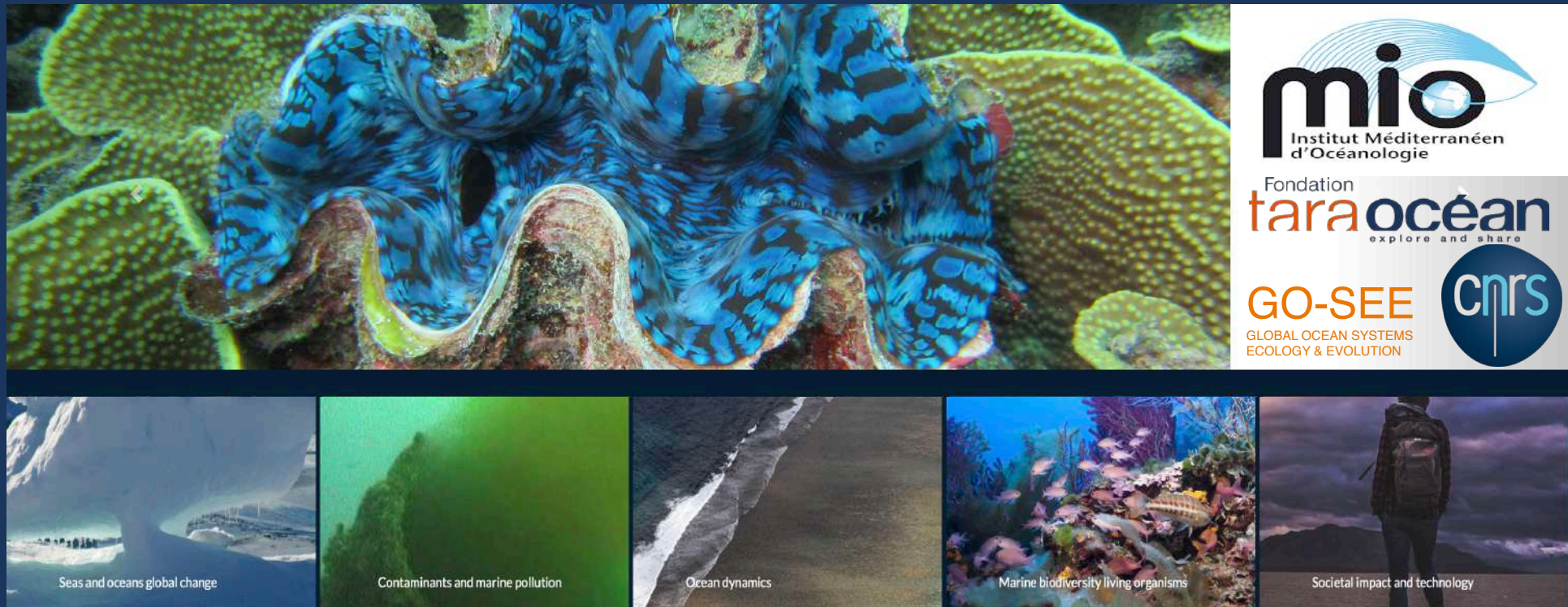


# Marine genomics web services



**Magali Lescot**

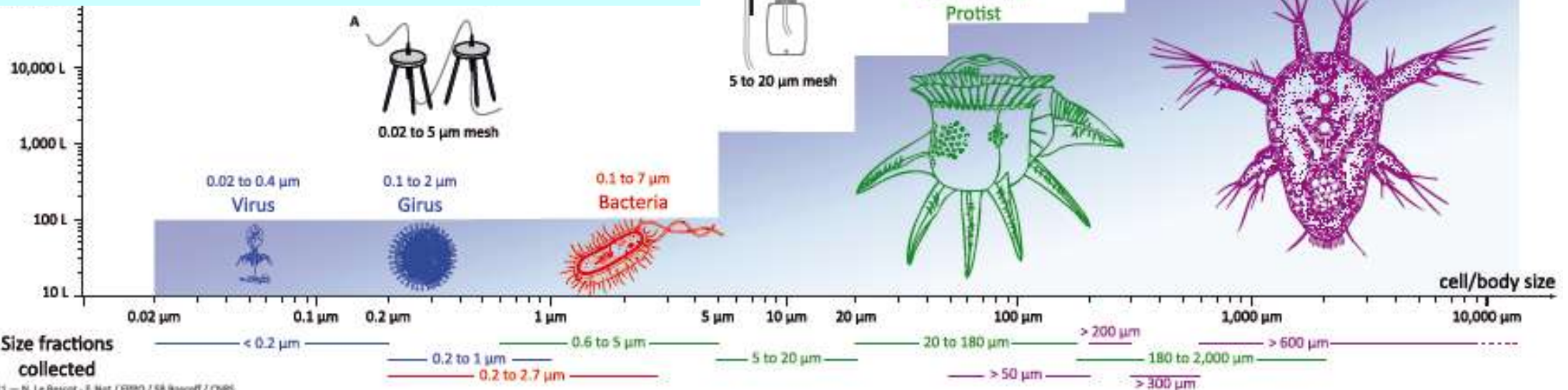
**MIO UMR 7294 – CNRS AMU IRD - Marseille**



