

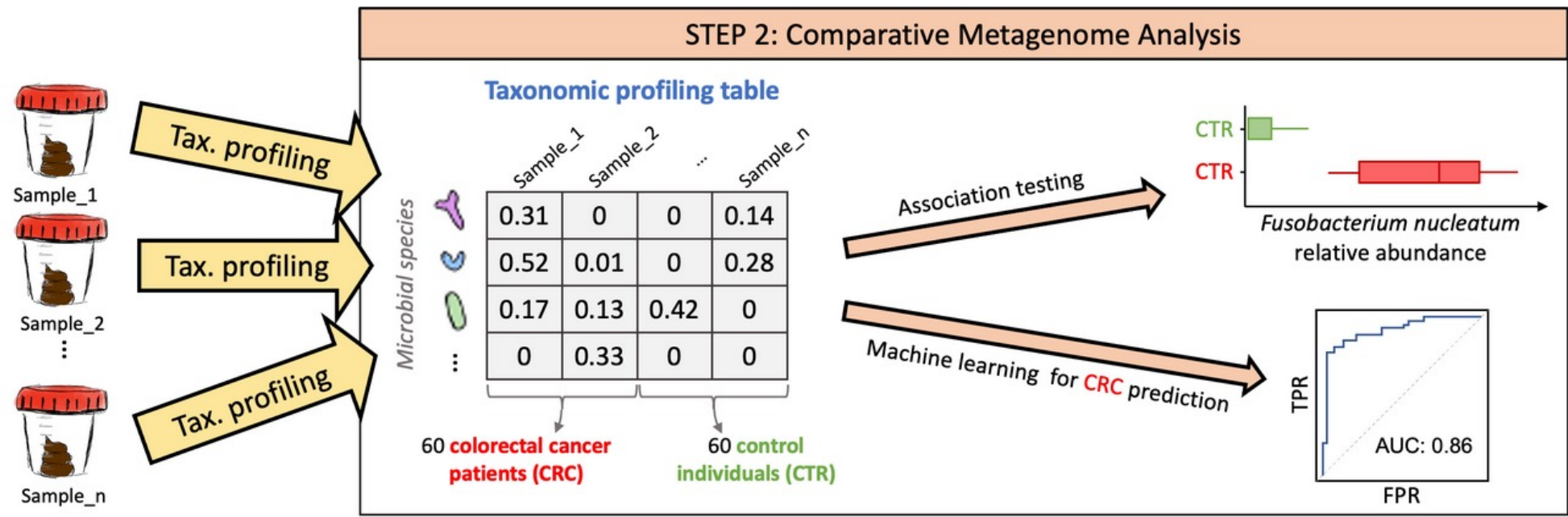
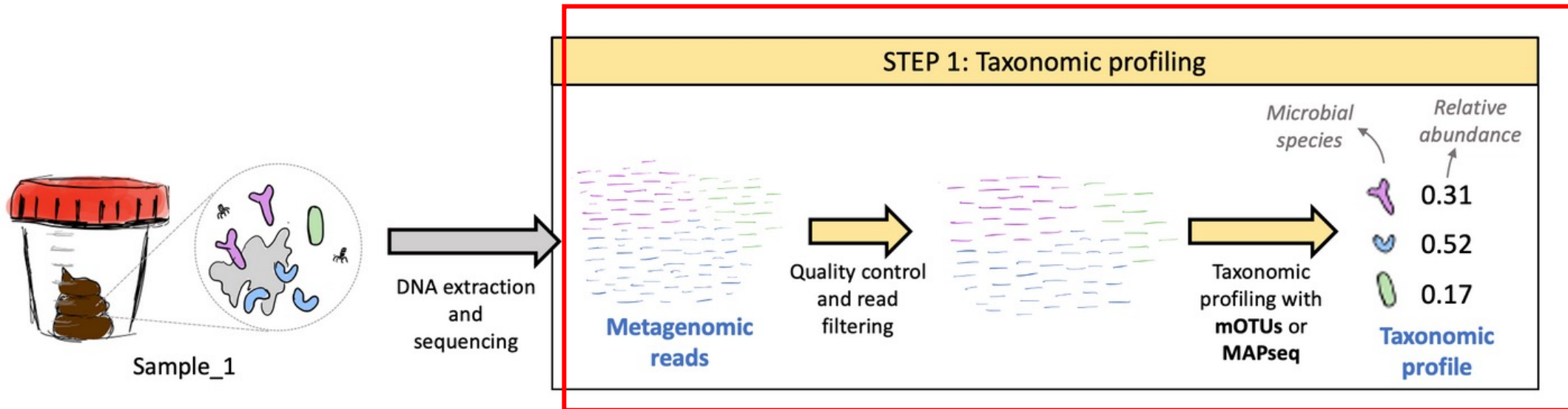


# Profiling and modeling the colorectal cancer microbiome

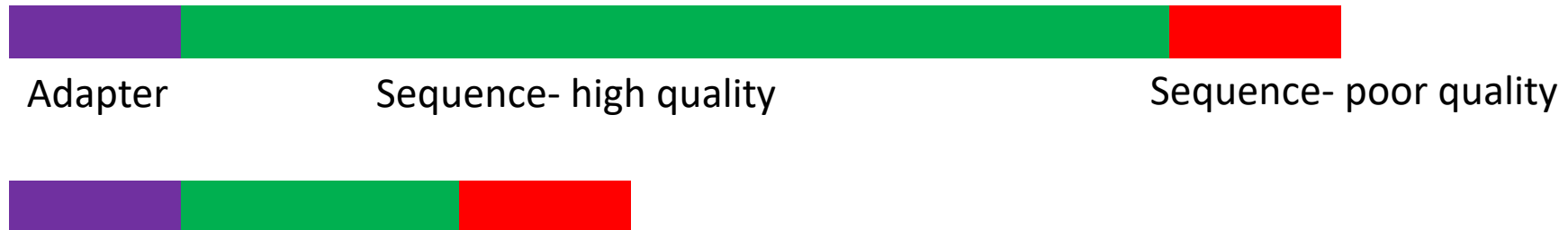
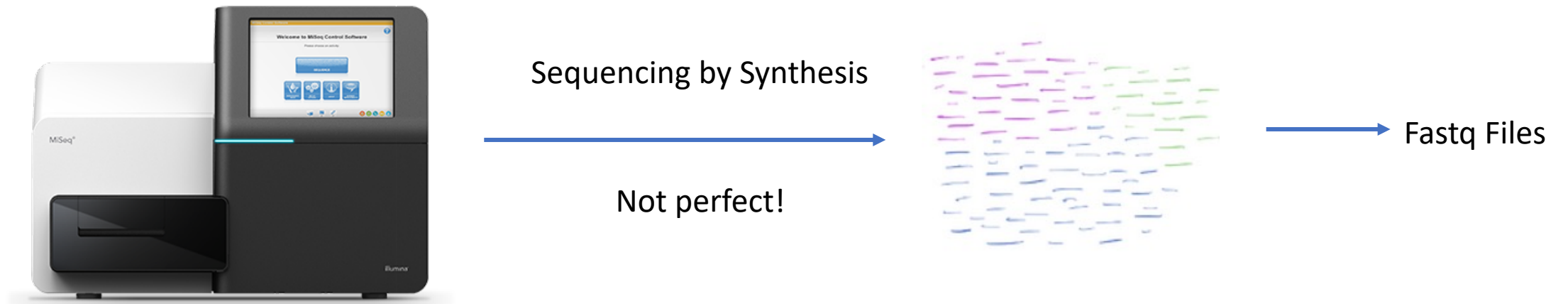
Alessio Milanese, Lukas Malfertheiner

Project 3

Spring School Bioinformatics and computational approaches in  
Microbiology



# Part 1: Quality Filtering and Trimming



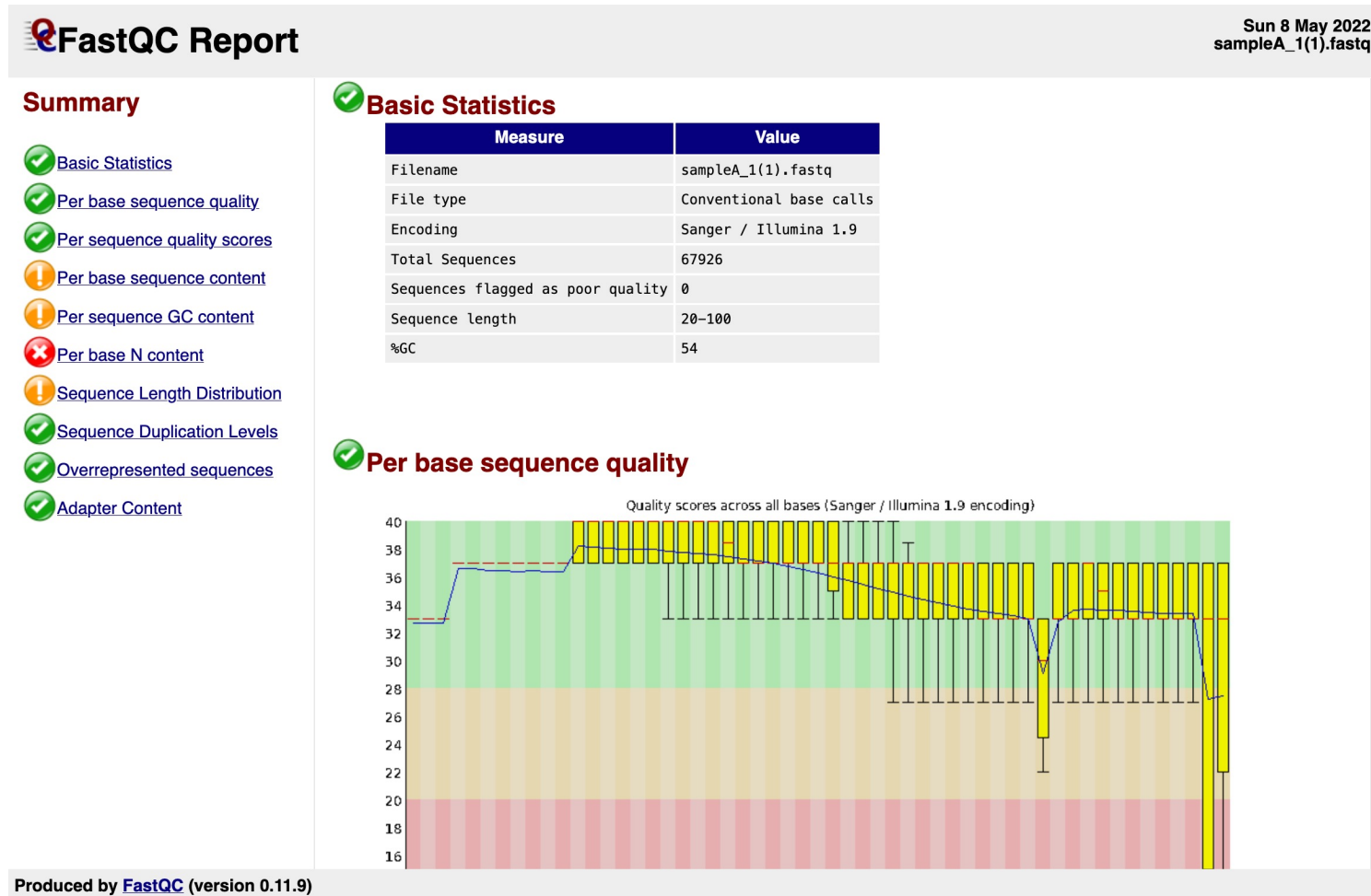


# Phred Score

Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%

# Fastqc

- Program to assess quality of sequences with a convenient web interface



---

*Genome analysis*

Advance Access publication April 1, 2014

## **Trimmomatic: a flexible trimmer for Illumina sequence data**

Anthony M. Bolger<sup>1,2</sup>, Marc Lohse<sup>1</sup> and Bjoern Usadel<sup>2,3,\*</sup>

