# What is metagenomics and why is it a big deal?

**Maliheh Mehrshad**

**Department of Aquatic Sciences and Assessment,** SLU

**BIONOMICS-MMlab**

MULTIPARTITE PARASITIC INTERACTIONS

# Who am I?

# BIONOMICS-MMlab
## MULTIPARTITE PARASITIC INTERACTIONS

## www.bionomics-mmlab.com
### Pushing the frontiers of multipartite parasitic interactions research

## Current bionomics-mmlab members

Zahra Goodarzi – PhD student
Dr. Vinicius Silva Kavagutti – PostDoc
Dr. Vesna Grujcic – PostDoc
Lauren Davies – PhD student

## Former lab members
Maryan Resaei-some – PostDoc
Fanny Persson – Field assistant
Armand Stoe – MSc student
Remco Hoogervorst – Internship student
Ezhilarasan Mani Ezhilan – MSc student

# What is metagenomics and why is it a big deal?

**Maliheh Mehrshad**

**Department of Aquatic Sciences and Assessment,** SLU

BIONOMICS-MMlab
MULTIPARTITE PARASITIC INTERACTIONS

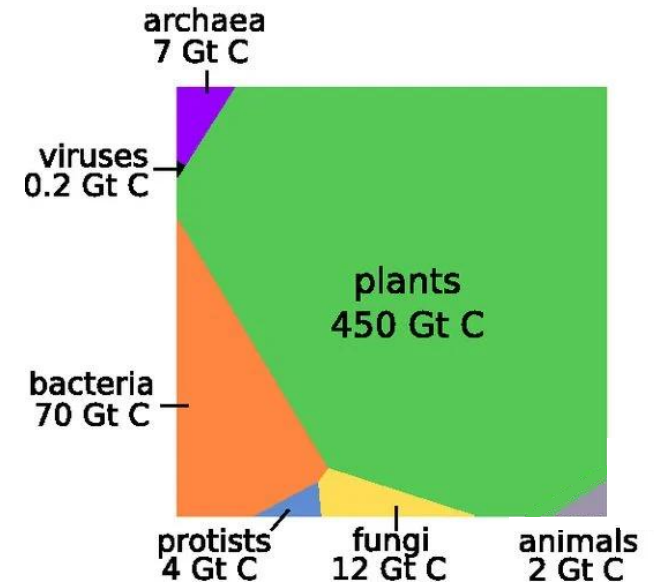# Let's answer some questions on menti!