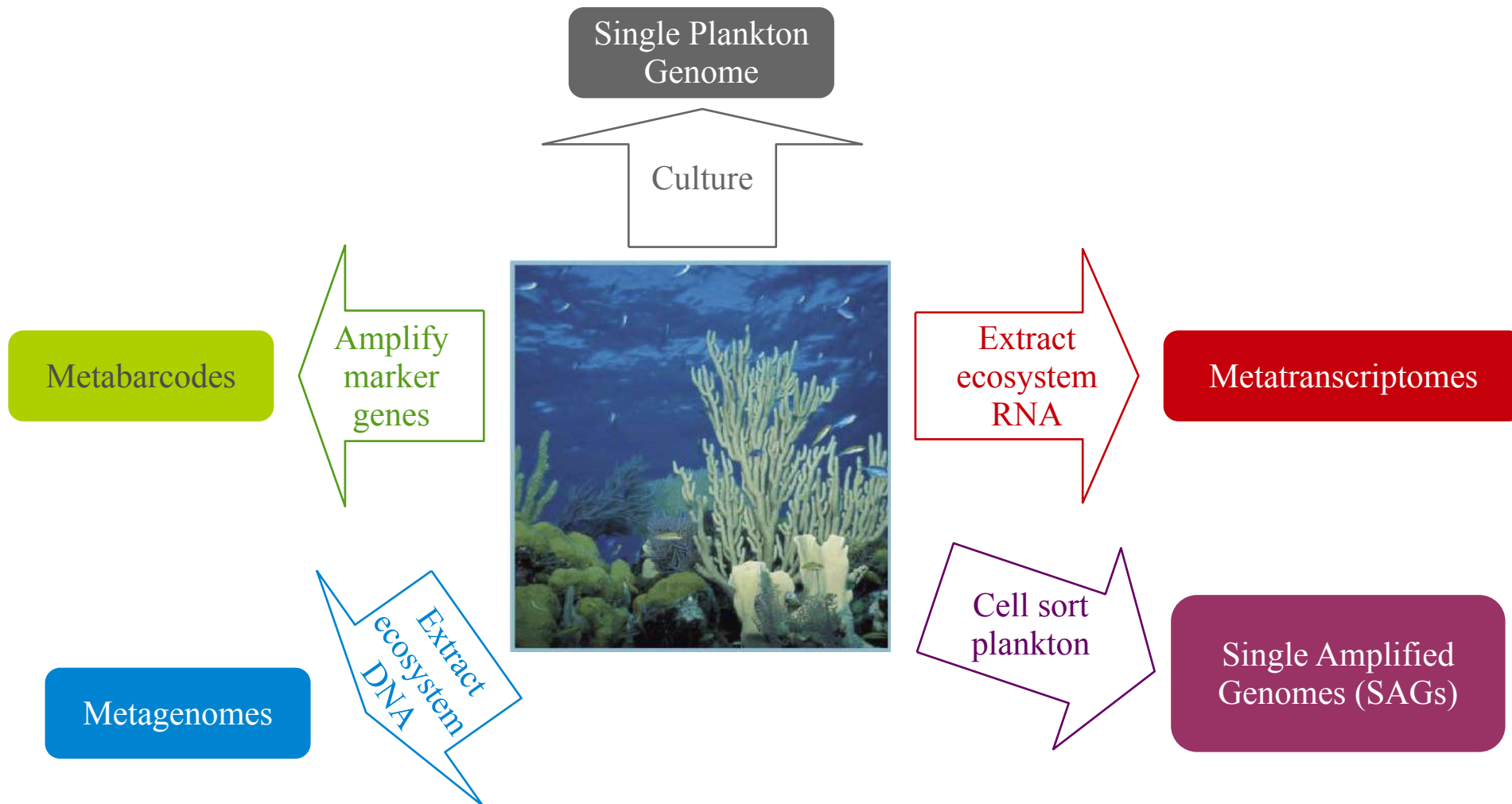


# Sampling strategy

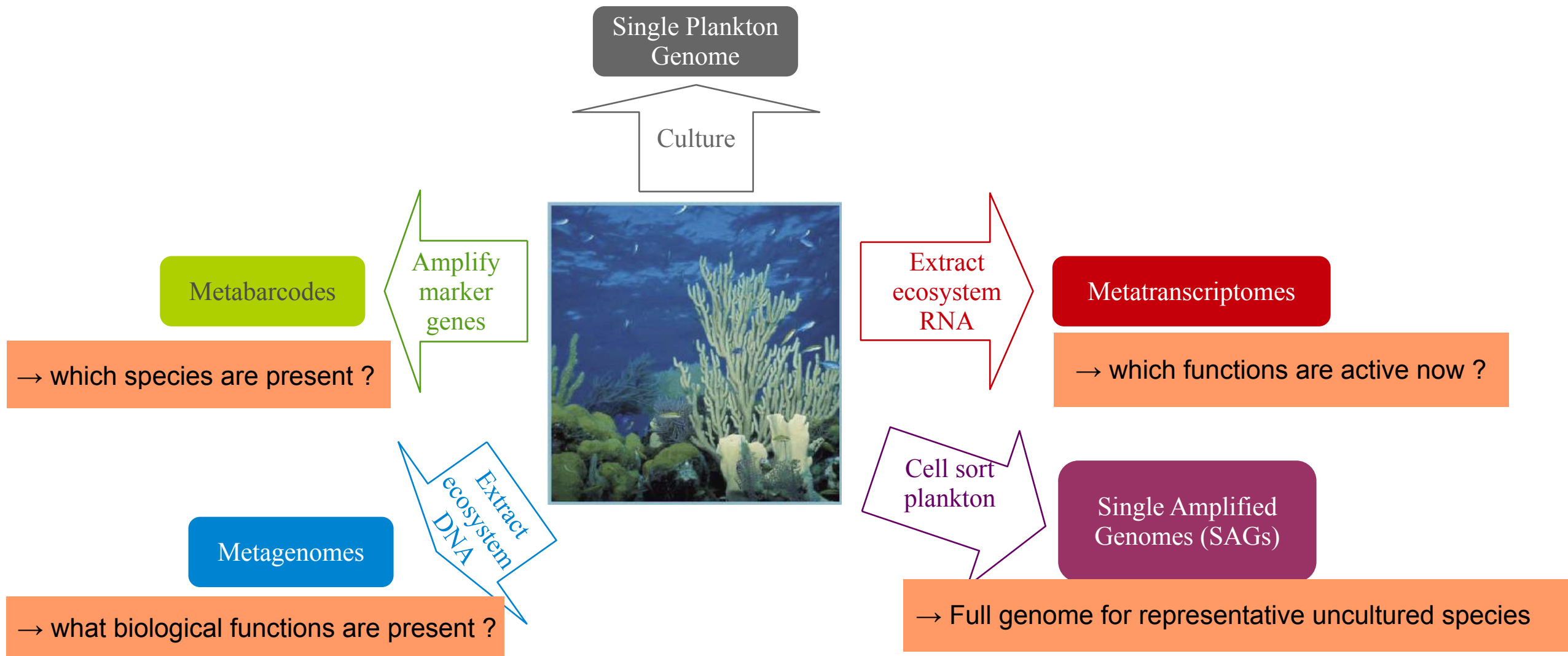
Filter pore sizes	Plankton fraction
$p < 0.2 \mu\text{m}$	1: Viruses (phages)
$0.2 < p < 1.6 \mu\text{m}$	2: Bacteria & Giruses
$0.8 < p < 5 \mu\text{m}$	3: Bacteria & Protists I
$5 < p < 20 \mu\text{m}$	4: Protists II
$20 < p < 180 \mu\text{m}$	5: Protists III
$180 < p < 2000 \mu\text{m}$	6: Protists IV



# Environmental genomics



# Environmental genomics







# OCEAN ATLAS

ONE CLICK MARINE BIOGEOGRAPHY



<http://tara-oceans.mio.osupytheas.fr/>

Explore the biogeography of a gene/protein sequence in a pan-oceanic collection of plankton metagenomes :

**OGA**  
Ocean Gene Atlas

Explore the distribution of plankton taxa across a pan-oceanic collection of metabarcodes :

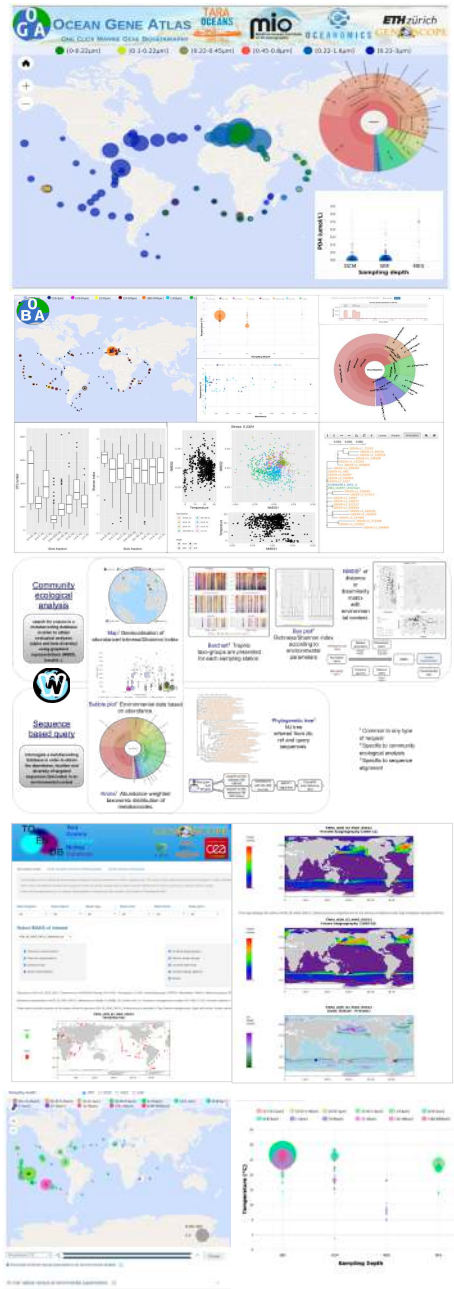
**OBA**  
Ocean Barcode Atlas

Explore the biogeography of shared k-mer ratios for one or multiple DNA sequence(s) :

**ORA**  
Ocean Read Atlas

Select below the SMAG of interest and visualize the effect of climate change :

**TOENDB**  
Tara Oceans Ecological Niches Database



**Ocean Gene Atlas => Ocean Gene Atlas V2.0**

- Marine metagenomic and metatranscriptomic datasets
- Tara Oceans* Prokaryotes and Eukaryotes genes catalogues
- 8 trillion of read sequences for 228 million genes**

**Ocean Barcode Atlas**

- Marine metabarcoding datasets
- 500 000 barcode sequences

**Ocean Read Atlas**

**First indexing method capable of processing terabyte-sized genomic sequence in one dataset!**

- index of 1,393 metagenome samples of raw sequences from virus to mesozooplankton from *Tara Oceans*
- 36.7TB raw data = one index (4.63TB)

**TOENDB**

- *Tara Oceans* Ecological Niches
- 713 Single Cell and Metagenomes
- Assembled Genomes