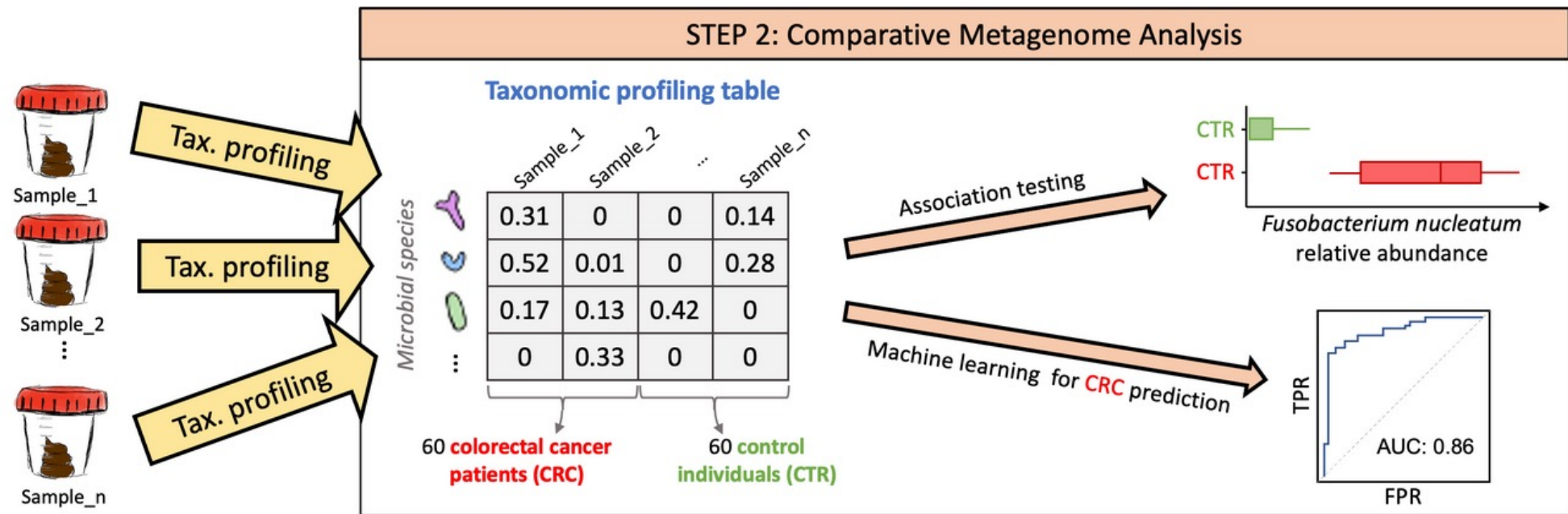
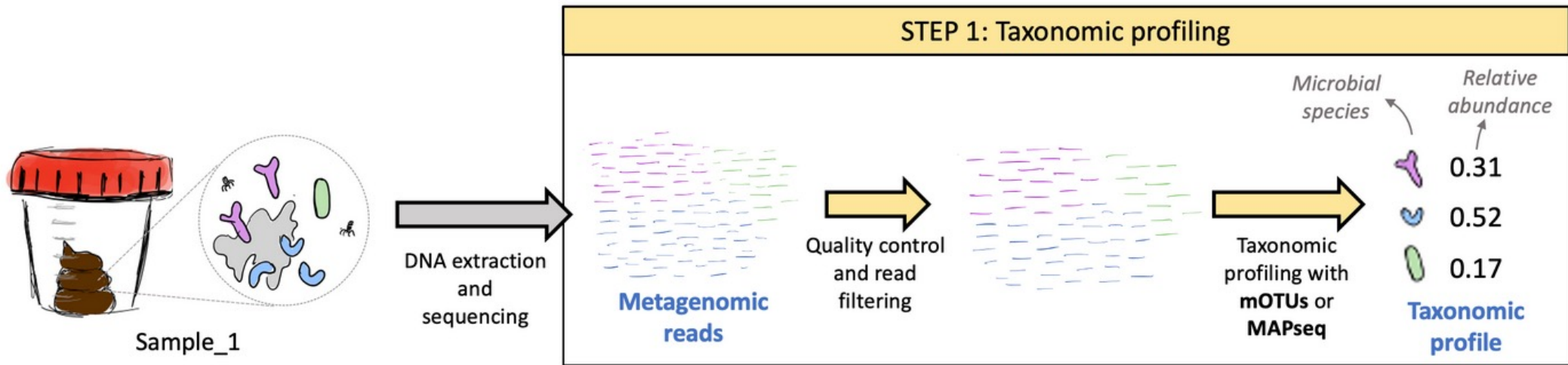
<sup>1</sup>Department Metabolic Networks, Max Planck Institute of Molecular Plant Physiology, Am Mühlberg 1, 14476 Golm, <sup>2</sup>Institut für Biologie I, RWTH Aachen, Worringer Weg 3, 52074 Aachen and <sup>3</sup>Institute of Bio- and Geosciences: Plant Sciences, Forschungszentrum Jülich, Leo-Brandt-Straße, 52425 Jülich, Germany

Associate Editor: Inanc Birol

---

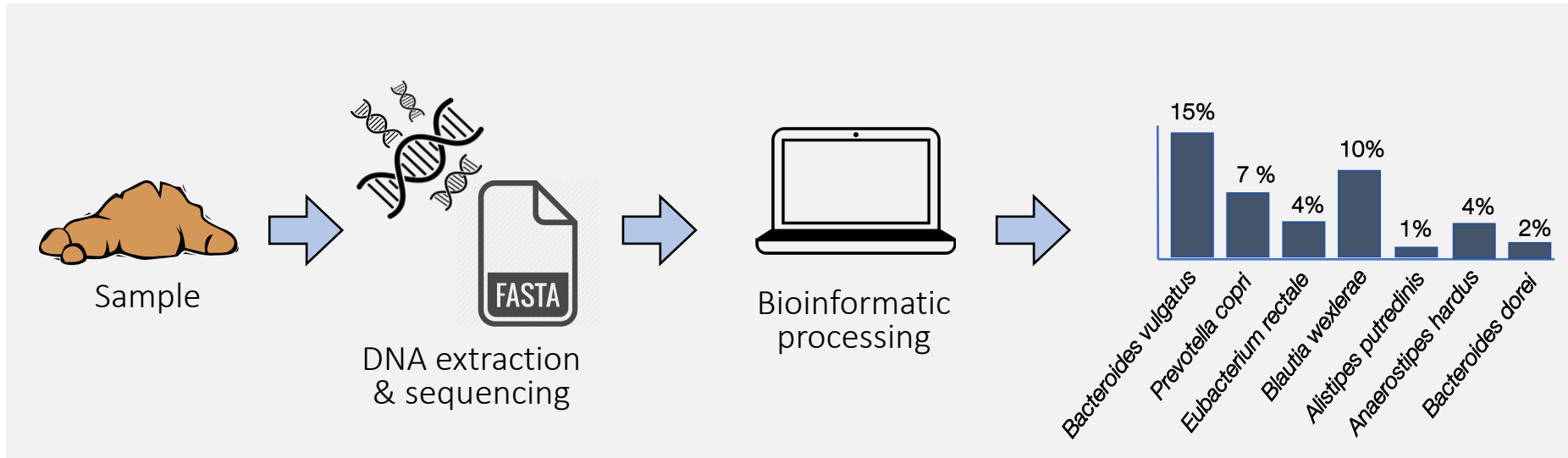
- Use quality information obtained by fastqc in order to trim our fastq files
- Important to not get misleading results!







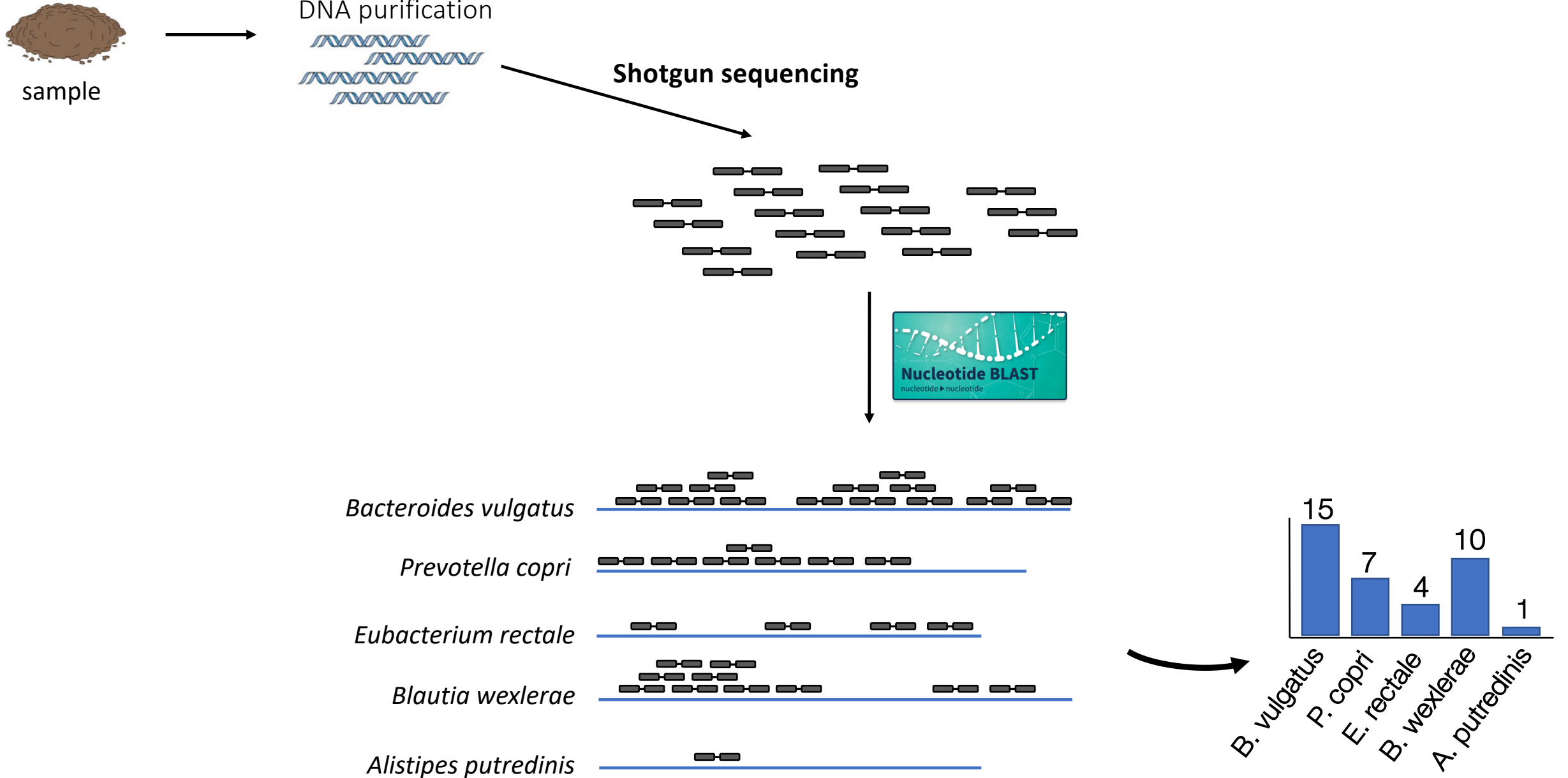
# Taxonomic profiling – what is it?