- Join at menti.com
- Use code 74440277

# Microbes Matter!

Life on Earth has been microscopic for much of its ~4 billion year history.
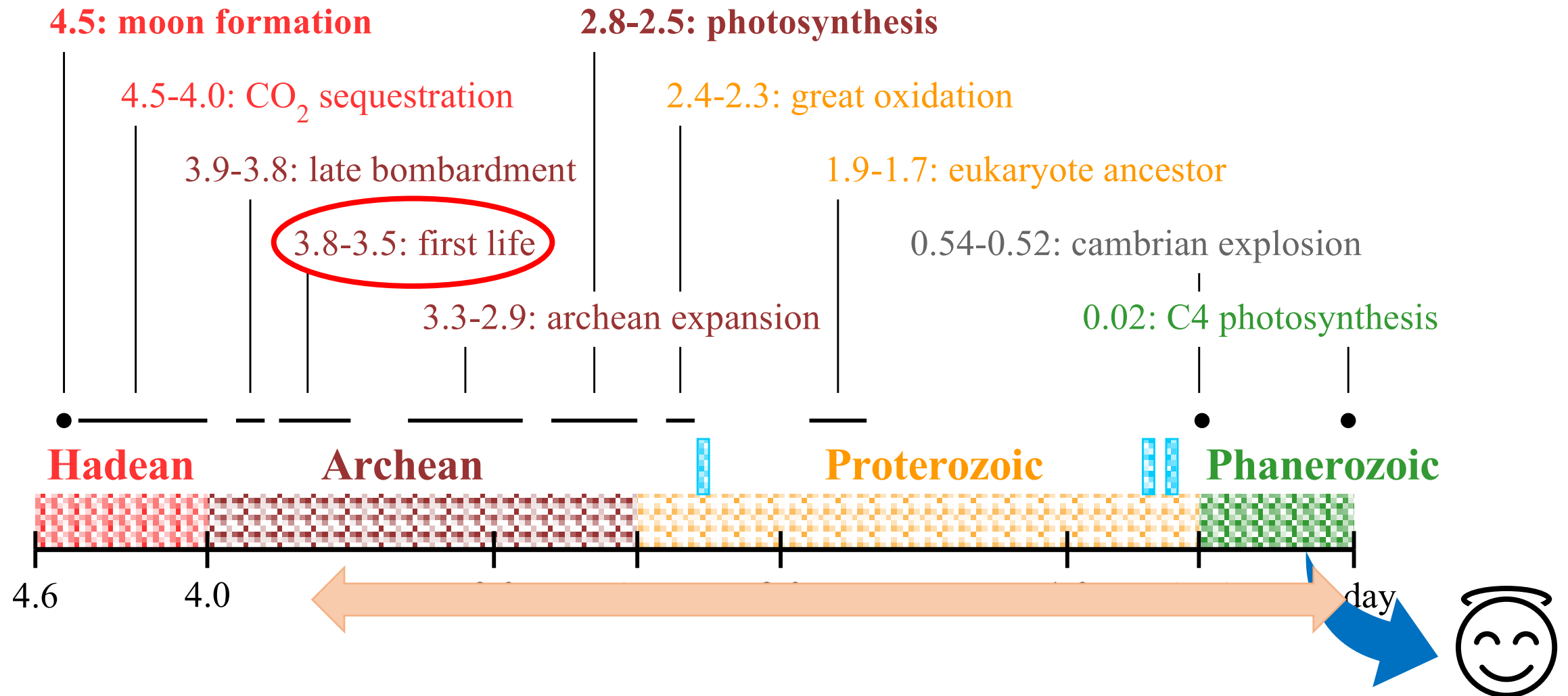
The metabolic activity of these organisms has left its mark.

- Great oxygenation event
- Photosynthesis
- Lignin and cellulose degradation
- Cycling of elements (C, N, S, Fe, …)
- Greenhouse gas sink/emission
- Interact with plants and animals
- …