<http://tara-oceans.mio.osupytheas.fr/>



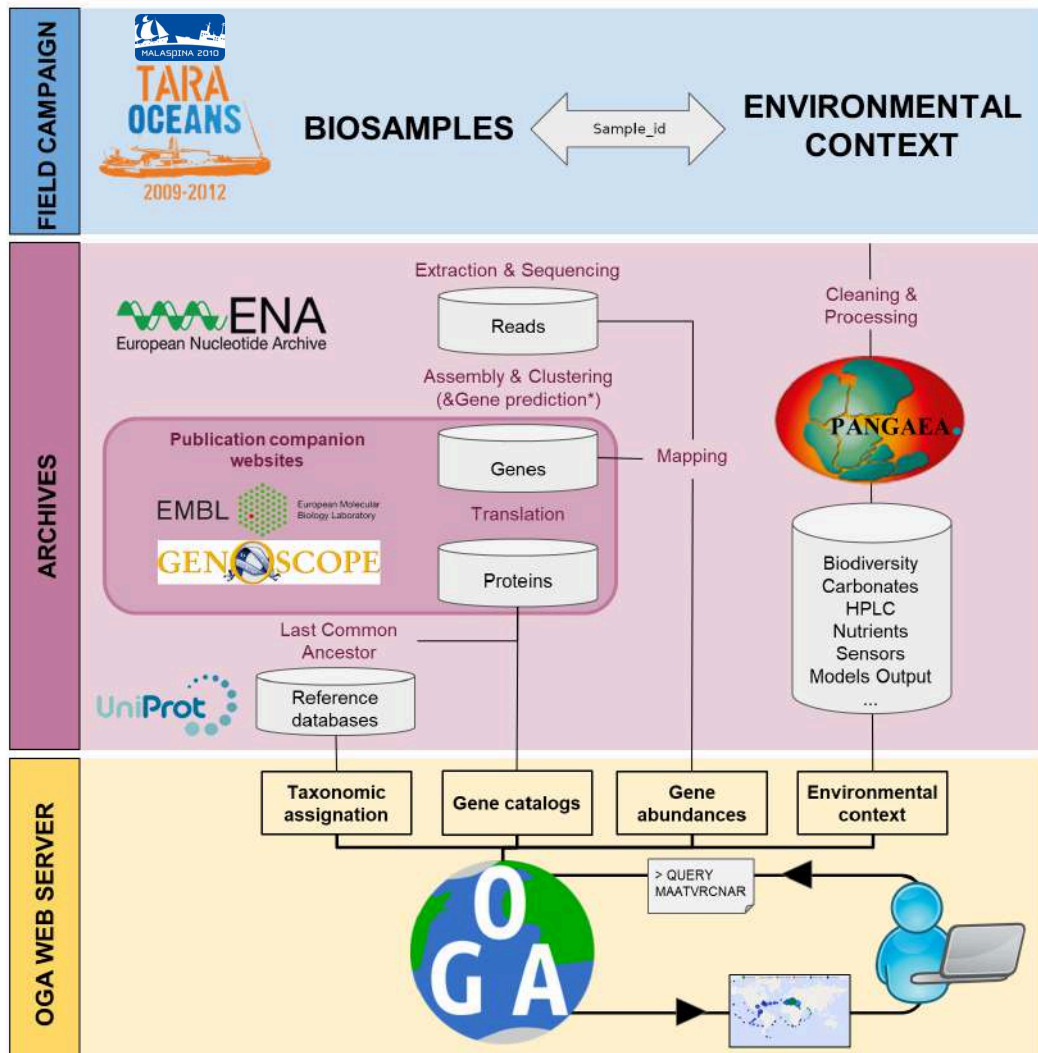




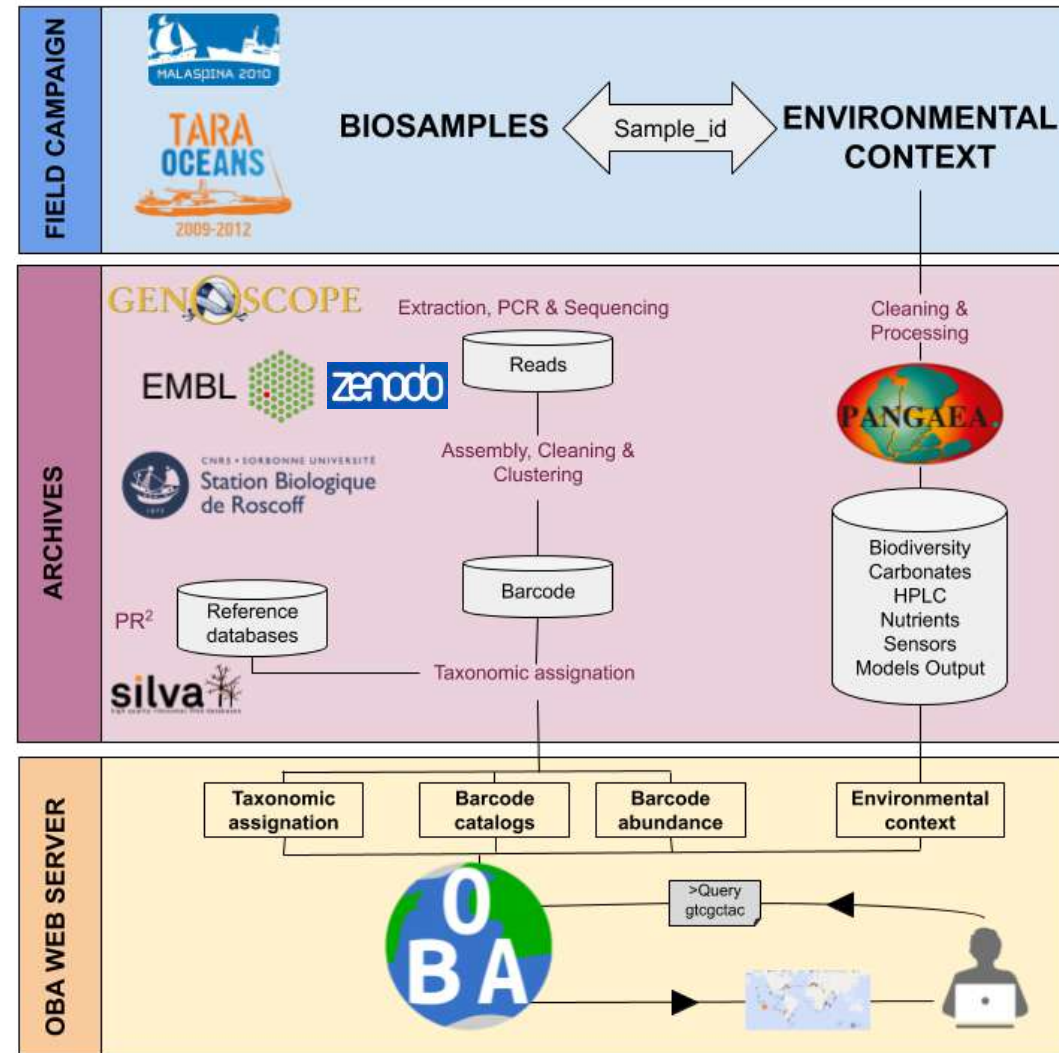
# Ocean Gene Atlas

# Ocean barcode Atlas

## Metagenomic & Metatranscriptomic datasets



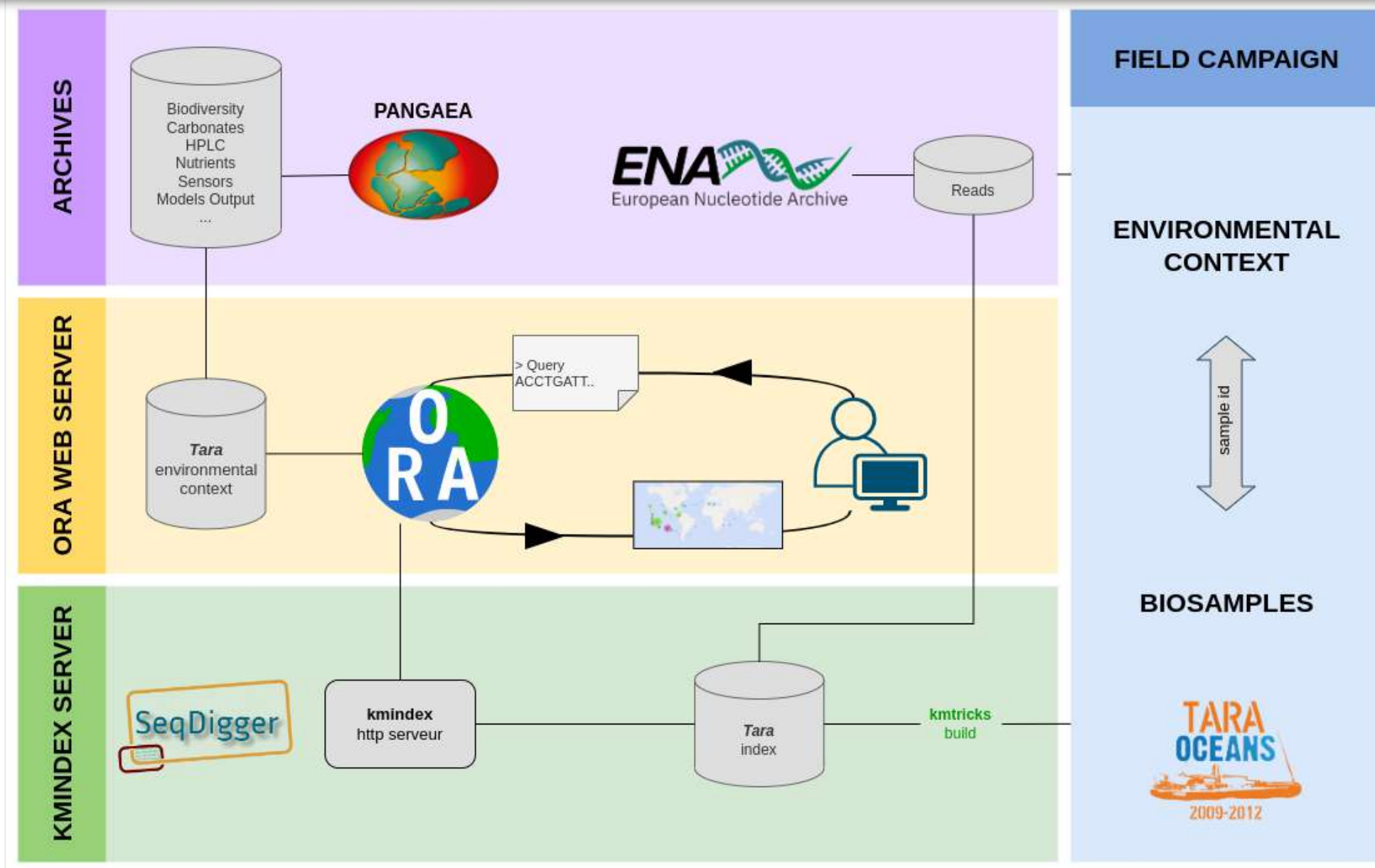
## Metabarcoding datasets







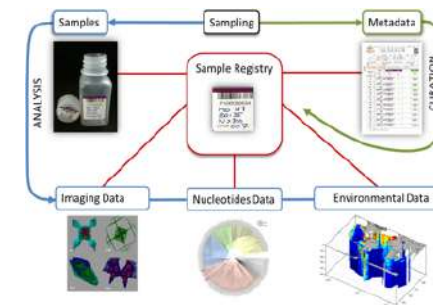
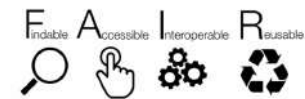
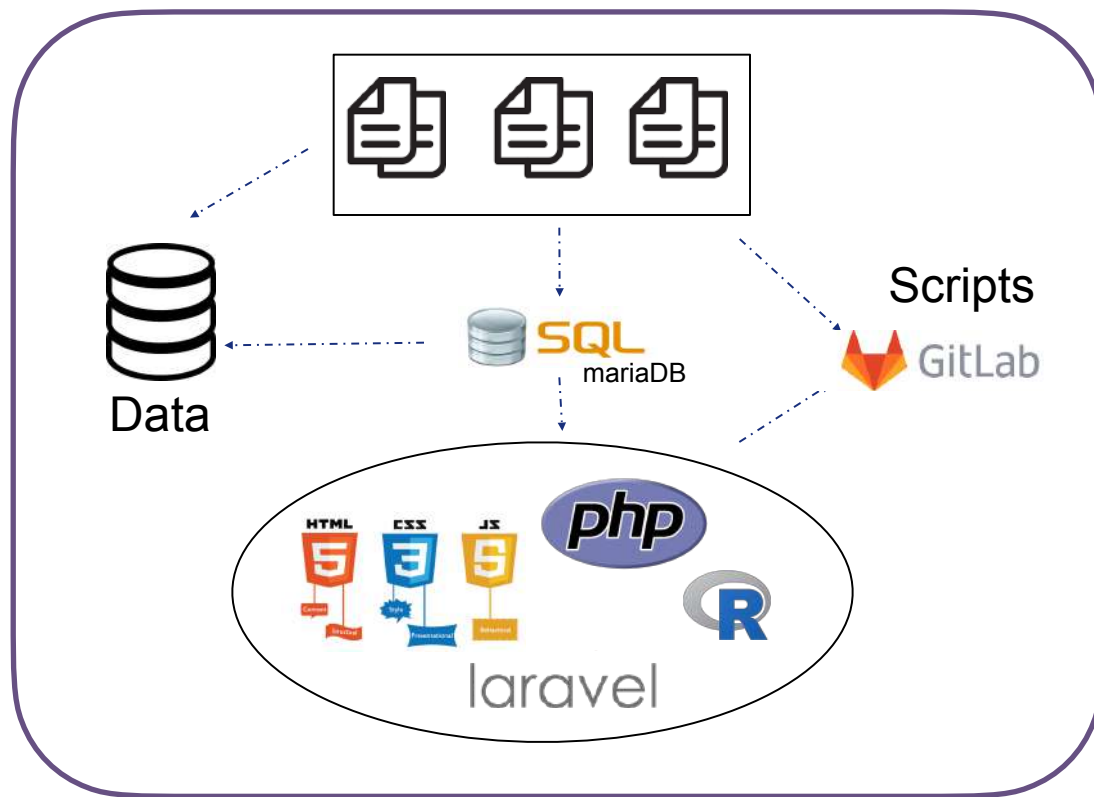
# The Ocean Read Atlas





Nolan Lezzoche

Services web



Metadata

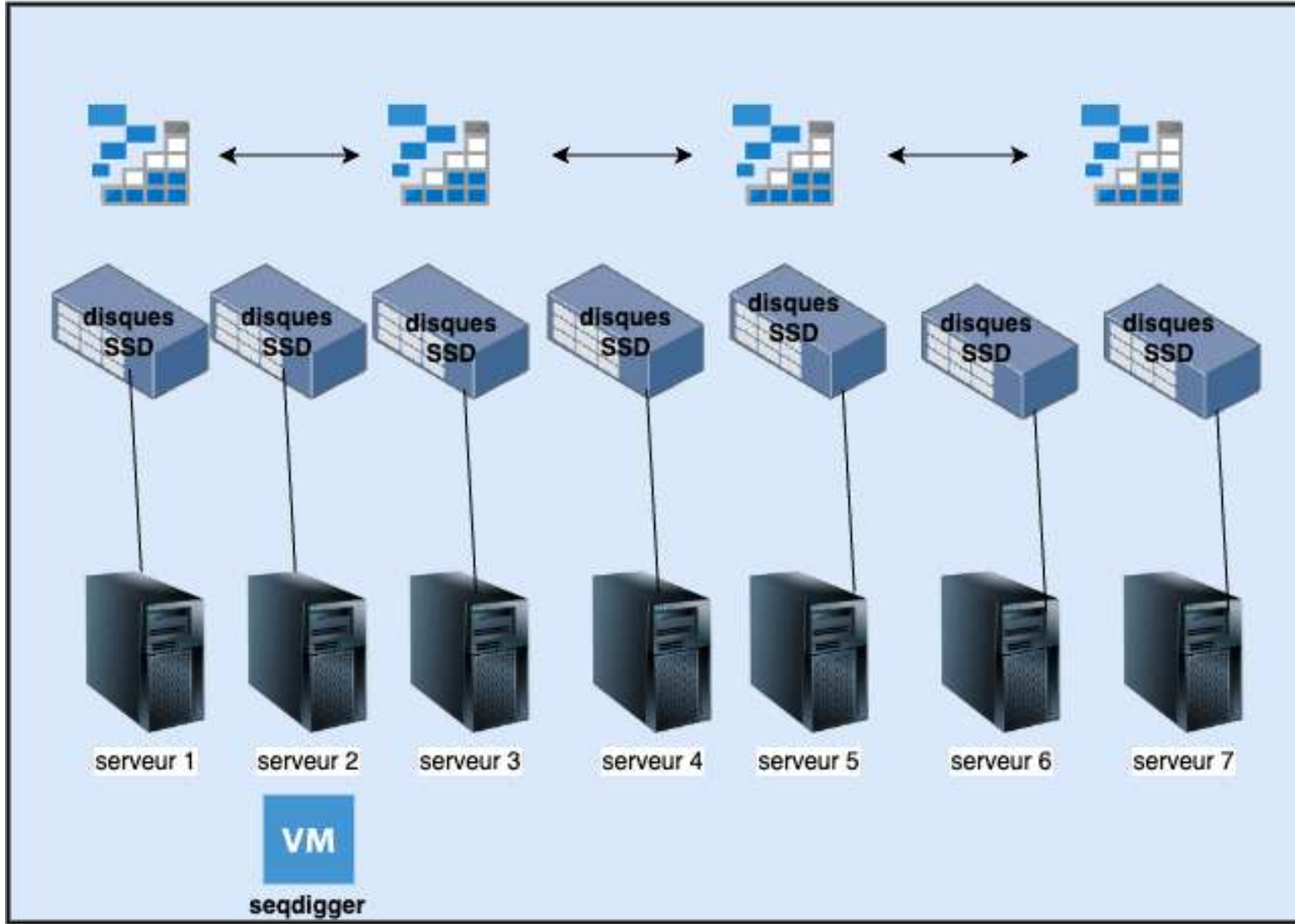




# Infrastructure SeqDigger



Julien Lecubin



Distributed storage system Ceph

virtualisation PROXMOX



# The Ocean Gene Atlas v2.0

Metagenomic (MetaG)  
Metatranscriptomic (MetaT)

ONE CLICK MARINE GENE BIOGEOGRAPHY

**OCEAN GENE ATLAS**  
ONE CLICK MARINE GENE BIOGEOGRAPHY

**if you use this web service, please cite:**  
 The Ocean Gene Atlas v2.0: online exploration of the biogeography and phylogeny of plankton genes. 10.1093/nar/gkac420  
 The Ocean Gene Atlas: exploring the biogeography of plankton genes online. 10.1093/nar/gky376.

Submit your gene or protein sequence below to:

- display its abundance on ocean maps
- observe its co-variation with marine environmental features (\* CO<sub>2</sub>, nutrients etc.)