## **Taxonomic Profiling:**

Estimate relative cell counts in a microbiome sample from metagenomic sequencing

# Taxonomic profiling – how it is done?



# Taxonomic profiling approaches – whole-genome mapping

Environmental sample



Shotgun sequencing



- DNA extraction bias
- sequencing biases
- sampling noise

# Taxonomic profiling approaches – whole-genome mapping

Environmental sample



Shotgun sequencing



True taxonomic annotation



Estimated by whole-genome mapping



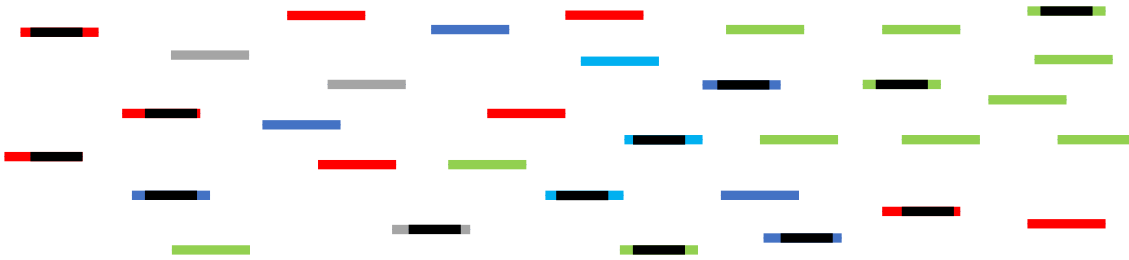
- genome size issue

# Taxonomic profiling approaches – marker gene mapping

Environmental sample



Shotgun sequencing



# Taxonomic profiling approaches – marker gene mapping

Environmental sample



Shotgun sequencing





# Taxonomic profiling approaches – marker gene mapping

Environmental sample



Shotgun sequencing



# Taxonomic profiling approaches – marker gene mapping

Environmental sample



Shotgun sequencing



True taxonomic annotation



Estimated by whole-genome mapping



- genome size issue

Estimated by universal marker



# Taxonomic profiling – mapping reads to genomes

Environmental sample

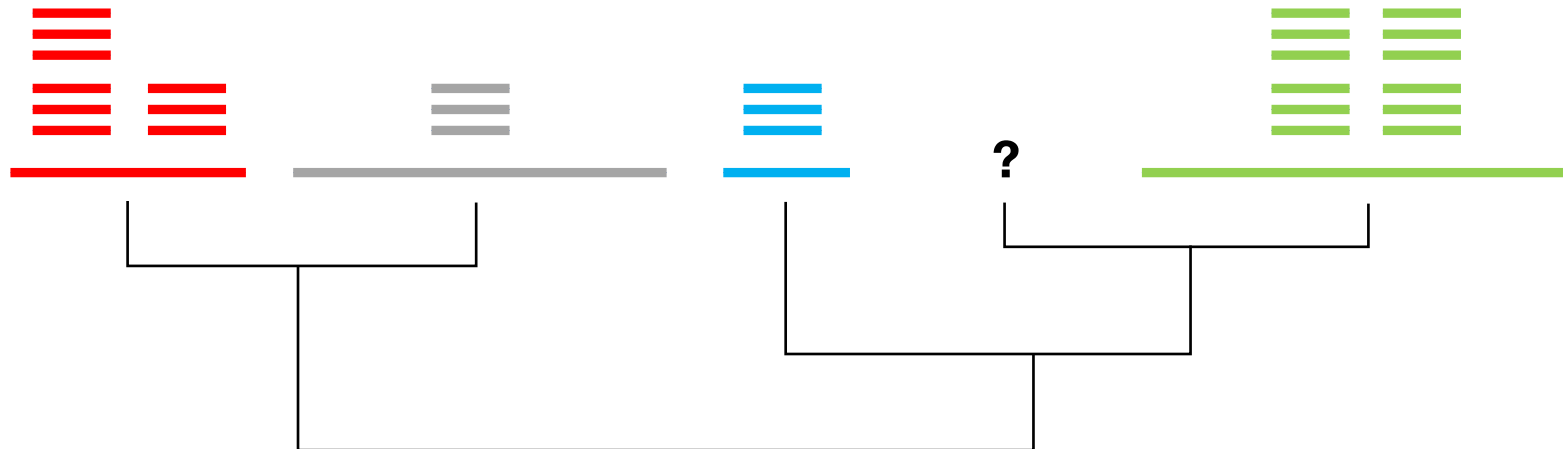


Shotgun sequencing



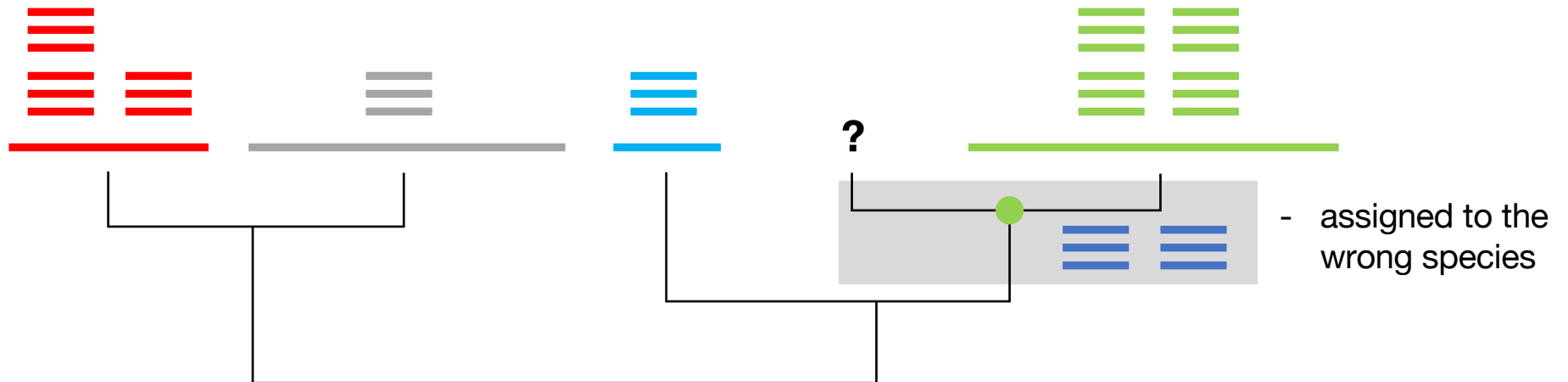
# Taxonomic profiling – incomplete reference databases

Environmental sample



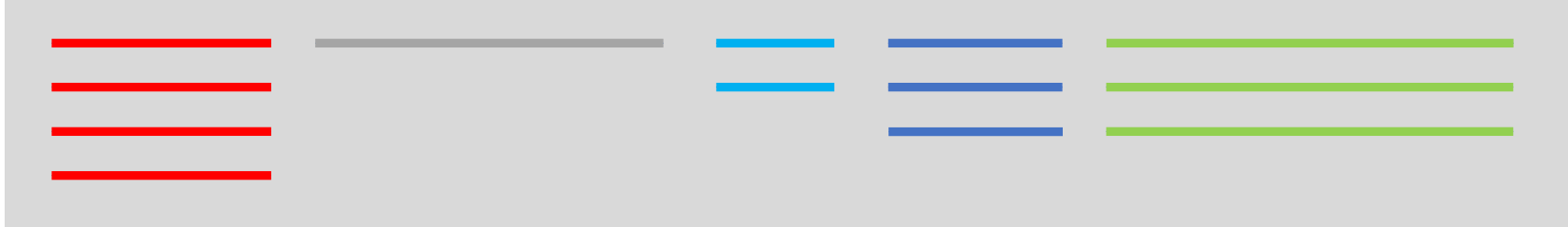
# Taxonomic profiling – incomplete reference databases

