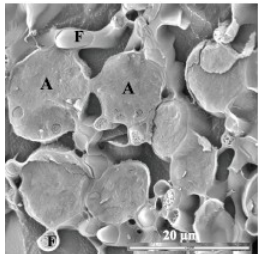
Protected : more than 6 years

## Others microorganisms :

Phage T1, *Synechococcus*

*Haloarcula*, *Deinococcus*

Recently : surprising **survival of lichen**

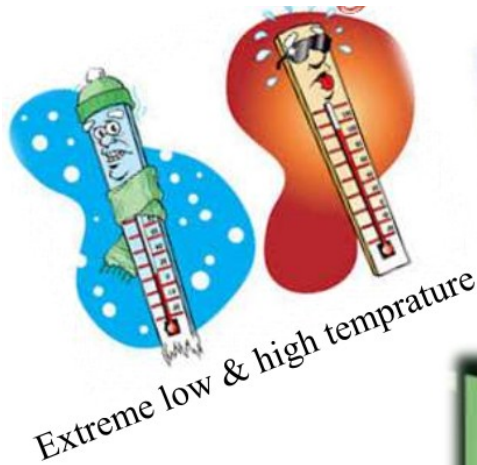




# Where ?

Almost everywhere !

=> Microbiome (gut, skin,..), ocean (TARA project) , soils and deep ground, space ! **And more**



volcano



Soil



Havey metal composition



Waste water



Acidic



alkaline

Many other metagenomics hot spots



# Methods

## \* Microarrays

Requires knowledge of the community in advance

- PhyloChip (taxonomic)
- Geochip (metabolic)

## \* (Meta) Barcoding sequencing

Amplicon based analysis through High throughput sequencing of a given gene (or part of) after amplification

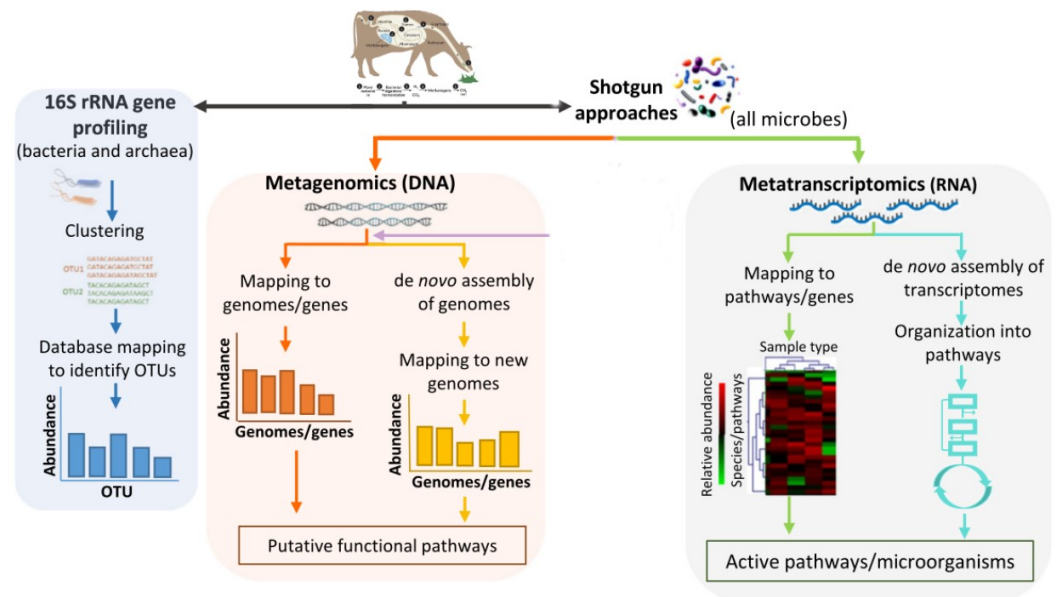
- SSU rRNA (=16S/18S)
- protein coding genes: rpoB, nifH, IRS, cytC, RecA,...
- ITS (internal transcribed spacer)

## \* (WGS / Total / Shotgun) MetaGenomics sequencing

Sequencing of the whole DNA in a sample. Complete community analysis, characterization of the pangenome

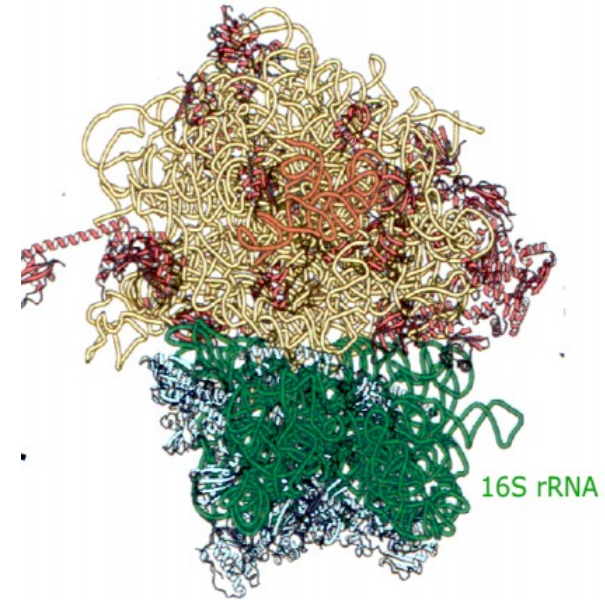
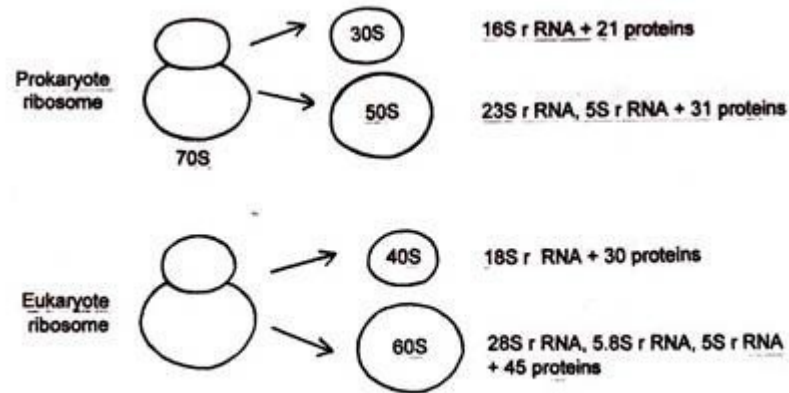
## \* MetaTranscriptomics

*Idem TotalMetagenomics, but from RNA*



# (Meta) Barcoding sequencing

\* Usually based on rRNA **16S** for **bacteria** (~1540nt), rRNA **18S** for **eucaryotes**  
=> current taxonomic classification for prokaryotes & eukaryotes. Species definition !



