

---

# Novel bioinformatics tools to assess microbial diversity in life support systems

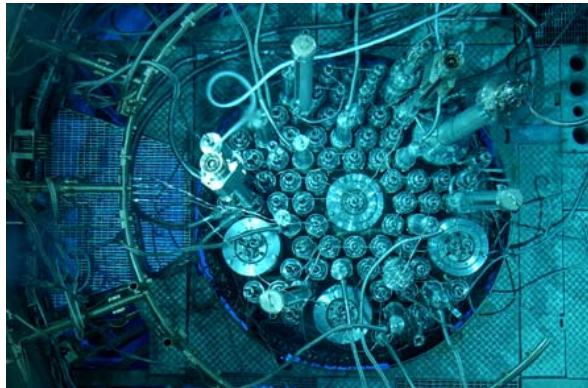
Mohamed Mysara, Natalie Leys, Pieter Monsieurs

[mahemd@sckcen.be](mailto:mahemd@sckcen.be); [pmonsieu@sckcen.be](mailto:pmonsieu@sckcen.be)

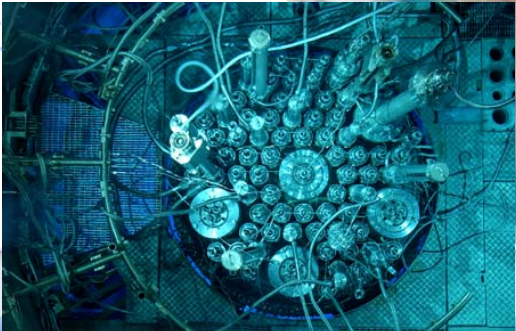
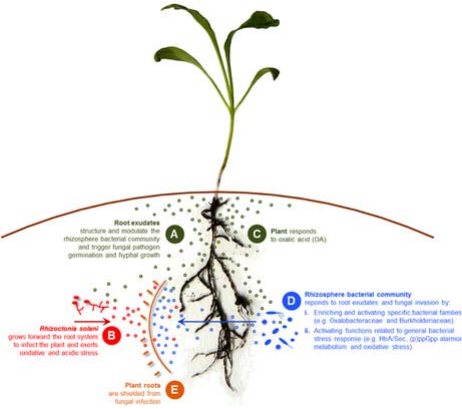
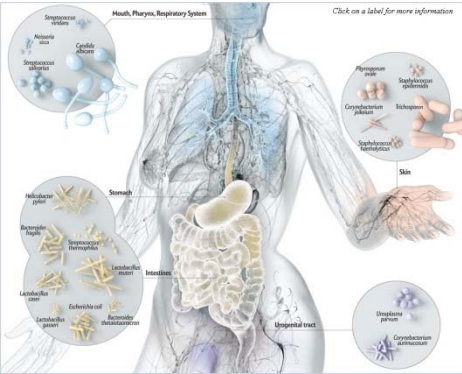


# Bacterial communities

- Bacterial species live in communities, rather than individual species
- They interact, depend and talk to each other
- These dynamic communities referred to as **microbiome**



# Microbiome studies

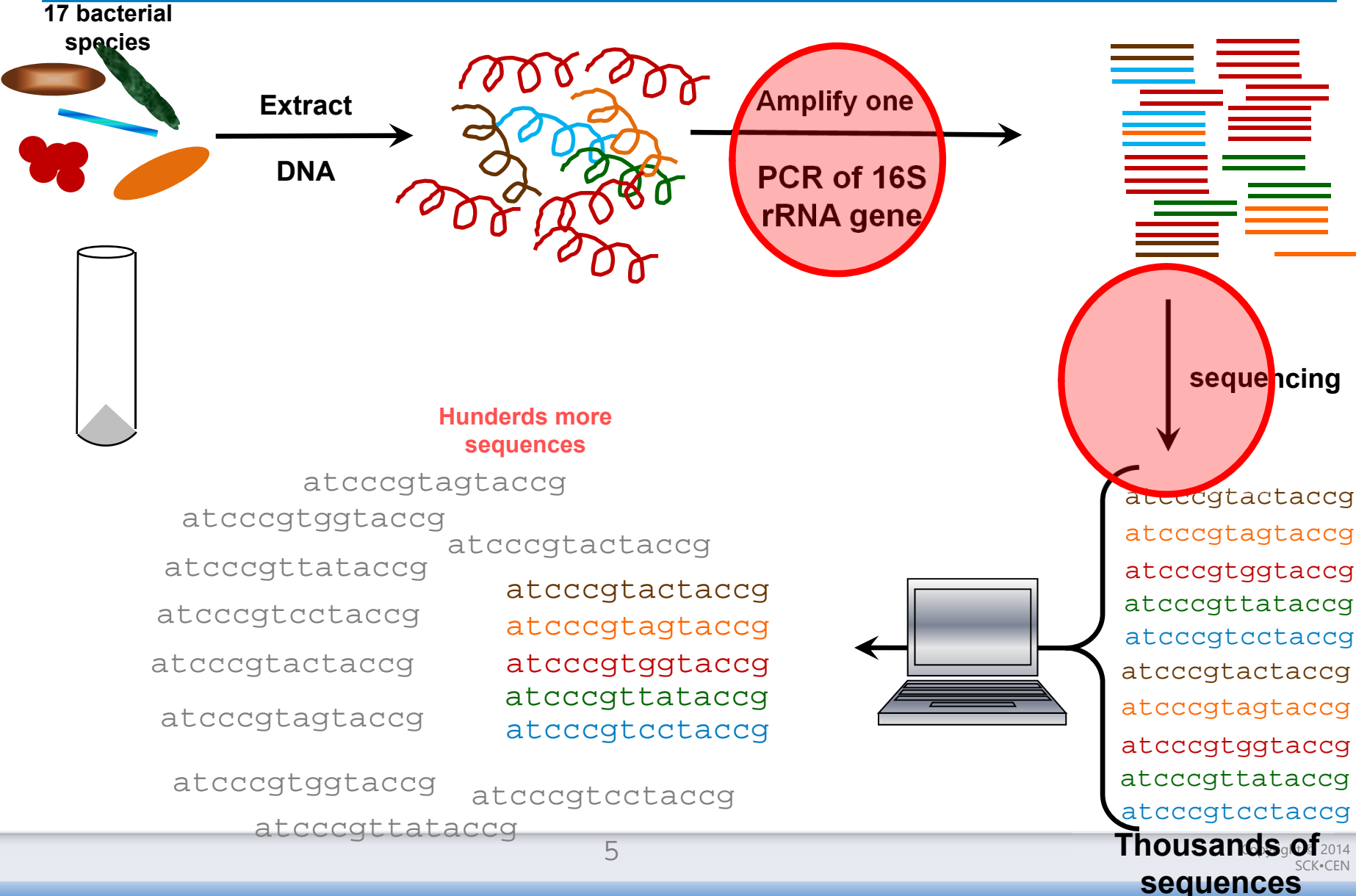


However, culture based approach exhibit various disadvantages:

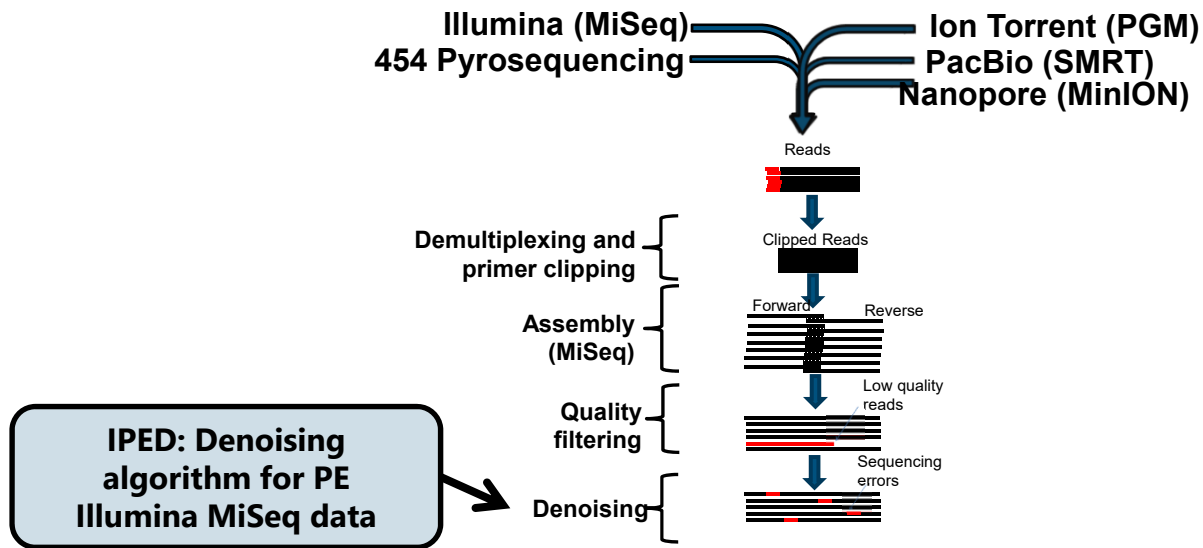
- Time and labor consuming
- Most of the species can not be cultured in lab conditions



# Microbiome & 16S rRNA metagenomics



# 16S rRNA Metagenomics Analysis Pipeline



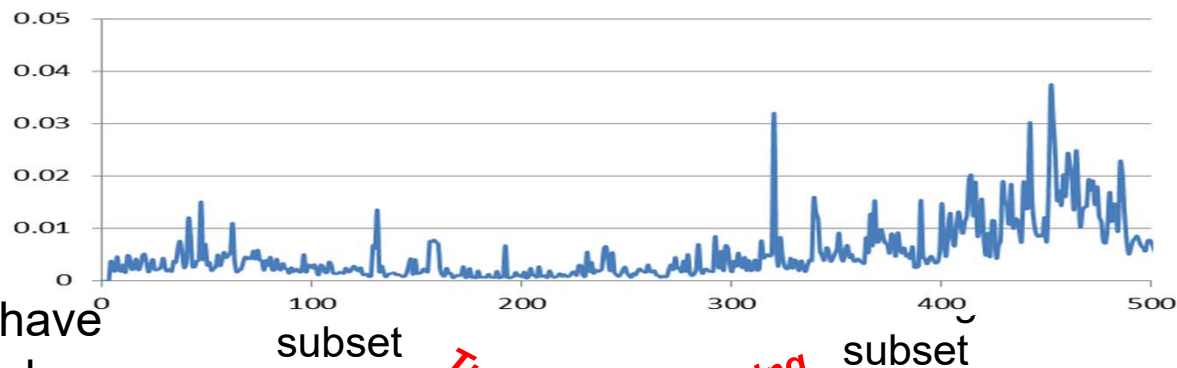
# Chapter 2 & 3: Training artificial intelligent model to handle this problem

I) Selection of datasets where the ground truth is known (Mock dataset)

What does the machine report	A	G	-	A	G	T	C	G	A	T
What should the machine report	A	G	C	A	G	G	C	-	A	T
<b>Status</b>	<b>T</b>	<b>T</b>	<b>D</b>	<b>T</b>	<b>T</b>	<b>S</b>	<b>T</b>	<b>I</b>	<b>T</b>	<b>T</b>



II) Identification of features contributing to the sequencing errors  
e.g. Position in the read



For this purpose, we have developed two artificial intelligence tools:  
A) NoDe for 454  
B) IPED for MiSeq

IV) Selection of the best performing model

# Denoising algorithms concept

ATCCC-TACTACCGA-CCCGTACTACC-G ← Correct (Count = 100)  
ATCCC-TACTACCGA-C~~G~~CGTACTACC-G ← Substitution (Count = 5)  
ATCCC-TACTACCGA-CCCGTACT-CC-G ← Deletion (Count = 3)  
ATCCC-TACTACCGA-CCCGTACTACC~~C~~G ← Insertion (Count = 2)

## The Classifier

- (i) Extracting quality-features for each position (Perl)
- (ii) Running a pre-trained classifier (WEKA using JAVA)
- (iii) Marking the potentially erroneous positions (Perl)

ATCCC-TACTACCGA-CCCGTACTACC-G  
ATCCC-TACTACCGA-C~~X~~CGTACTACC-G  
ATCCC-TACTACCGA-CCCGTACT~~X~~CC-G  
ATCCC-TACTACCGA-CCCGTACTACC~~X~~G