Environmental sample



# Taxonomic profiling – incomplete reference databases

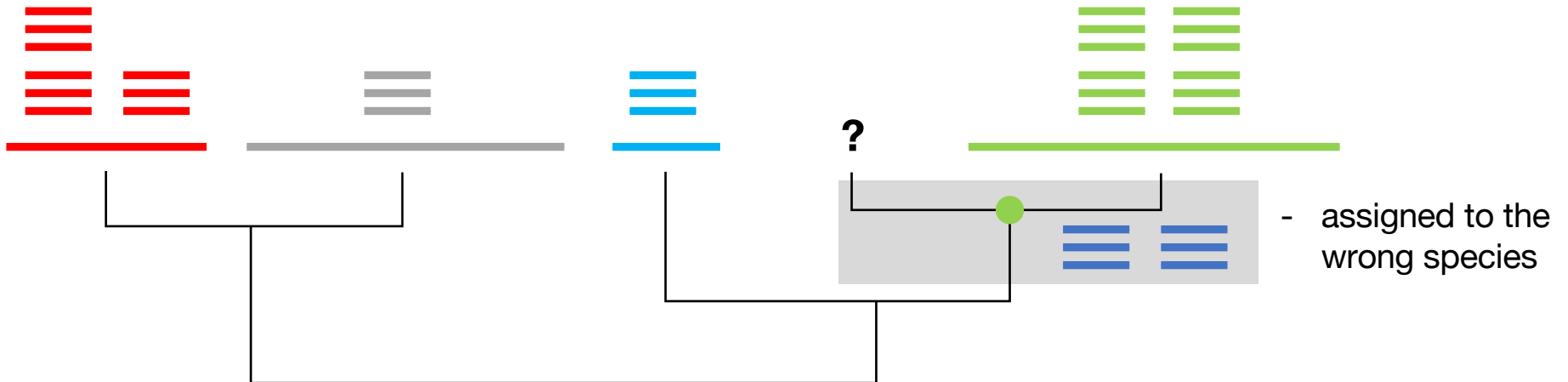
Environmental sample



True taxonomic annotation



Estimated when dark blue is missing



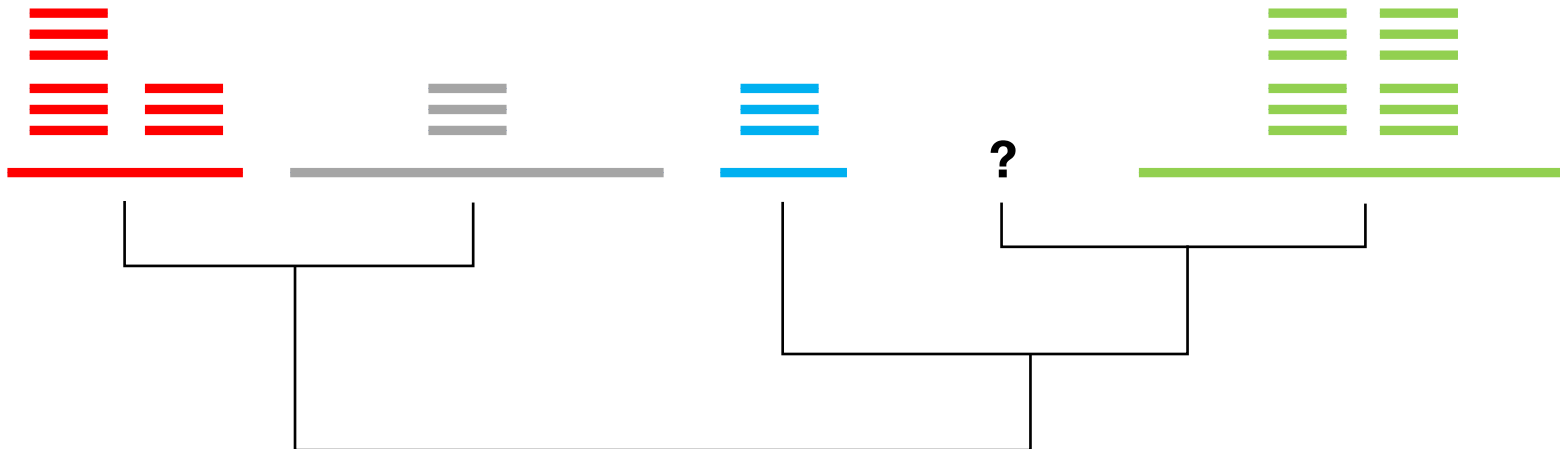


# Taxonomic profiling – incomplete reference databases

Environmental sample



- ignore the reads



# Taxonomic profiling – incomplete reference databases

Environmental sample



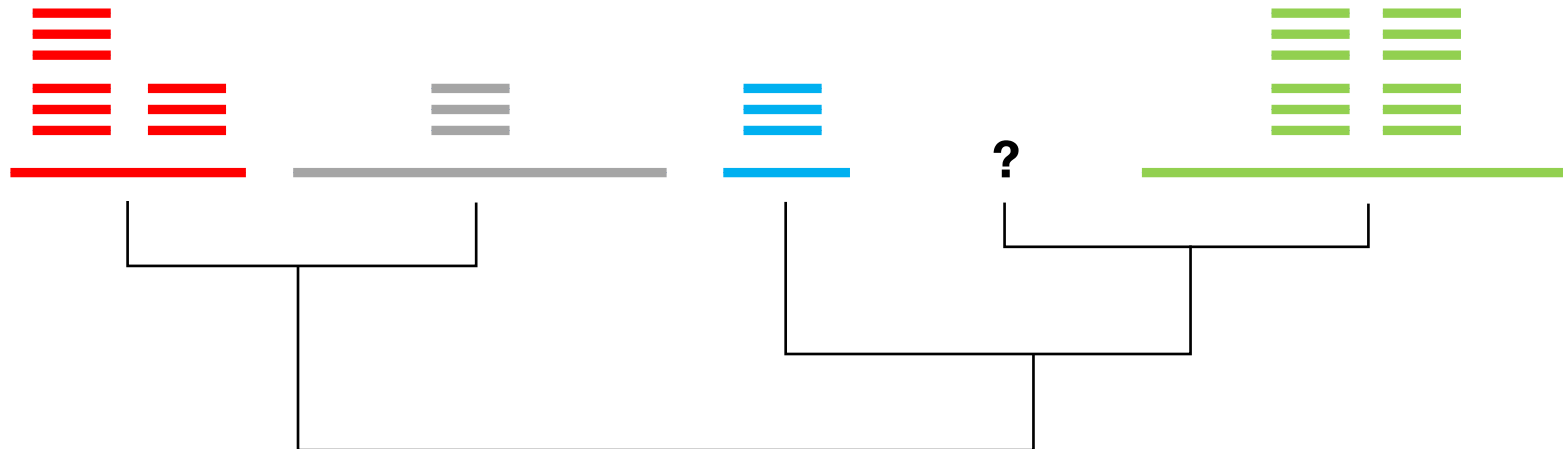
True taxonomic annotation



Estimated when dark green is missing

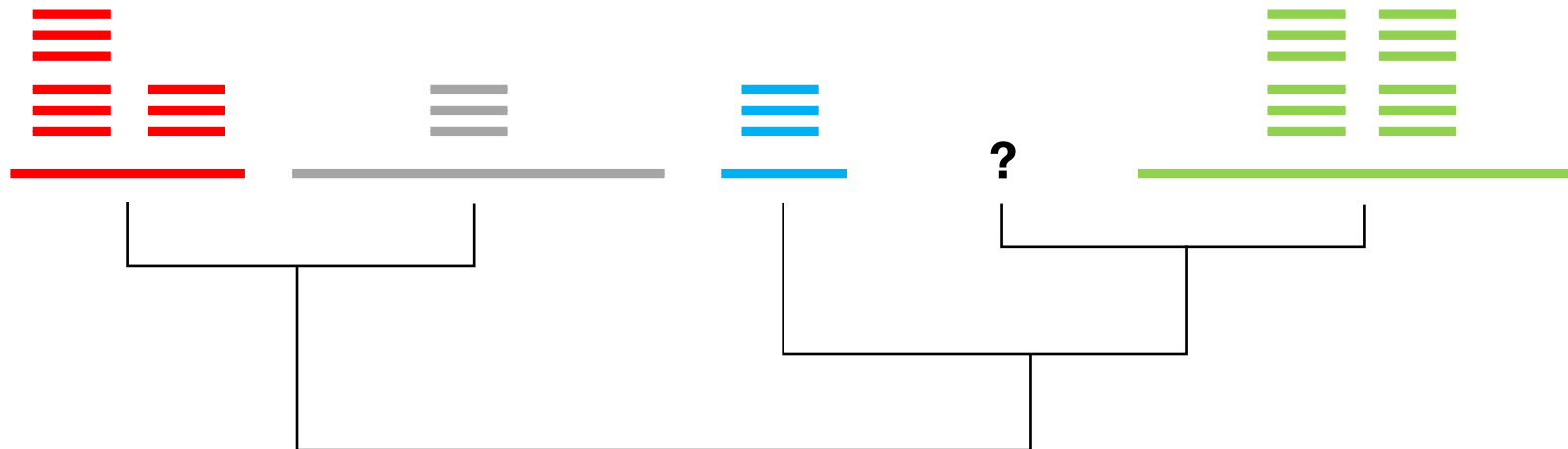


- ignore the reads



# Taxonomic profiling – incomplete reference databases

Environmental sample



Globally  
unassigned



# Taxonomic profiling – incomplete reference databases

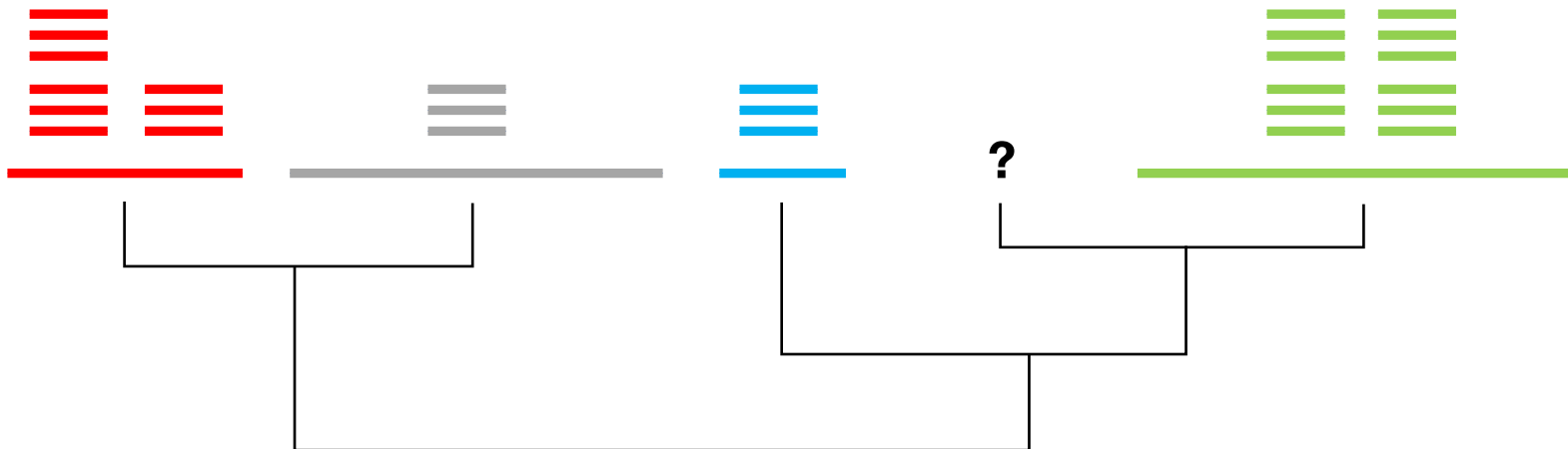
Environmental sample



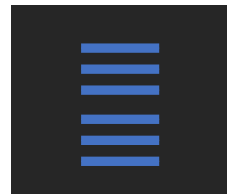
True taxonomic annotation



Estimated when dark blue is missing



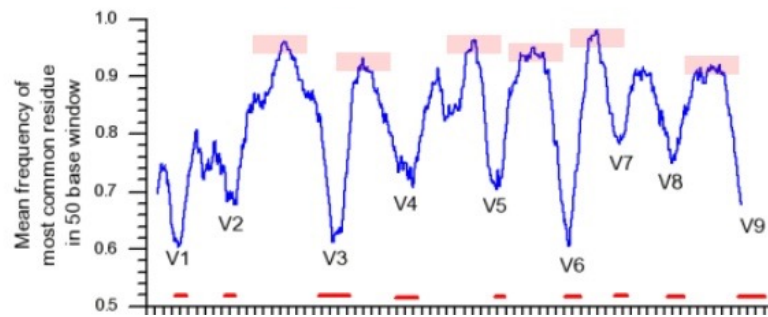
Globally unassigned



Trimmed, filtered reads



16S rRNA



10 universal protein marker genes



[nature](#) > [nature communications](#) > [articles](#) > [article](#)

Article | [Open Access](#) | [Published: 04 March 2019](#)

## Microbial abundance, activity and population genomic profiling with mOTUs2

[Alessio Milanese](#), [Daniel R Mende](#), [Lucas Paoli](#), [Guillem Salazar](#), [Hans-Joachim Ruscheweyh](#), [Miguelangel Cuenca](#), [Pascal Hingamp](#), [Renato Alves](#), [Paul I Costea](#), [Luis Pedro Coelho](#), [Thomas S. B. Schmidt](#), [Alexandre Almeida](#), [Alex L Mitchell](#), [Robert D. Finn](#), [Jaime Huerta-Cepas](#), [Peer Bork](#), [Georg Zeller](#) & [Shinichi Sunagawa](#)

[Nature Communications](#) **10**, Article number: 1014 (2019) | [Cite this article](#)

**25k** Accesses | **107** Citations | **78** Altmetric | [Metrics](#)

> [Bioinformatics](#). 2017 Dec 1;33(23):3808–3810. doi: 10.1093/bioinformatics/btx517.

## MAPseq: highly efficient k-mer search with confidence estimates, for rRNA sequence analysis

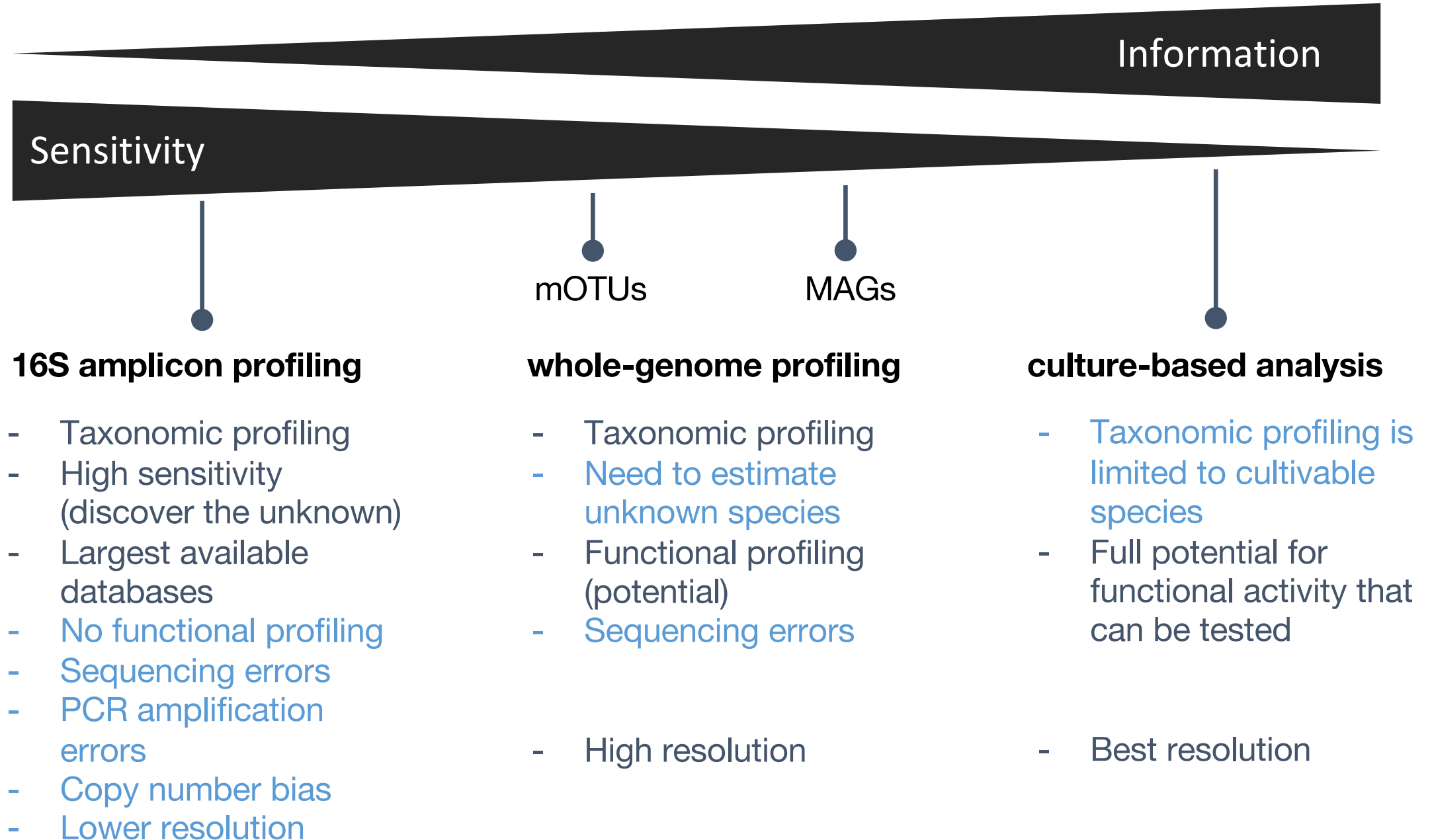
[João F Matias Rodrigues](#)<sup>1</sup>, [Thomas S B Schmidt](#)<sup>1</sup>, [Janko Tackmann](#)<sup>1</sup>, [Christian von Mering](#)<sup>1</sup>

Affiliations + expand

PMID: 28961926 PMCID: [PMC5860325](#) DOI: [10.1093/bioinformatics/btx517](#)

[Free PMC article](#)