\* Most commonly used molecular marker  
– essential function  
– Ubiquity  
– evolutionary properties

\* OTU (operational taxonomic units)  
definition based on 16S rRNA gene

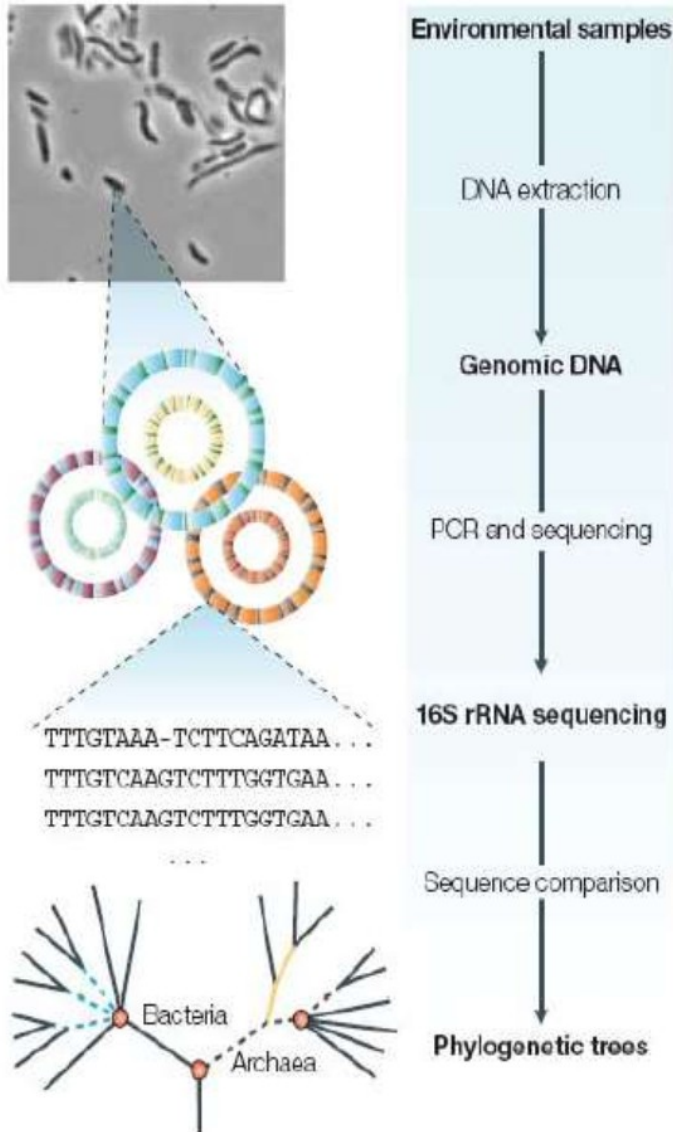
=> organisms displaying 97 to 98% identity  
in this gene to be part of the same OTU

\* Rapid and cost-effective approaches for  
assessing diversity and abundance.

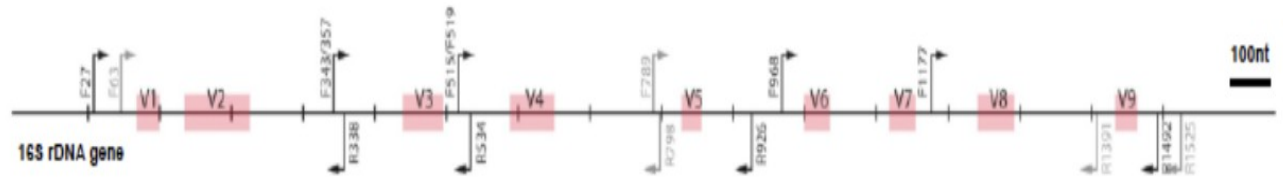
\* Highly conserved gene, 9 hypervariable  
domains interspersed with conserved  
fragments

# (Meta) Barcoding sequencing

16s rRNA structure. Design primers for PCR (polymerase chain reaction) for specific regions of the rRNA, not the whole molecule (ie. V4 and V6 region for 16S)



\* PCR amplification with primers that hybridize to highly conserved regions in bacterial or archaeal 16S rRNA, followed by cloning and sequencing



\* Phylogenetic analysis of 16S rRNA helps to reveal the species diversity in a community

# (Meta) Barcoding sequencing

## Limits of rRNA use

- Sampling challenges : rich species Vs sparse species. Rare species could not be sequenced
- Do not tell much about the functional abilities of a community
- Based on the assumption that the level of interspecies rRNA variation is homogeneous among genera
- PCR bias: not all rRNA genes amplify equally well with the same “universal” primers
- Multiple copies of rRNA genes in some species (which may artificially lead to the over-representation of some species)
- Speed of evolution of rRNA genes may vary according to the phylum
- Homology cutoff not applicable to all genera => may underestimate, or overestimate the diversity
- The discriminatory power may be insufficient at the species level, especially for closely related species

ig.

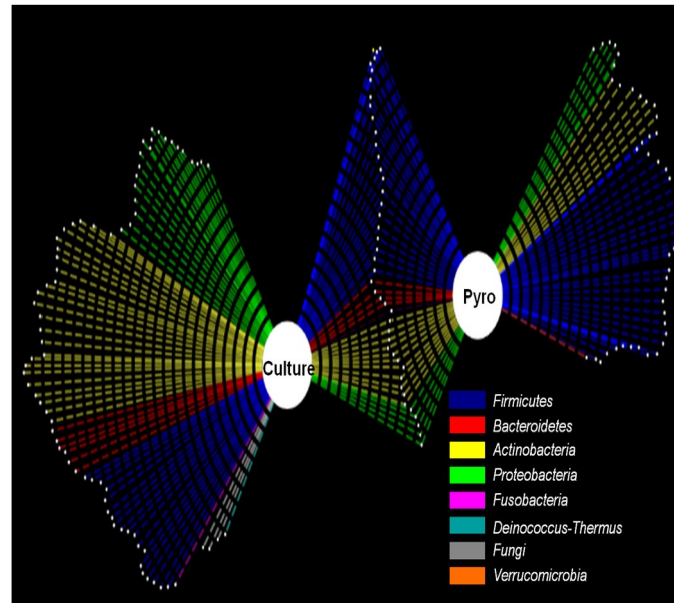
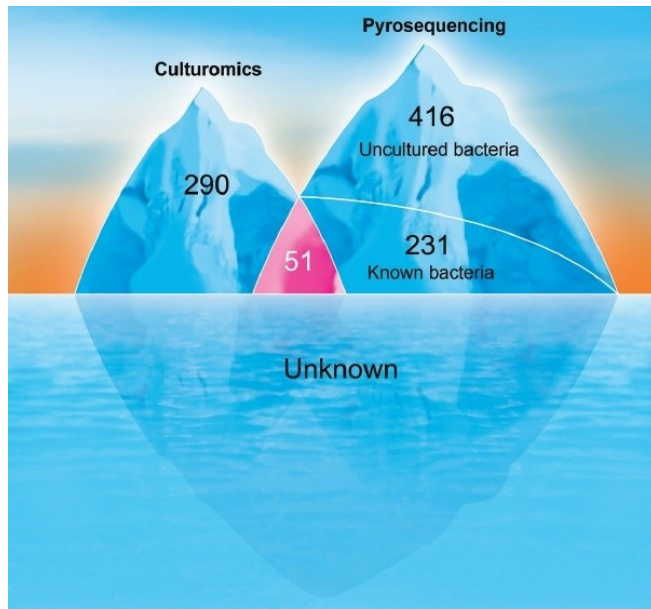
- *Pantoea agglomerans* strains may exhibit up to 27 bp differences, which does not validate the 98.7% cutoff)
- *Clostridium tetani* and *C. innocuum* exhibit 104 bp differences, which does not validate the 95% threshold => classification in distinct genera?