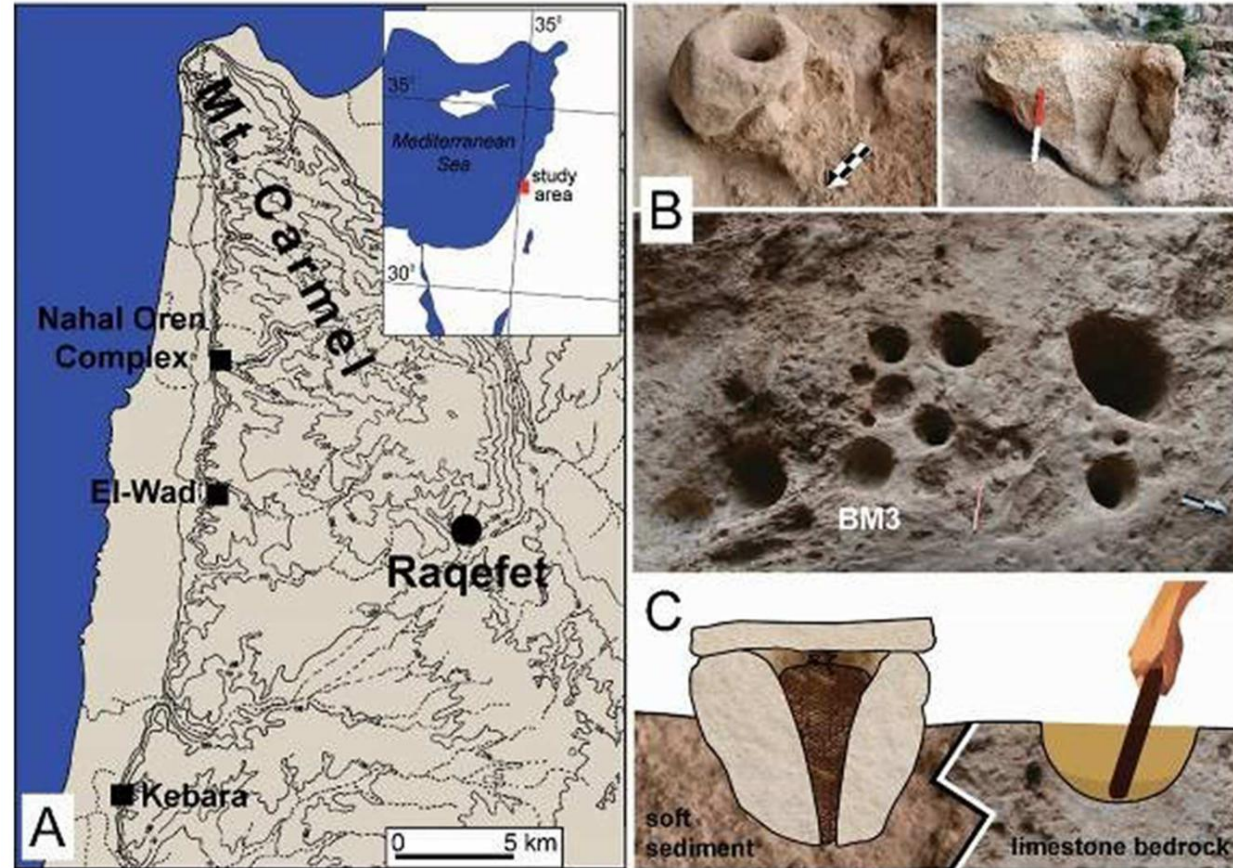archaea
7 Gt C

viruses
0.2 Gt C

plants
450 Gt C

bacteria
70 Gt C

protists
4 Gt C

fungi
12 Gt C

animals
2 Gt C

Bar-On *et.al. PNAS,* 2018

# Life on the planetary timescale



**4.5: moon formation**

4.5-4.0: $CO_2$ sequestration

3.9-3.8: late bombardment

3.8-3.5: first life

3.3-2.9: archean expansion

**2.8-2.5: photosynthesis**

2.4-2.3: great oxidation

1.9-1.7: eukaryote ancestor