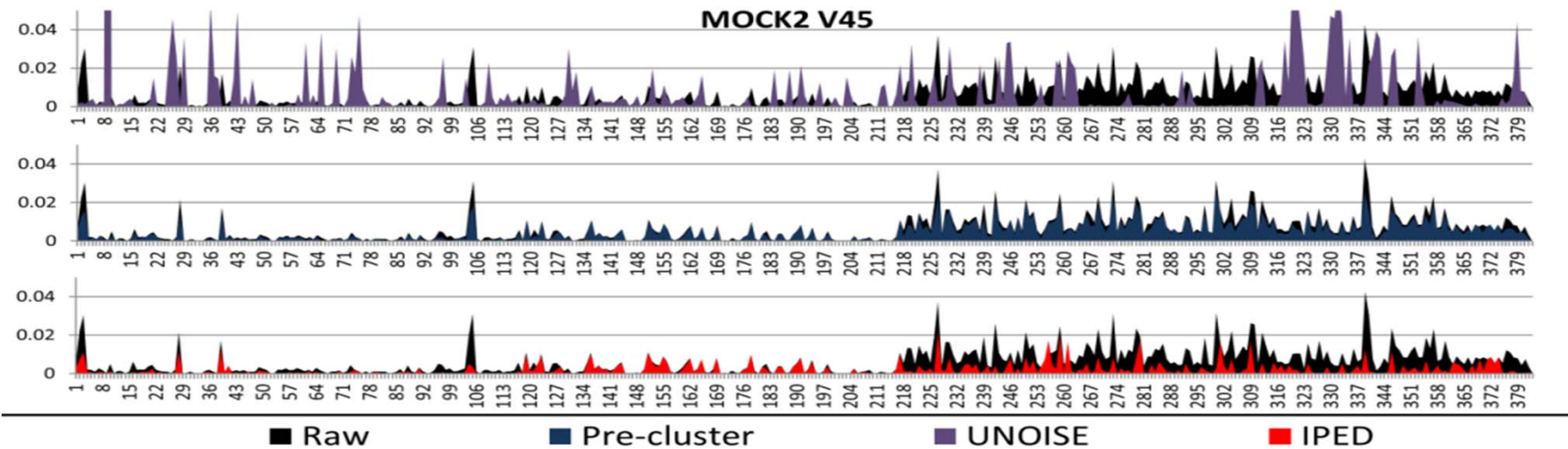
X = correctly marked erroneous positions

## Modified Pre-cluster (mothur using C++)

ATCCC-TACTACCGA-CCCGTACTACC-G ← Representative Read (Count = 110)

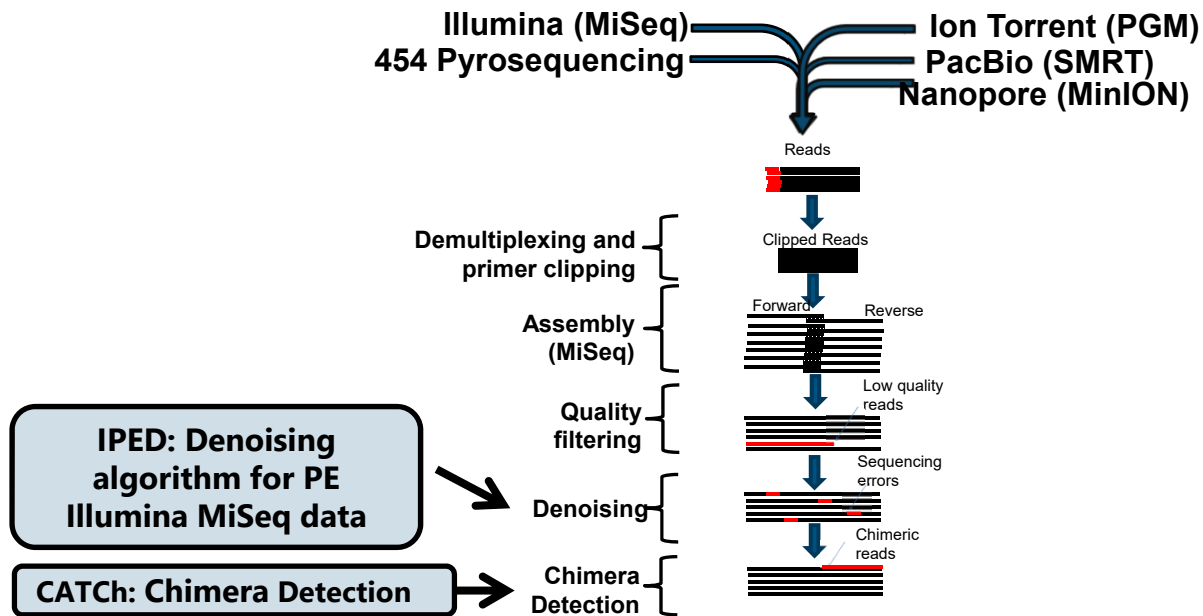


# Chapter 3: IPED comparative analysis



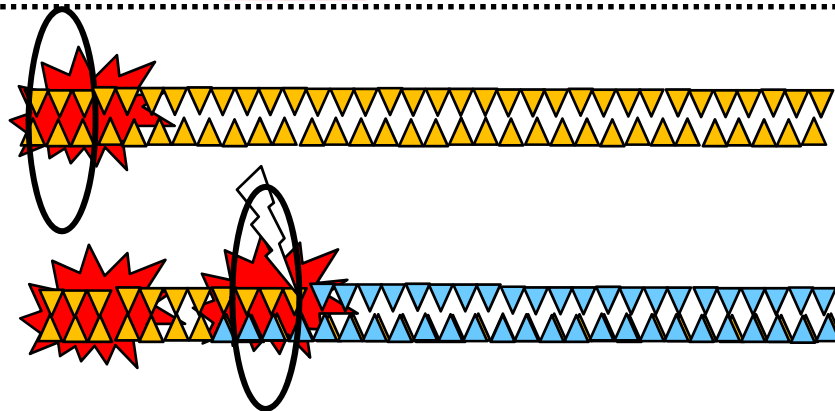
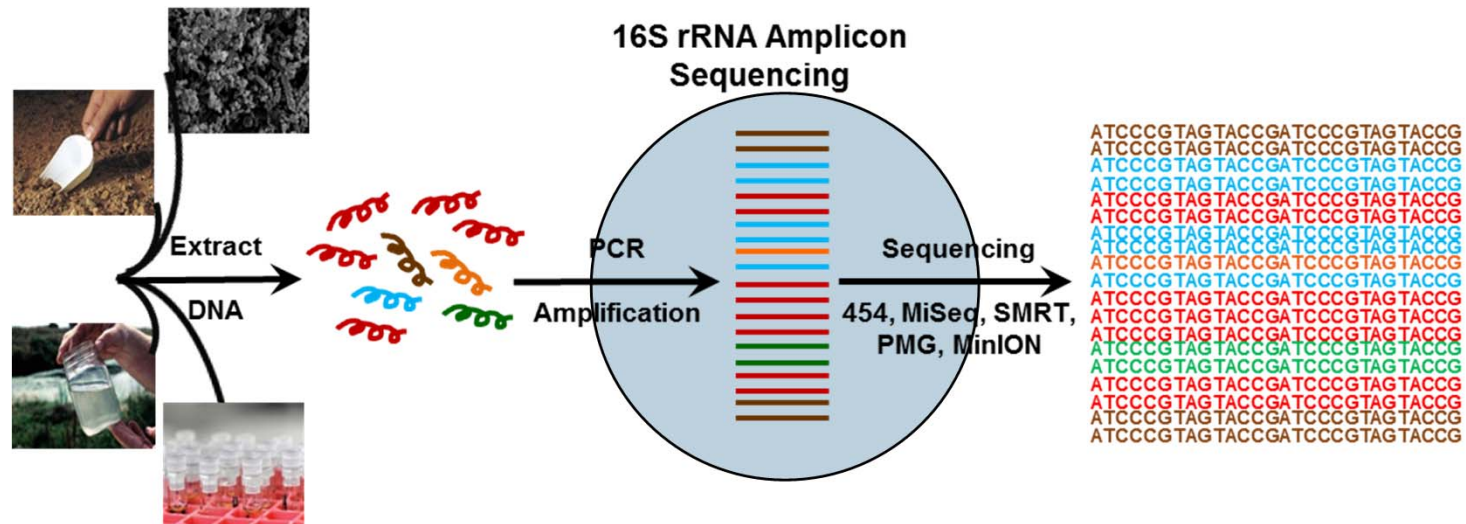
	Error	CPU Cost
UNOISE	0.18%	14 sec
Pre-cluster	0.18%	12 sec
<b>IPED</b>	<b>0.10%</b>	<b>70 sec</b>

# 16S rRNA Metagenomics Analysis Pipeline





# Chimeric problem



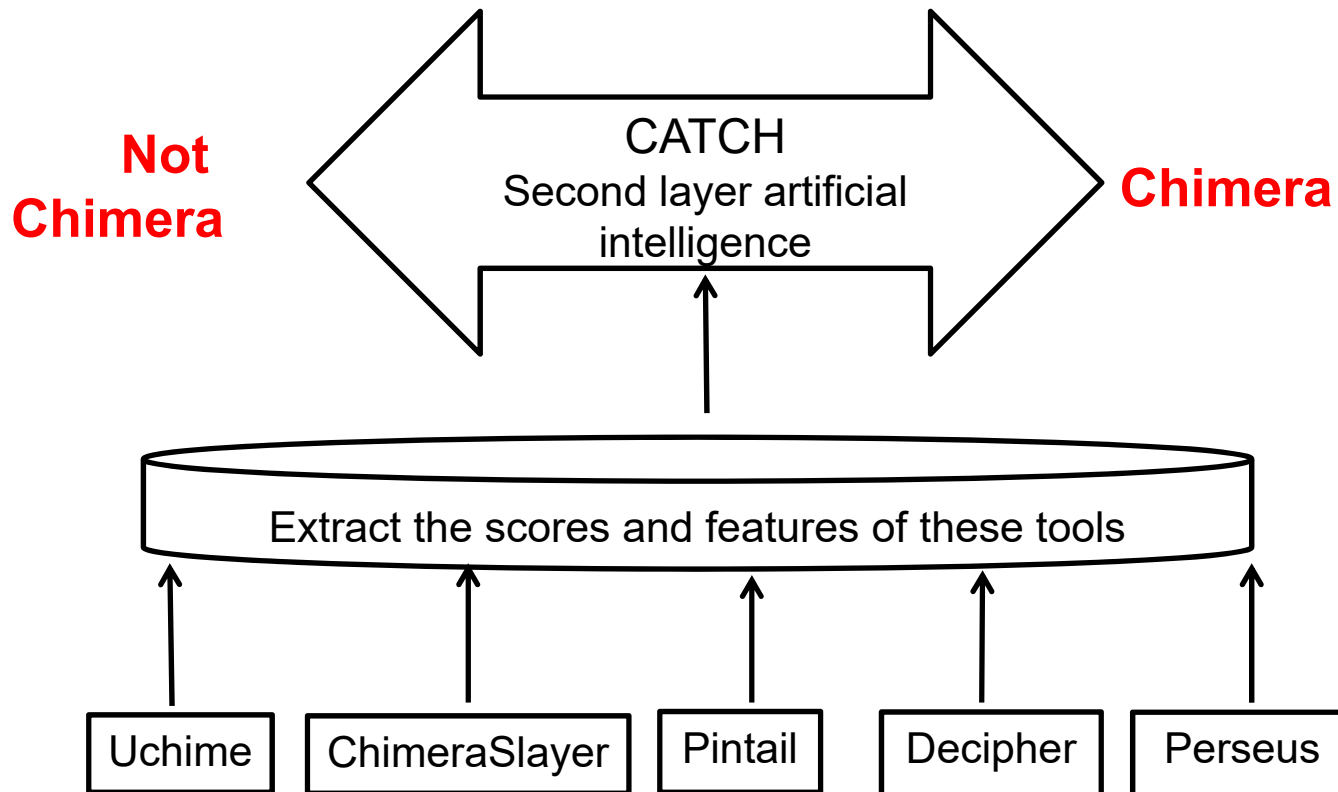
Sequence1  
atcccgtagtaccg

Sequence 2  
tagctacgtacgat

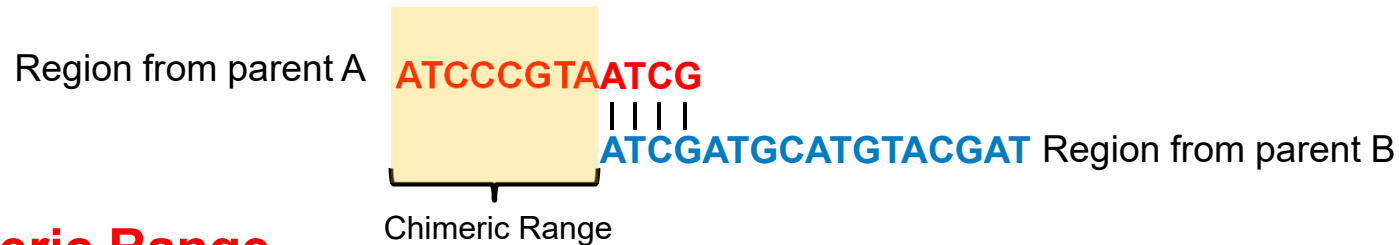
**Chimera**

This hybrid sequence is not real, it could be mistaken as a **false NOVEL** species  
Chimeric rate in Next generation sequencing run can reach up to 45% of the reads.