# (Meta) Barcoding sequencing

## Limits of rRNA use

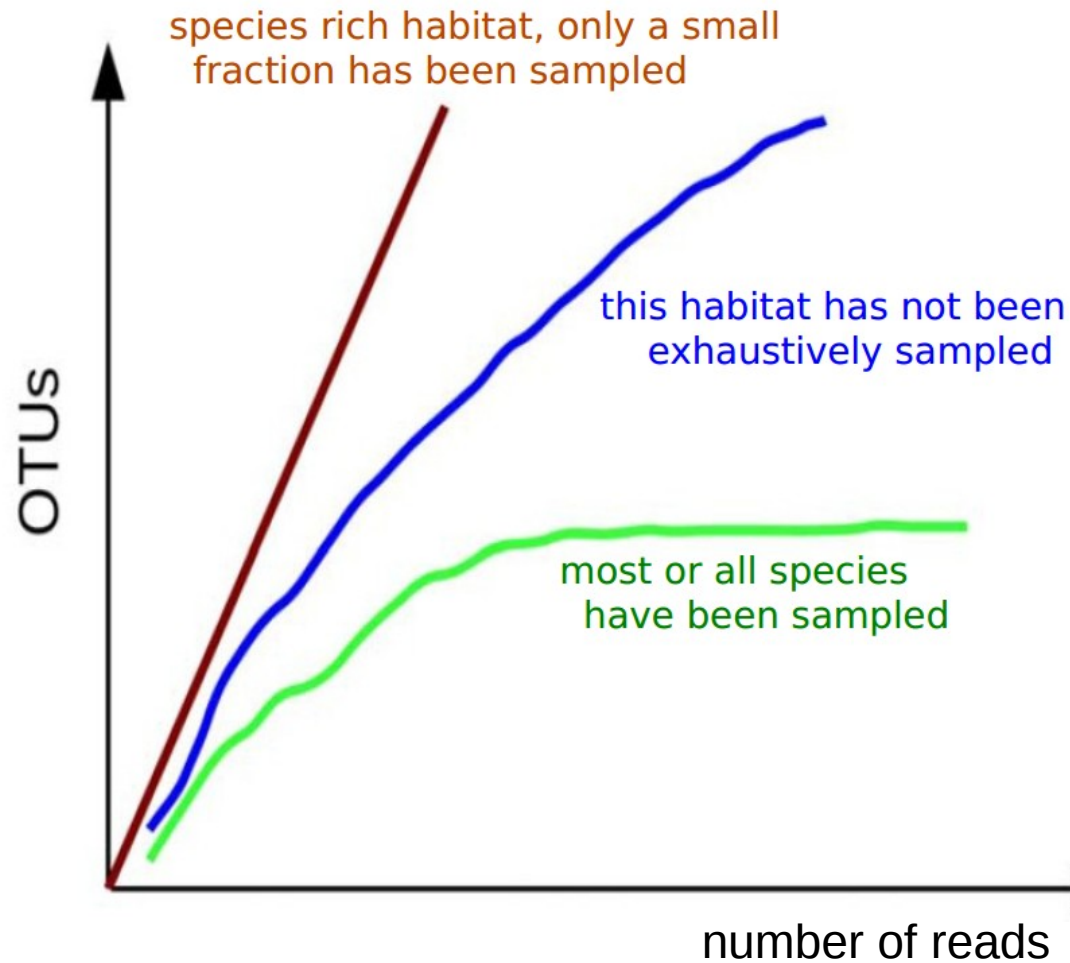
- Possibility of 16S rRNA genes acquired by HGT (Jain et al. Horizontal gene transfer among genomes: the complexity hypothesis. Proc Natl Acad Sci USA 1999;96:3801-6)
- Missing sequences or variable quality of available sequences, including from validly published species



**For procaryotes** : NGS could be complementary to culturomics. 90-95% microorganisms remain uncultivable in laboratory.

**For eucaryotes** : cultures is not even possible (plankton, fungi..), direct observation (microscope) is exhausting => NGS is widely uses (rRNA 18S, or other genes)

# (Meta) Barcoding sequencing



**Rarefaction** allows the calculation of species richness for a given number of individual samples, based on the construction of so-called rarefaction curves. This curve is a plot of the number of species as a function of the number of samples.

# (Total) MetaGenomics sequencing

=> Environmental sample : all DNA sequencing, no specific amplification

Who is here ? => Biodiversity characterization

Who does what ? => Physiological characterization

## RESEARCH ARTICLE

### Environmental Genome Shotgun Sequencing of the Sargasso Sea

J. Craig Venter,<sup>1\*</sup> Karin Remington,<sup>1</sup> John F. Heidelberg,<sup>3</sup>  
Aaron L. Halpern,<sup>2</sup> Doug Rusch,<sup>2</sup> Jonathan A. Eisen,<sup>3</sup>  
Dongying Wu,<sup>3</sup> Ian Paulsen,<sup>3</sup> Karen E. Nelson,<sup>3</sup> William Nelson,<sup>3</sup>  
Derrick E. Fouts,<sup>3</sup> Samuel Levy,<sup>2</sup> Anthony H. Knap,<sup>6</sup>  
Michael W. Lomas,<sup>6</sup> Ken Nealson,<sup>5</sup> Owen White,<sup>3</sup>  
Jeremy Peterson,<sup>3</sup> Jeff Hoffman,<sup>1</sup> Rachel Parsons,<sup>6</sup>  
Holly Baden-Tillson,<sup>1</sup> Cynthia Pfannkoch,<sup>1</sup> Yu-Hui Rogers,<sup>4</sup>  
Hamilton O. Smith<sup>1</sup>

We have applied "whole-genome shotgun sequencing" to microbial populations collected en masse on tangential flow and impact filters from seawater samples collected from the Sargasso Sea near Bermuda. A total of 1.045 billion base pairs of nonredundant sequence was generated, annotated, and analyzed to elucidate the gene content, diversity, and relative abundance of the organisms within these environmental samples. These data are estimated to derive from at least 1800 genomic species based on sequence relatedness, including 148 previously unknown bacterial phylotypes. We have identified over 1.2 million previously unknown genes represented in these samples, including more than 782 new rhodopsin-like photoreceptors. Variation in species present and stoichiometry suggests substantial oceanic microbial diversity.

1.2 million unknown genes  
(Venter et al., 2004)

#### INTRODUCTION TO SPECIAL ISSUE

### Tara Oceans studies plankton at planetary scale

P. Bork<sup>1</sup>, C. Bowler<sup>2</sup>, C. de Vargas<sup>3,4</sup>, G. Gorsky<sup>5,6</sup>, E. Karsenti<sup>2,7</sup>, P. Wincker<sup>8</sup>

+ See all authors and affiliations

Science 22 May 2015:  
Vol. 348, Issue 6237, pp. 873  
DOI: 10.1126/science.aac5605

Article    Figures & Data    Info & Metrics    eLetters    PDF

The ocean is the largest ecosystem on Earth, and yet we know very little about it. This is particularly true for the plankton that inhabit the ocean. Although these organisms are at least as important for the Earth system as the rainforests and form the base of marine food webs, most plankton are invisible to the naked eye and thus are largely uncharacterized. To study this invisible world, the multinational *Tara* Oceans consortium, with use of the 110-foot research schooner *Tara*, sampled microscopic plankton at 210 sites and depths up to 2000 m in all the major oceanic regions during expeditions from 2009 through 2013 (1).

