



Abims

04/10/2014

Bioinformatique microbienne *IFB & FG variants*

Marine microbial
and metagenomic
development : *Data & Tools*

Erwan Corre
ABiMS





IFB-Grand Ouest Platform

Team leader : C. Caron

- E-infrastructure : HPC environment (600 core cluster, 500 TB storage) and Bio analysis tools, Galaxy
- Expertise in the analysis of omics data
- Software engineering (information systems, databases, etc.)
- Training modules

Research team

Marine and environmental Biology

- Tools and methodology developments ,
- Research projects in marine genomics and metagenomics



Metagenomics

- Carnivorous sponge microflora (col. Museum)
- Abalone stomach microflora (PIA Idealg)
- Macroalgae epibiont microflora (PIA Idealg)
- Meta-barcoding , -genomic, -transcriptomic analysis "Post TARA" (PIA OCEANOMICS).
- Metagemonic on thermophilic hydrothermal vent flora (LM2E : Lois Maignien)
- Metagenomics of microbial flora associated with the methane cycle in marine sediments on continental margins (IFREMER).

Comparative Microbial genomics

- Marine Cyanobacteria (L. Garczareck)
- Marine vibrio (F. Le Roux)
- Marine Flavobacteria (G. Michel)



Developments

DATA

- Comparative genomic database

TOOLS

- Genome scaffolding tool
- Metabarcoding pipeline

***Cyanorak v.2*, an information system dedicated to the expert curation and annotation of clusters of orthologous sequences from marine picocyanobacteria**

Antoine Bisch, Loraine Guéguen, Gregory Farrant, Mark Hoebeke, Gildas Le Corguillé, Erwan Corre, Wilfrid Carré, Christophe Caron, Laurence Garczarek and Frédéric Partensky

Marine Phototrophic Prokaryotes (MaPP) Team, Phytoplankton Group, UMR7144 CNRS-UPMC, Station Biologique, Roscoff
Analysis and Bioinformatics for Marine Science (ABiMS), FR2424 CNRS, Station Biologique, Roscoff

Cyanorak

- Storage and visualization of picocyanobacteria sequences
- Clusters of orthologous groups
- Aims: Curation, annotation and Export

History

- **Cyanorak v.1 (2005) : 14 genomes (11 Syn., 3 Proc.)**
 - Manual annotation of a fixed set of genomes
- **Cyanorak v.2 (2013): 57 genomes (40 Syn., 3 Cya., 14 Proc.)**
 - Genome sequence origin:
 - Cyanorak v.1: 11 Syn, 3 Proc.
 - NCBI/RefSeq: 4 Syn., 2 Cya., 11 Proc.
 - Newly sequenced genomes : 25 Syn., 1 Cya.
 - Import and semi-automatic annotation of new sets of genomes

Database Contents


Organisms	14
Genes	35,993
Clusters	7,847

Database Contents

Organisms	57
Genes	154,723
Clusters	24,415

Organisms page

Gene page

<div> <div>  <div> <div>CyanoLogos</div> <div>Roscoff</div> </div> </div> <div> <div>Small Station Biologique Roscoff</div> </div> </div>		Cyanoarak - Organism List																			
<div> <div>HOME</div> <div>ORGANISMS</div> <div>SEARCH</div> <div>MAPPINGS</div> <div>HISTORY</div> <div>ABOUT US</div> <div>REFERENCES</div> <div>LINKS</div> <div>BLAST</div> </div>																					
<div> <div>Cyanorak Release 2.0 (restricted)</div> <div>Cyanorak Release 1.0</div> </div>																					
<div> <div>Create new gene</div> </div>																					
<div> <div> <div>Prochlorococcus (14)</div> <div>Synechococcus / Cyanobium (43)</div> </div> </div>																					
<div> <div> <div> <div>Name</div> <div>Species</div> <div>SubCluster</div> <div>Clade</div> <div>SubClade</div> <div>Pigment type</div> <div>Sequencing center</div> <div>Genome status</div> <div>Genome size</div> <div>Contigs</div> <div>CGS</div> <div>Chapters</div> <div>GC %</div> </div> </div> </div>																					
Cy1_N501	Cyanobium sp.	5.2			1	Genoscope	WGS	2 727 770	47	2,902	2,777	67.3									
CyA_C030302	Cyanobium gracile	5.2			1	JGI	Complete	3 342 364	1	3,280	3,061	68.7									
CyA_C03091	Cyanobium sp.	5.2			1	JCVI	WGS	2 832 412	1	2,771	2,953	68.7									
Syn_A15-122	Synechococcus sp.	5.1		WPC1	3c	Genoscope	WGS	2 537 768	18	2,791	2,686	60.5									
Syn_A15-28	Synechococcus sp.	5.1		III	IIIb	3c	Genoscope	2 336 026	30	2,690	2,603	60.2									
Syn_A15-44	Synechococcus sp.	5.1		II	IIa	2	Genoscope	2 586 506	77	3,118	2,975	52.7									
Syn_A15-60	Synechococcus sp.	5.1		VII		3c	Genoscope	2 537 019	25	3,097	3,001	59.1									
Syn_A15-52	Synechococcus sp.	5.1		II	RC	3a/b	Genoscope	2 237 102	37	2,777	2,699	61.3									
Syn_A15-64.1	Synechococcus sp.	5.1		CRD1	CRD1b	3a/b	Genoscope	1 596 041	96	4,268	4,091	61.3									
Syn_BI05-14.1-3	Synechococcus sp.	5.1		CRD1	CRD1a	3a/a	Genoscope	2 693 428	45	3,308	3,178	53.2									
Syn_B1-107	Synechococcus sp.	5.1		IV	IVa	3a/a	JCVI	2 285 035	1	2,525	2,438	55.0									
Syn_BMK_MC-1	Synechococcus sp.	5.1		V		2	Genoscope	2 619 408	75	3,053	2,920	58.9									
Syn_BIOM18	Synechococcus sp.	5.1		IIIa		3c	Genoscope	2 314 805	186	2,832	2,717	61.9									
Syn_CB0104	Synechococcus sp.	5.2		CB4		1	TIGR	6 991 045	1	3,010	2,800	64.2									
Syn_C030205	Synechococcus sp.	5.2		CB5		2	TIGR	2 411 159	1	2,748	2,617	60.0									
Syn_C03011	Synechococcus sp.	5.1		I	Ia	3a/a	Complete	6 606 749	1	2,931	2,789	52.5									
Syn_CC0905	Synechococcus sp.	5.1		II	IIc	3c	JGI	2 510 660	1	2,713	2,559	59.2									
Syn_CC0902	Synechococcus sp.	5.1		IV	IVa	3a/a	Complete	2 234 829	1	2,387	2,326	54.2									
Syn_M48.1	Synechococcus sp.	5.1		II	IIa	3a	Genoscope	2 109 294	1	2,467	2,411	60.3									
Syn_M02045	Synechococcus sp.	5.1		VI	Vla	3c	Genoscope	2 422 988	40	2,701	2,613	61.8									
Syn_M05014.1	Synechococcus sp.	5.1			3a/b		Genoscope	2 290 219	200	2,767	2,704	46.8									
Syn_M05B_1	Synechococcus sp.	5.1		I	Ib	3a/a	Genoscope	2 433 234	35	3,043	2,964	52.8									
Syn_NDUM7013	Synechococcus sp.	5.1		VII		3a	Genoscope	2 541 016	28	2,893	2,810	57.4									

Search tool

[Fast Search](#)
[Advanced Search](#)
[One-Click Search](#)

CHOOSE SEARCH TYPE:

[Search for an identifier / annotation](#)
[Search for common / specific clusters](#)

Clusters

Search for

Search type ☐ exact search ☒ pattern search

Identifiers: ☐ [do/iselect all](#) ☐ cluster number ☐ Cyanorak V1 cluster number

Annotations: ☐ [do/iselect all](#) ☐ gene name ☐ product ☐ COG / NOG ☐ batNOG ☐ CyANO ☐ CyOG ☐ EC number ☐ TIGR role ☐ MaPP role ☐ comment

☐ qualifiers: ☐ [mapping_merge_rule](#) ☐ [GO terms](#)

☐ [ALL](#) ☐ [biological process](#) ☐ [cellular component](#) ☐ [molecular function](#)

Genes

Identifiers: ☐ [do/iselect all](#) ☐ [orf id](#) ☐ Cyanorak V1 locus tag ☐ Cyanorak V2 locus tag ☐ NCBI GI ☐ RefSeq ☐ JGI ☐ JCVI ☐ MicrobesOnline