# CATCh Training/Running

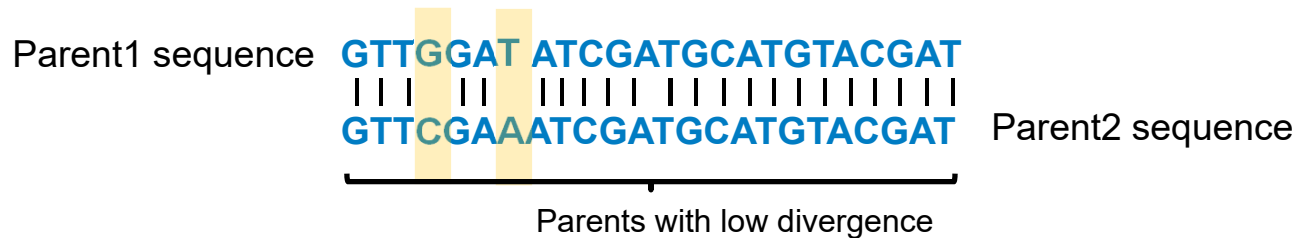


# Chapter 4: Chimera Detection Challenges



## Chimeric Range

Length added by the smaller parent



## Divergence

A measure of the differences between parents

Bimera **ATCCCGTAATGCATGTACGAT**

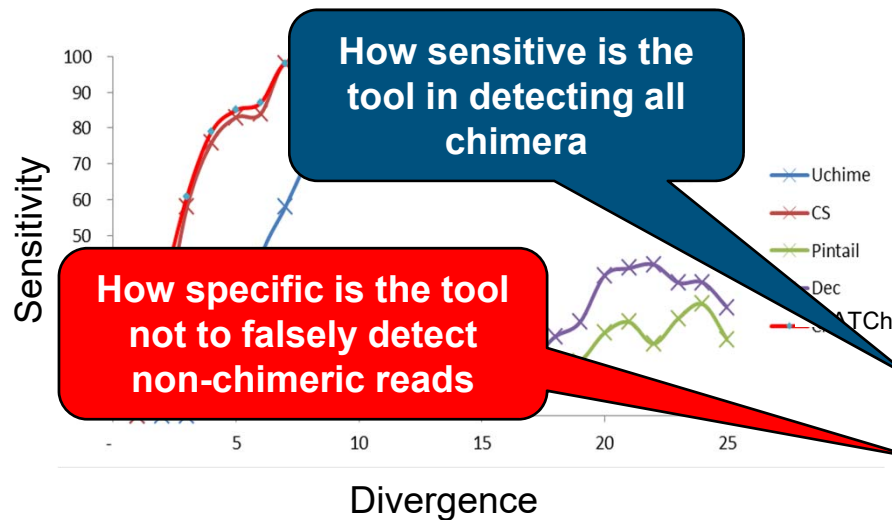
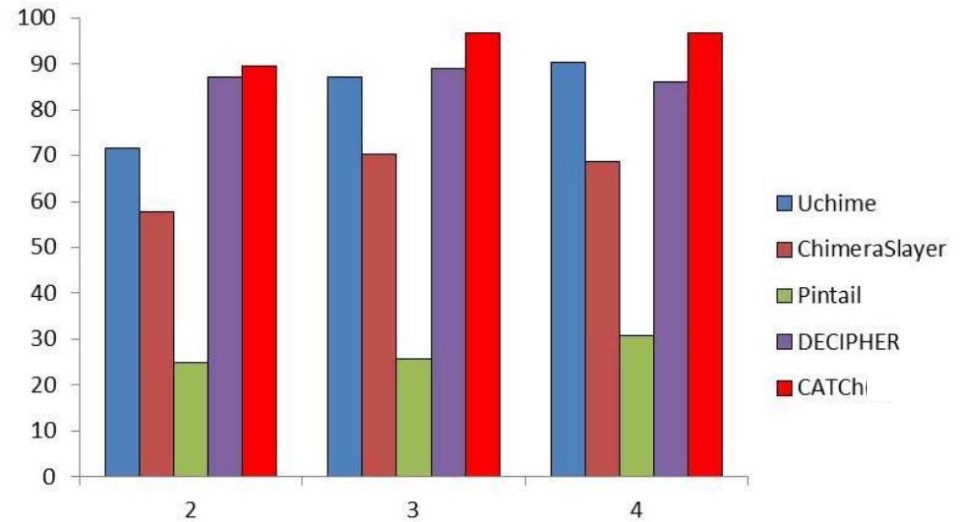
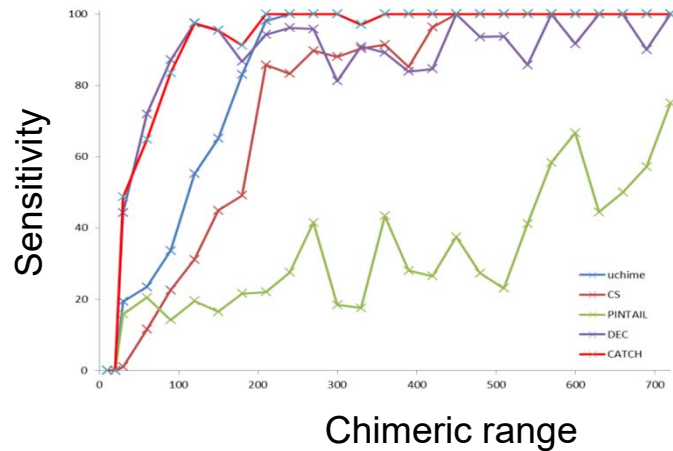
Trimera **ATCCCGTA ATGCATGTTCTAGCTAGC**

Tetramera **ATCCGTA ATGCATCTAGCTAGCATGCAT**

## Number of parents

Number of parent read forming the chimeras

# CATCh Comparative Analysis

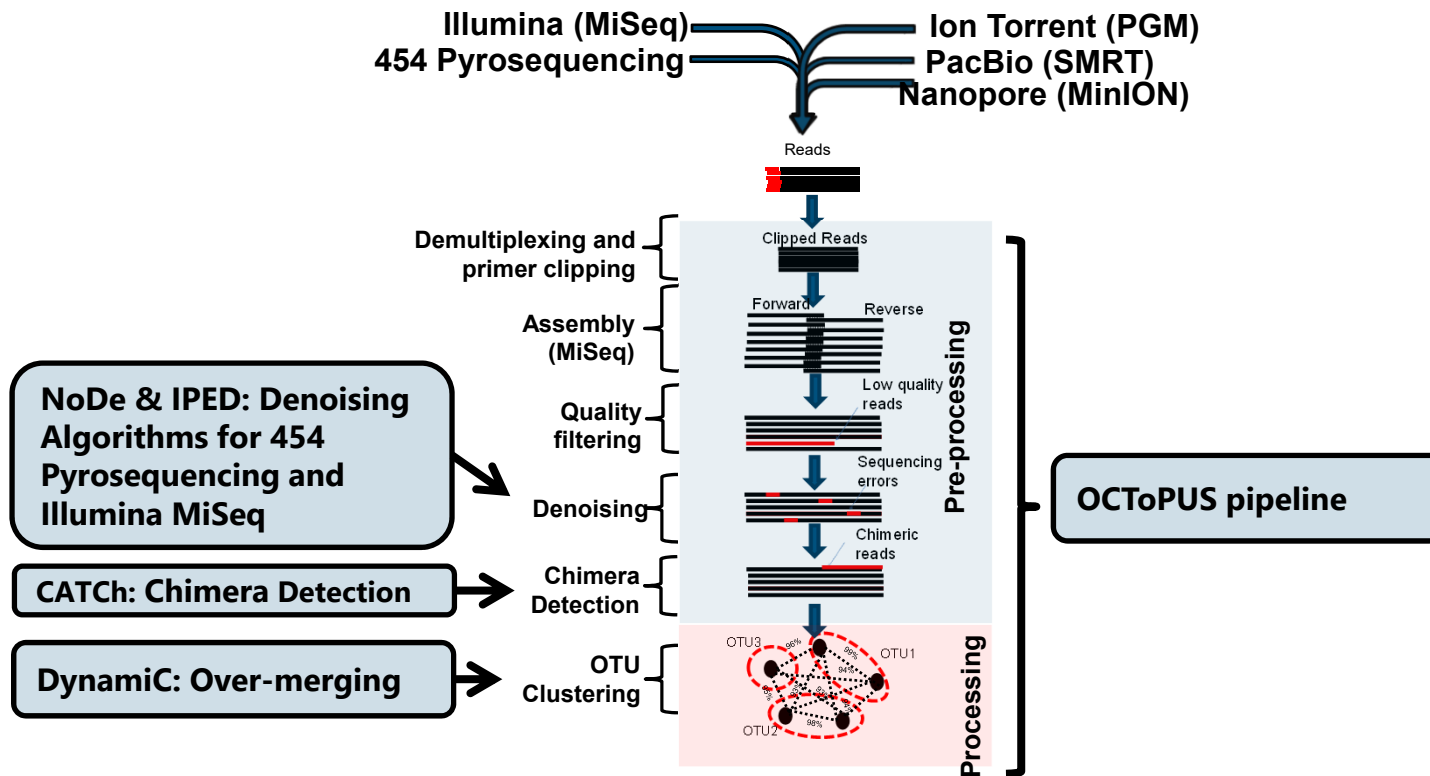


How sensitive is the tool in detecting all chimera

How specific is the tool not to falsely detect non-chimeric reads

Evaluation	Reference tools					<i>De novo</i> tools			
	UCHIME	ChimeraSlayer	Pintail	DECIPHER	CATCh	UCHIME	ChimeraSlayer	Perseus	CATCh
Sensitivity	78	67	29	57	<b>85</b>	60	53	62	<b>70</b>
Specificity	97	98	75	97	<b>96</b>	97	96	96	<b>95</b>

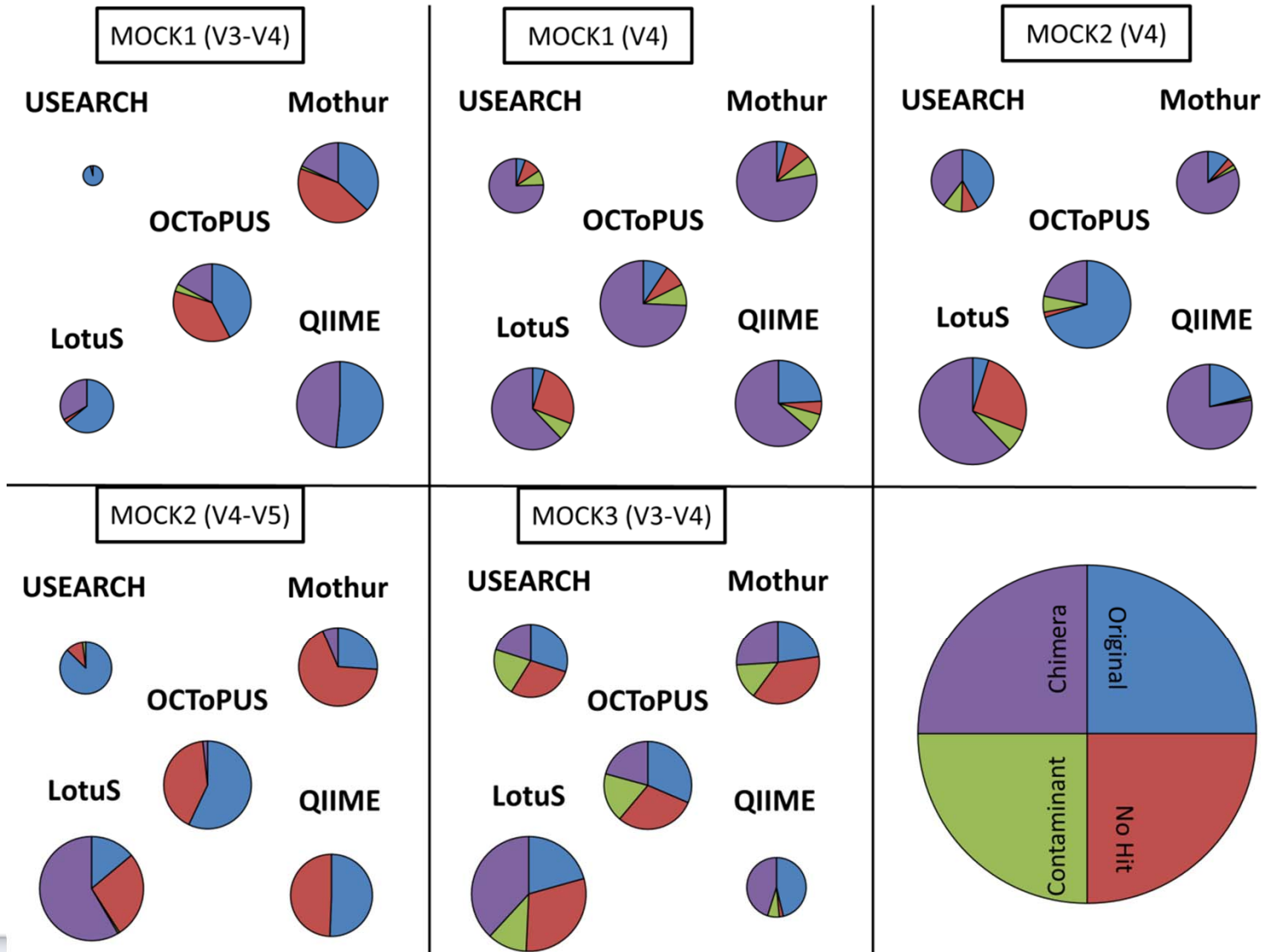
# 16S rRNA Metagenomics Analysis Pipeline







# Benchmark - accuracy



MOCK2 (V4)