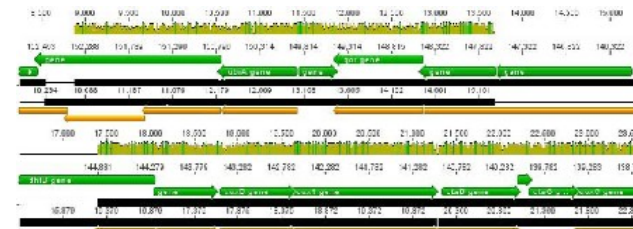
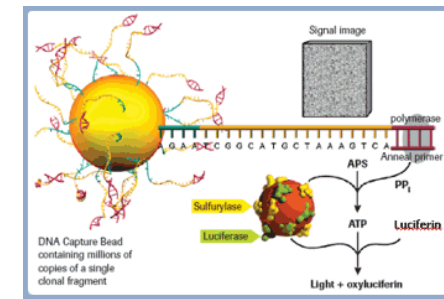
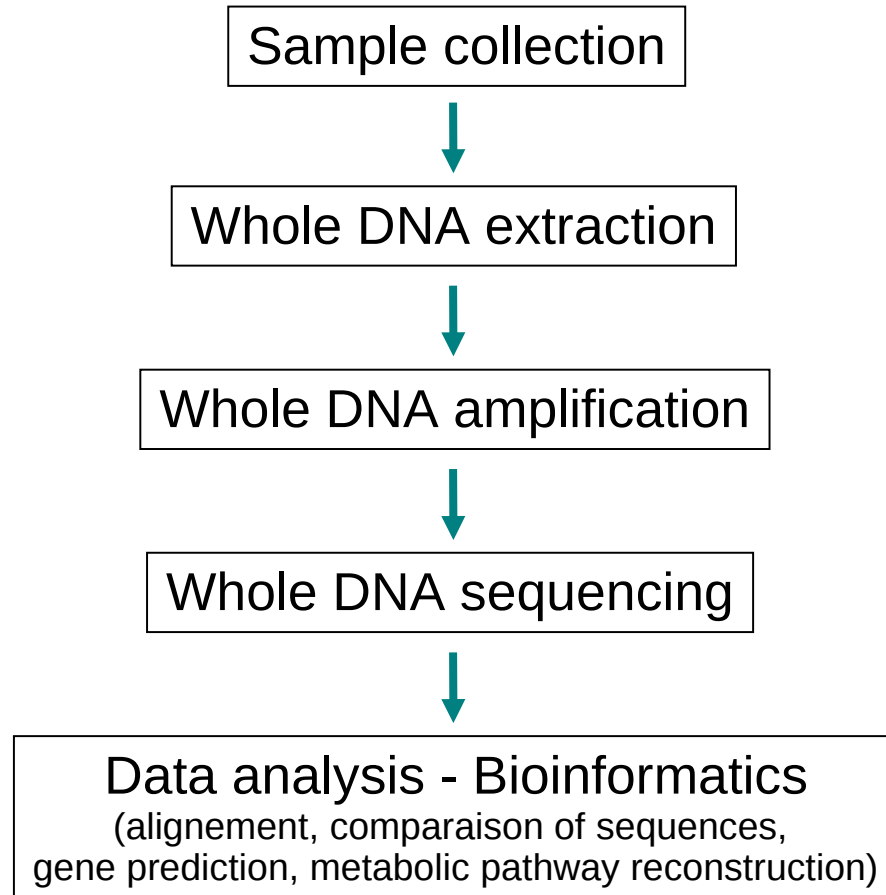
Success depended on collaboration between scientists and the *Tara* Expeditions logistics team. The journey involved not only science but also outreach and education as well as negotiation through the shoals of legal and political regulations, funding uncertainties, threats from pirates, and unpredictable weather (2). At various times, journalists, artists, and teachers were also on board. Visitors included Ban Ki-moon (Secretary-General of the United Nations) and numerous youngsters, including



Tara Oceans : 117 millions of oceanic genes  
(Bork and al., 2015)

# (Total) MetaGenomics sequencing

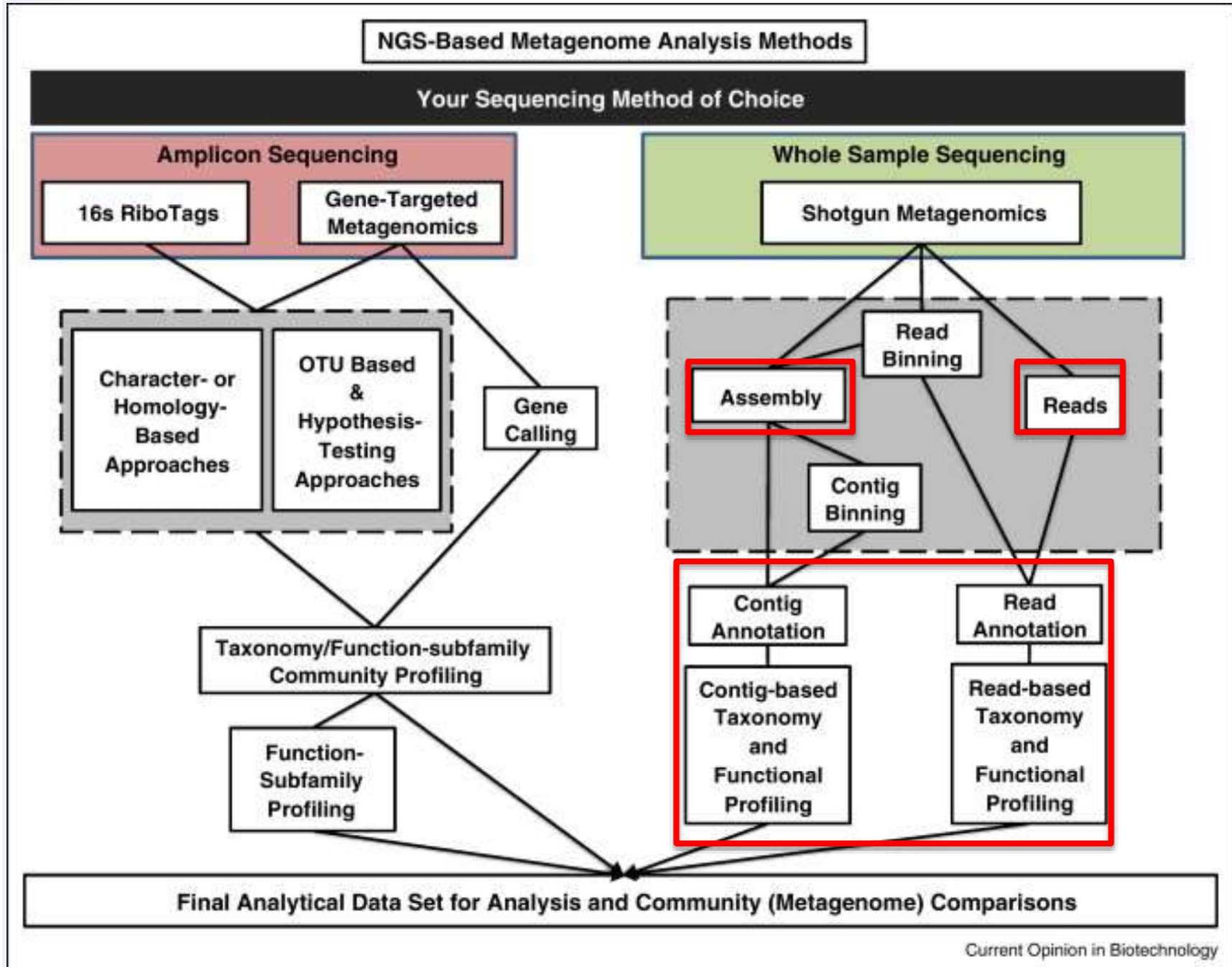
## Principle of metagenomics



=> Information about : Biodiversity, but also physiology, metabolic pathways...