# Strengths and weaknesses of different approaches



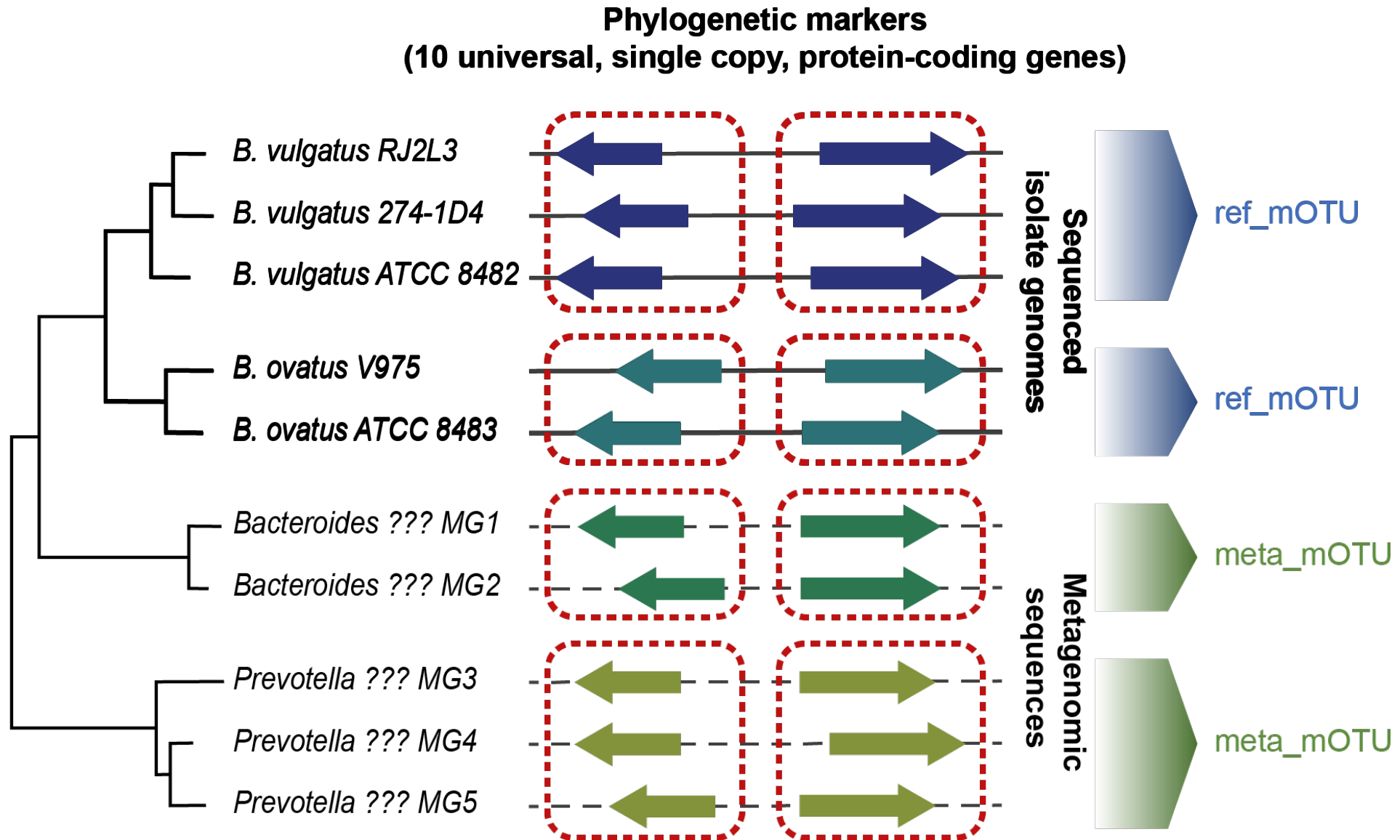


# The mOTUs framework – DB construction

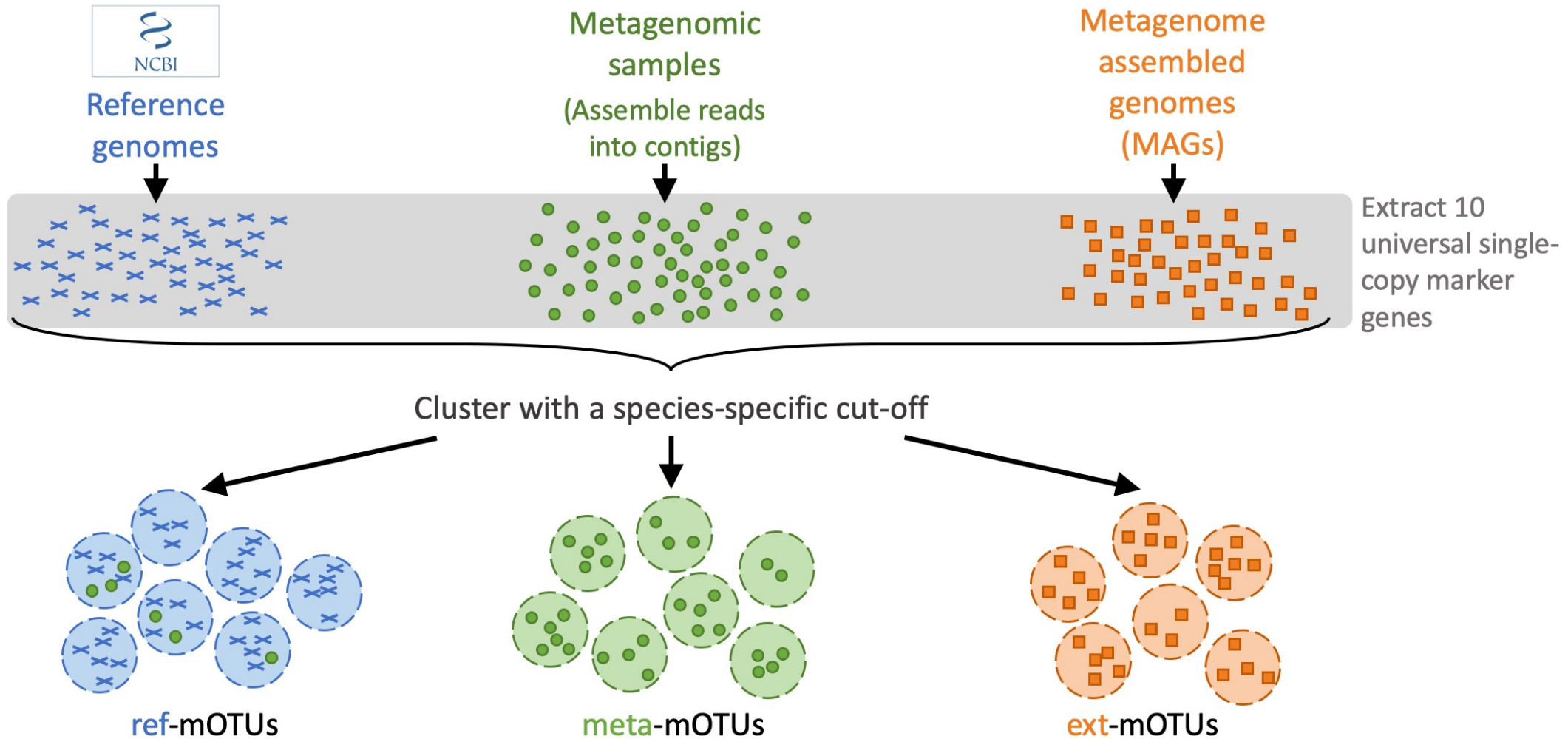
[Ciccarelli et al. *Science* 2006]

[Sunagawa et al. *Nat. Methods* 2013]

[Milanese et al. *Nat. Commun.* 2019]

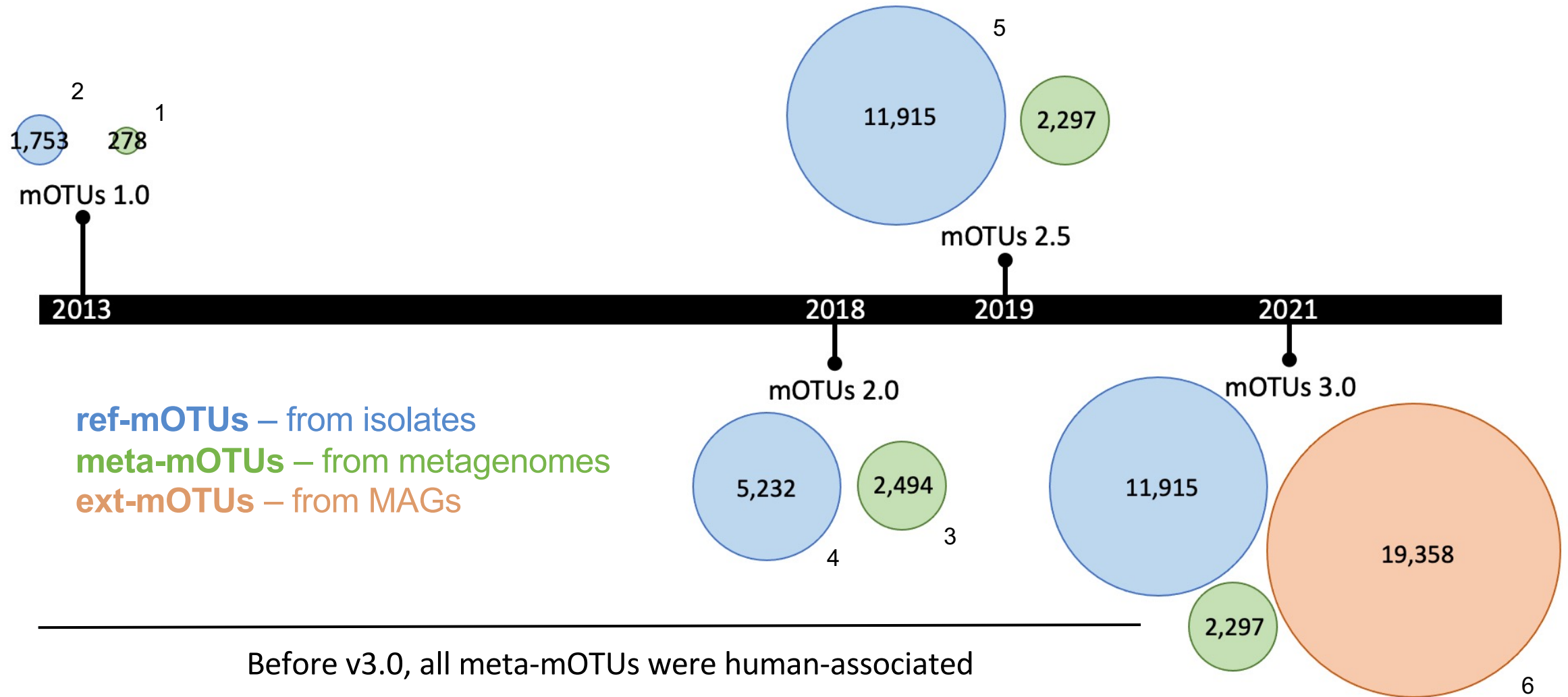


# Incorporation of MAGs into the mOTUs3 database



MAG-derived mOTUs are called ext\_mOTUs

# Improvement of scope in mOTUs since first version



1. [Sunagawa et al., *Nat. Methods* 2013]

2. [Mende et al. *Nat. Methods* 2013]

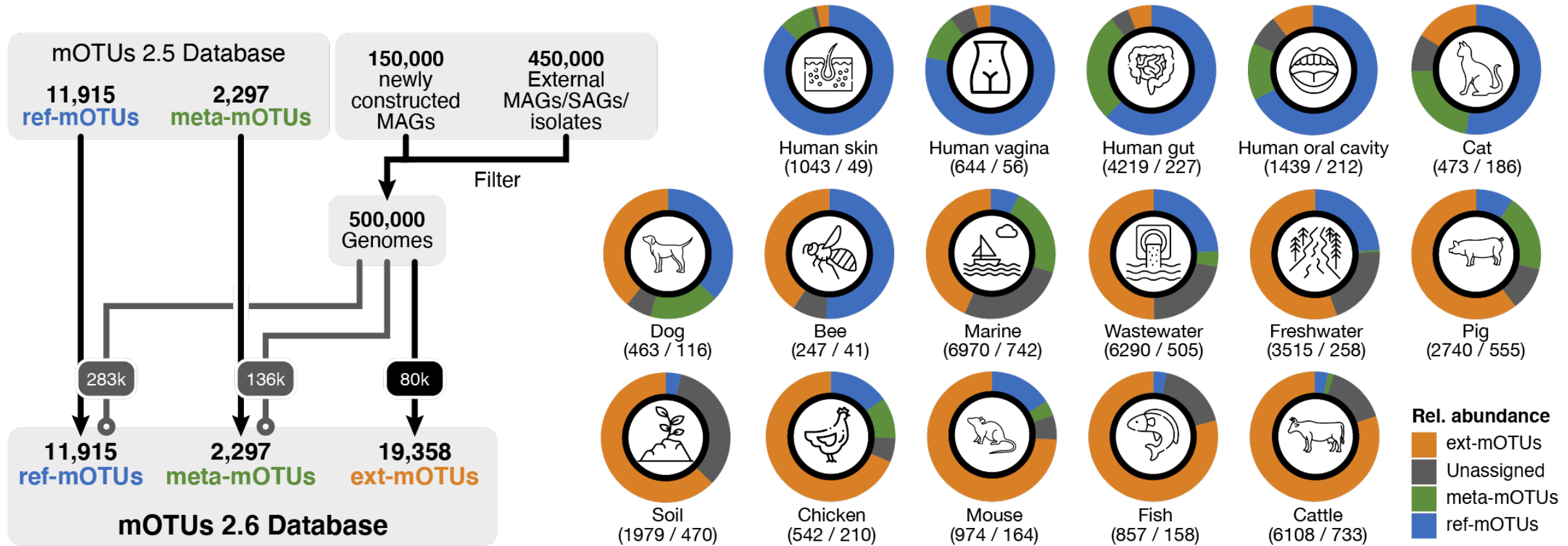
3. [Milanese et al., *Nat. Commun.* 2019]

4. [Mende et al., *Nucleic Acids Res.* 2017]

5. [Mende et al., *Nucleic Acids Res.* 2020]