0.54-0.52: cambrian explosion

0.02: C4 photosynthesis

**Hadean**    **Archean**    **Proterozoic**    **Phanerozoic**

4.6    4.0    day

# We knew microbial ecology before knowing microbes



Fermented beverage and food storage in 13,000 year-old stone mortars

# The Unseen Majority



1676

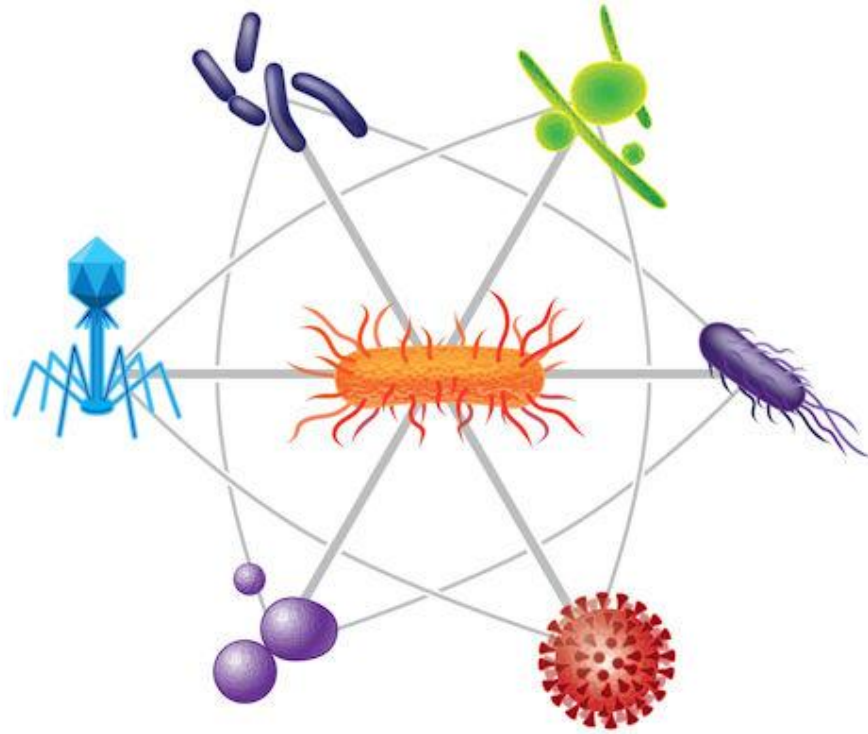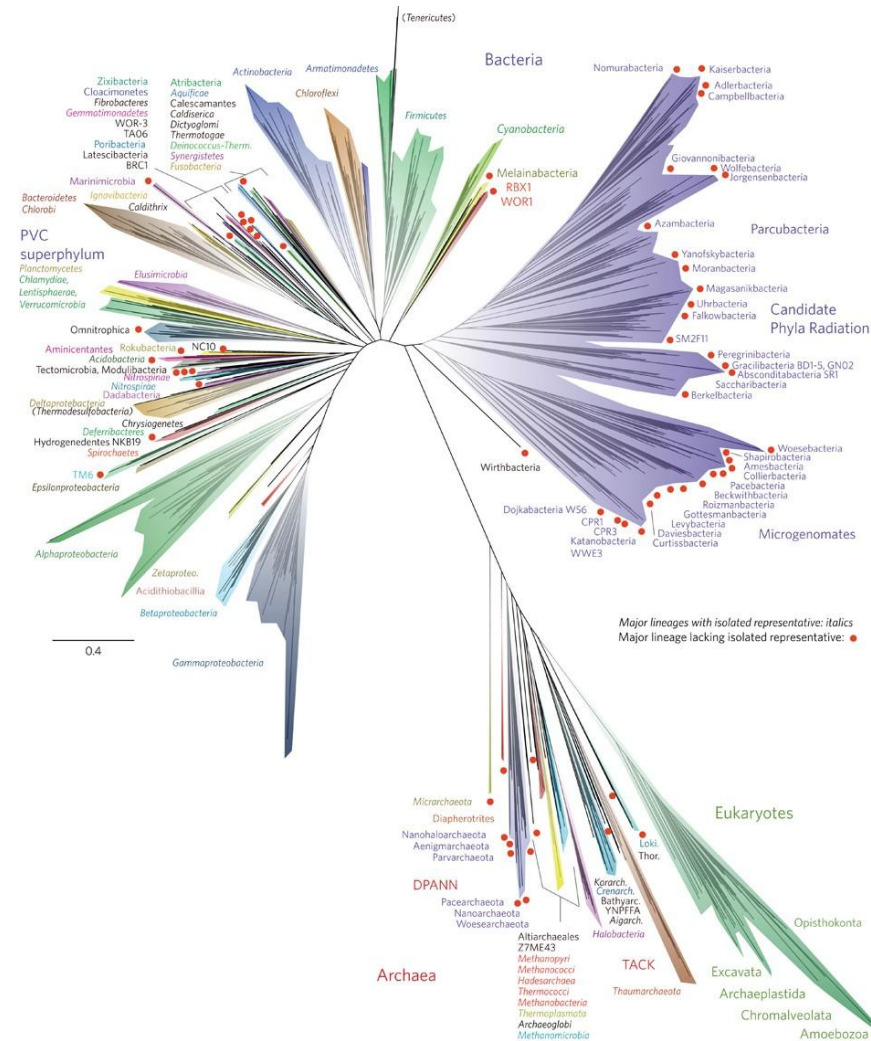# The Unseen Majority: Great plate count anomaly