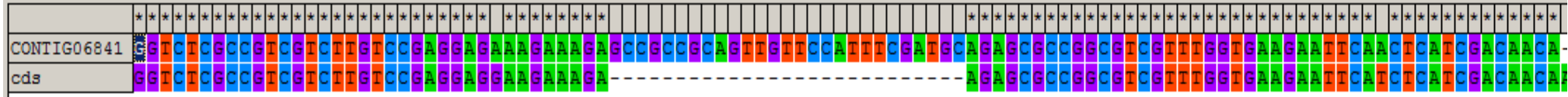
# (Total) MetaGenomics sequencing



# (Total) MetaGenomics sequencing

## Analysis and expected results

- Direct results : DNA sequences and contigs



- After analysis

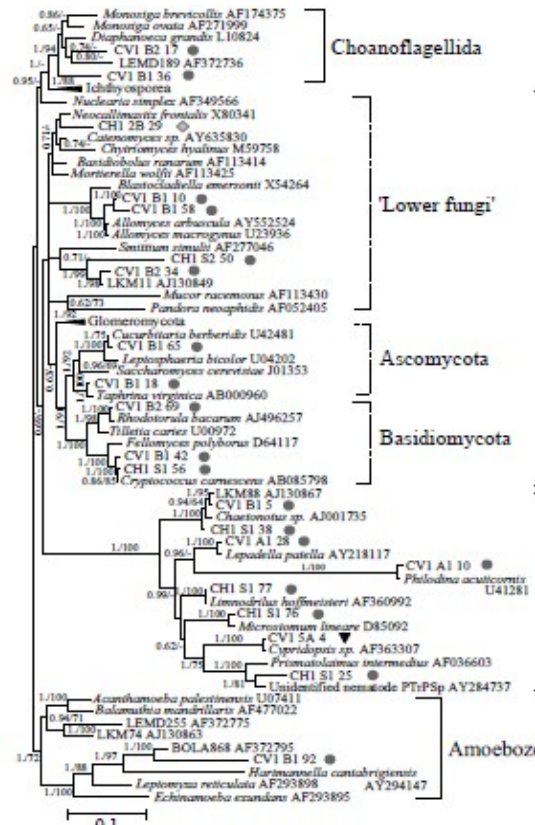
### (1) Biodiversity

- Phylogenetic trees :

Known species and

Unknown microorganisms

=> Who are they closest to ?



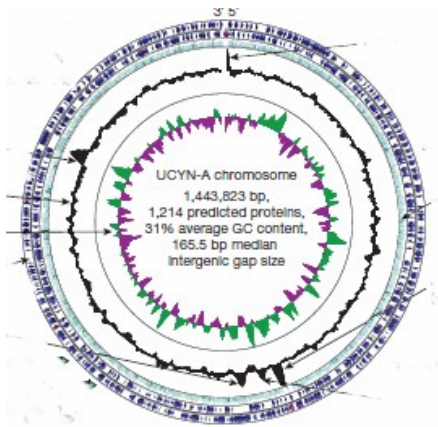
Discovery of :

- new species ?
- new phyla ?

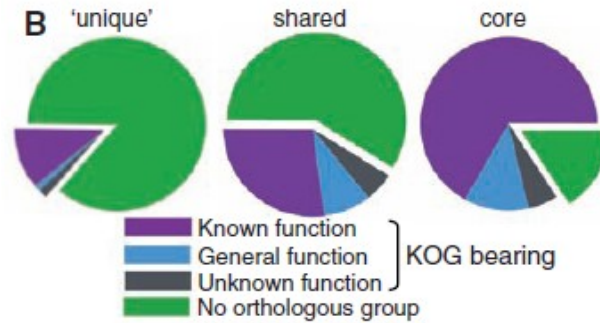
# (Total) MetaGenomics sequencing

## (2) Physiology and genomic

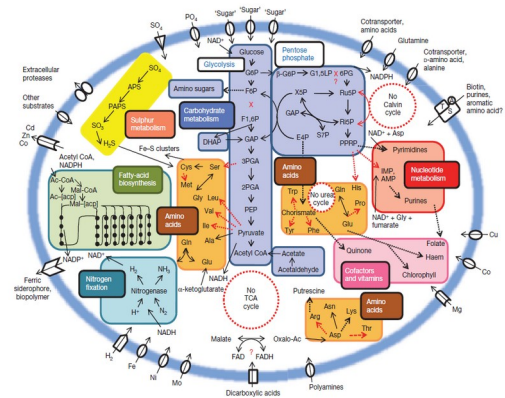
### Genomes



### Genes functional groups



### Pathways metabolism, physiology



Informations as :

- G+C contents
- Genome sizes
- DNA repair mechanisms
- Pathways of excretion, polysaccharides secretion
- ...

# (Total) MetaGenomics sequencing

## Sequence Classification

More difficult than for barcoding sequencing => creating bins

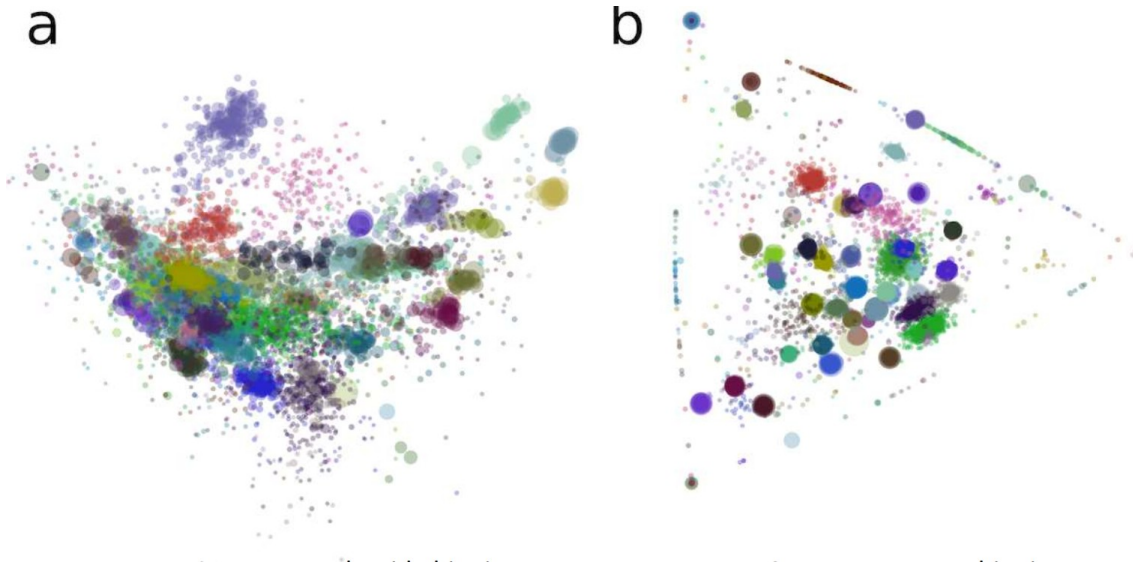
Sequence classification (binning) is the process of separating sequence data using specific information. Sequence classification by:

### **(1) Sequences composition**

- \* Tetranucleotide frequency (kmer counting)
- \* Clustering of reads. (e.g. swarm, cd-hit)
- \* Sequence (co-) assembly (MetaHit, Metavelvet)
- \* Differential coverage of contigs (GroopM, Concoct)

Advantage : read with unknown origin can be classified into a bin

Disadvantage: impossible to determine taxonomy or function of the reads



Input data: 1159  
(genomesImelfort et al., PeerJ, 2014)