- explore its taxonomic distribution in seawater

Data mined from Tara Oceans and Malaspina metagenomes & metatranscriptomes (user manual) - Credits

- Try an [Example](#) with OM-RGC dataset (prokaryotes)
- Try an [Example](#) with MATOU dataset (eukaryotes)
- Try an [Example](#) with EUK\_SMAGs dataset (eukaryotes)
- Try an [Example](#) with MGT dataset (eukaryotes)
- Try an [Example](#) with Arctic\_MAGs dataset (prokaryotes)
- Try an [Example](#) with BAC\_ARC\_MAGs dataset (prokaryotes)
- Try an [Example](#) with MDeep-MAGs dataset (prokaryotes)

Dataset: OM-RGCv1 - Tara Oceans Microbiome Reference

Job title:

Sequence type: Protein Nucleotide

Either, query sequence: Paste your fasta sequence here

Phylogenetic tree

or HMM file: Choisir le fichier: aucun fichier sélectionné

or results file: Choisir le fichier: aucun fichier sélectionné

or Pfam ID: PF00111

or (un)genes name list: OM-RGC.v1.009423385

Search method: blastp

Expect threshold: 1E-10

Abundance as: percent of total reads

Maps: 2

Bubble plots: 2

Email: Optional

Reset Submit

Nucleic Acids Research, 2022, 1  
<https://doi.org/10.1093/nar/gkac420>

Nucleic Acids Research, 2018, 1  
[doi: 10.1093/nar/gky376](https://doi.org/10.1093/nar/gky376)