Annotations: ☐ [do/iselect all](#) ☐ gene name ☐ product ☐ COG / NOG ☐ batNOG ☐ CyANO ☐ CyOG ☐ EC number ☐ TIGR role ☐ comment

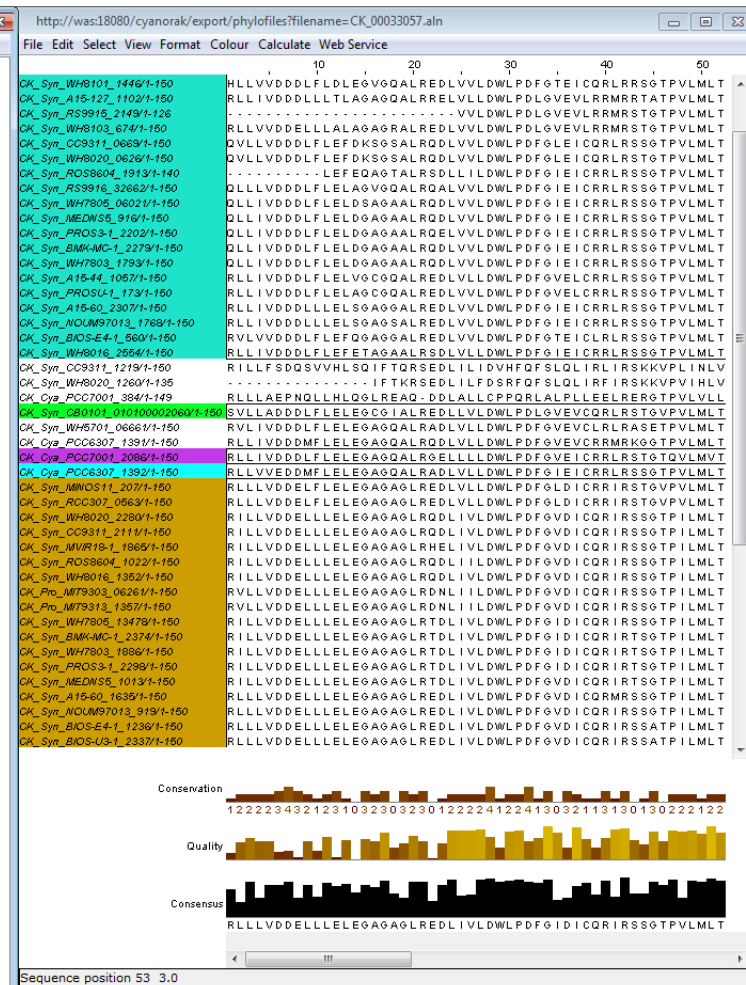
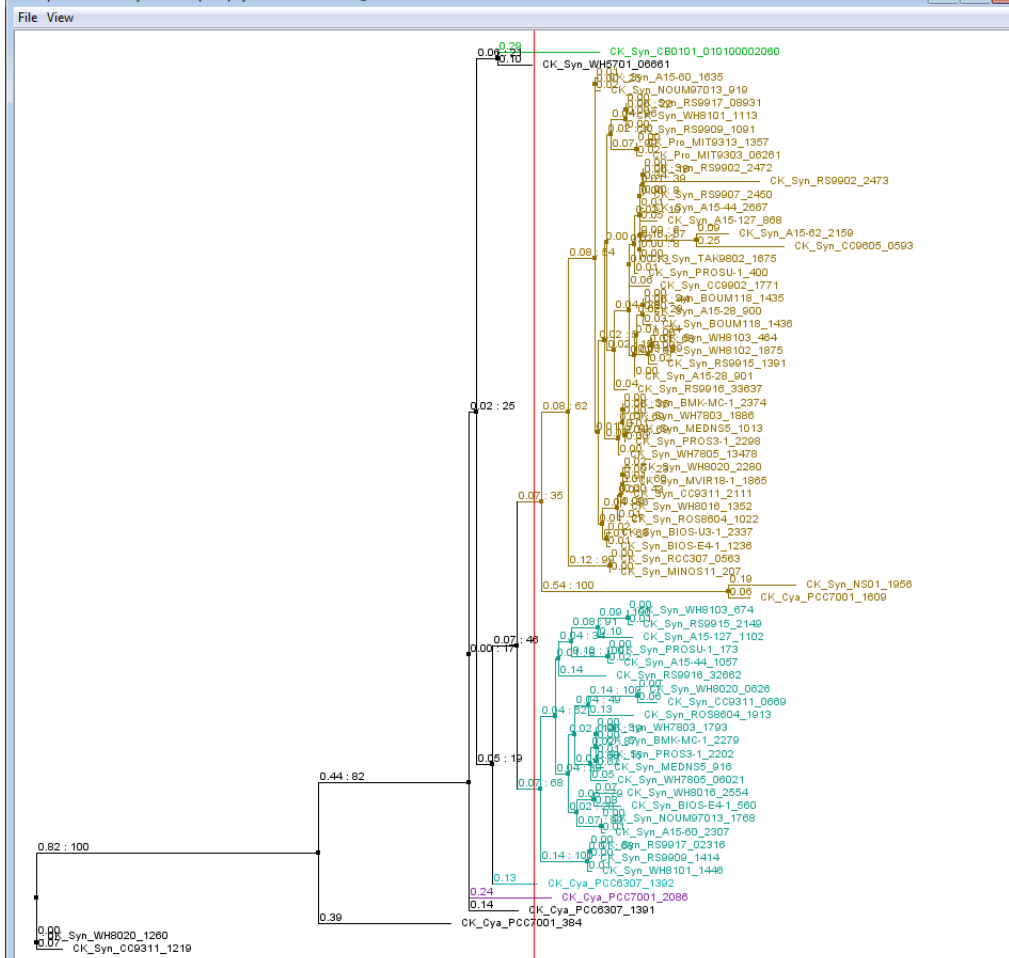
☐ qualifiers: ☐ [artificial_location](#) ☐ [codon_start](#) ☐ [color](#)

[illegible][illegible]

Cluster page

Identifiers	
Cytosine V1 cluster number 579	
Annotated features	
Gene name pcdB	
Product cytochrome b5	
Functional categories	
EggNOG	COG2 DGG COG1290 : Energy production and conversion METABOLISM
	hsr(NOG) hsr(COG)5234 : Energy production and conversion METABOLISM
	cyt(NOG) cyo(NOG)5314 : Energy production and conversion METABOLISM cyo(NOG)1014 : Energy production and conversion METABOLISM
	EC number 1.10.99.1
	TEIG rule
	MAPP rule
	Biological processes 0005767 0022000
	Cellular component 0005514
GO terms	Molecular function 0015138
	To be classified
Other qualifiers	
mapping_merge_rule	NLS_SINRA_SUPERSET_nbf_SINRA_SUPERSET
Comment	
None.	

http://was:18080/cyanorak/export/phylofiles?filename=CK_00033057.tree



PhyML tree

MAFFT Alignment

○ Genome context

Genome context of gene **CK_Cya_NS01_1303** (cluster: **CK_00000009**)



○ Local BLAST

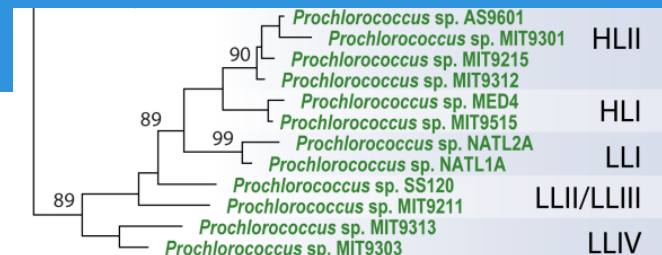

BLAST
[BLAST](#)
[Entrez](#)
[?](#)

Marine Cyanobacteria BLAST Server

Choose program to use and database to search:

Program
 Genome
☒ Syne blastp vs proteins CA4) genome
☒ Syne blastx vs proteins , CA4) genome
☒ Syne tblastn vs genome Ib, 3a) genome
☒ Syne tblastn vs genes b, 3a) genome
☒ Syne tblastx vs genome , 3a) genome
☒ Syne tblastx vs genes 2) genome
☒ Syne tblastx vs genes 3a) genome
☒ Syne blastn vs rRNA or tRNA
☒ Synechococcus sp. TAK9802 (IIa, 3a) genome
☒ Synechococcus sp. WH8109 (IIa, 3b) genome
☒ Synechococcus sp. RS9902 (IIa, 3c) genome
☒ Synechococcus sp. CC9605 (IIc, 3c) genome
☒ Synechococcus sp. A15-62 (IIc, CA4) genome
☒ Synechococcus sp. PROSU-1 (IIIa, CA4) genome
☒ Synechococcus sp. WH8103 (IIIa, 3b') genome
☒ Synechococcus sp. BOUM118 (IIIa, 3c) genome
☒ Synechococcus sp. WH8102 (IIIa, 3c) genome
☒ Synechococcus sp. RS9915 (IIIa, CA4) genome
☒ Synechococcus sp. A15-28 (IIIb, 3c) genome