1985

# The Unseen Majority: Great plate count anomaly
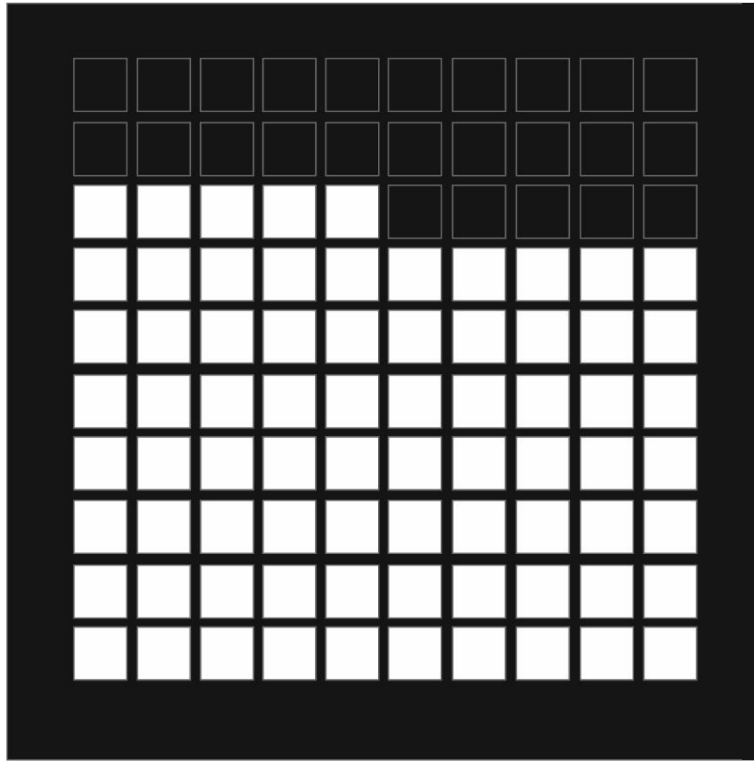


1985

# The Unseen Majority: metagenomics revolution

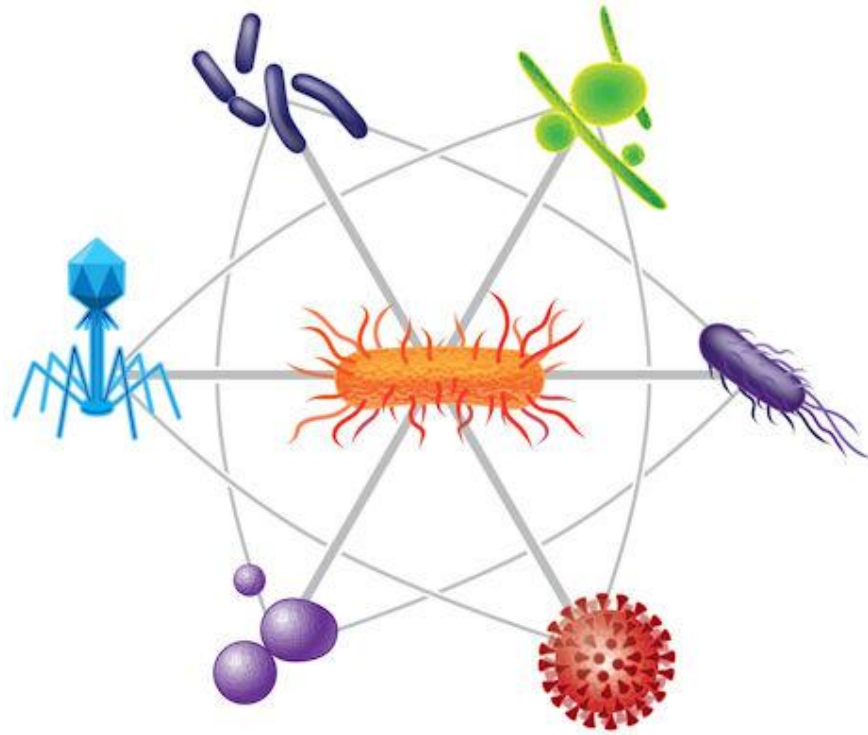# The Unseen Majority: metagenomics revolution