- a) PCA - Tetranucleotide binning
- b) GroopM coverage binning



# (Total) MetaGenomics sequencing

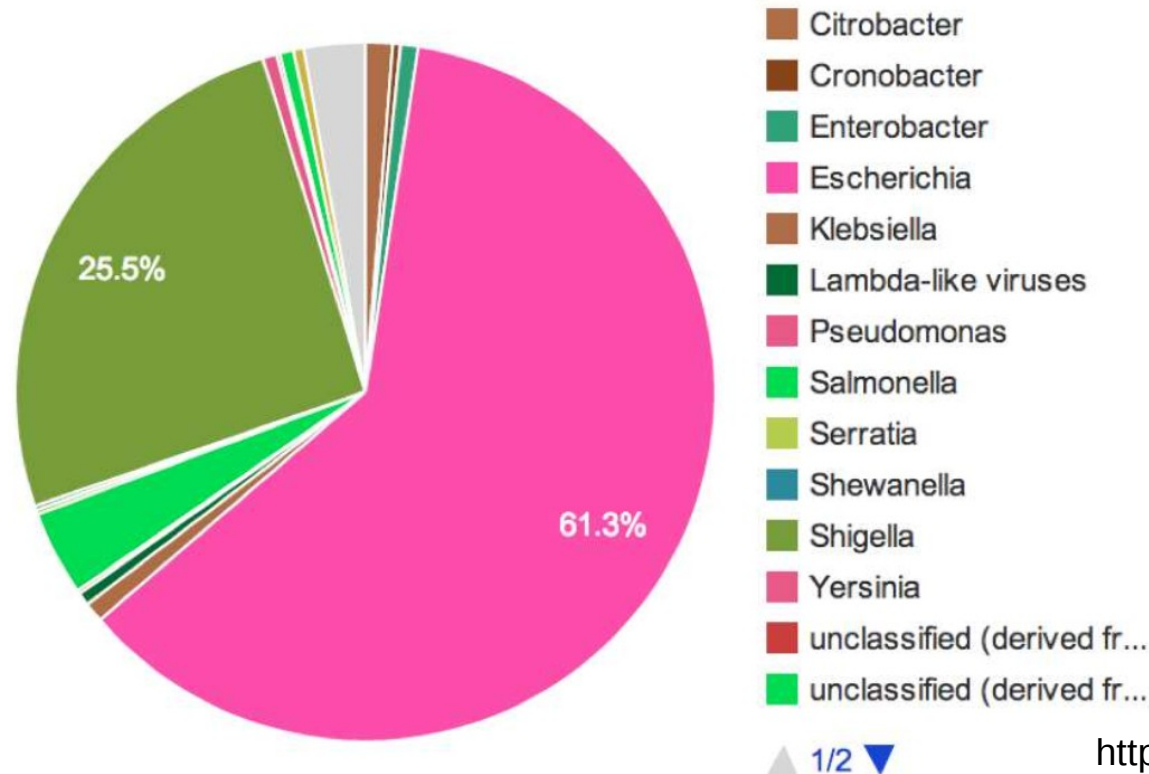
## (2) Sequences similarity

- \* Compare sequences to reference database (e.g. Blast, bwa, bowtie)
- \* Use phylogenetics to classify sequences.

Advantage: One can determine taxonomy and function of reads.

Disadvantage: reads with no similarity to databases sequences, can not be classified

Usual example : Using the best blast hit



# (Total) MetaGenomics sequencing - tools

## Tools for sequences classification

### **Nucleotide composition:**

CompostBin , PCA-analysis of k-mer, frequencies, Self-Organizing Maps (different variants), MetaCluster, PhyloPythia, Naïve Bayes classifier (NBC), etc

### **Sequence similarity:**

MEGAN, SorT-Items, Threephyler, COMET, Metaphlan, PhyloSift, Kraken, etc

### **Both:**

Phymm / PhymmBL, Phylopythia, RAphy, Metaxa2\*(rRNA), PhyloOTU (rRNA), MLTreeMap, RITA, STAMP, WGSQuikr.

### **Differential Coverage:**

GroopM, Concoct, Blobology

# (Total) MetaGenomics sequencing - tools

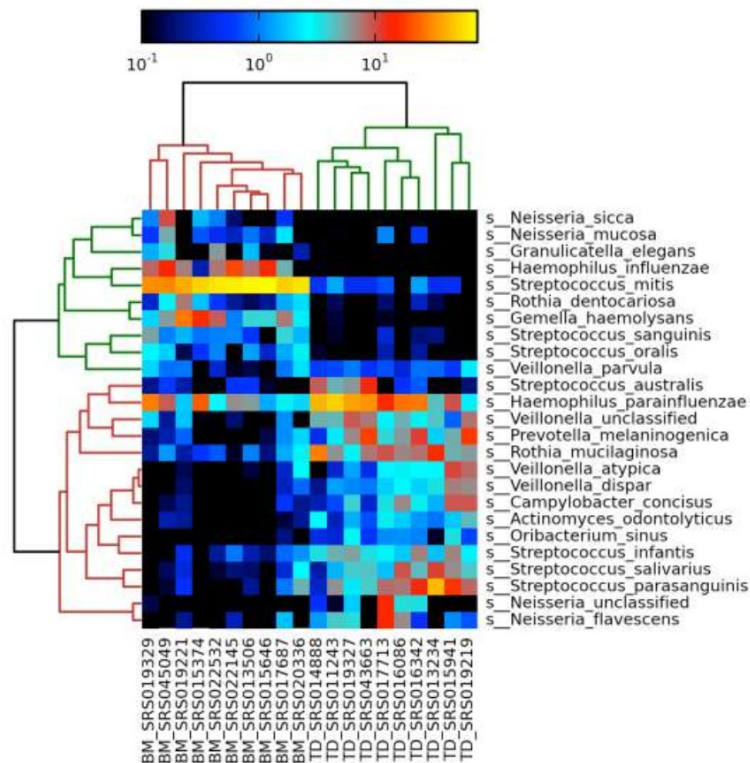
## MetaPhlan: Metagenomic Phylogenetic Analysis

- \* Uses a database of taxon specific marker genes
- \* Works well with known ecosystems: e.g. gut communities

## PhyloSift:

- \* Uses a database of 37 universal proteins & rRNA genes.
- \* Designed to classify using phylogenies

Both databases are smaller than NCBI NR, depending on your ecosystem, one will work better



<https://bitbucket.org/nsegata/metaphlan/>

<http://sourceforge.net/p/krona>

# (Total) MetaGenomics sequencing - tools

**MEGAN** (Huson et al., Genome Research, 2007)

- Developed for characterization of metagenomic shotgun reads
- LCA assignment based on BLAST hit scores
- Support for paired-end reads and comparison of datasets.
- Latest version can analyze RDP files / QIIME OTU files
- Analysis of metabolism via SEED, KEGG or COG maps
- Comparison of multiple metagenomes (> 2)

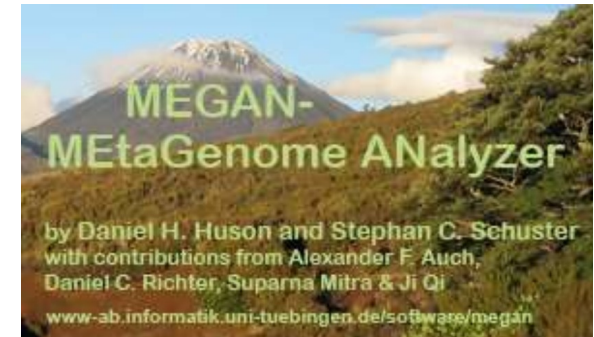
<https://github.com/husonlab/megan-ce>  
<http://megan.informatik.uni-tuebingen.de/>

*Why to use MEGAN ?*

Easy to work with on a desktop / laptop computer:  
Extra things needed: Java, a BLAST server

MEGAN gives a visualization of BLAST results

- Study diversity
- Compare samples
- Contamination filtering
- Special gene of interest
- Extraction of sequences based on taxonomic / metabolic information.