**The Ocean Gene Atlas v2.0: online exploration of the biogeography and phylogeny of plankton genes**

Caroline Vernet<sup>1,2</sup>, Julien Lecubin<sup>3</sup>, Pablo Sánchez<sup>4</sup>, Tara Oceans Coordinators, Shinichi Sunagawa<sup>5</sup>, Tom O. Delmont<sup>2,4</sup>, Silvia G. Acinas<sup>3,4</sup>, Eric Pelletier<sup>2,4</sup>, Pascal Hingamp<sup>1</sup> and Magali Lescot<sup>1,2,\*</sup>

**The Ocean Gene Atlas: exploring the biogeography of plankton genes online**

Emilie Villar<sup>1,2,\*</sup>, Thomas Vannier<sup>2</sup>, Caroline Vernet<sup>2</sup>, Magali Lescot<sup>2</sup>, Miguelangel Cuenca<sup>2</sup>, Aurélien Alexandre<sup>2</sup>, Paul Bachelier<sup>2</sup>, Thomas Rosnet<sup>2</sup>, Eric Pelletier<sup>1</sup>, Shinichi Sunagawa<sup>2</sup> and Pascal Hingamp<sup>2,\*</sup>

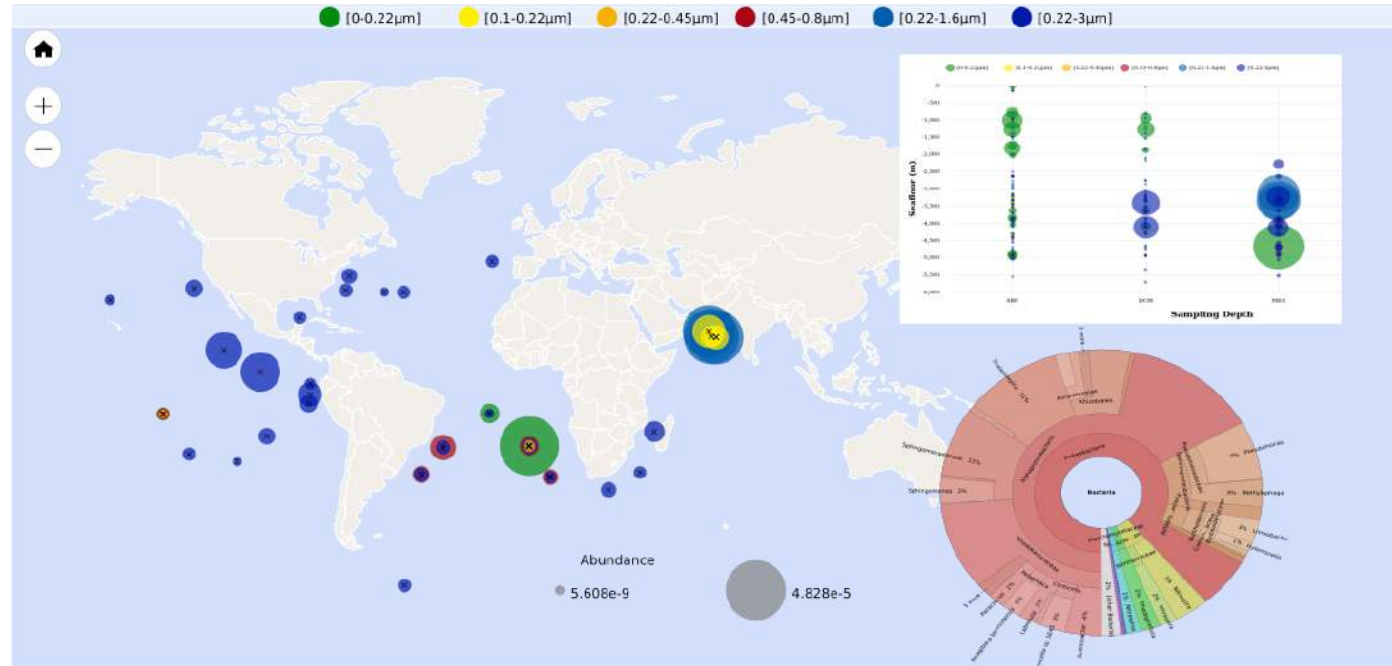
Dataset	Read number	MAG number	Gene number	Sample number
OM-RGCv1	7.20E+12	-	40,154,822	243
OM-RGCv2_metaG	1.13E+11	370	46,775,154	180
OM-RGCv2_metaT	5.00E+09	187		187
MATOU_metaG	1.85E+11	-	116,849,350	445
MATOU_metaT	8.70E+10			440
MGT	5.80E+07	924	6,946,068	364
EUK_SMAGs	2.80E+11	713	10,207,435	939
BAC_ARC_MAGs	2.80E+11	1,888	4,567,982	922
Arctic_MAGs_metaG	1.40E+08			68
Arctic_MAGs_metaT	4.50E+07	530	1,033,381	53
MDeep-MAGs	6.49E+08	317	867,795	58
<b>Total</b>	<b>8,15E+12</b>	<b>4,929</b>	<b>227,401,987</b>	



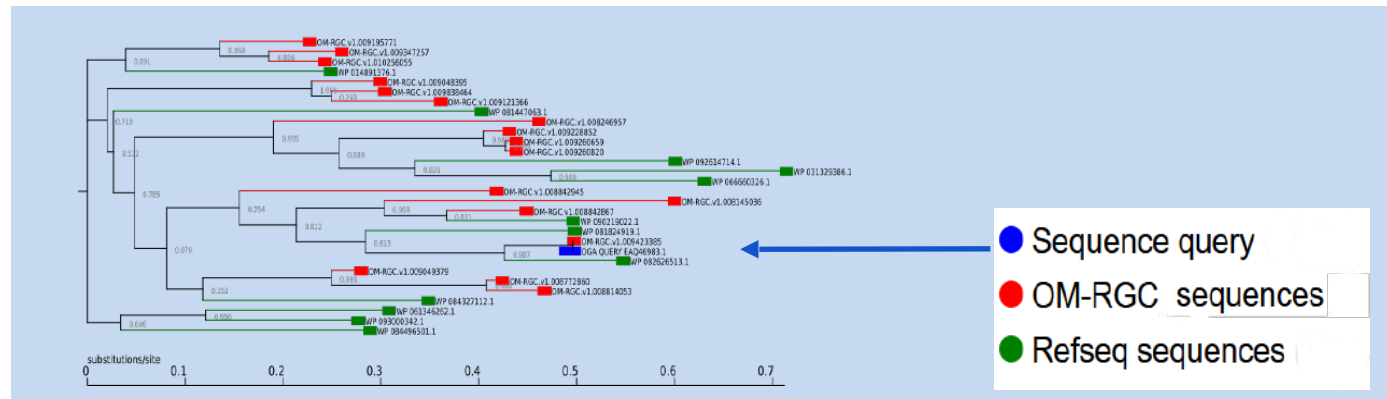


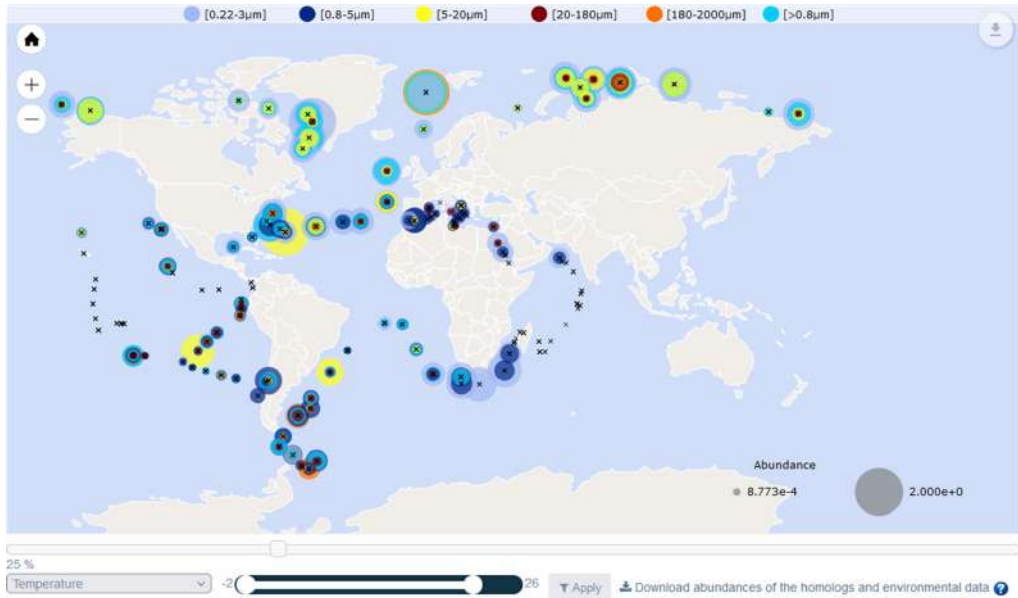


# Biogeography of homologous sequences



- Search for a gene or protein in the different samples (sequence alignment or hmm search)
- Computation of the homologous genes abundance in each sample (/ station and / filter)
- Correlation with environmental parameters
- Taxonomic assignment of identified sequences
- Phylogenetic tree of the homologous sequences





Phylogenetic tree : marine environmental genes in context of reference sequences ?

View multiple alignment ?