☒ Check All ☐ Uncheck All
☒ Prochlorococcus sp. MED4 (HLI) genome
☒ Prochlorococcus sp. MIT9515 (HLI) genome
☒ Prochlorococcus sp. AS9601 (HLII) genome
☒ Prochlorococcus sp. MIT9202 (HLII) genome
☒ Prochlorococcus sp. MIT9215 (HLII) genome
☒ Prochlorococcus sp. MIT9301 (HLII) genome
☒ Prochlorococcus sp. MIT9312 (HLII) genome
☒ Prochlorococcus sp. UH18301 (HLII) genome
☒ Prochlorococcus sp. NATL1A (LLI) genome
☒ Prochlorococcus sp. NATL2A (LLI) genome
☒ Prochlorococcus sp. SS120 (LLII) genome
☒ Prochlorococcus sp. MIT9211 (LLIII) genome
☒ Prochlorococcus sp. MIT9303 (LLIV) genome
☒ Prochlorococcus sp. MIT9313 (LLIV) genome



Scanlan et al., MMBR, 2009

Phyletic pattern:

- core/accessory/specific genes
- genes involved in adaptation to specific environmental niches

Phyletic pattern

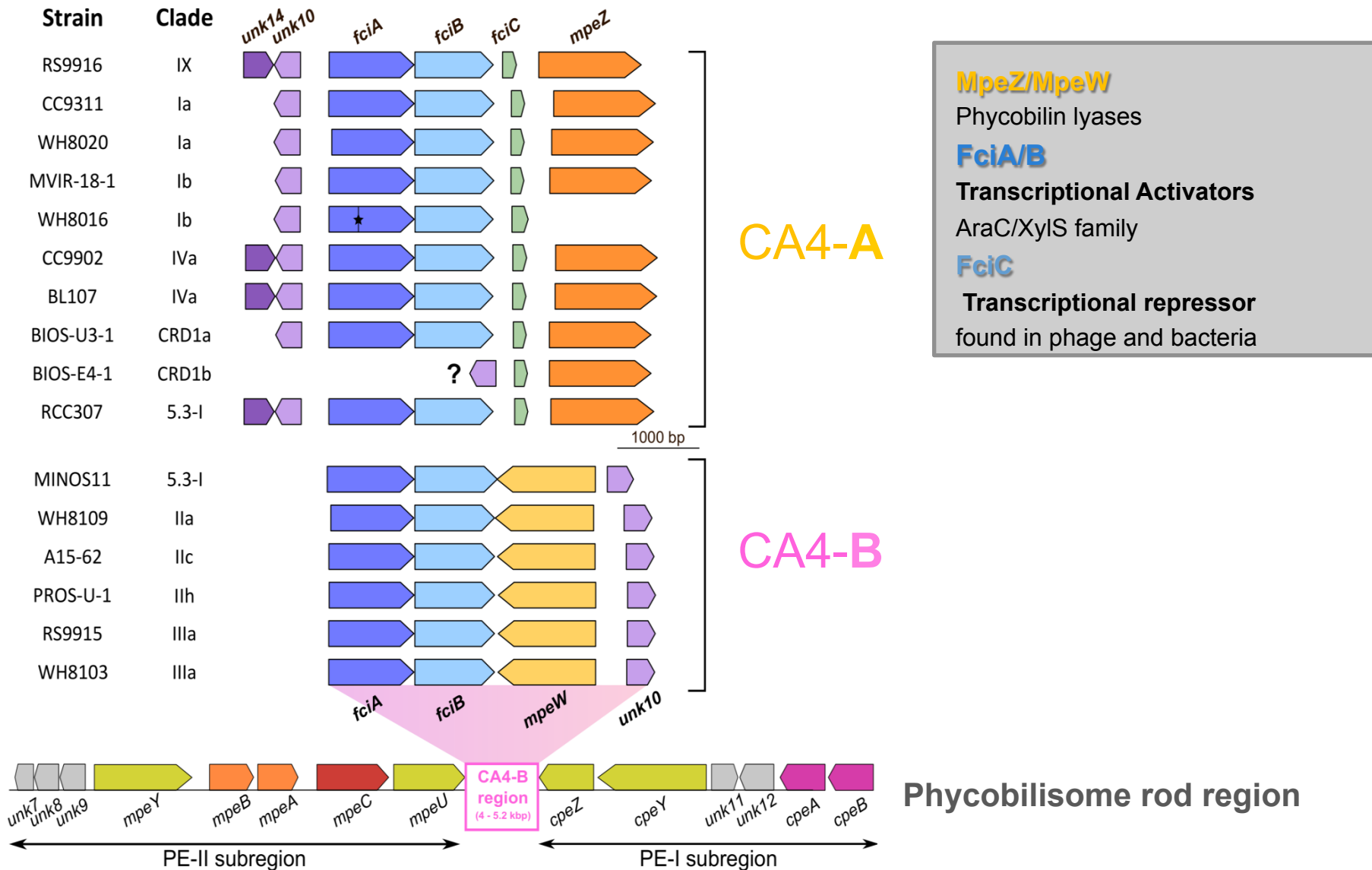
Prochlorococcus															
Subcluster HL								Subcluster LL							
HLI		HLII						LLI		LLII	LLIII	LLIV			
MED4	MIT9515	AS9601	MIT9202	MIT9215	MIT9301	MIT9312	UH18301	NATL1A	NATL2A	SS120	MIT9211	MIT9303	MIT9313		
Low b/a	Low b/a	Low b/a	Low b/a	Low b/a	Low b/a	Low b/a	Low b/a	High b/a	High b/a	High b/a	High b/a	High b/a	High b/a		
1	1	1	1	1	1	1	1	1	1	0	0	0	0		
Synechococcus															
Subcluster 5.1															
Ia		Ib			IIa				IIc		IIh	IIIa			
CC9311	WH8020	MVIR-18-1	ROS8604	WH8016	A15-44	M16.1	RS9902	RS9907	TAK9802	WH8109	A15-62	CC9605	PROS-U-1	BOUM118	RS9915
3dA	3dA	3aA	3a	3aA	2	3a	3c	3a	3a	3bB	3dB	3c	3dB	3c	3dB
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Synechococcus															
Subcluster 5.1															
V		VIa		VIb	VII		VIII		IX	CRD1a	CRD1b	WPC1			
BMK-MC-1	WH7803	MEDNS5	WH7805	PROS-7-1	A15-60	NOUM97013	RS9909	RS9917	WH8101	RS9916	BIOS-U3-1	BIOS-E4-1	A15-127		
2	3a	3c	2	2	3c	3a	1	1	1	3dA	3dA	3cB	3c		
1	1	1	1	1	1	1	1	1	1	1	1	1	1		
Synechococcus				Cyanobium				Cyanobium				Cyanobium			
Subcluster 5.2		Subcluster 5.3		Subcluster 5.2		Subcluster 5.2		Subcluster 5.2		Subcluster 5.2		Subcluster 5.2		Subcluster 5.2	
CB4	CB5														
CB0101	CB0205	WH5701	MINOS11	RCC307	NS01	PCC6307	PCC7001								
1	2	1	3dB	3eA	1	1	1								
1	1	1	1	1	1	1	1								

DNA photolyase *phrA* gene absent
in strictly LL *Prochlorococcus* strains

List of genes (53)

25/09/2014

The specific island of chromatic adaptors (CA4)