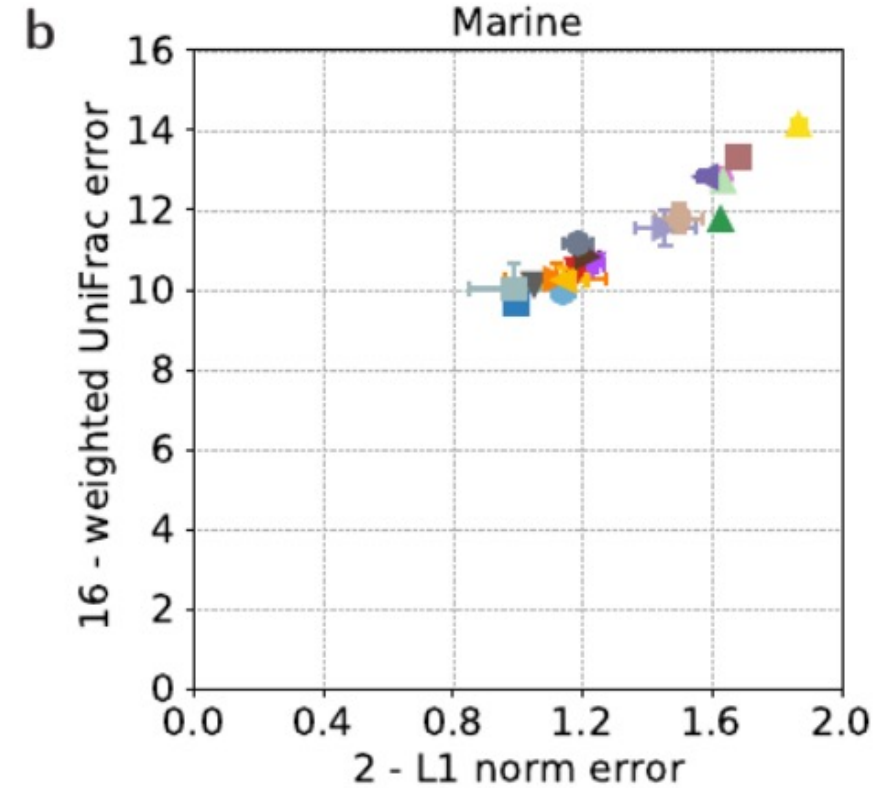
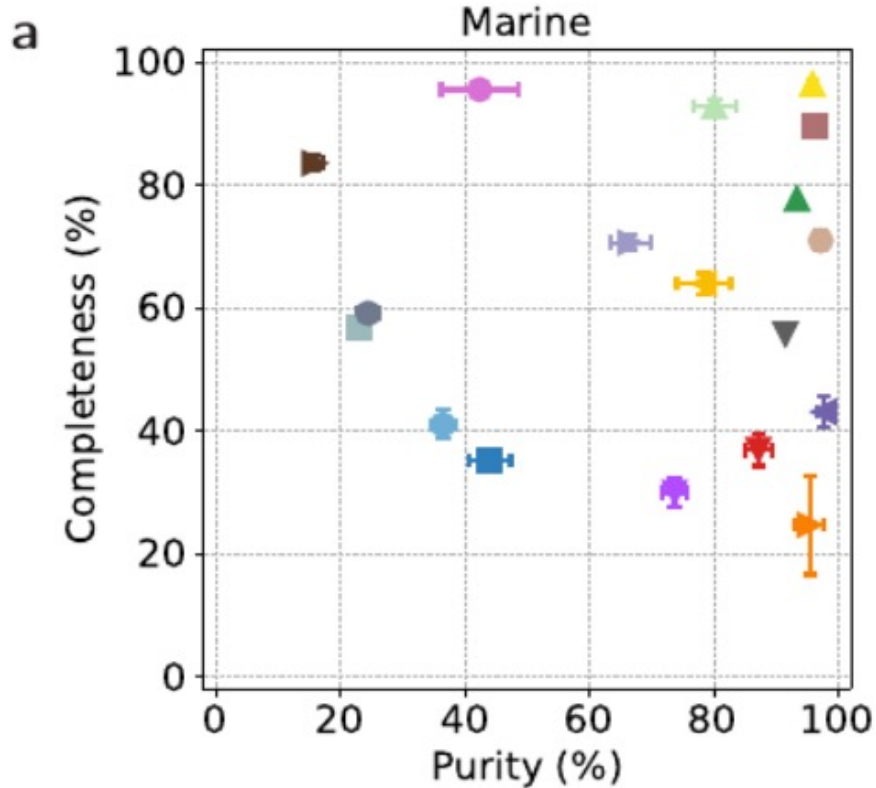
6. [Ruscheweyh, Milanese et al. *bioRxiv* 2021]

# mOTUs3 – database extension by marker genes from metagenome-assembled genomes (>500,000 MAGs)

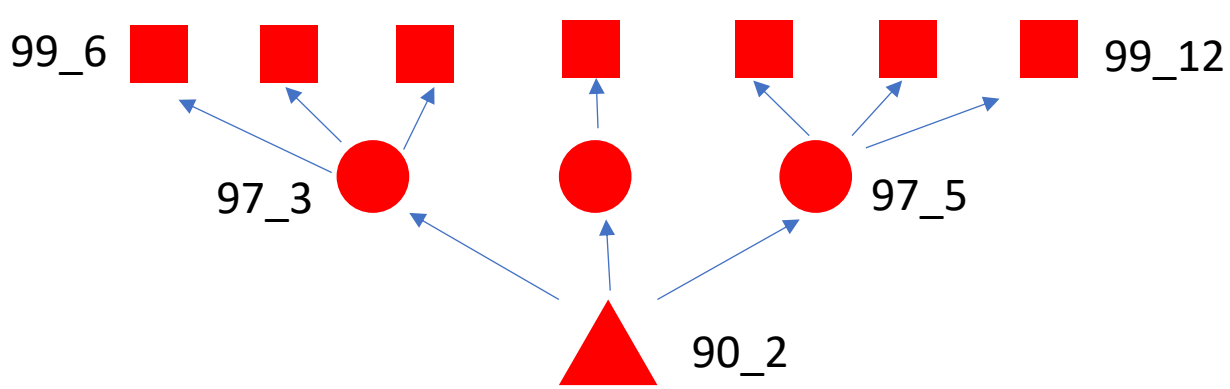
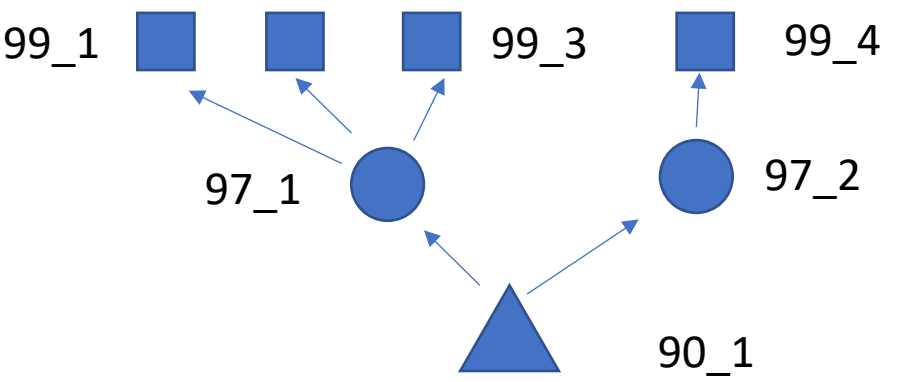
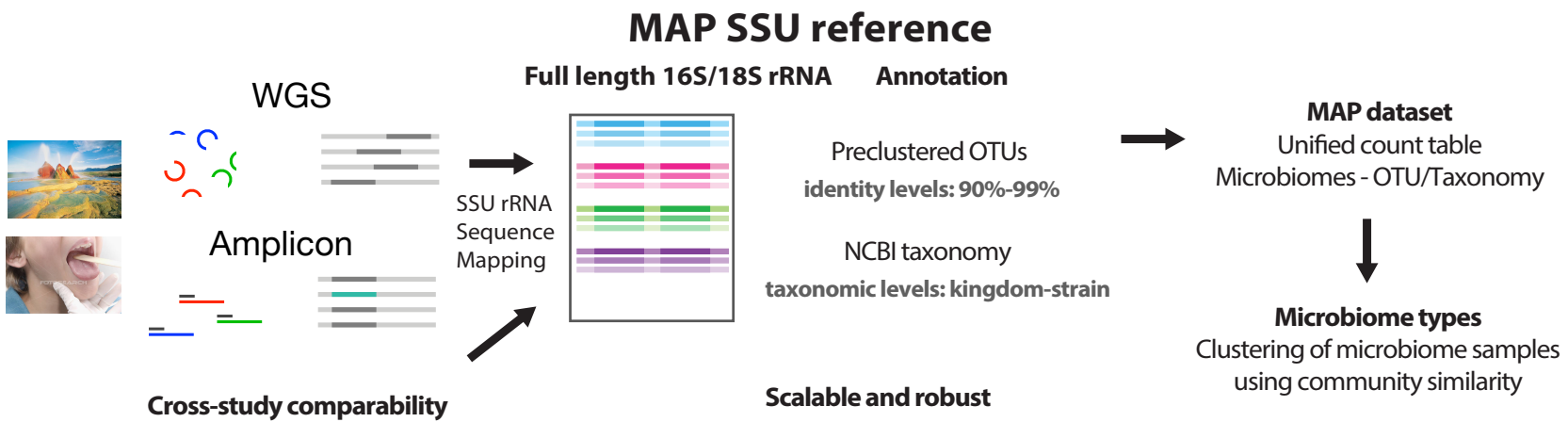


Enables profiling an unprecedented diversity of prokaryotes (33,570 species) across many environments.