- Radial / Linear tree
- Root on the longest branch
- Remove tree root
- Add branch length values
- Increase / Decrease leaf size, Hide / Show leaf labels
- Grow / Shrink tree size
- Reset (cancel all changes)

- Download tree in svg format
- Download sequences in fasta format
- Download full fasta alignment
- Download the intermediate column cleaned alignment (Trimal) in fasta format
- Download the final cleaned alignment (MaxAlign) in fasta format
- Download output from second alignment curation step (MaxAlign)
- Download tree in newick format

Root the tree on the following leaf:  Root

FastTree option's :

WAG substitution model  JTT substitution mode

Gamma20 distribution  no Gamma distribution  ?

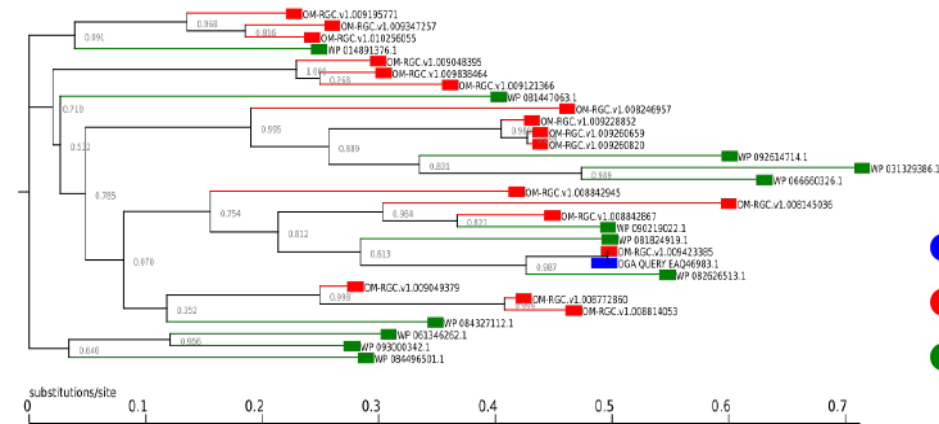
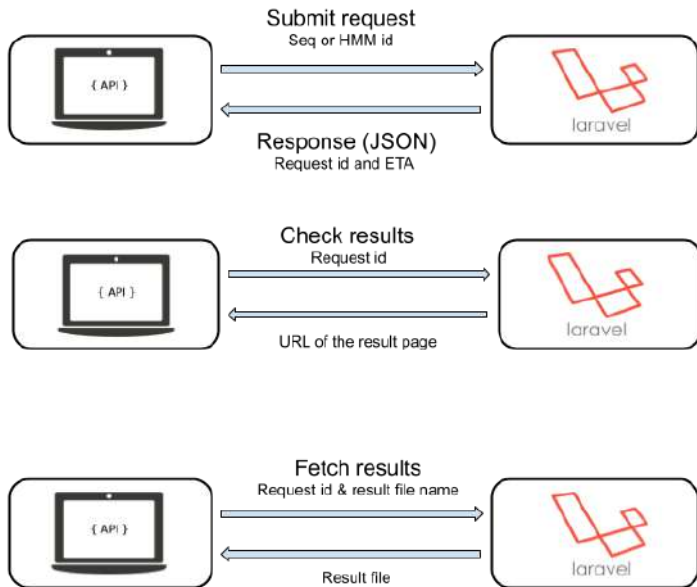
(It may take several minutes)

Number of sequences from **SMAGs** : 77 (4 sequence(s) were/was excluded during the maxalign step)

Number of sequences from **RefSeq** after clustering: (8 sequence(s) were/was excluded during the maxalign step)

Clustering at identity:	-	100%	95%	90%	85%	80%	75%	70%	65%	60%
Number of sequences:	1254	1234	891	685	545	399	280	189	112	64

## API: Application and Programming Interface



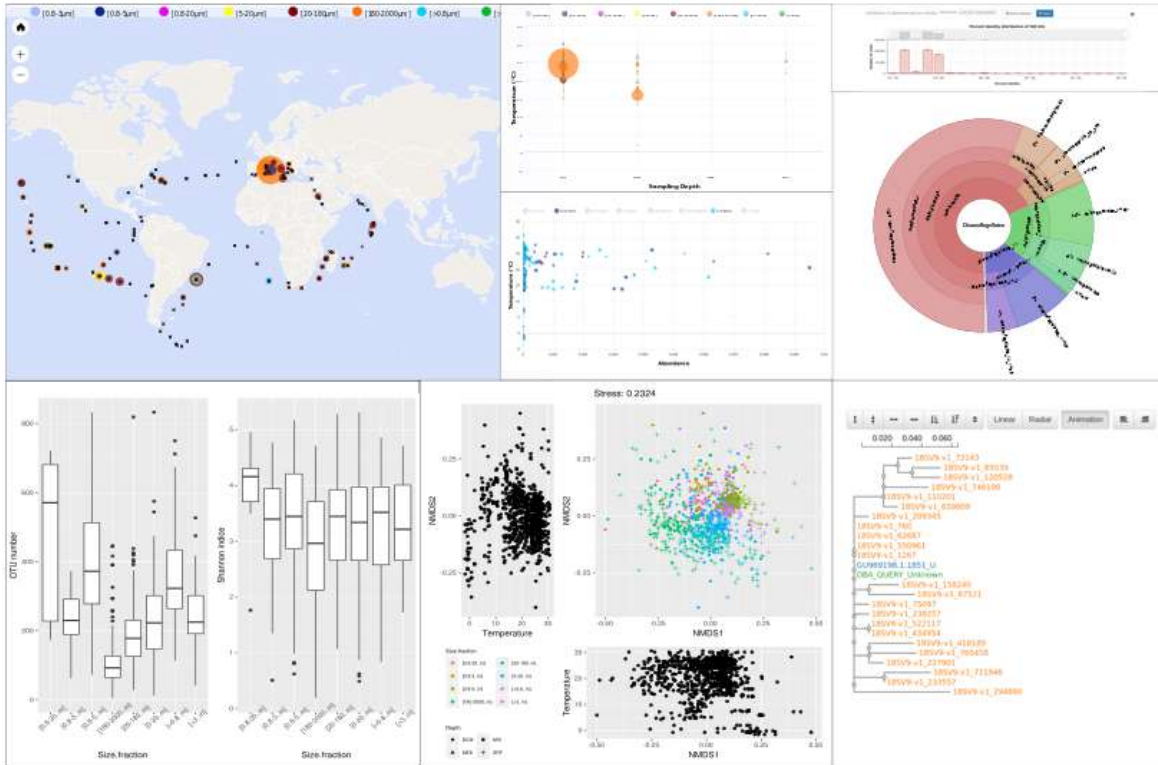
- Sequence query
- OM-RGC sequences
- Refseq sequences





# The Ocean Barcode Atlas

## Metabarcoding datasets



<http://oba.mio.osupytheas.fr/ocean-atlas/>

*The Ocean Barcode Atlas: a web service to explore the biodiversity and biogeography of marine organisms.*

C. Verne, N. Henry, J. Lecubin, C. de Vargas, P. Hingamp, M. Lescot (2020)  
*Molecular Ecology Resources*. Link: <https://doi.org/10.22541/au.160193452.23998228/v1>