MpeZ/MpeW

Phycobilin lyases

FciA/B

Transcriptional Activators

AraC/XylS family

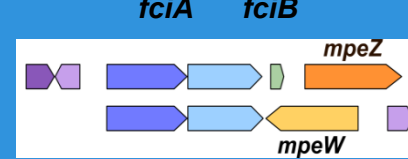
FciC

Transcriptional repressor

found in phage and bacteria

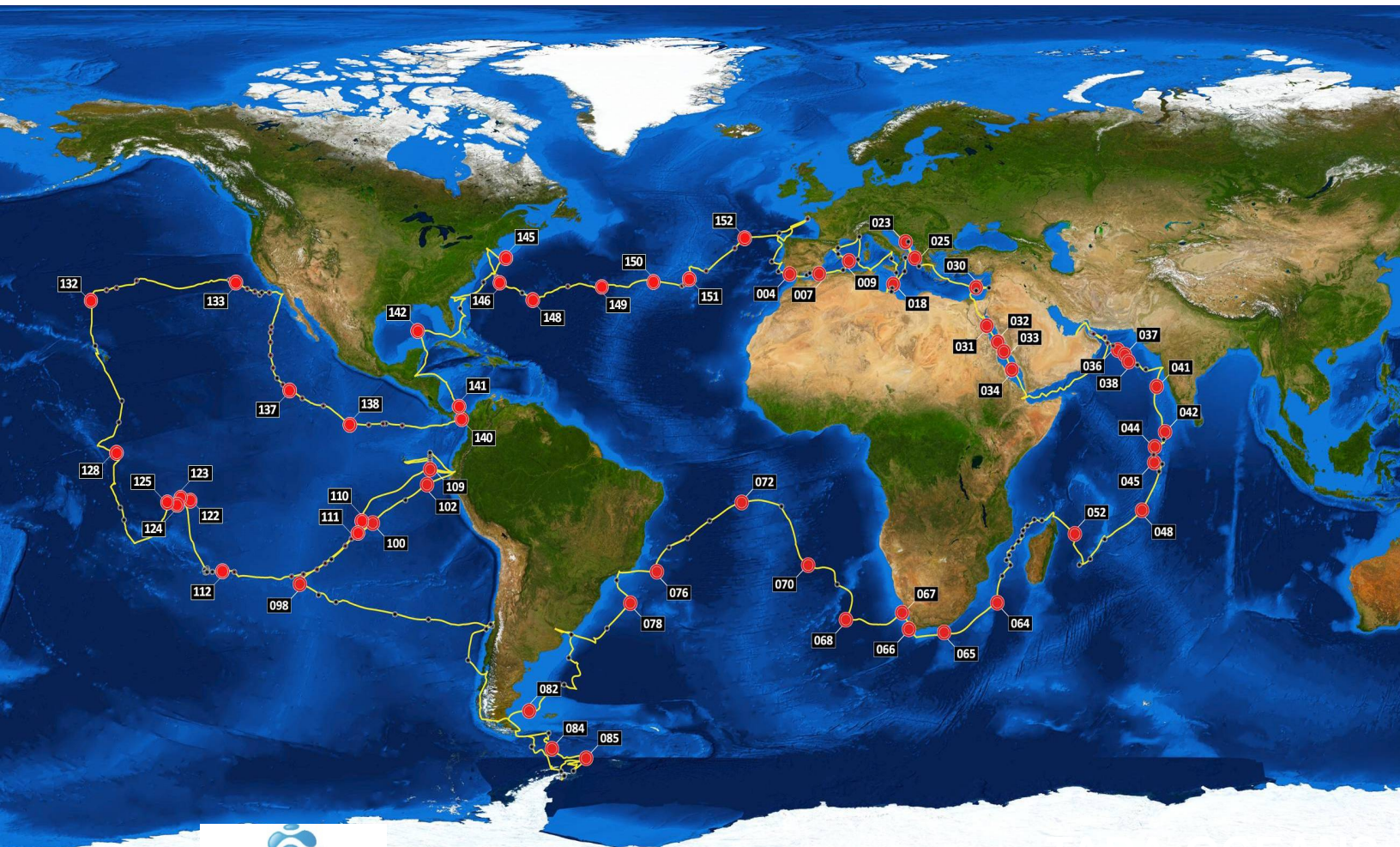
Geographical Distribution of CA4

CA4-A
CA4-B



- *mpeW* / *mpeZ* recruitment -> Abundance of CA4-A and B
- *petB* recruitment -> Abundance of *Synechococcus*

=> Relative proportion of CA4-A and B



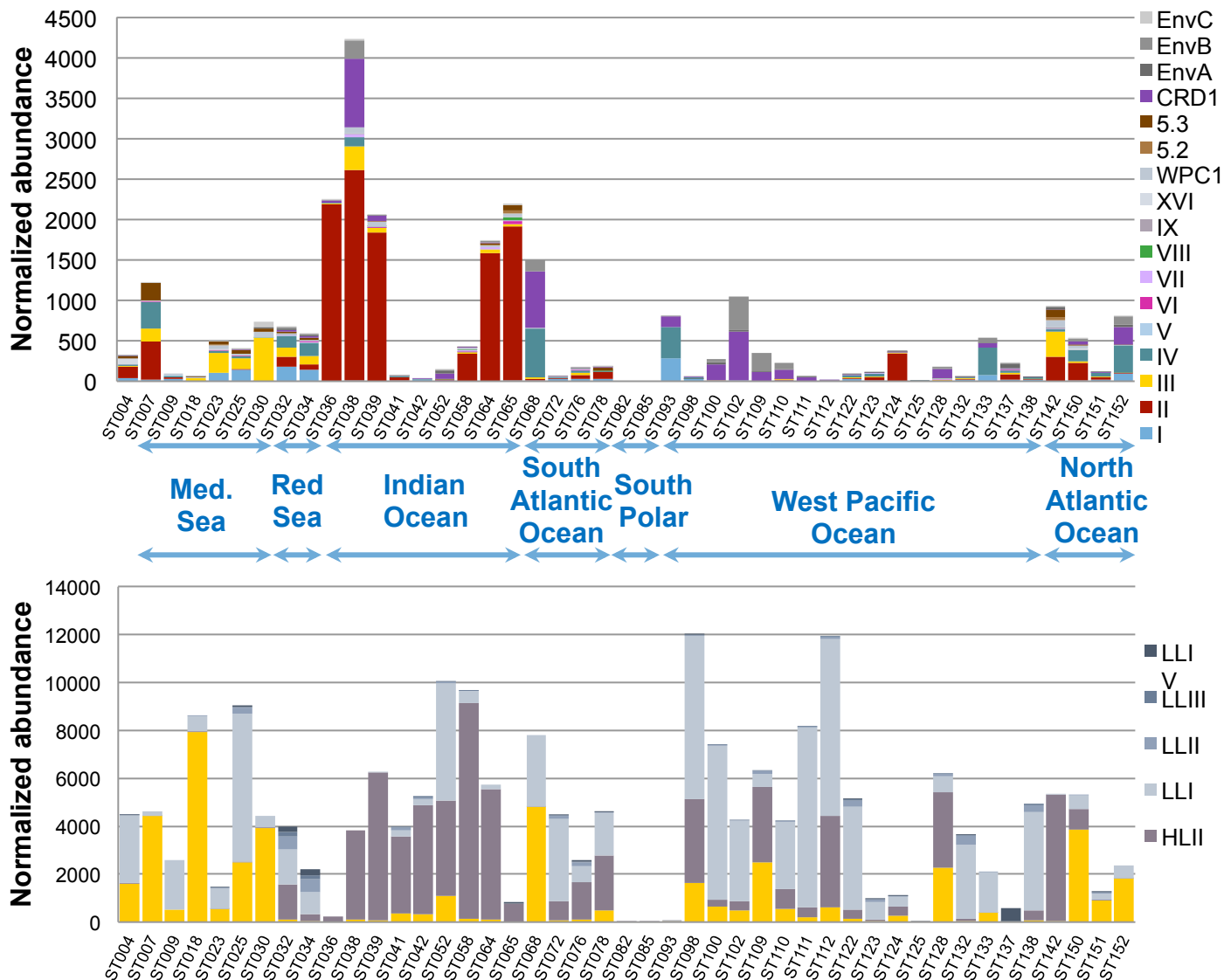
Metagenomes
~5.10⁶ reads
(~175 bp)