# (Total) MetaGenomics sequencing - tools

## The basics of MEGAN

MEGAN uses BLAST, a database and a taxonomy file

- BLAST N : nucleotides against a nucleotide database.
- BLAST X : Translated nucleotide against a protein database.

- Which database?

one of the many available database like the NCBI-nonredundant database (nr), or a your own custom database.

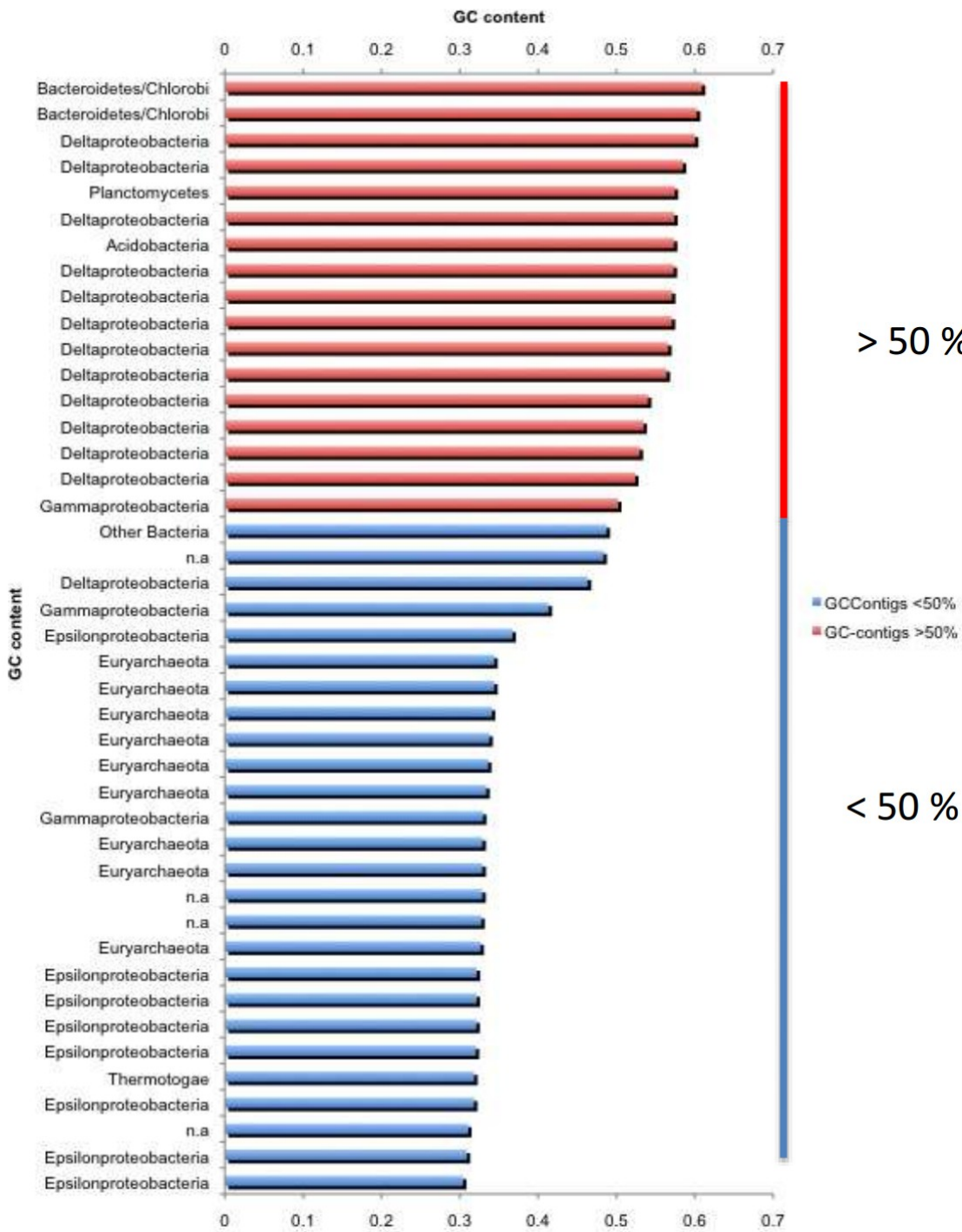
- Taxonomy: NCBI taxonomy, or your own custom taxonomy

- LCA clustering

BLAST output file is used to bin sequences using the LCA (“Lowest Common Ancestor” \*\*) assignment algorithm into specific taxons

*“ In this approach, every read is assigned to some taxon. If the read aligns very specifically only to a single taxon, then it is assigned to that taxon. The less specifically a read hits taxa, the higher up in the taxonomy it is placed. Reads that hit ubiquitously may even be assigned to the root node of the NCBI taxonomy” (MEGAN manual)*

# (Total) MetaGenomics sequencing - tools



MEGAN classifications

Major taxa are separated by different GC content

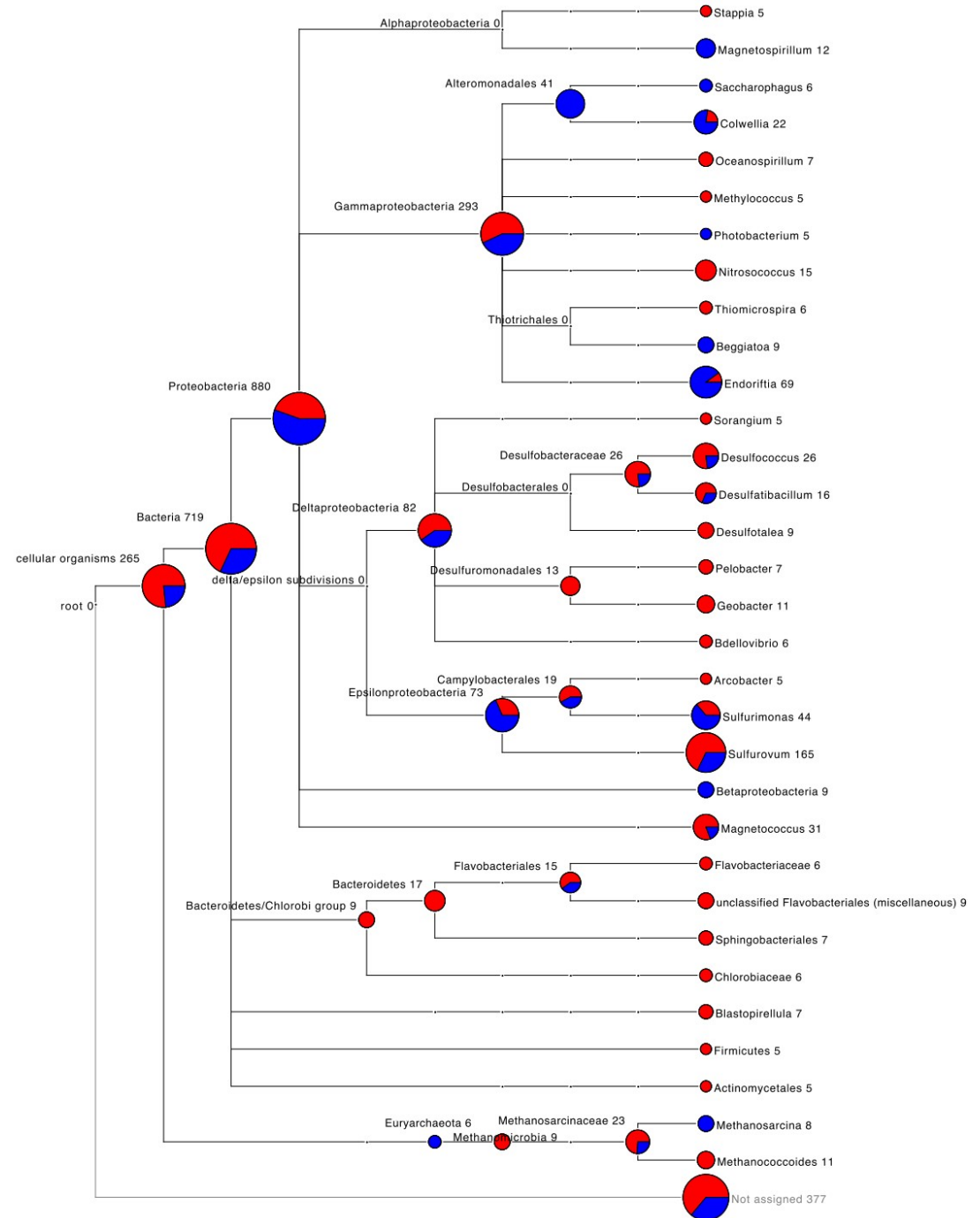
Adaptation to environment

(Kotlar et al., 2011)