## Diversity analysis of planktonic communities

Public

Private

**18S V9**  
474,303 OTUs  
1,046 samples

**ASV 16S V4V5**  
1,361,502 ASVs  
1,134 samples

**18S V4**  
156,648 OTUs  
1,191 samples

**MOOSE 18S V4**  
15,743 OTUs  
271 samples

**16S V4**  
3,902 OTUs  
60 samples

**ASV 18S V4**  
1,011 samples

**ASV 18S V9**  
1,069 samples

**miTAGs**  
23,987 sequences  
180 samples





# Request types

Community ecological analysis

Sequence based query

Submission interface

User friendly web service

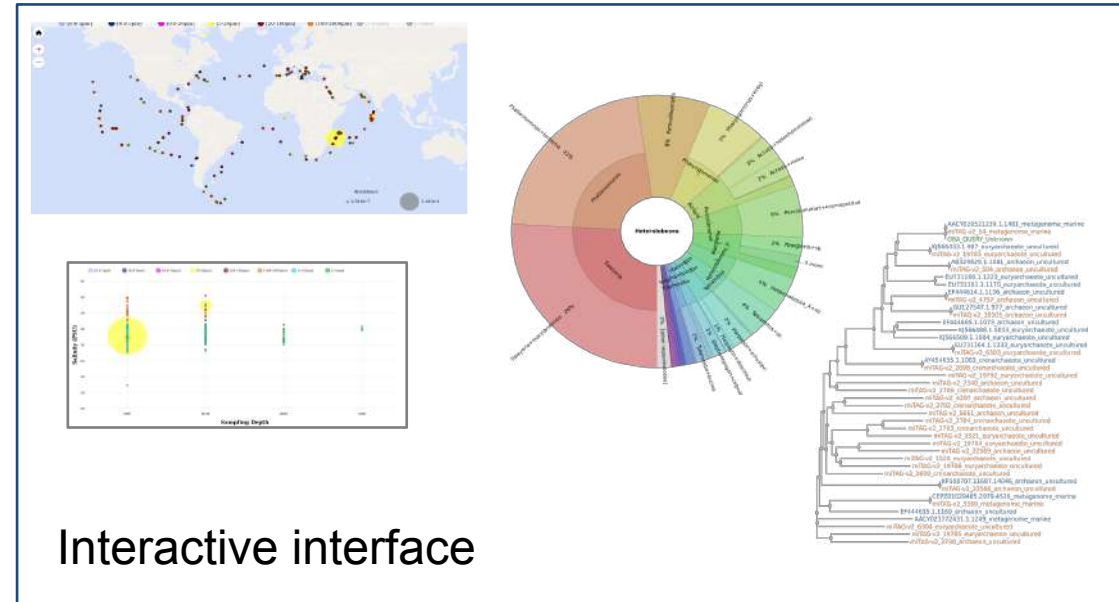
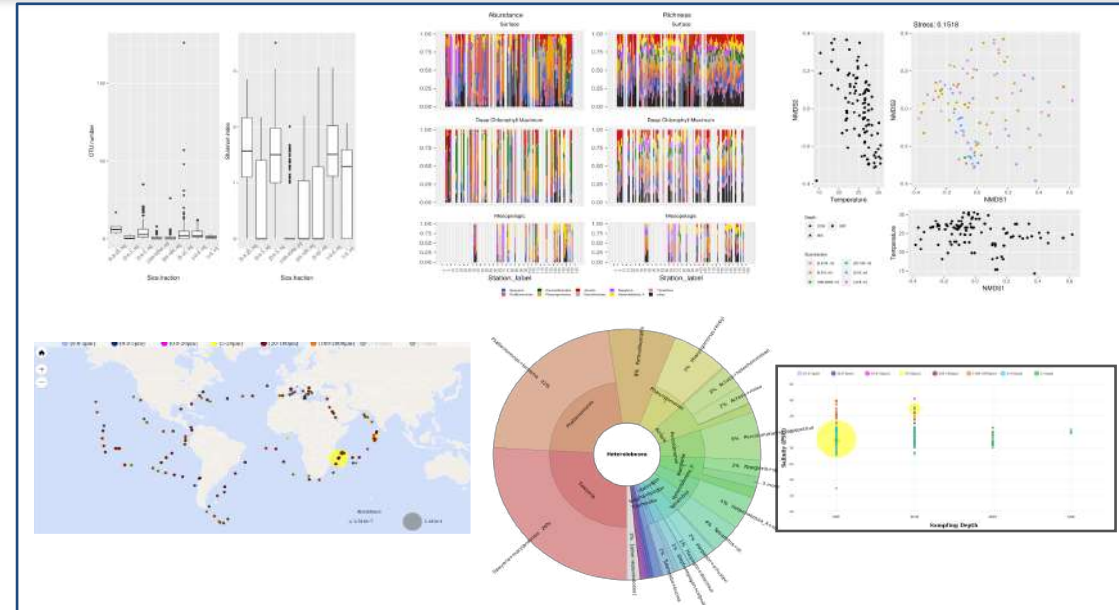
Search by taxonomy

Submit your sequence

Submit from a ref db

Search from an ID

Alignment



Interactive interface



# Data

## Metagenome sequencing (soil, seawater, waste water, air, gut, ...)



### Raw sequences

- **Redundant** & fragmented data
- Error-prone (0.1% to 10% error rate)
- Important **background noise**
- Heterogenous
  - Quality and quantity
- **Volumes:**
  - hundreds millions fragments / experiment
  - Millions of experiments



- **Archived**

**ENA** (European Nucleotide Archive) contains almost **50 Petabytes** of raw data (raw sequences)

=> equivalent to **100,000 desktop computer**

Query this data with a sequence « ATGAGAAAAGTAGCAATTTACGGAAAAGGC » is IMPOSSIBLE!

The cumulative total size of the compressed files representing these genomes is 125 Gigabytes.

That's around **500 thousand times smaller than the 50 Petabytes** of raw sequencing data (also compressed).

=> querying these genomes would take **several hundred days**.

We know how to query assembled genomes.

**But we don't know how to query the raw data 2 by 2.**

=> **2 weeks at TGCC on Tara Oceans datasets**

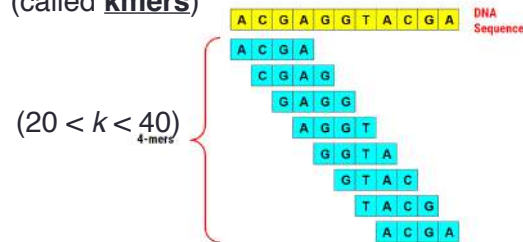
=> so we can extend to 2000 weeks = **38 years of queries on "The Sequence Read Archive"**

---

# kmers

## Words

- **No word in DNA**
- Split to subsequences of fixed length  $k$  (called **kmers**)



- Thousand billions distinct kmers
  - (*google indexes millions*)

## Compare sequences

- Sequence similarity  $\sim$  shared kmers count

ACGAGG <u>T</u> ACGA	ACGAG <u>I</u> TACGA
ACGA	ACGA
CGAG	CGAG
GAGG	GAGT
AGGT	AGTT
GGTA	GTTA
GTAC	TTAC
TACG	TACG
ACGA	ACGA

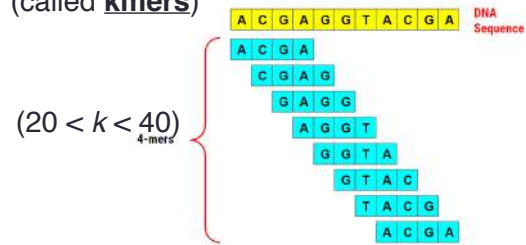
- 4 over 8 kmers shared



# kmers

## Words

- **No word in DNA**
- Split to subsequences of fixed length  $k$  (called **kmers**)



- Thousand billions distinct kmers
  - (*google indexes millions*)

## Query vs Bank

- Sequence similarity  $\sim$  shared kmers count

ACGAGGTACGA

BANK

ACGA

CGAG

GAGG

AGGT

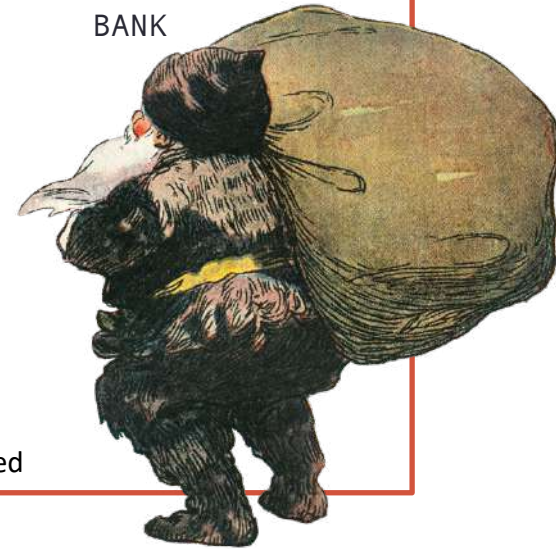
GGTA

GTAC

TACG

ACGA

- 6 over 8 kmers shared



# Indexing: conceptual view

## One read set:

- Extract & count kmers
- Filter kmers
- Generate a [counting] bloom filter with kmtricks

## N read sets:

- Create N [counting] bloom filters
- This is the index

```
Reads  
>read1  
ACGAG...ACGTA  
>read2  
ACGGC...GGACT  
...  
>read1000000  
GGCGA...AGATA
```

```
Counted  
kmers  
AAAAAC 12  
ACCATA 4  
AGGTAT 1  
...  
TCGGAT 5
```

```
cBloom  
Filter  
0  
12  
4  
...  
0
```

```
Reads  
Reads  
Reads  
Reads  
>  
Z >  
... C >  
>read1  
ACGAG...ACGT  
...  
Z >  
7 >  
>read1000000  
GGCGA...AGAT
```

```
cBloom  
Filters  
0 8 3 8  
12 0 13 0  
4 7 6 0  
... ..  
0 24 2 9
```

Probabilistic data structure

# Kmindex & Ocean Read Atlas

**kmindex** and **ORA**: indexing and real-time user-friendly queries in terabyte-sized complex genomic datasets

Téo Lemane<sup>1,\*</sup>, Nolan Lezsoche<sup>2</sup>, Julien Lecubin<sup>3</sup>, Eric Pelletier<sup>4,5</sup>, Magali Lescot<sup>2,5</sup>, Rayan Chikhi<sup>6</sup>, and Pierre Peterlongo<sup>1,\*</sup>



**Input:** *Tara Oceans* Metagenomic

Compressed fastq.gz files:  
36.7 TB, 1,393 samples

Each sample:  
Position  
Species fraction sizes  
Physico-chemical env:  
Ph, salinity, T°, ...