BLASTX
alignment
vs 2 *petB*

BLASTN
alignment
vs 280 *petB*

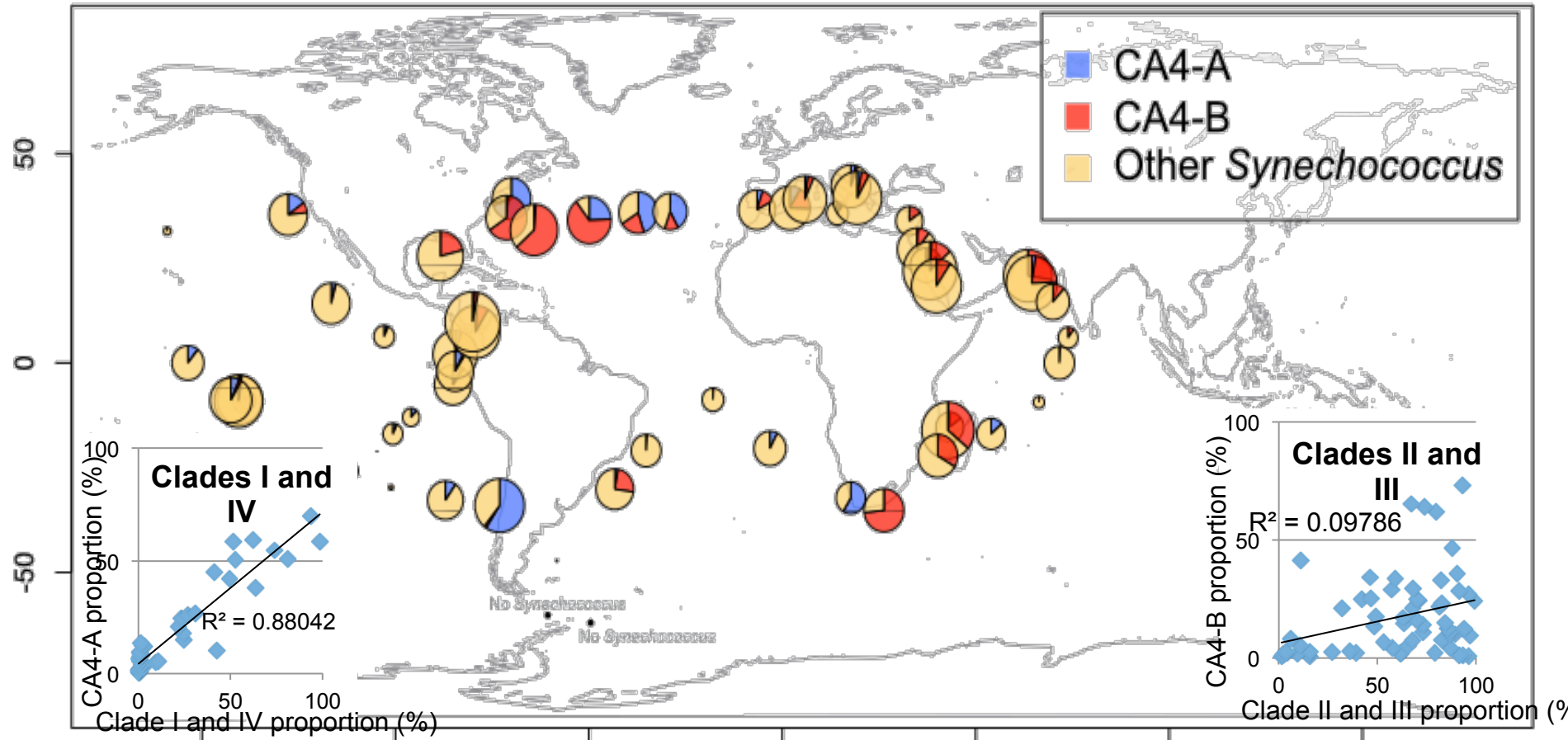
Taxonomic
Assignment
Abundance
profiles

Abundance of *Synechococcus* (A) and *Prochlorococcus* (B) clades in DCM stations of the TARA Oceans metagenomes (small size fraction) based on *petB* miTAG recruitments



Geographical Distribution of CA4

- Abundant in coastal region / Temperate oceanic waters
 - Until 90% of *Synechococcus* populations
- Most clades I and IV are CA4
- Only some clades II and III display this capacity



The future of Cyanorak (v.3):

○ More data type stored/visualized

- Additional sequences types (rRNA, ncRNA, tRNA, UTR)
- Transcriptomic data (Arrays and RNAseq)

Transcriptomes from 4 *Proc.* and 7 *Syn.* (Genoscope)

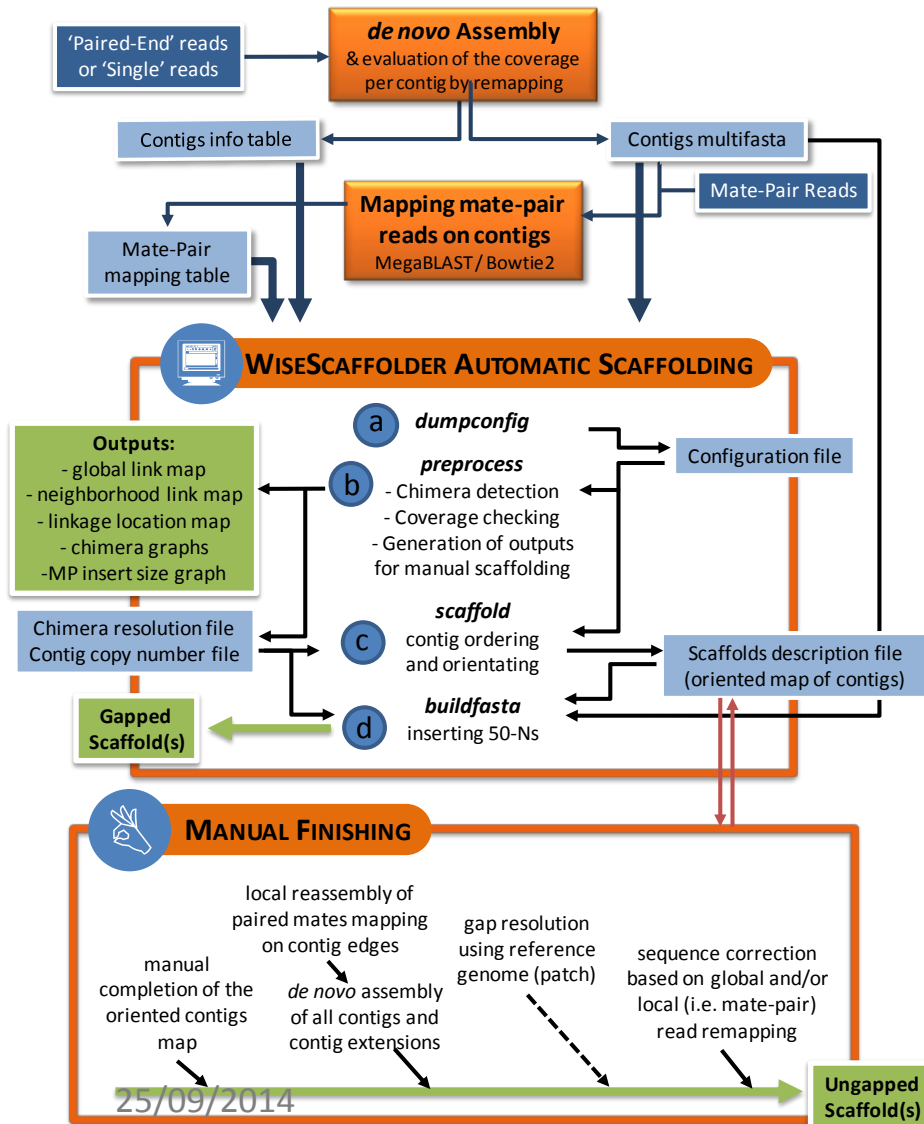
Transcriptomes from stress experiments (HL/UV, stress ox, T°C, etc.) on 4 *Syn.* strains
(~ 334 transcriptomes, ANR SAMOSA project)

- Metagenomic and metatranscriptomic data
Recruitment plots of environmental reads
- Information about mutants

○ New modules

- Genome Browser (JBrowse)
- Genomic island viewer
- Comparative analysis tools (core/accessory/unique genes)

TOOLS : WiseScaffolder Workflow



Inputs:

- Contigs produced by CLC
- Contigs description (coverage, length)
- Mapping of mate-pairs

4 subcommands, 4 steps of tuning

- dumpconfig
- preprocess
- scaffold
- buildfasta

Chimera & contig copy number detection

Outputs:

- scaffolds description
- fasta of scaffolds and unscaffolded contigs
- various tables and graphs to help with manual scaffolding