USEARCH



Mothur



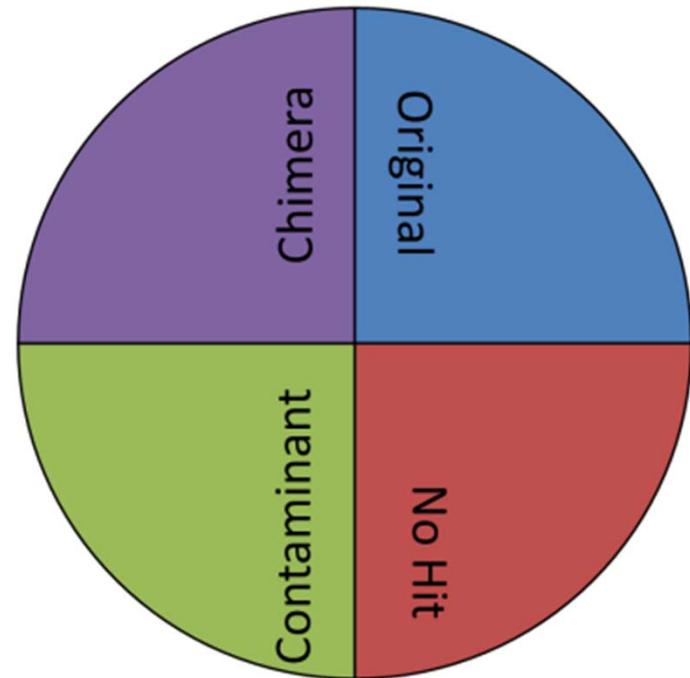
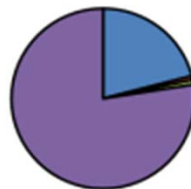
OCToPUS



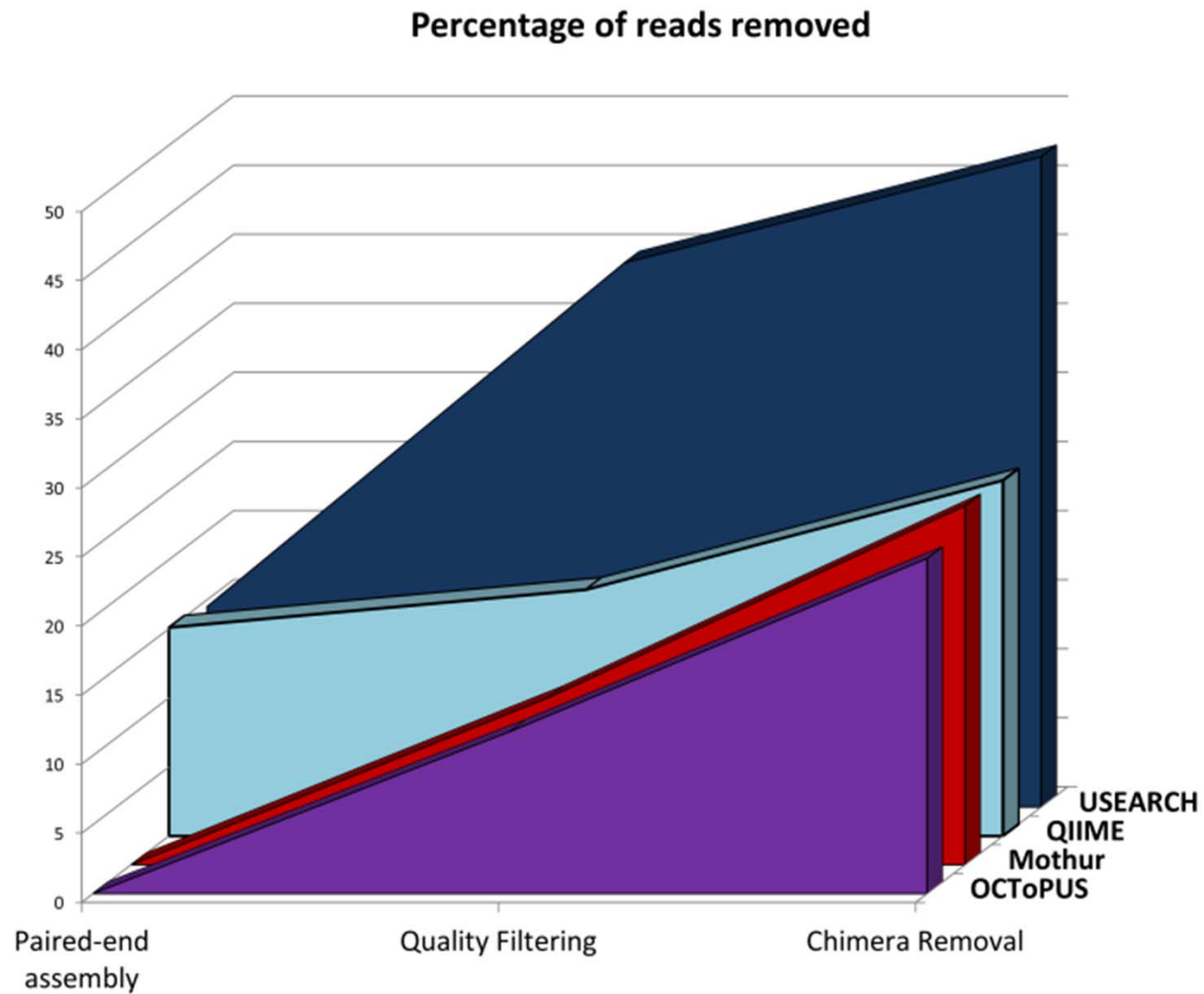
LotuS



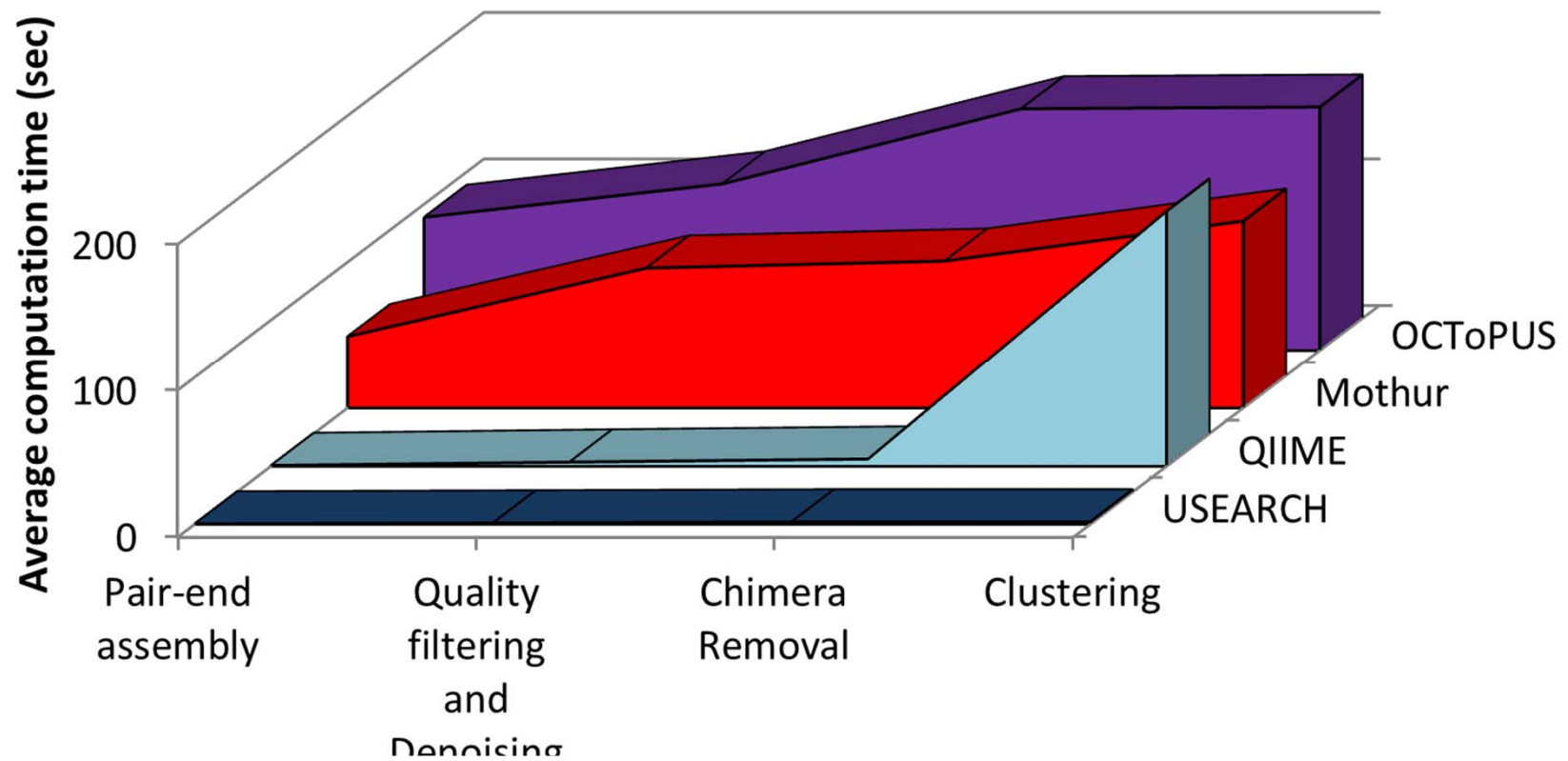
QIIME



# Benchmark – data retrieval



# Benchmark – Computational cost

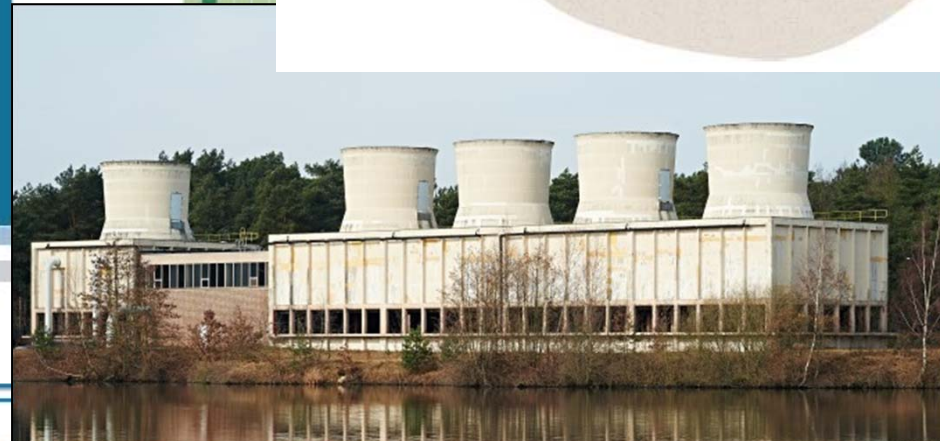
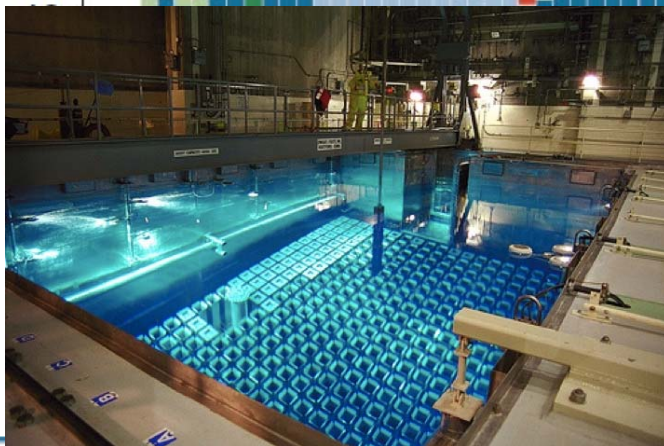
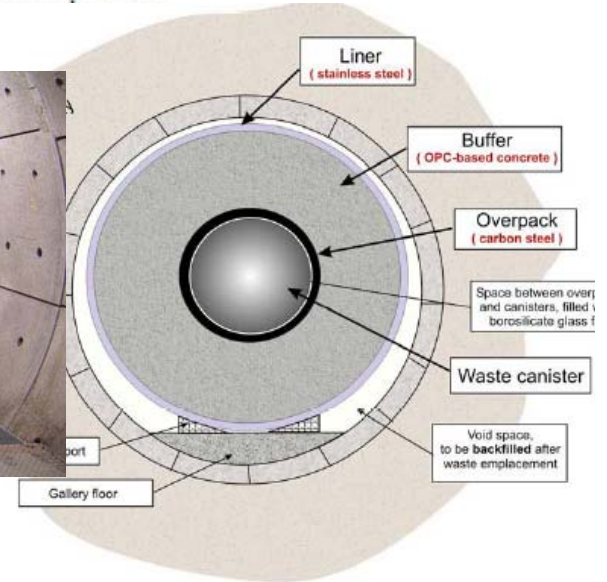


Introduction

Pelvic radiotherapy is a commonly used treatment to treat specific types of cancer (e.g. colon cancer). After exposure to ionizing radiation, the intestine is always affected. The intestinal epithelium is very sensitive to the increased oxidative stress after exposure to ionizing radiation (Riley, 1994), resulting in a loss

- Chitinophagaceae
- Comamonadaceae
- Chitinophagaceae

- Acidobacteria subdivision 3 (incertae sedis)



Conclusion

In future experiments: first the dose of the irradiation has to be optimized for our mouse model in order to exhibit relevant clinical symptoms. Second, a more specific examination of apoptosis and examination of oxidative stress will be performed. And third, a larger subset of inflammation markers will also be used. In addition, the formulation and dose of *Arthrospira* sp. as food supplement will also be further optimized.