# High-accuracy profiling as evaluated by an independent benchmark - CAMI



# MAPseq

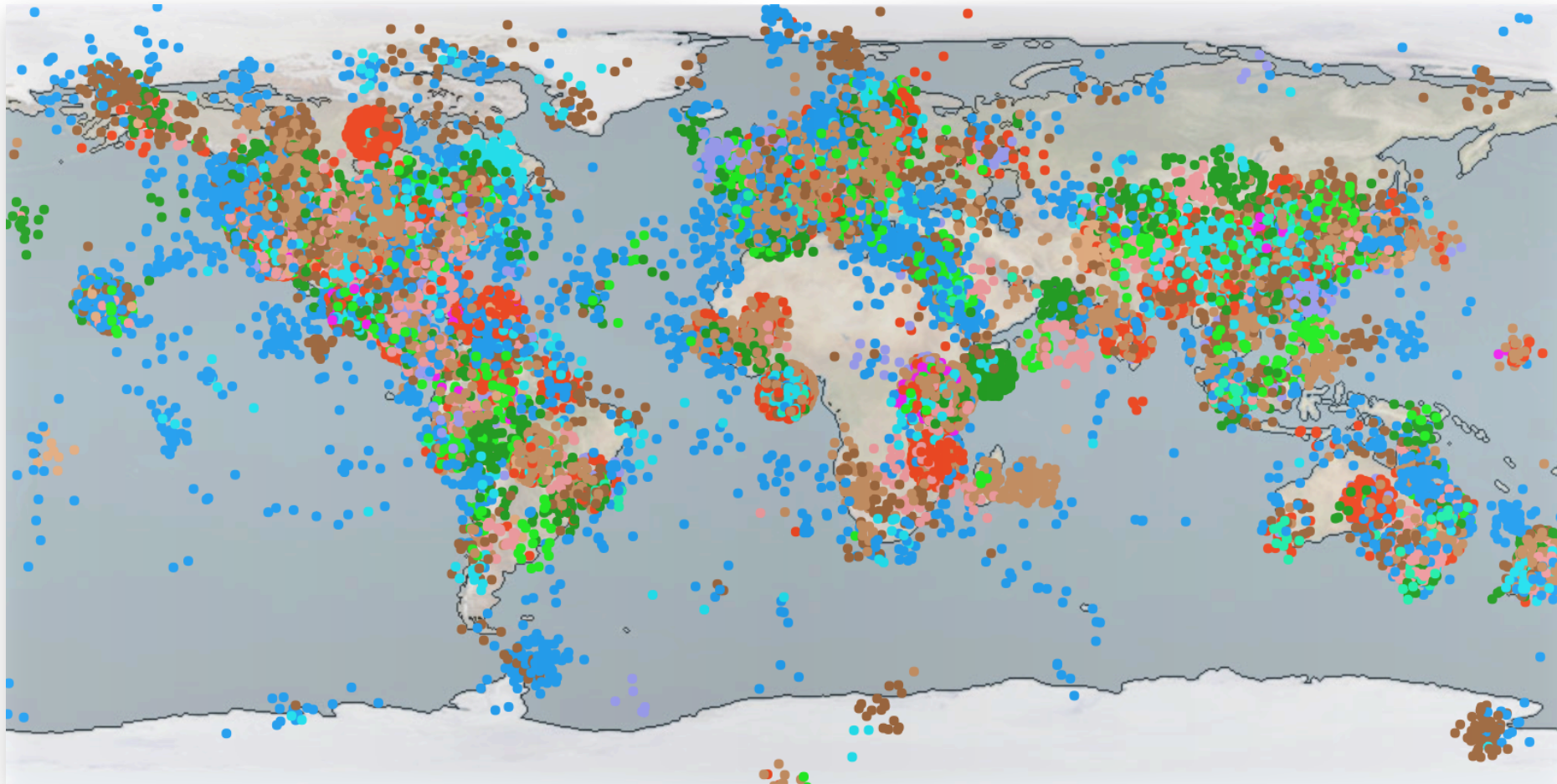






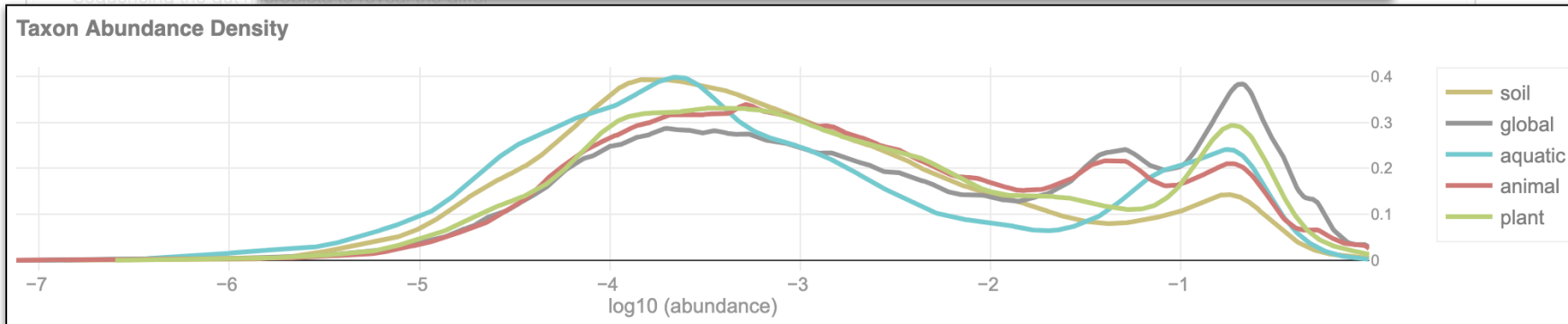
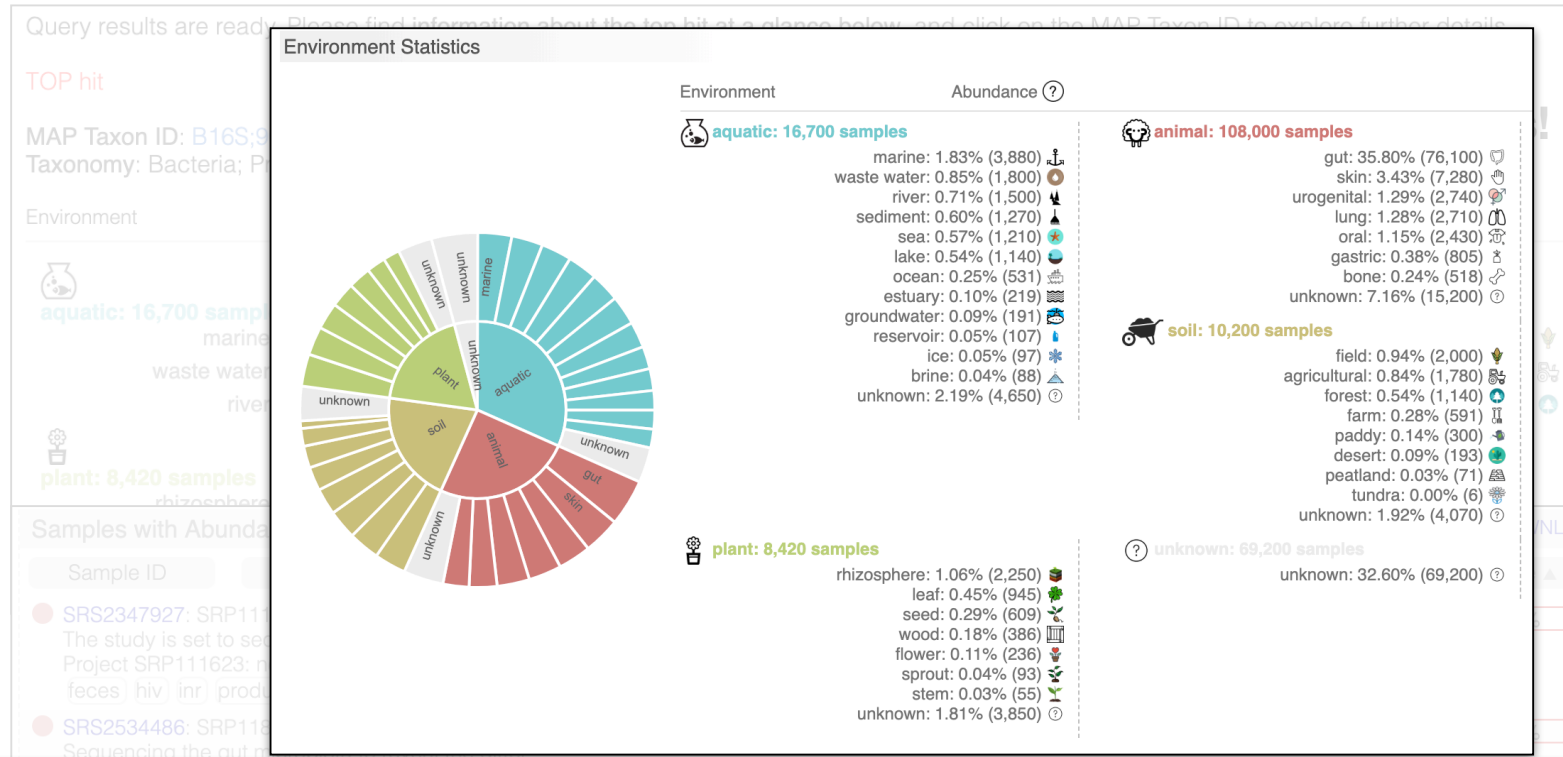
[microbeatlas.org](http://microbeatlas.org)

compare your metagenomic data to a global reference set of a million microbiome samples

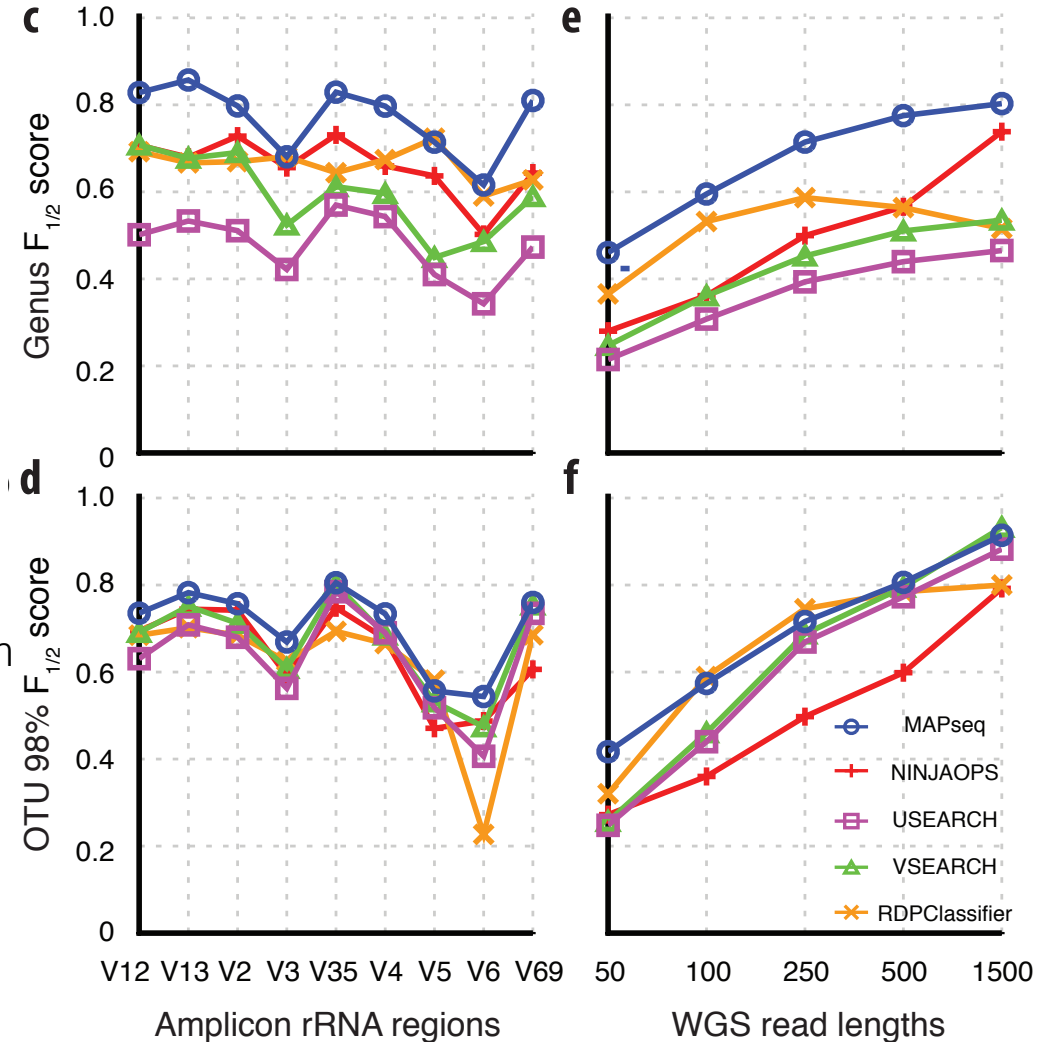
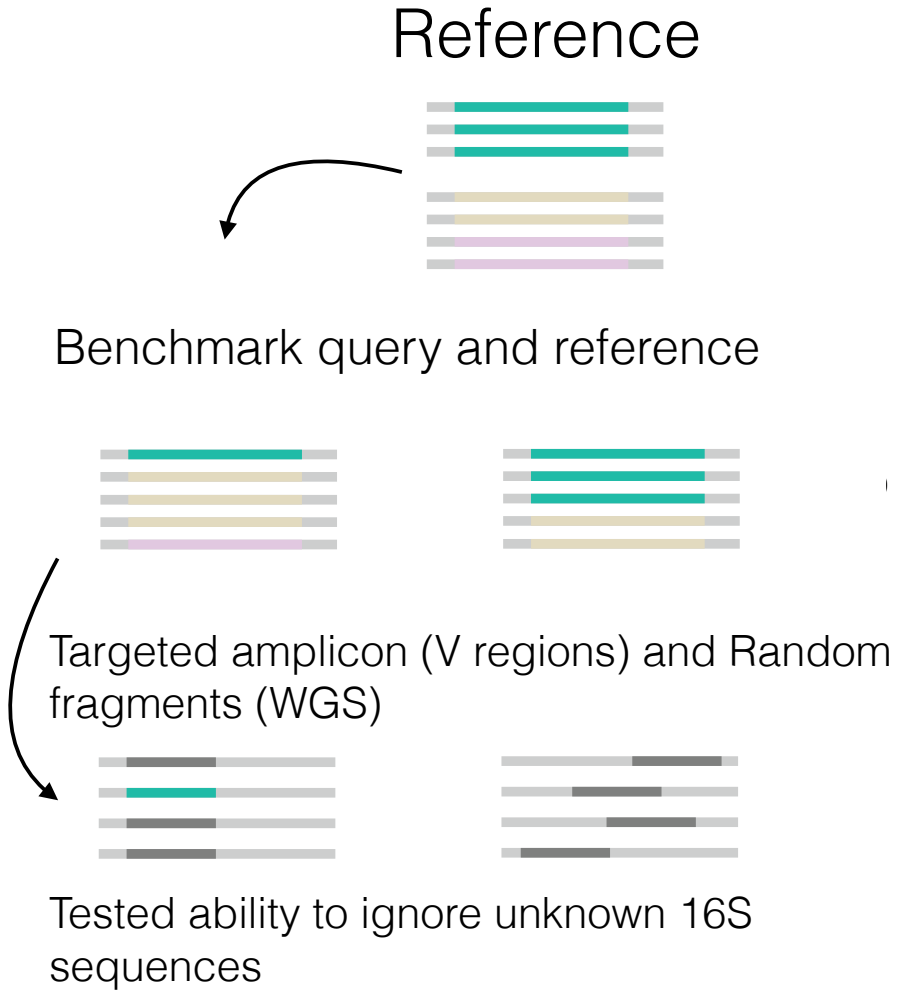


# Microbe Atlas Project Website

Search for any microbial taxa by sequence or name



# Benchmark of accuracy on known taxonomy

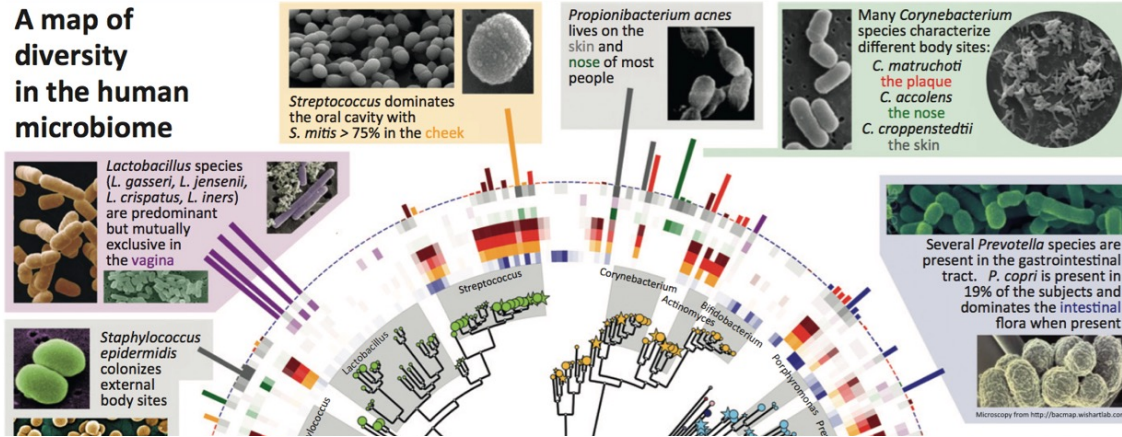


# Taxonomic profiling – why it is important?

Taxonomic analysis is fundamental to the analysis of microbial communities

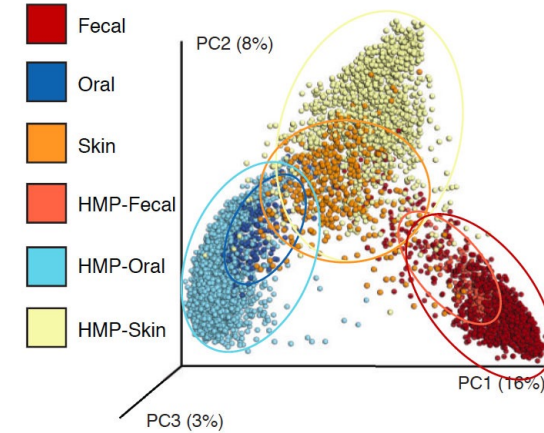
## Describing the microbial community under study