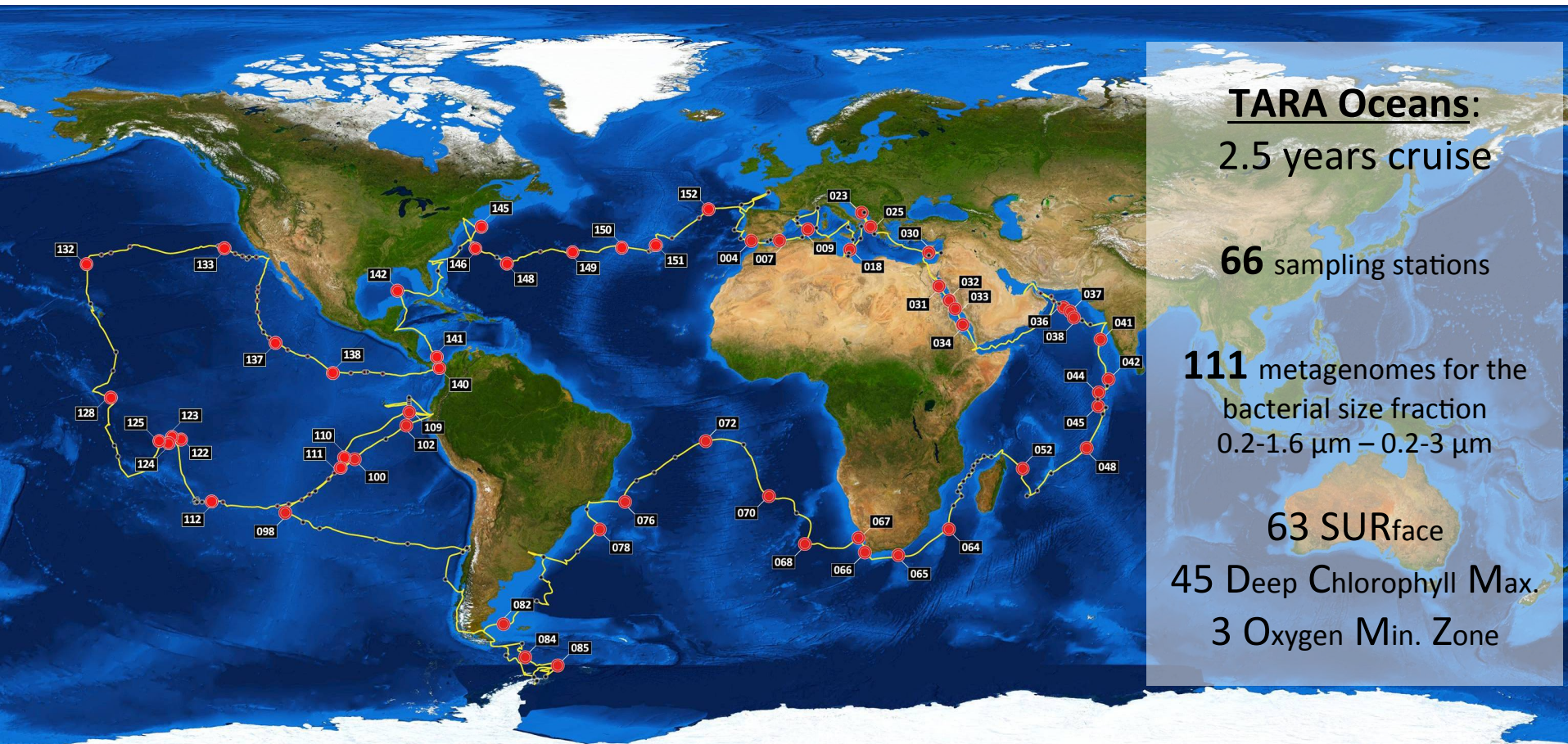
Stéphane Audic (Évolution du Plancton et Écosystèmes Pélagiques,
Adaptation et Diversité en Milieu Marin, Station biologique de
Roscoff)

stephane.audic@sb-roscoff.fr

Collaboration Abims (Analyses and bioinformatics for Marine
Science) <http://abims.sb-roscoff.fr>

METAGENOMICS

TARA Oceans Metagenomic Dataset



Metagenic Analysis

Environment



Metagenomic

Sequencing of the full DNA content of an environment. Works for prokaryotic and small eukaryotic genomes

Metatranscriptomic

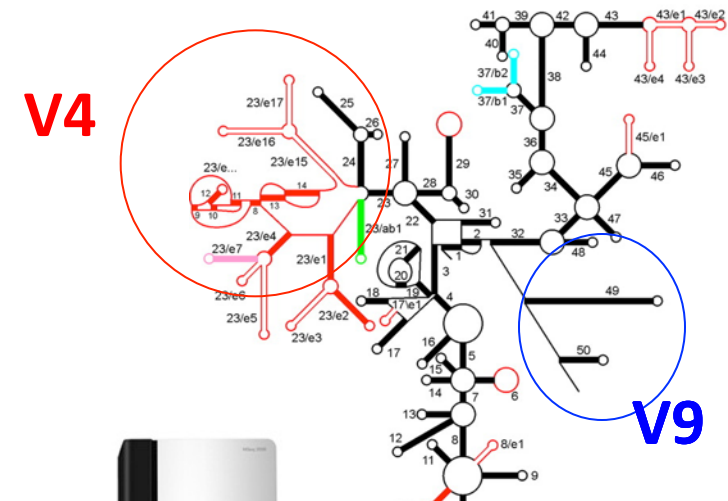
Sequencing of the full RNA content of an environment. Selection of organism categories : ex - eucaryotes poly-A selection

metabarcodingMetagenic

Sequencing of all « instances » of a gene in an environment.

Selected gene by PCR with primers with a great taxonomic coverage but specific enough to amplify only the gene of interest

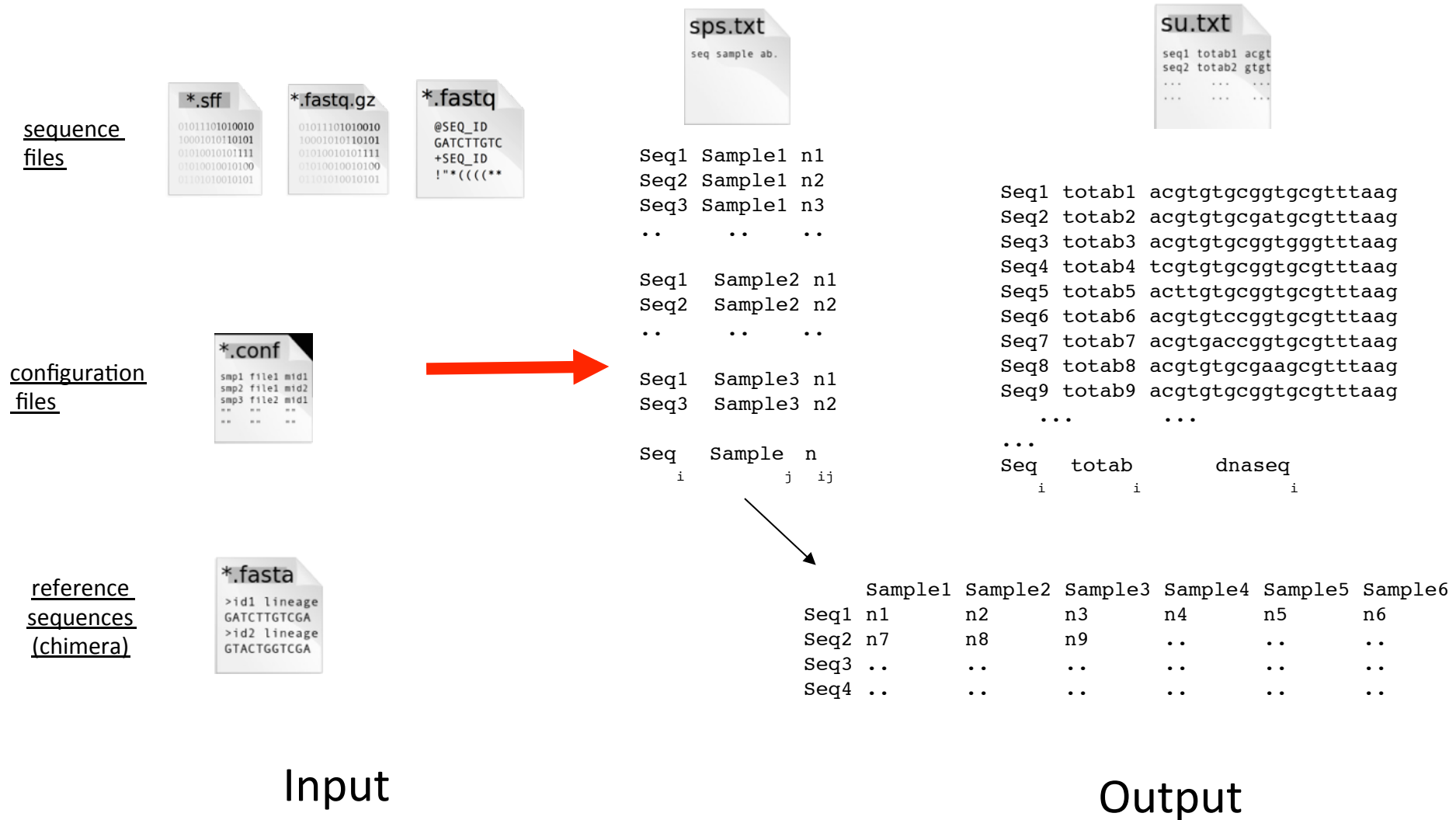
Qualitative and semi-quantitative approach.