# Summary

Illumina (MiSeq)  
 454 Pyrosequencing  
 Ion Torrent (PGM)  
 PacBio (SMRT)  
 Nanopore (MinION)

Reads

Demultiplexing and primer clipping

Assembly (MiSeq)

Quality filtering

Denoising

Chimera Detection

OTU Clustering

Pre-processing

Processing

Post-processing

Biodiversity analysis

NoDe & IPED: Denoising Algorithms for 454 Pyrosequencing and Illumina MiSeq

CATCh: Chimera Detection

DynamiC: Over-merging

## Publications

*CATCh: Mysara et al., AEM (2015)*

*NoDe: Mysara et al. BMC Bioinformatics (2015)*

*iPED: Mysara et al. BMC Bioinformatics (2016)*

*DynamiC: Mysara et al. GigaScience (2017)*

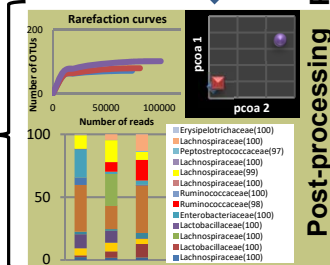
*OCToPUS: Mysara et al. FEMS Ecology (2017)*

OCToPUS pipeline

<https://github.com/M-Mysara/OCToPUS>

## Future work – Extend expertise

- New long-read sequencing technologies (NanoPore, PacBio)
- Complex mock communities (n > 200)
- Shotgun metagenomics



- New long-read sequencing technologies

- NanoPore



- PacBio



- Challenging 16S rRNA amplicon sequencing pipelines

- Complex mock communities ( $n > 200$  strains)

- Strain / subspecies variation detection

- Shotgun metagenomics

- Taxonomy versus metabolic potential



# Acknowledgements

Dr. Mohamed Mysara  
Ruben Props  
Dr. Natalie Leys

Prof. Dr. Peter Vandamme  
Prof. Dr. Yvan Saeys

Prof. Dr. Jeroen Raes

Prof. Dr. Daniel Charlier



**KU LEUVEN**



**Copyright © 2014 - SCK•CEN**

PLEASE NOTE!

This presentation contains data, information and formats for dedicated use ONLY and may not be copied, distributed or cited without the explicit permission of the SCK•CEN. If this has been obtained, please reference it as a "personal communication. By courtesy of SCK•CEN".

**SCK•CEN**

Studiecentrum voor Kernenergie  
Centre d'Etude de l'Energie Nucléaire  
Belgian Nuclear Research Centre

Stichting van Openbaar Nut  
Fondation d'Utilité Publique  
Foundation of Public Utility

Registered Office: Avenue Herrmann-Debrouxlaan 40 – BE-1160 BRUSSELS  
Operational Office: Boeretang 200 – BE-2400 MOL



STUDIECENTRUM VOOR KERNENERGIE  
CENTRE D'ETUDE DE L'ENERGIE NUCLEAIRE

# Microarrays

- Types

- Two-color arrays: *Rhodospirillum rubrum*, *Cupriavidus metallidurans*
- Affymetrix: *Pseudomonas aeruginosa*
- Nimblegen: *Arthrospira sp.* PCC8005
- Agilent: *Cupriavidus metallidurans*

- In-house facilities

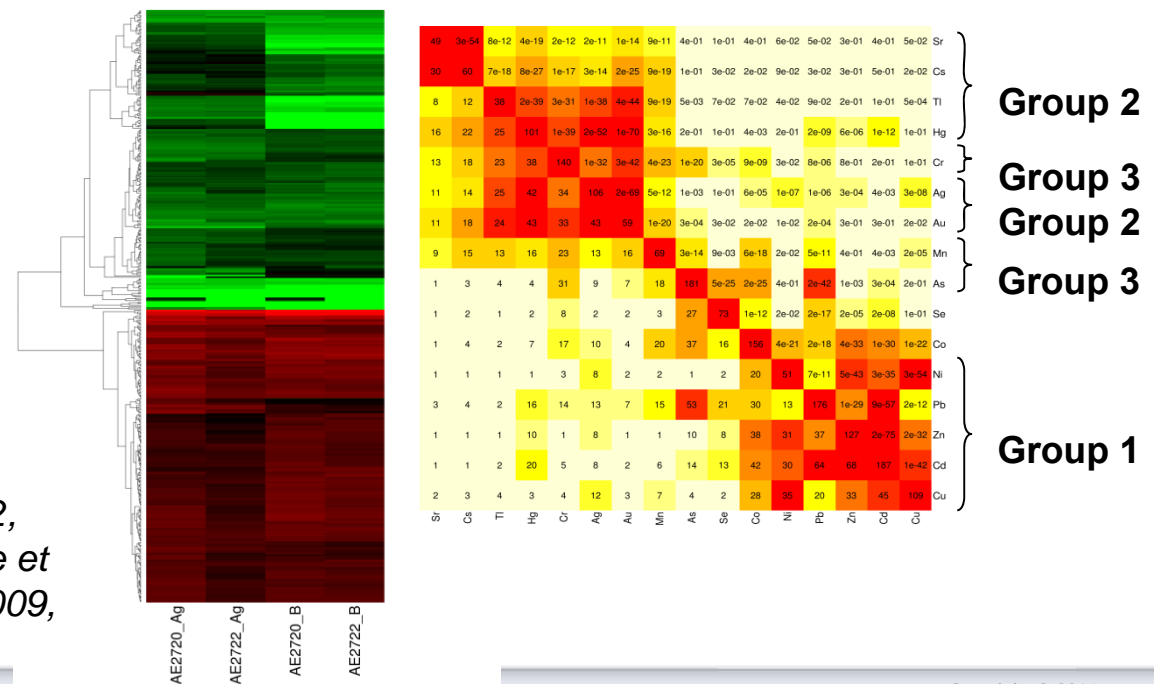
- Two-color microarrays
- Affymetrix arrays

- Analysis

- MIC: BioConductor
- RDB: Partek

Monsieurs et al., 2011, Van Houdt et al. 2012, Crabbe et al. 2011, Crabbe et al. 2010, Pycke et al, 2010, Badri et al. 2014, Mastroleo et al. 2009,

...



# Transcriptomics via RNA-seq

- Application area

- MIC

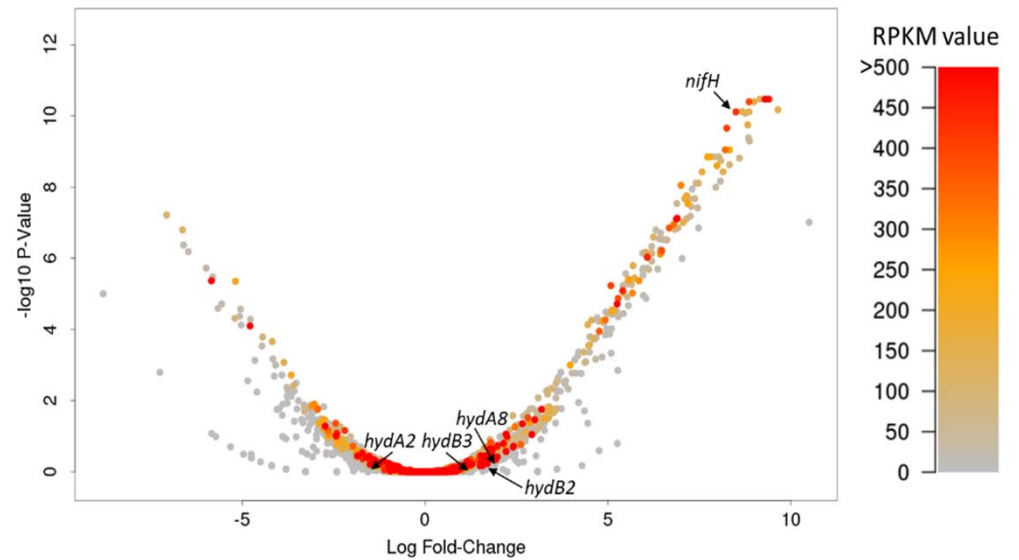
- *Arthrospira* sp. PCC8005
    - *Clostridium butyricum*
    - *Pseudomonas aeruginosa*

- BIS

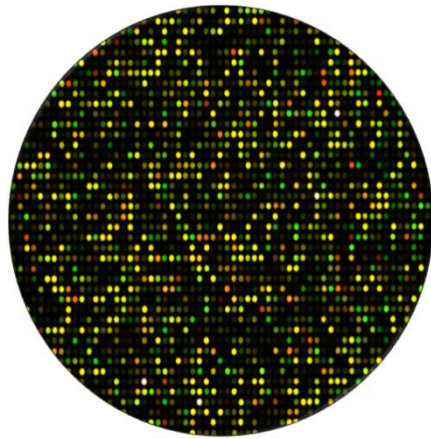
- To be started

- Analysis

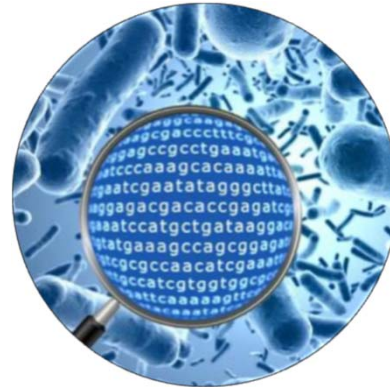
- BWA / Bowtie
  - EdgeR



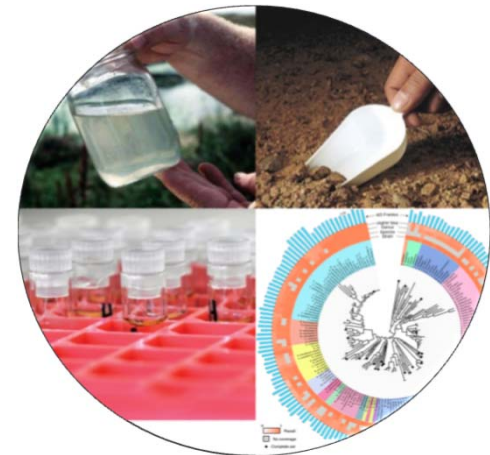
# Overview



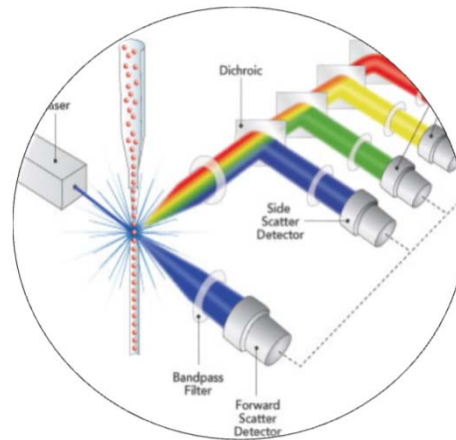
**Transcriptomics**



**Genomics**



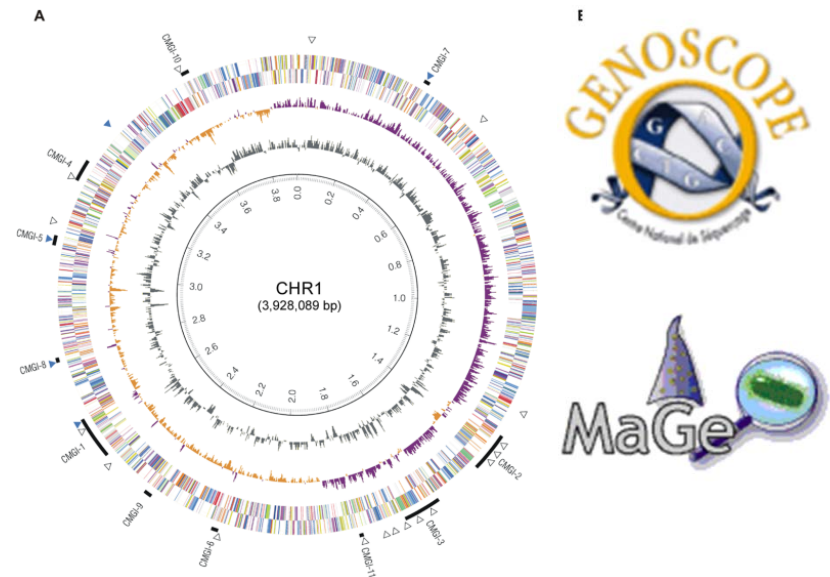
**Metagenomics**



**Flow cytometry**

# Genome assembly

- *De novo* genome assembly
  - Based on 454 pyrosequencing data
  - Based on Illumina HiSeq data
  - Annotation → GenoScope - MaGe
    - CmetScope
    - ArthroScope
- Applications
  - MIC: Metal resistant bacteria
  - BIS: *Lemna minor* (~ 400 Mbp genome)