### A map of diversity in the human microbiome



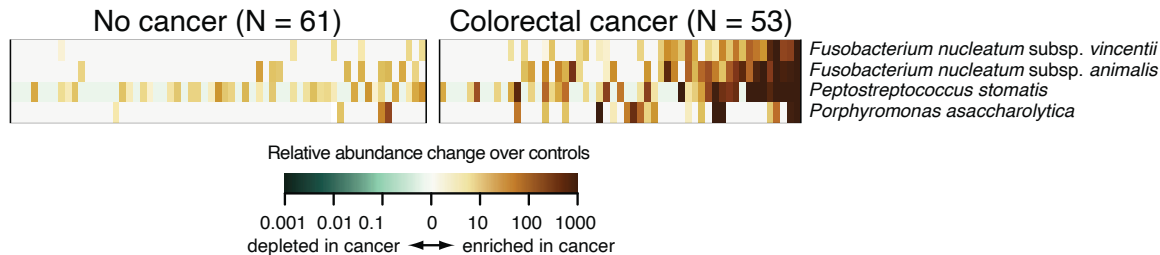
[Morgan et al., Trends in Genetics, 2013]

## Comparing different microbial communities



[McDonald et al., mSystems, 2018]

## Correlating environm. or host features to microbes



[Zeller et al., MSB, 2014]

## Comparing findings to literature

*Fusobacterium nucleatum* Contributes to the Carcinogenesis of Colorectal Cancer by Inducing Inflammation and Suppressing Host Immunity

**RESEARCH**  
**CANCER**

**Analysis of *Fusobacterium* persistence and antibiotic response in colorectal cancer**

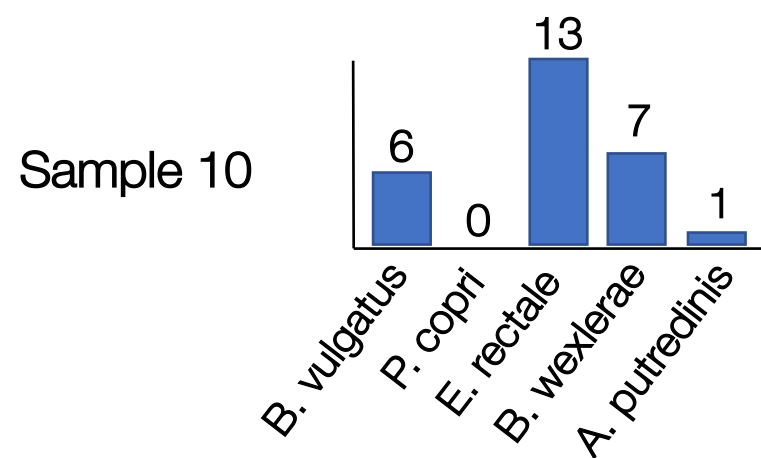
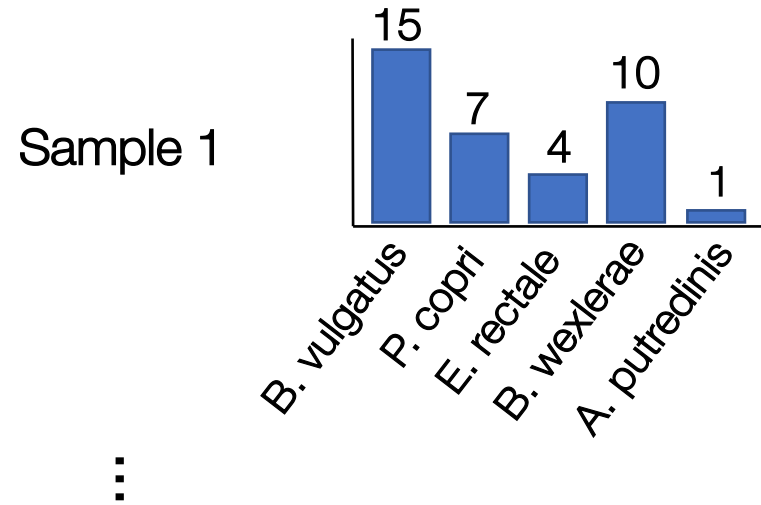
Susan Bullman,<sup>1,2</sup> Chandra S. Pedamallu,<sup>1,2</sup> Ewa Sidinska,<sup>1</sup> Thomas E. Clancy,<sup>2</sup> Xiaoyang Zhang,<sup>1,2</sup> Diana Cai,<sup>1,2</sup> Donna Neuberg,<sup>1</sup> Katherine Huang,<sup>2</sup> Fatima Guevara,<sup>1</sup> Timothy Nelson,<sup>1</sup> Otari Chilpashvili,<sup>1</sup> Timothy Hagan,<sup>1</sup> Mark Walker,<sup>1</sup> Aruna Ramachandran,<sup>1,2</sup> Begonia Diez-Ordaz,<sup>1,2</sup> Gamal Serres,<sup>1</sup> Nuria Mulet,<sup>1</sup> Stefania Landolfi,<sup>1</sup> Santiago Ramon y Cajal,<sup>1</sup> Roberta Fasani,<sup>1</sup> Andrew J. Aguirre,<sup>1,2,3</sup> Kimmie Ng,<sup>1</sup> Elena Elez,<sup>1</sup> Shuji Ogino,<sup>1,2,3</sup> Josep Taberner,<sup>1</sup> Charles S. Fuchs,<sup>1</sup> William C. Hahn,<sup>1,2,3</sup> Paolo Nucci,<sup>1,2</sup> Matthew Meyerson<sup>1,2,3</sup>

***Fusobacterium nucleatum* Promotes Colorectal Carcinogenesis by Modulating E-Cadherin/ $\beta$ -Catenin Signaling via its FadA Adhesin**

Mara Roxana Rubinstain,<sup>1,2</sup> Xiaowei Wang,<sup>1,2</sup> Wendy Liu,<sup>2</sup> Yujun Hao,<sup>2,3</sup> Guifang Cai,<sup>2</sup> and Yiping W. Han<sup>1,2,4\*</sup>

<sup>1</sup>Department of Periodontics

# Profiling multiple samples



	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
<i>B. vulgatus</i>	15	0	9	6	9	21	3	0	45	6
<i>P. copri</i>	7	11	0	0	12	0	6	0	0	0
<i>E. rectale</i>	4	4	0	4	0	7	0	0	0	13
<i>B. wexlerae</i>	10	0	2	0	0	5	0	0	4	7
<i>A. putredinis</i>	1	0	0	0	0	3	0	0	0	1
<i>E. coli</i>	0	3	12	0	0	5	0	4	1	0
<i>C. innocuum</i>	0	2	0	0	0	1	2	8	0	6
<i>R. intestinalis</i>	12	0	0	6	4	0	5	2	0	0
<i>A. finegoldii</i>	6	1	1	0	0	0	2	0	0	23



# Profiling multiple samples – Library size

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	4	6	12	2	4	0	2	3	0	9
P. copri	6	2	4	1	4	8	6	1	5	0
E. rectale	3	0	0	8	1	2	0	0	3	3
B. wexlerae	0	3	6	0	8	4	4	3	4	0
A. putredinis	0	0	0	6	0	14	1	0	0	7
E. coli	0	0	0	0	0	12	6	8	21	4
C. innocuum	0	1	2	0	4	2	1	1	0	5
R. intestinalis	5	1	2	0	2	3	9	0	2	0
A. finegoldii	0	0	0	1	0	0	0	0	0	0

# Profiling multiple samples – Library size

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	4	6	12	2	4	0	2	3	0	9
P. copri	6	2	4	1	4	8	6	1	5	0
E. rectale	3	0	0	8	1	2	0	0	3	3
B. wexlerae	0	3	6	0	8	4	4	3	4	0
A. putredinis	0	0	0	6	0	14	1	0	0	7
E. coli	0	0	0	0	0	12	6	8	21	4
C. innocuum	0	1	2	0	4	2	1	1	0	5
R. intestinalis	5	1	2	0	2	3	9	0	2	0
A. finegoldii	0	0	0	1	0	0	0	0	0	0
SUM	18	13	26	18	23	45	29	16	35	28

# Profiling multiple samples – Library size

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	4	<b>6</b>	<b>12</b>	2	4	0	2	3	0	9
P. copri	6	<b>2</b>	<b>4</b>	1	4	8	6	1	5	0
E. rectale	3	<b>0</b>	<b>0</b>	8	1	2	0	0	3	3
B. wexlerae	0	<b>3</b>	<b>6</b>	0	8	4	4	3	4	0
A. putredinis	0	<b>0</b>	<b>0</b>	6	0	14	1	0	0	7
E. coli	0	<b>0</b>	<b>0</b>	0	0	12	6	8	21	4
C. innocuum	0	<b>1</b>	<b>2</b>	0	4	2	1	1	0	5
R. intestinalis	5	<b>1</b>	<b>2</b>	0	2	3	9	0	2	0
A. finegoldii	0	<b>0</b>	<b>0</b>	1	0	0	0	0	0	0
SUM	18	<b>13</b>	<b>26</b>	18	23	45	29	16	35	28



# Profiling multiple samples – Relative abundance

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	0.2	<b>0.5</b>	<b>0.5</b>	0.1	0.2	0	0.1	0.2	0	0.3
P. copri	0.3	<b>0.2</b>	<b>0.2</b>	0.1	0.2	0.2	0.2	0.1	0.1	0
E. rectale	0.2	<b>0</b>	<b>0</b>	0.4	0	0	0	0	0.1	0.1
B. wexlerae	0	<b>0.2</b>	<b>0.2</b>	0	0.3	0.1	0.1	0.2	0.1	0
A. putredinis	0	<b>0</b>	<b>0</b>	0.3	0	0.3	0	0	0	0.3
E. coli	0	<b>0</b>	<b>0</b>	0	0	0.3	0.2	0.5	0.6	0.1
C. innocuum	0	<b>0.1</b>	<b>0.1</b>	0	0.2	0	0	0.1	0	0.2
R. intestinalis	0.3	<b>0.1</b>	<b>0.1</b>	0	0.1	0.1	0.3	0	0.1	0
A. finegoldii	0	<b>0</b>	<b>0</b>	0.1	0	0	0	0	0	0

# Profiling multiple samples – Richness

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	0.2	0.5	0.5	0.1	0.2	0	0.1	0.2	0	0.3
P. copri	0.3	0.2	0.2	0.1	0.2	0.2	0.2	0.1	0.1	0
E. rectale	0.2	0	0	0.4	0	0	0	0	0.1	0.1
B. wexlerae	0	0.2	0.2	0	0.3	0.1	0.1	0.2	0.1	0
A. putredinis	0	0	0	0.3	0	0.3	0	0	0	0.3
E. coli	0	0	0	0	0	0.3	0.2	0.5	0.6	0.1
C. innocuum	0	0.1	0.1	0	0.2	0	0	0.1	0	0.2
R. intestinalis	0.3	0.1	0.1	0	0.1	0.1	0.3	0	0.1	0
A. finegoldii	0	0	0	0.1	0	0	0	0	0	0

↓  
4

↓  
5

- The richness is calculated per sample
- It represents the total number of species observed in a sample

# Profiling multiple samples – Prevalence

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	0.2	0.5	0.5	0.1	0.2	0	0.1	0.2	0	0.3
P. copri	0.3	0.2	0.2	0.1	0.2	0.2	0.2	0.1	0.1	0
E. rectale	0.2	0	0	0.4	0	0	0	0	0.1	0.1
B. wexlerae	0	0.2	0.2	0	0.3	0.1	0.1	0.2	0.1	0
A. putredinis	0	0	0	0.3	0	0.3	0	0	0	0.3
E. coli	0	0	0	0	0	0.3	0.2	0.5	0.6	0.1
C. innocuum	0	0.1	0.1	0	0.2	0	0	0.1	0	0.2
R. intestinalis	0.3	0.1	0.1	0	0.1	0.1	0.3	0	0.1	0
A. finegoldii	0	0	0	0.1	0	0	0	0	0	0

→ 8

- The prevalence is calculated per species
- It measure the number of sample where the species is detected

→ 1

# Profiling multiple samples – Prevalence

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10
B. vulgatus	0.2	0.5	0.5	0.1	0.2	0	0.1	0.2	0	0.3
P. copri	0.3	0.2	0.2	0.1	0.2	0.2	0.2	0.1	0.1	0
E. rectale	0.2	0	0	0.4	0	0	0	0	0.1	0.1
B. wexlerae	0	0.2	0.2	0	0.3	0.1	0.1	0.2	0.1	0
A. putredinis	0	0	0	0.3	0	0.3	0	0	0	0.3
E. coli	0	0	0	0	0	0.3	0.2	0.5	0.6	0.1
C. innocuum	0	0.1	0.1	0	0.2	0	0	0.1	0	0.2
R. intestinalis	0.3	0.1	0.1	0	0.1	0.1	0.3	0	0.1	0
A. finegoldii	0	0	0	0.1	0	0	0	0	0	0

→ 8 (0.8)

- The prevalence is calculated per species
- It measure the number of sample where the species is detected
- It can also be represented as fraction of the total amount of samples

→ 1 (0.1)