Wuyts J et al. Nucl. Acids Res.
2001;29:5017-5028

Step 1: raw reads-> table of occurrences

From sequence files and configuration files, get number of occurrences for **each detected sequence in each sample**.



Creation of "Configuration file"

Separate samples:
In which file?
Wich MID?
Wich primerset? (= Which marker)
Wich primer direction

Index description for sample
separation (MID)

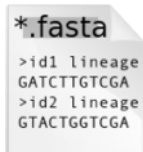
sequence
files



configuration
files



reference
sequences
(chimera)



Sample file

```
SMP1 fichier1 MID1 pset Dir
SMP2 fichier1 MID2 pset BOTH
SMP1 fichier2 MID1 pset CMP
etc...
```

MID file

```
MID1 ACGCGTG
MID2 GCTAGTG
MID3 CCGTGTA
MID4 GCTGGTC
etc...
```

Pset file

```
pset1 primerF primerR base_de_ref
\ error-rate size-range
pset2 primerF primerR base_de_ref
\ error-rate size-range
etc...
```

Input

Description of Primers, reference
database, etc.

Raw data extraction

- 1- convert sff to fastq (if necessary)
- 2- 'scaling' of quality values (33;64; ...)
- 3- search in raw data the exact motif
`"MID-PrimerF- ([acgt]+)primerR(rev) "`
 -> output the expected error raw in the corresponding
 sequence (base on the worst 50 nuc)

One file per sample

sequence
files



configuration
files



```
np_complete_S*

>name1 E=0.2 L=300 S=S001
AGCGTGCGTTGGTGTCCGTAGTC
>name2 E=0.7 L=310 S=S001
AGCGTGCGAACGTGTCCGTAGTC
>etc...
```

4- dereplication
Number of
occurence

reference
sequences
(chimera)



5- collect of common
identifiers in several
samples

md5.several

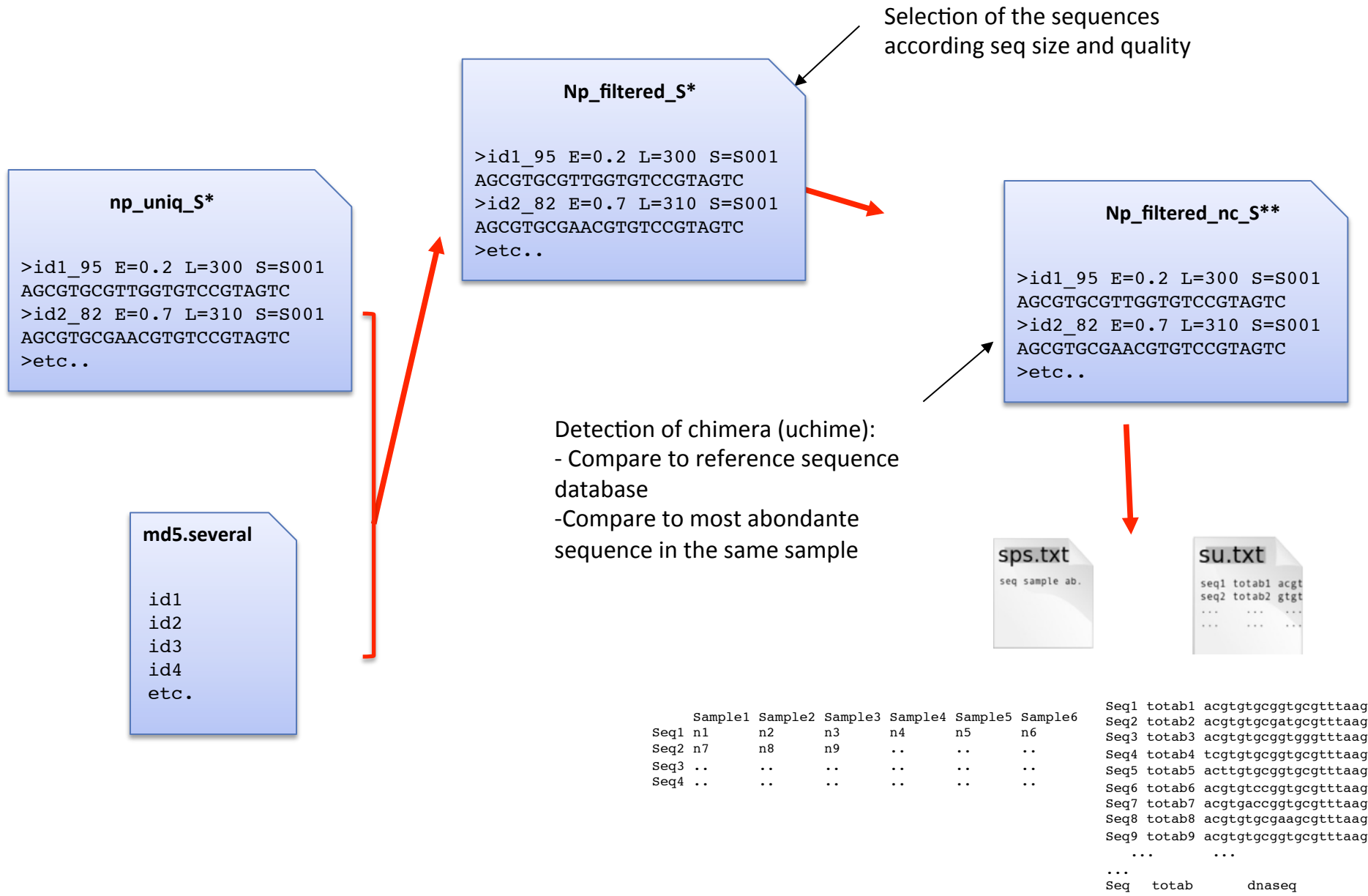
```
id1
id2
id3
id4
etc.
```

np_uniq_S*

```
>id1_95 E=0.2 L=300 S=S001
AGCGTGCGTTGGTGTCCGTAGTC
>id2_82 E=0.7 L=310 S=S001
AGCGTGCGAACGTGTCCGTAGTC
>etc..
```

Input

Cleaning and compilation of results



Graphics



- 

Analyses and Bioinformatics for Marine Science

24

22



Graphics

Execute

```
6 lines
format: txt, database: 2
adding: EukV4/md5.several (stored
0%) adding:
EukV4/np_uniq_S0006 (deflated
83%) Creates the input file list per
/w/galaxy/prod/galaxy-dist/tools
/abims/metabarcoding/scripts
/Perl/build_config.pl -p /w/galaxy
/prod/galaxy-dist/database/files
/052/dataset
```

Step 2: occurrence table -> assignation

Node

Lineage

```
Eukaryota|Opisthokonta|Mesomycetozoa|Ichthyosporea|Rhynosporida|Rhynosporidae
Eukaryota|Archaeplastida|Chlorophyta|Trebouxiophyceae|Chlorellales|Chlorellales_X|Meyerella|Meyerella+sp.
Eukaryota|Opisthokonta|Metazoa|Arthropoda|Hexapoda|Insecta|Coccobius
Eukaryota|Opisthokonta|Metazoa|Arthropoda|Hexapoda|Insecta|Saldula|Saldula+sp.
Eukaryota|Alveolata|Ciliophora|Litostomatea|Trichostomatia|Amylavorax|Amylavorax+dogieli
Eukaryota|Opisthokonta|Metazoa|Arthropoda|Crustacea|Ostracoda|Parasterope|Parasterope+pollex
Eukaryota|Archaeplastida|Chlorophyta|Chlorophyceae|Chlorophyceae_X|Cw-Chlamydomonadales|Pteromonas|Pteromonas+angulosa
Eukaryota|Opisthokonta|Metazoa|Mollusca|Gastropoda|Caenogastropoda|Crepidula|Crepidula+fornicata
Eukaryota|Stramenopiles|Stramenopiles_X|Bacillariophyta|Bacillariophyta_X|Raphid-pennate|Cymbella|Cymbella+minuta
Eukaryota|Opisthokonta|Metazoa|Arthropoda|Hexapoda|Insecta|Hemiptera|Hemiptera+cyaneipennis
```