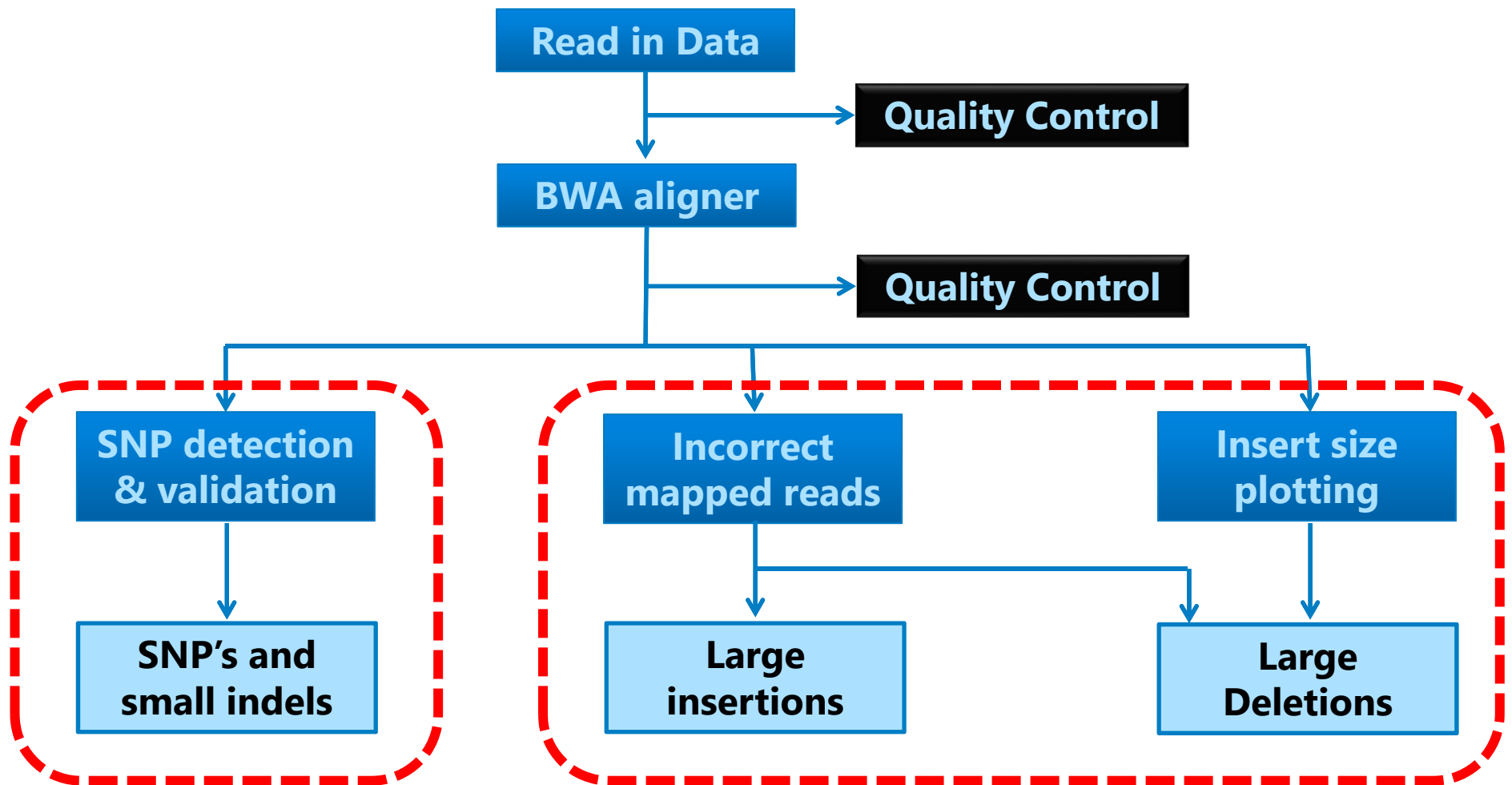
Platform	Insert Size	Read Length	No. Reads	No. Nucleotides	Genome Coverage
Illumina HiSeq 2000	200 bp	2*100 bp	207.985.822	40.730.561.447	100X

Statistics	SOAPdenovo2 <sup>2</sup>	CLC Bio	CLC Bio
Input	preprocessed data	preprocessed data	processed reads + flash data
K-mer size	63	53	53
No. Scaffolds	108607	116254	117403
Max scaffold length	110,000	110,000	110,000
genome length	401	388	410
N50	10194	8059	8423

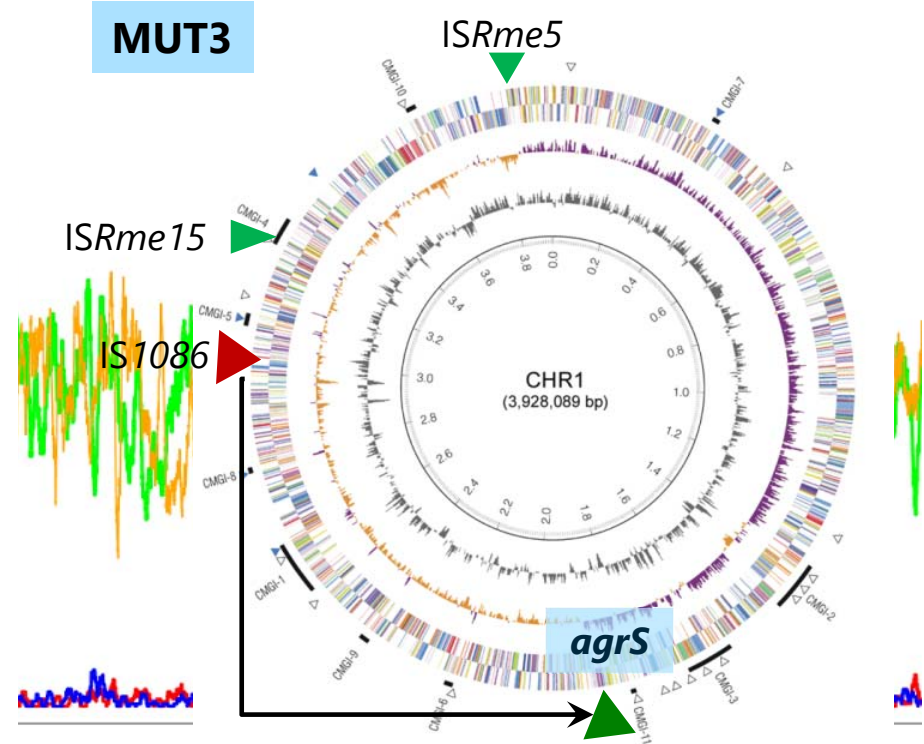
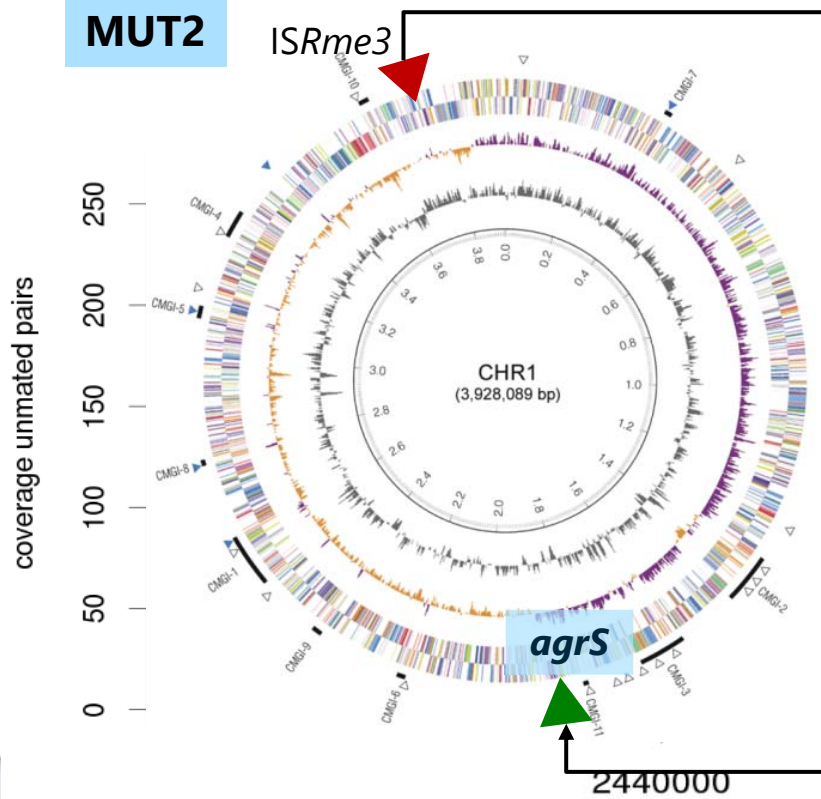
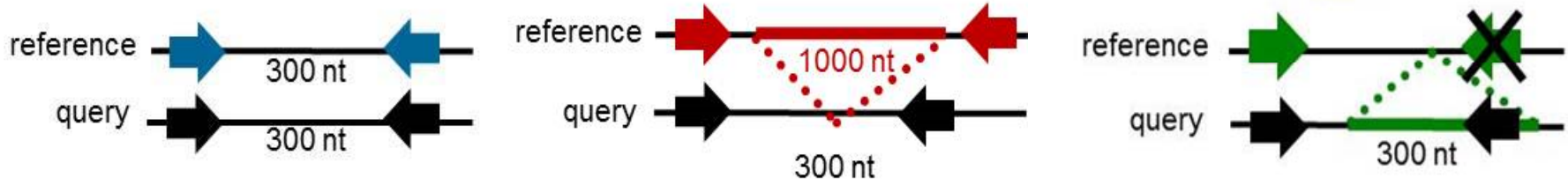
Janssen et al., 2010, Monsieurs et al., 2013, Monsieurs et al., 2014a, Monsieurs et al., 2014b

Arne Van Hoeseke

# Bacterial resequencing - workflow



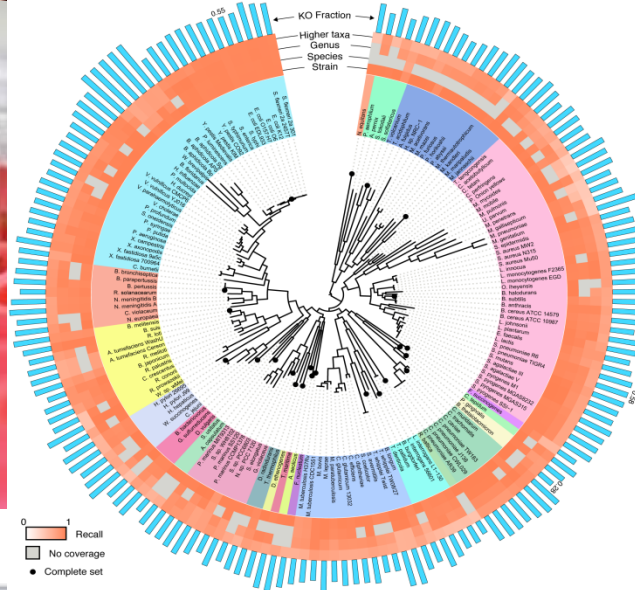
**BWA / Samtools /  
GATK genomic suite**





- Bio-informatics server – calculations
  - 24 processors – 3.2 GHz
  - 96 Gb memory
  - 6.0 Tb hard drive (RAID5 configuration)
- Data server - storage
  - 1 Tb hard drive in RAID10 configuration
  - Automatic backup to EqualLogic tapes
  - Storage of all raw microarray and sequencing data
- FERMI cluster
  - **Details to be added**

- To be added?: Slides Roel on Affymetrix and Partek



# 16S rRNA metagenomics algorithms

## Article Contents

Abstract

Supplementary data

ACCEPTED MANUSCRIPT

## From reads to operational taxonomic units: an ensemble processing pipeline for MiSeq amplicon sequencing data

Mohamed Mysara; Mercy Njima; Natalie Leys; Jeroen Raes; Pieter Monsieus

Gigascience giw017. DOI: <https://doi.org/10.1093/gigascience/giw017>

Published: 18 January 2017

Views PDF Cite Share Tools

### Abstract

**Introduction:** The development of high-throughput sequencing technologies has provided microbial ecologists with an efficient approach to assess bacterial diversity at an unseen depth, particularly with the recent advances in the Illumina MiSeq sequencing platform. However, analysing such high-throughput data is posing important computational challenges, requiring specialized bioinformatics solutions at different stages during the processing pipeline, such as