# (Total) MetaGenomics sequencing - tools

## multiple samples

Comparison between reads assigned to **Phosphorus** metabolism and **Nitrogen** metabolism



# (Total) MetaGenomics sequencing - tools

## Other tools

- MG-RAST\* (<http://metagenomics.anl.gov/>) (1,2)
- IMG/M (<http://img.jgi.doe.gov/>) (1)
- WebMGA (<http://weizhong-lab.ucsd.edu/metagenomic-analysis/>) (1)
- METAgen assist\* (<http://www.metagenassist.ca/METAGENassist/faces/Home.jsp>) (1,2)
- Real-Time metagenomics (<https://edwards.sdsu.edu/RTMg/>) (1)
- Ribosomal Database Project (RDP) ([rdp.cme.msu.edu](http://rdp.cme.msu.edu)) (2)
- Qiime (Quantitative Insights Into Microbial Ecology) ([www.qiime.org](http://www.qiime.org)) (1)
- Mega (and note “Megane”), more focus on phylogeny (<https://www.megasoftware.net/>) (2)
- Mothur (<https://www.mothur.org/>) (1,2)

(1) *MetaG*

(2) *Barcoding / amplicon sequences*



# (Total) MetaGenomics sequencing - tools

MG -RAST (online tool)

<http://metagenomics.anl.gov/>

**MG-RAST**  
metagenomics analysis server

Browse Metagenomes

Register Contact Help Upload News

**About**

MG-RAST (the Metagenomics RAST) server is an automated analysis platform for metagenomes providing quantitative insights into microbial populations based on sequence data.

# of metagenomes	74,462
# base pairs	23.64 Tbp
# of sequences	218.27 billion
# of public metagenomes	12,322

Updates | MG-RAST Version 3.2 released [May 30, 2012]

\* login required

This project has been funded in part with Federal funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN272200900040C.  
This work was supported in part by the Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, under Contract DE-AC02-09CH11357.

cite MG-RAST

2012

**MG-RAST**  
metagenomics analysis server  
version 4.0.3

398,227 metagenomes containing 1,534 billion sequences and 215.12 Tbp processed for 29,881 registered users.

for programmatic access visit our API site

search string e.g. mgp128 or mgm4447970.3 search

upload download analyze

Report

Turn your raw sequence into analyzed data.

November 2019

# (Total) MetaGenomics sequencing - tools

## MG -RAST

### Metagenome Analysis

**1 Data Type**

ORGANISM ABUNDANCE

Representative Hit Classification

» Best Hit Classification

Lowest Common Ancestor

FUNCTIONAL ABUNDANCE

Hierarchical Classification

All Annotations

OTHER

Recruitment Plot

**2 Data Selection**

Metagenomes 4441147.3

Annotation Sources M5NR

Max. e-Value Cutoff 1e-5

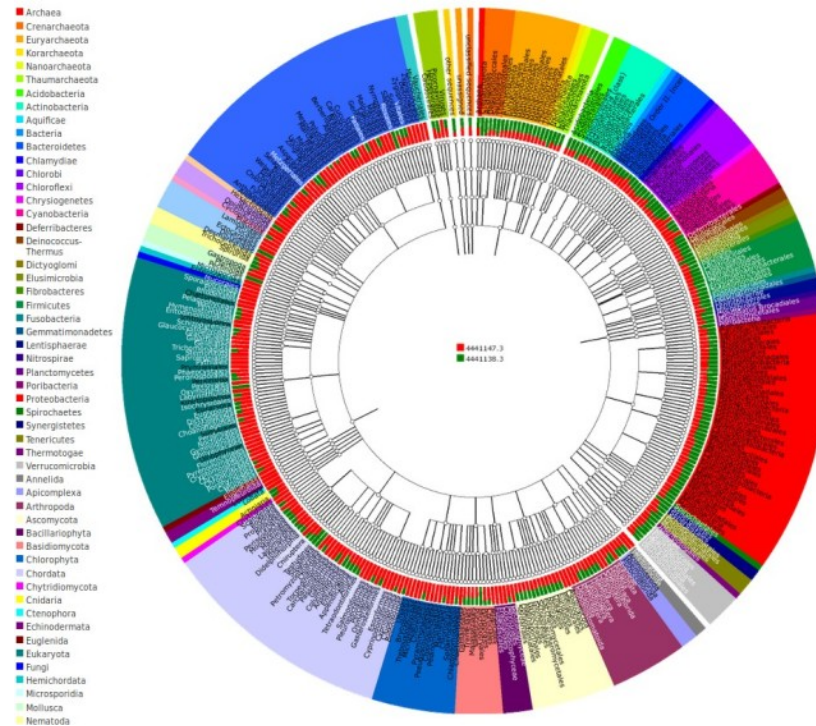
Min. % Identity Cutoff 60 %

Min. Alignment Length Cutoff 15

Workbench  use features from workbench

**3 Data Visualization**

barchart  tree  table  heatmap  PCoA  rarefaction



# (Meta) Barcoding sequencing - tools

## Ribosomal Database Project (RDP)

<https://rdp.cme.msu.edu/>



RDP Release 11, Update 5 :: September 30, 2016

3,356,809 16S rRNAs :: 125,525 Fungal 28S rRNAs  
Find out what's new in RDP Release 11.5 [here](#).



[Cite RDP's latest tool articles.](#)

RDP provides quality-controlled, aligned and annotated Bacterial and Archaeal 16S rRNA sequences, and Fungal 28S rRNA sequences, and a suite of analysis tools to the scientific community. New to RDP release 11:

- RDP tools have been updated to work with the new fungal 28S rRNA sequence collection.
- A new Fungal 28S Aligner and updated Bacterial and Archaeal 16S Aligner. We optimized the parameters for these secondary-structure based Infernal aligners to provide improved handling for partial sequences.
- Updated RDPipeline offers extended processing and analysis tools to process high-throughput sequencing data, including single-strand and paired-end reads.
- Most of the RDP tools are now available as open source packages for users to incorporate in their local workflow.



<https://rdp.cme.msu.edu/help/tutorial.jsp>

[http://rdp.cme.msu.edu/tutorials/init\\_process/RDPtutorial\\_INITIAL-PROCESS.html](http://rdp.cme.msu.edu/tutorials/init_process/RDPtutorial_INITIAL-PROCESS.html)

[https://rdp.cme.msu.edu/tutorials/classifier/classifer\\_cover\\_page.html](https://rdp.cme.msu.edu/tutorials/classifier/classifer_cover_page.html)