Experimental Data

```
Seq1 totab1 acgtgtgcggtgcggtttaag
Seq2 totab2 acgtgtgcgatgcggtttaag
Seq3 totab3 acgtgtgcggtgggtttaag
Seq4 totab4 tcgtgtgcggtgcggtttaag
Seq5 totab5 acttgtgcggtgcggtttaag
Seq6 totab6 acgtgtccggtgcggtttaag
Seq7 totab7 acgtgaccggtgcggtttaag
Seq8 totab8 acgtgtgcgaagcggtttaag
Seq9 totab9 acgtgtgcggtgcggtttaag

...
Seq totab dnaseq
i i i
```

Reference sequences

Reference	Node	Node	Parent
14.1.1837_U	17059	58684	58175
82.1.1888_U	41847		
85.1.1715_U	28767	57573	57566
75.1.1781_U	4702		
90.1.1711_U	46762		
59.1.1955_U	43646	58017	57997
54.1.1769_U	39467		
80.1.2094_U	3238	44919	44918
82.1.2144_U	1923		
90.1.1695_U	2505	8563	8562
		28336	1194
		28801	28800
		5006	2228
		55548	55546
		29429	2292

ggsearch
alignment global
multi-threaded

Hit List

```
Seq1 %id ref1
Seq1 %id ref2
Seq2 %id ref3
Seq3 %id ref4
Seq3 %id ref5
... ..
```

Assignation

```
Seq1 %id ref1,ref2 lineage
Seq2 %id ref3 lineage2
Seq3 %id ref4,ref5 lineage3
... ..
... ..
```

Step 3:co-occurrence table -> OTUisation

Experimental data

```
Seq1 totab1 acgtgtgcggtgcgtttaag
Seq2 totab2 acgtgtgcatgcgtttaag
Seq3 totab3 acgtgtgcggtgggtttaag
Seq4 totab4 tcgtgtgcggtgcgtttaag
Seq5 totab5 acttggtcggtgcgtttaag
Seq6 totab6 acgtgtccggtgcgtttaag
Seq7 totab7 acgtgaccggtgcgtttaag
Seq8 totab8 acgtgtgcgaagcgtttaag
Seq9 totab9 acgtgtgcggtgcgtttaag
...
Seq totab dnaseq
i i i
```

Swarms

```
Seq1 totab1 acgtgtgcggtgcgtttaag swarm_id1
Seq2 totab2 acgtgtgcatgcgtttaag swarm_id2
Seq3 totab3 acgtgtgcggtgggtttaag swarm_id3
Seq4 totab4 tcgtgtgcggtgcgtttaag swarm_id3
Seq5 totab5 acttggtcggtgcgtttaag .....
Seq6 totab6 acgtgtccggtgcgtttaag .....
Seq7 totab7 acgtgaccggtgcgtttaag .....
Seq8 totab8 acgtgtgcgaagcgtttaag .....
Seq9 totab9 acgtgtgcggtgcgtttaag .....
...
Seq totab dnaseq
i i i
```

Swarm fast and accurate OTU construction

Swarm: robust and fast clustering method for amplicon-based studies

Frédéric Mahé^{1,2,3}, Torbjørn Rognes^{4,5}, Christopher Quince⁶, Colombran de Vargas^{1,2}, and Micah Dunthorn³

PeerJ PrePrints 2:e386v1 <http://dx.doi.org/10.7287/peerj.preprints.386v1>

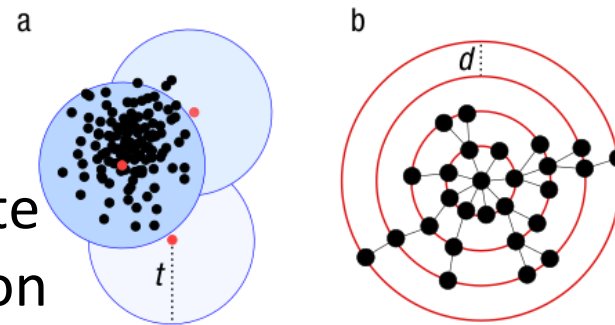


Figure 1. Visualization of the widely used greedy clustering approach based on centroid selection and a global clustering threshold, t , where closely related amplicons can be placed into different OTUs. (b) By contrast, Swarm clusters iteratively by using a small local clustering threshold, d , allowing OTUs to reach their natural limits.

	Sample1	Sample2	Sample3	Sample4	Sample5	Sample6
Swarm1	n1	n2	n3	n4	n5	n6
Swarm2	n7	n8	n9
Swarm3
Swarm4

Pipeline as a service

- Metagenetic is being developed, but it is an area that needs to offer STANDARDIZATION and dataset easily COMPARABLE.
- That's why we want to encourage the use of set of primers COMPATIBLE. (Developed as part of BioMarks / Tara Ocean project).
- Dissemination of protocols / primers
- Reference Database of specific marker: V4 specific database - V9 specific database - LSUD1D2 specific database (D. Grzebyk, lagunar ecosystemes) - Chloroplastic 16S specific database (J. Decelle, S. Romac)

Nucleic Acids Research, 2012, 1–8
doi:10.1093/nar/gks1160

The Protist Ribosomal Reference database (PR²): a catalog of unicellular eukaryote Small Sub-Unit rRNA sequences with curated taxonomy

Laure Guillou^{1,2,*}, Dipankar Bachar^{3,4}, Stéphane Audic^{1,2}, David Bass⁵, Cédric Berney⁶, Lucie Bittner^{1,2}, Christophe Boute^{1,2}, Gaetan Burgaud⁶, Colombar de Vargas^{1,2}, Johan Decelle^{1,2}, Javier del Campo⁷, John R. Dolan⁸, Micah Dunthorn⁹, Bente Edvardsen¹⁰, Maria Holzmänn¹¹, H.C.F. Kooistra Wiebe¹², Enrique Lara¹³, Noan Le Bescot^{1,2}, Ramiro Logares⁷, Frédéric Mahé^{1,2}, Ramon Massana⁷, Marina Montresor¹², Raphael Morard^{1,2}, Fabrice Not^{1,2}, Jan Pawlowski¹¹, Ian Probert^{14,15}, Anne-Laure Sauvadet^{1,2}, Raffaele Siano¹⁶, Thorsten Stoeck⁹, Daniel Vaultot^{1,2}, Pascal Zimmermann¹⁷ and Richard Christen^{3,4,*}

OPEN ACCESS Freely available online

Community Page

CBOL Protist Working Group: Barcoding Eukaryotic Richness beyond the Animal, Plant, and Fungal Kingdoms

Jan Pawlowski^{1*}, Stéphane Audic², Sina Adl³, David Bass⁴, Lassaad Belbahri⁵, Cédric Berney⁶, Samuel S. Bowser⁶, Ivan Cepicka⁷, Johan Decelle², Micah Dunthorn⁸, Anna Maria Fiore-Donno⁹, Gillian H. Gile¹⁰, Maria Holzmänn¹, Regine Jahn¹¹, Miloslav Jirků¹², Patrick J. Keeling¹³, Martin Kostka^{12,14}, Alexander Kudryavtsev^{1,15}, Enrique Lara⁵, Julius Lukeš^{12,14}, David G. Mann¹⁶, Edward A. D. Mitchell¹⁷, Frank Nitsche¹⁷, Maria Romeralo¹⁸, Gary W. Saunders¹⁹, Alastair G. B. Simpson²⁰, Alexey V. Smirnov¹⁵, John L. Spouge²¹, Rowena F. Stern²², Thorsten Stoeck⁹, Jonas Zimmermann^{11,23}, David Schindel²⁴, Colombar de Vargas^{2*}

Why Galaxy?

User point of view

- Complexes pipelines available to a larger community
- No complex installation (f.e. Qiime)
- Allows a transparent access to HPC facility

Developer point of view

- Code reviewing / clearly break the steps of an analysis / identify what can be parallelized and what can not.
- Lighten the distribution step
- Frees up time (to allow others to use its codes without having to run it by yourself)

Platform point of view

- Collaboration research/platform
- Integration of tools (workflow) / data
- Competences and visibility
- Spreading : New training subjects