Tara Schooner - Creative Commons Attribution 3.0





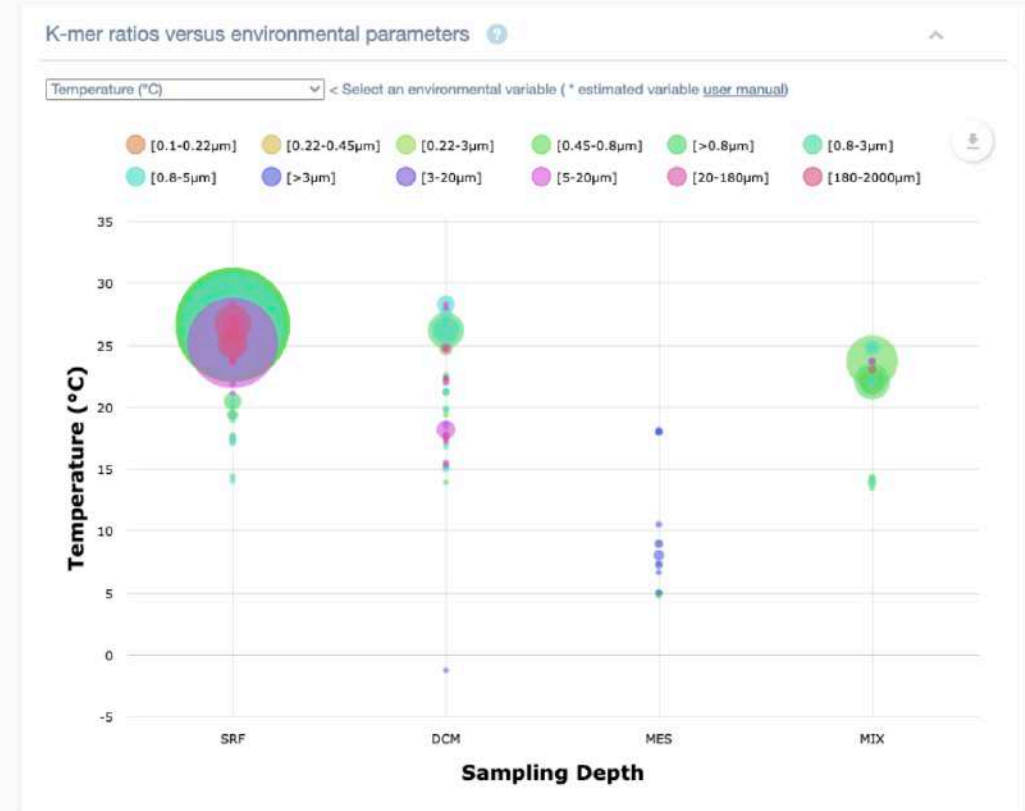
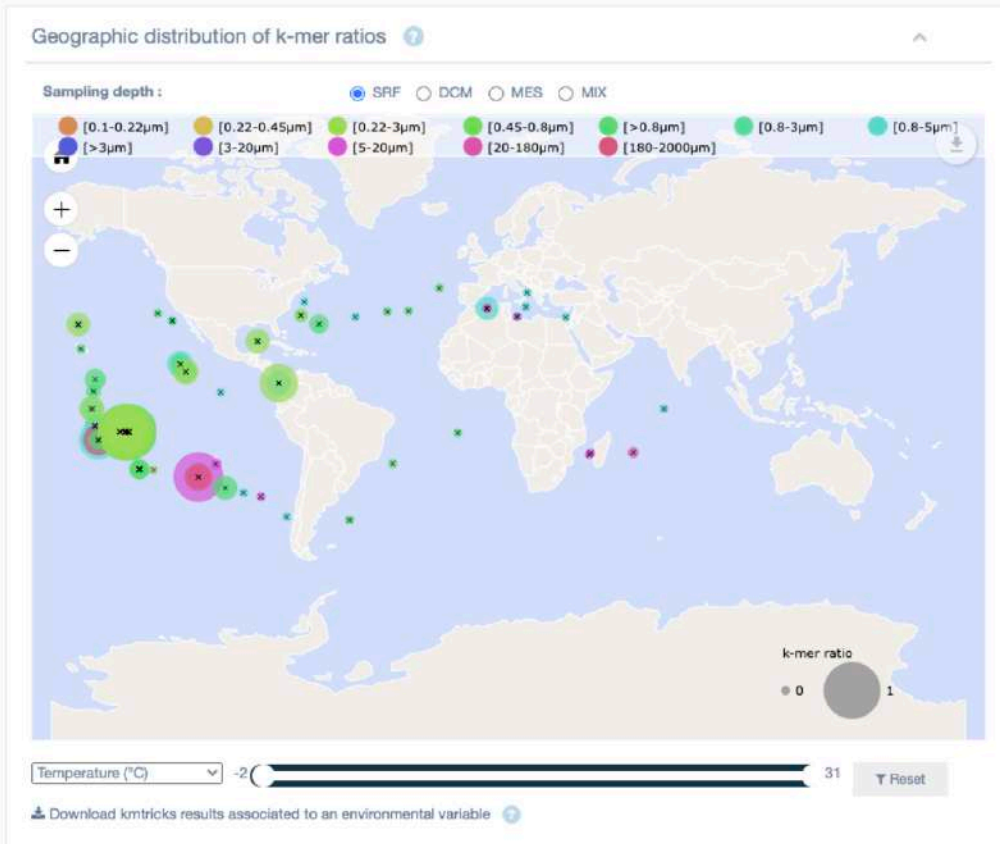
# The Ocean Read Atlas - Like OGA for metagenomic reads

	Build index				Query time		FP rate (%)	
	Time	RAM GB	Disk GB	Index size GB	Nb. queried reads 1	10 million	Average	Max
MetaProFi	30h15	278	5,684	226	12s72	1h29	11.18	21.55
COBS	26h30	278	5,684	184	1s51	15h56	13.29	24.60
<b>kmindex</b>	<b>2h56</b>	<b>107</b>	<b>878</b>	<b>164</b>	<b>0s06</b>	<b>4m21s</b>	<b>0.006</b>	<b>0.18</b>

Index construction and read query performance of kmindex on 50 *Tara* Oceans samples.

# The Ocean Read Atlas - Like OGA for metagenomic reads

- Exploring the biogeography of all sequences from *Tara* Oceans metagenomes (1393 samples)
- Limited to presence/absence query



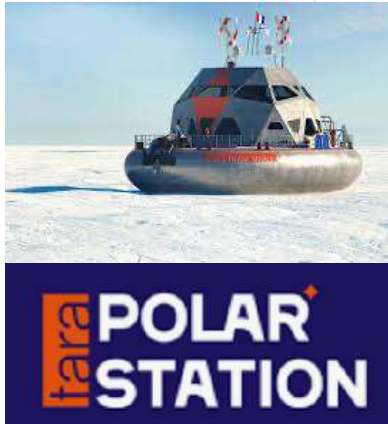
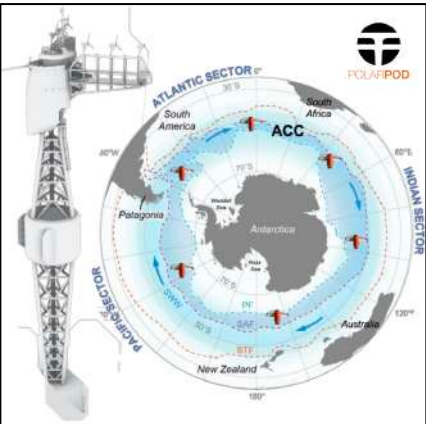




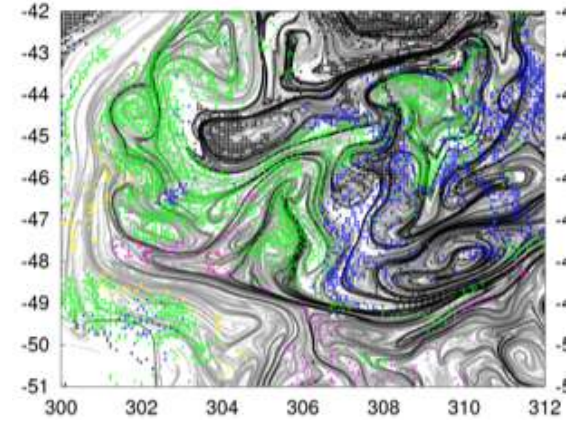
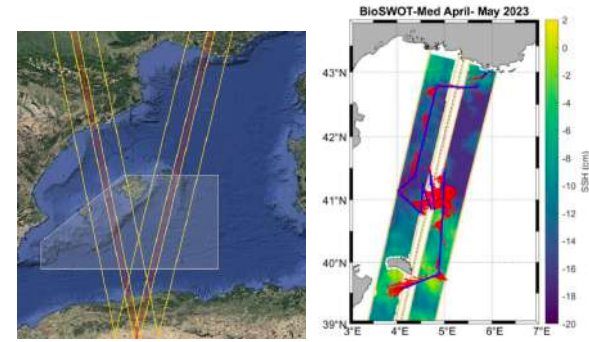
# TARA OCEANS



## BioSWOT Med



# tara POLAR STATION



- Sequence abundance included soon
- Scale (from TB to PB)
- Reduce (environmental impact)
- Smart (queries and answers)
- Deploy (distribute indexes)
- Apply to all sequencing projects (environment, health, agronomy)



**mio**  
Institut Méditerranéen  
d'Océanologie

Fondation  
**tara océan**  
explore and share

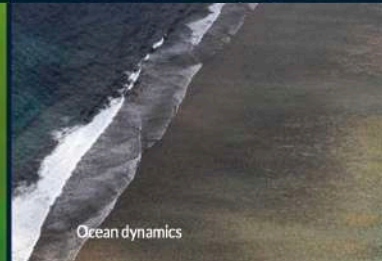
**GO-SEE**  
GLOBAL OCEAN SYSTEMS  
ECOLOGY & EVOLUTION



Seas and oceans global change



Contaminants and marine pollution



Ocean dynamics



Marine biodiversity living organisms



Societal impact and technology

<http://tara-oceans.mio.osupytheas.fr>