## Operational taxonomic units

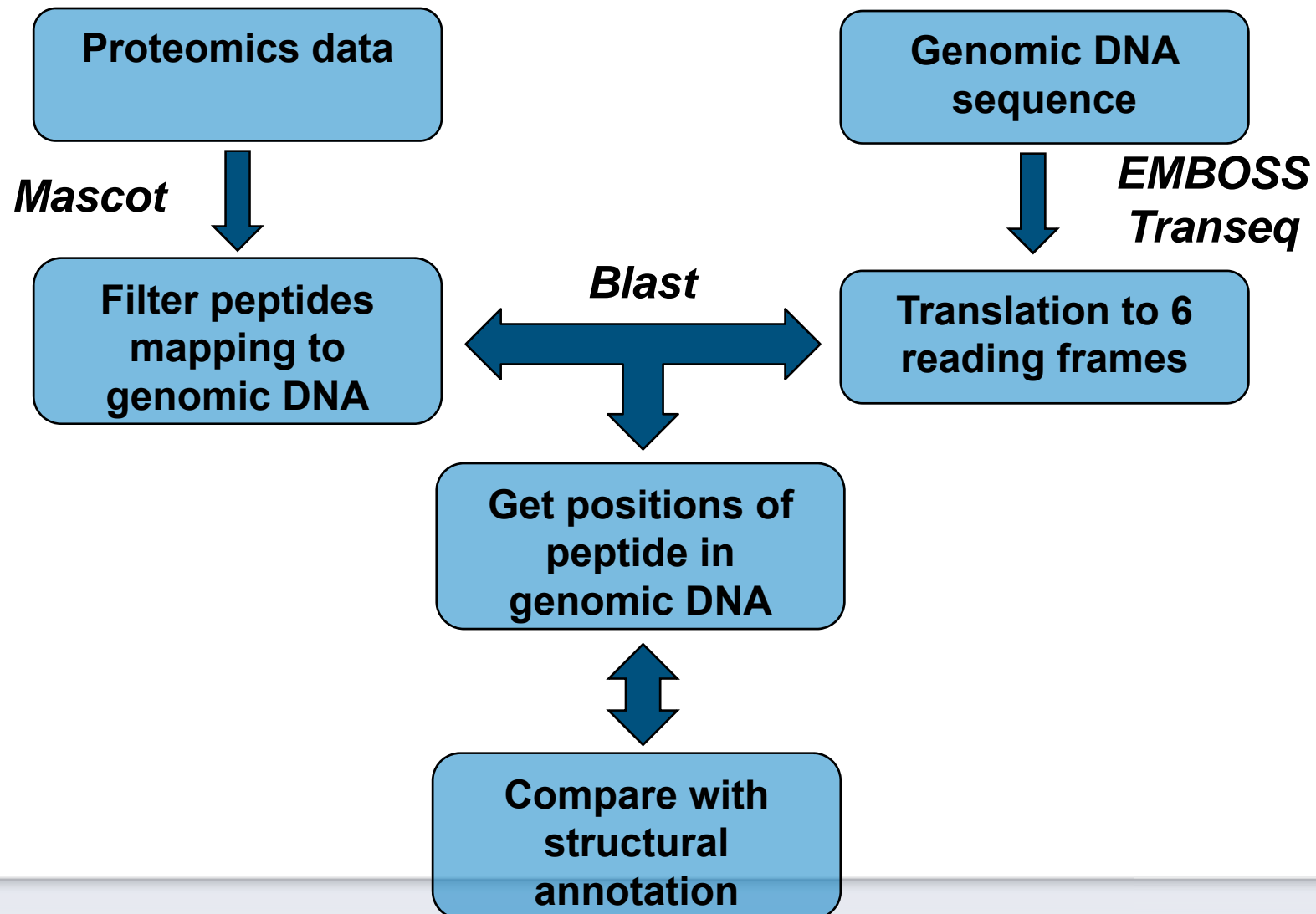
Ben Props<sup>6</sup>,  
Nico Boon<sup>6</sup>, Jeroen Raes<sup>3,4</sup>

Department of Microbiology, Ghent University, Ghent, Belgium, <sup>2</sup>Department of Microbiology, Ghent University, Ghent, Belgium, <sup>3</sup>VIB lab for Bioinformatics and Microbiology and Immunology, REGA institute, Ghent University, Ghent, Belgium and <sup>4</sup>Department of Microbiology, Ghent University, Ghent, Belgium

Our ensemble method (combining reference-based and *de novo* CATCH) were developed by integrating a powerful method. When comparing our classifiers with existing methods, the performance of our ensemble method was observed on a wide range of Illumina MiSeq data sets. Since our algorithm combines reference-based and *de novo* approaches, it produces more robust results when challenged with a wide range of chimeric sequences, and various numbers of parents. Adding this pipeline has a beneficial effect on the quality of the clustering results.

**background:** The development of high-throughput sequencing technologies has revolutionized the field of microbial ecology via the sequencing of phylogenetic marker genes (e.g. 16S rRNA gene amplicon sequencing). Denoising, the removal of sequencing errors, is an important step in preprocessing amplicon sequencing data. The increasing popularity of the Illumina MiSeq platform for these applications requires the development of appropriate

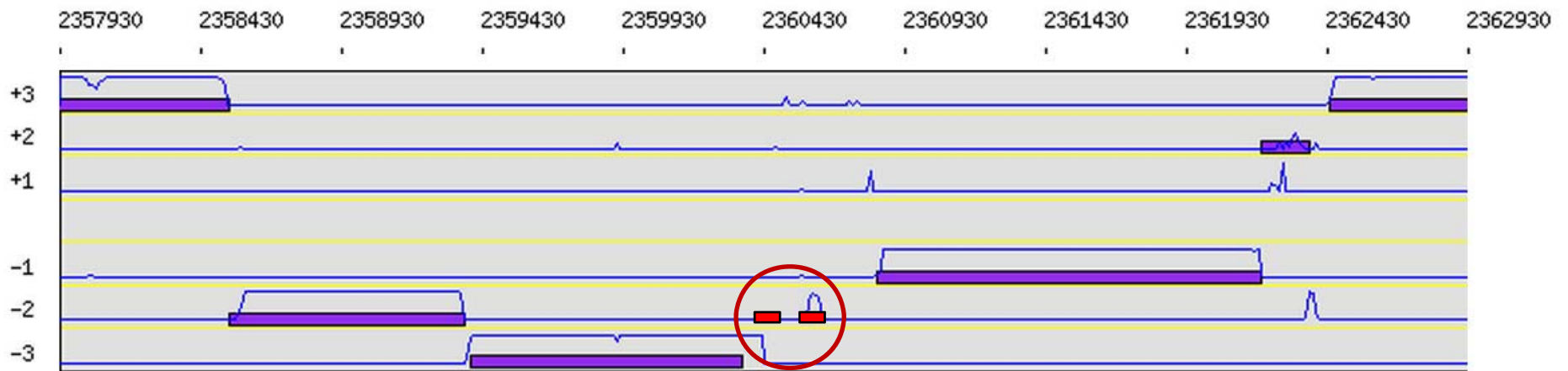
# Proteogenomics pipeline



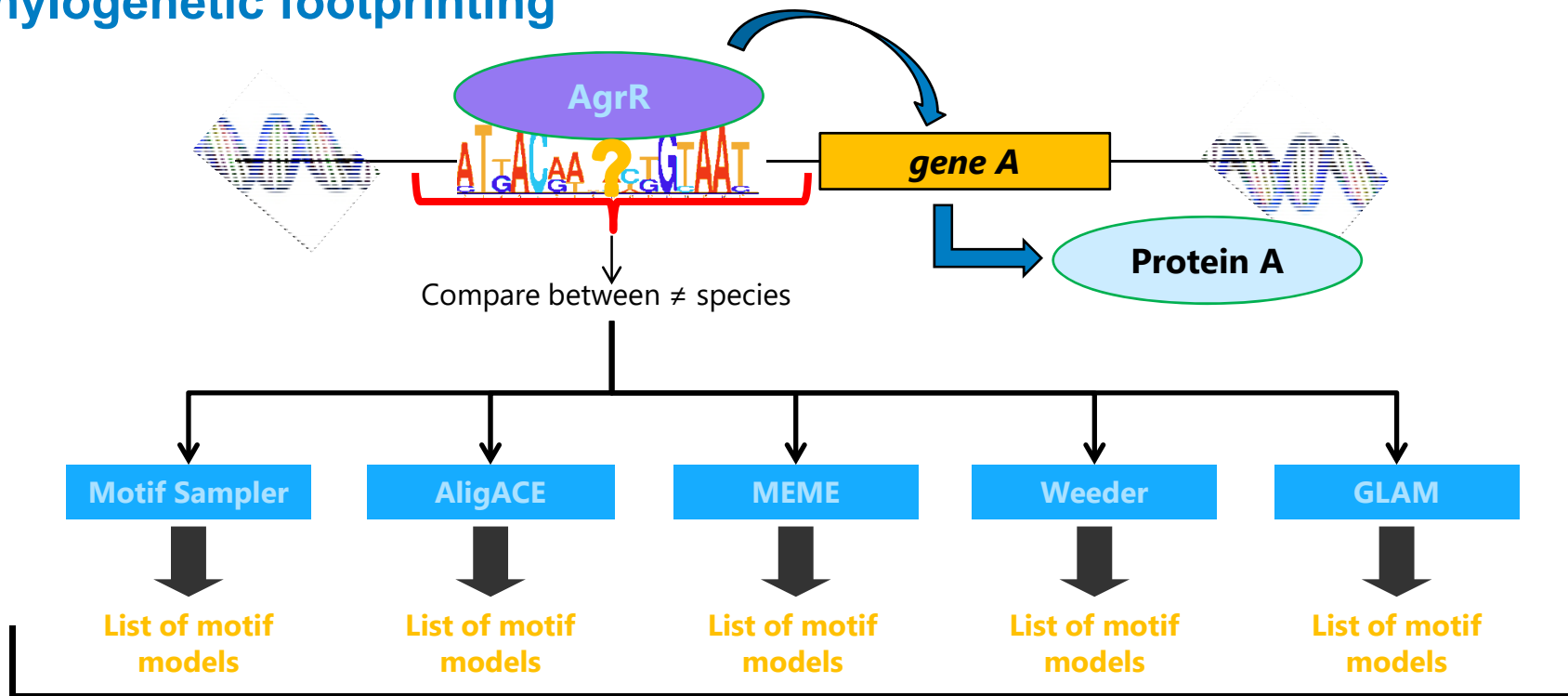


***Rhodospirillum rubrum* ATCC 11170 - chromosome Rru\_A NC\_007643**  
***Rhodospirillum rubrum* ATCC 11170 - chromosome Rru\_A NC\_007643**  
**2357930 -- 2362930**

( sequence length : 4352825 bases )

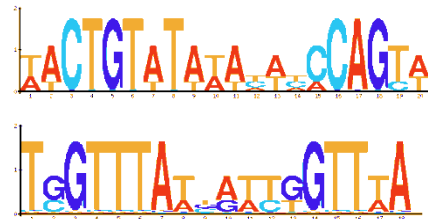


## Phylogenetic footprinting



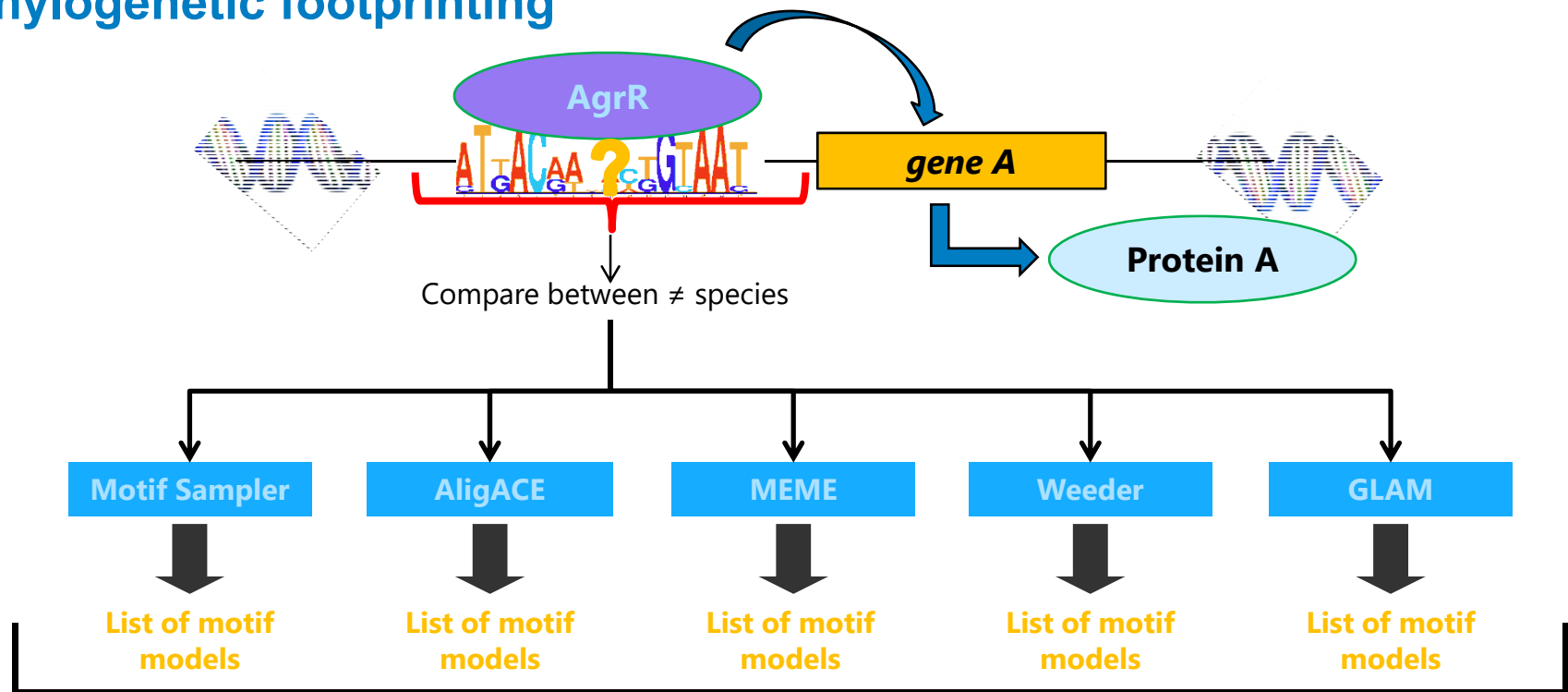
### Cluster motif models:

Markov clustering (MCL) of output motif models based on the Pearson Correlation Coefficient

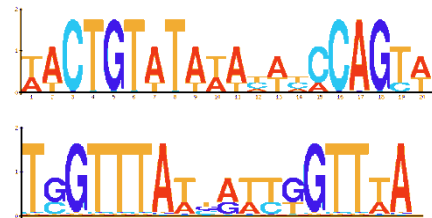


# Other pipelines

## Phylogenetic footprinting



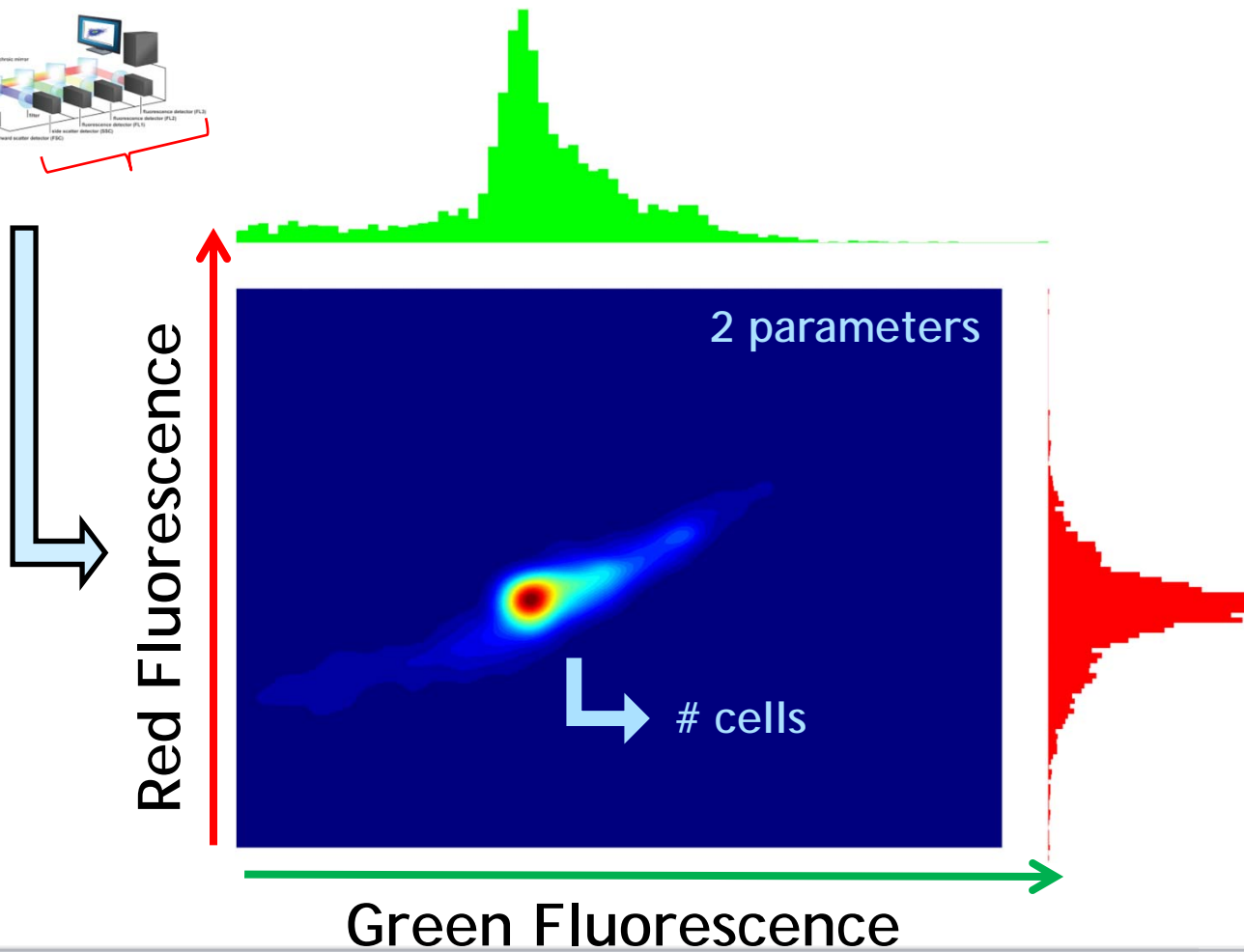
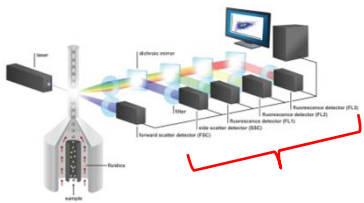
**Cluster motif models:**  
Markov clustering (MCL) of output motif models based on the Pearson Correlation Coefficient





# Phenotypic Biomarkers

Multivariate data

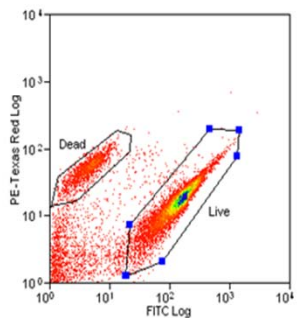




# Flow Cytometry

A

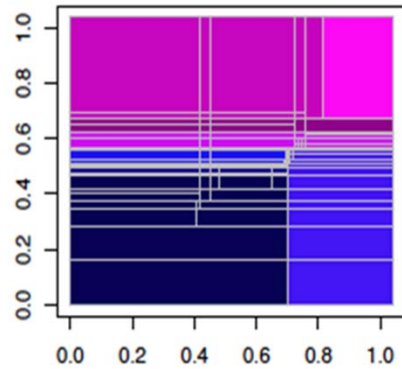
Flow cytometric analysis



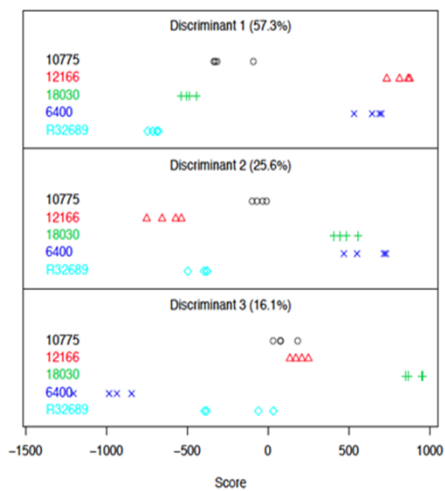
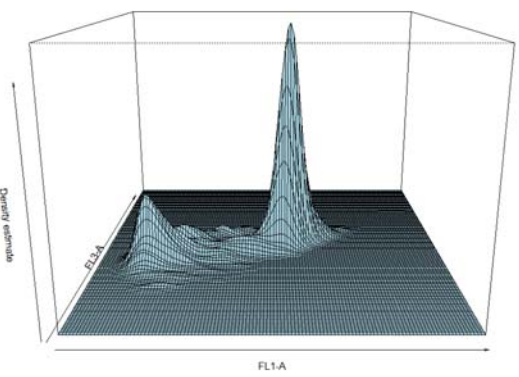
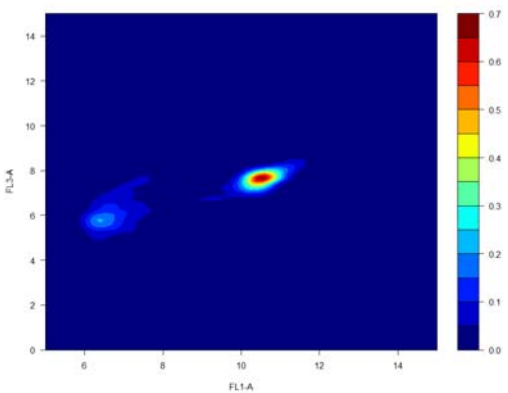
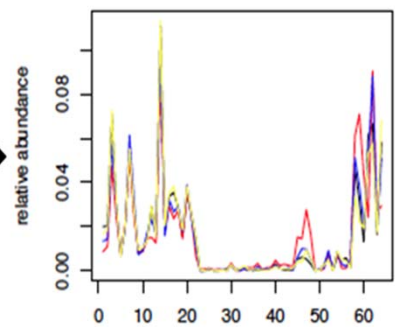
Data extraction

Cell	MSlow	MSlow	MSlow	MSlow	MSlow	MSlow	MSlow
1	6690	0	1420	0	1070	9396	0
2	0	0	842	481	11209	0	0
3	1028	0	471	424	9788	1891	0
4	18788	52278	29700	49890	49899	42888	28920
5	1220	0	828	2894	1291	8837	0
6	2879	0	697	2408	1428	804	0
7	0	0	1892	1288	1239	17	0
8	18798	20892	82825	82467	82284	82112	42812
9	14447	28817	28082	47781	48191	28788	28820
10	0	0	1788	2229	8712	7848	0
11	18678	28888	82147	82224	81717	82024	42888
12	0	0	2482	12071	9818	10184	0
13	8848	0	28890	28007	21888	22072	14248
14	14248	28871	4221	12887	14727	4181	0
15	0	0	8298	8888	12804	20222	8888
16	0	0	2892	8200	18228	4924	0
17	4884	0	12784	42284	28918	20222	2204
18	0	0	1717	848	8248	0	0
19	0	0	2821	18818	17284	11782	0
20	0	0	1788	12288	12828	12880	0

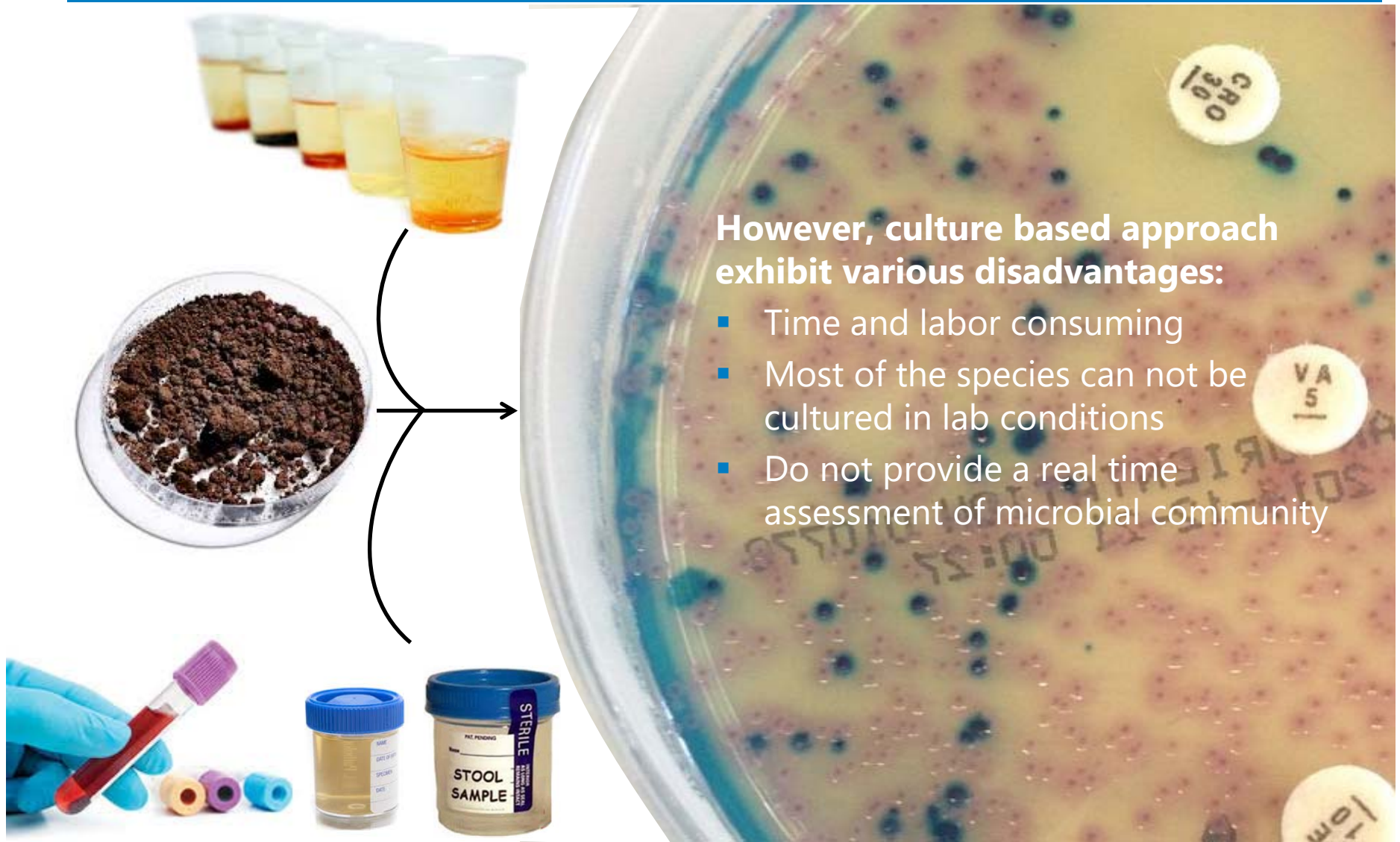
Create Fingerprint model



Create Fingerprints



# Culture-based approach



# Flow cytometry – test case