75 %

Now

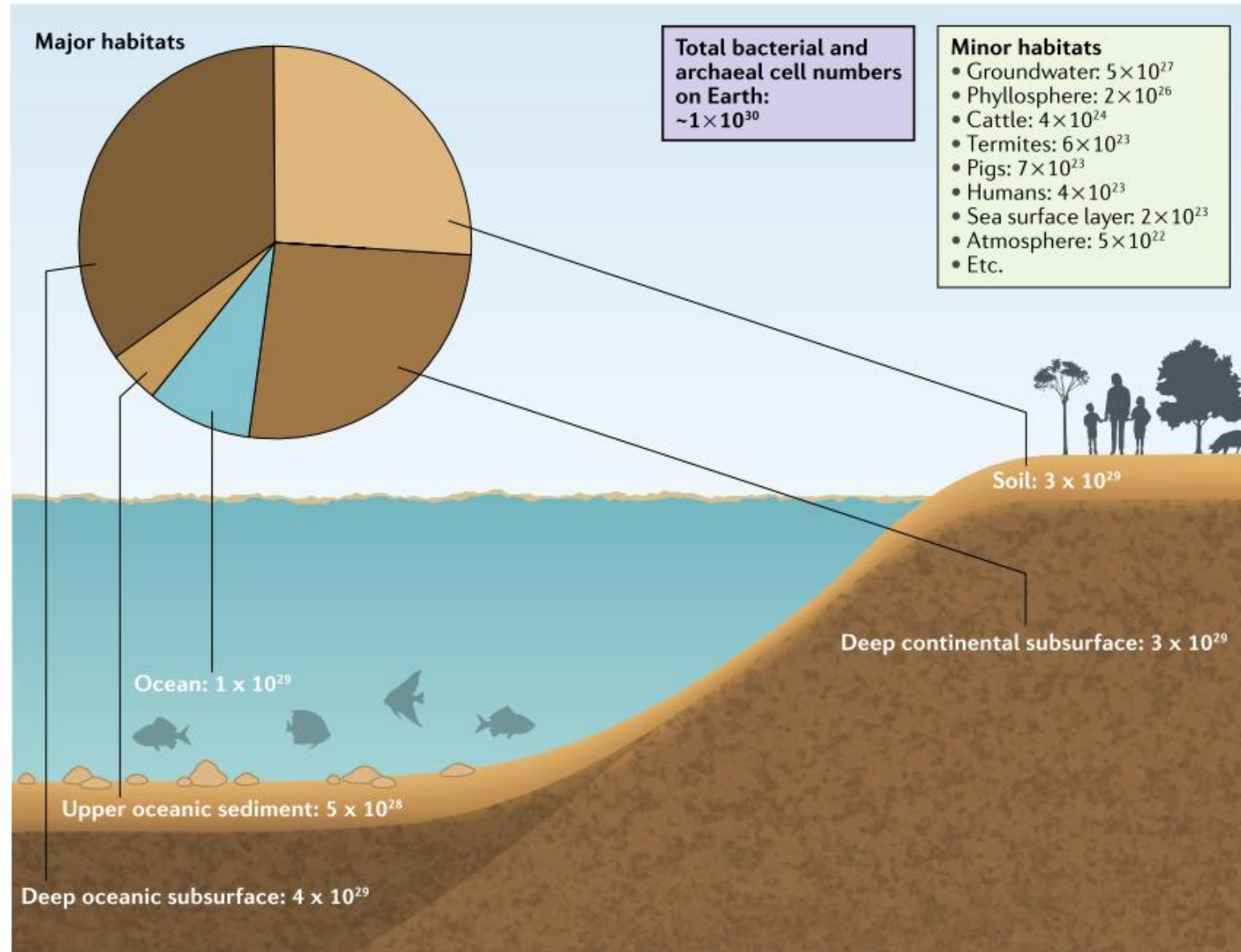# The Unseen Majority: metagenomics revolution



Now

This enormous biomass is distributed in microscopic cells

$\sim 1.2 \times 10^{30}$ bacterial/archaeal cells exist in the "big five" habitats of Earth
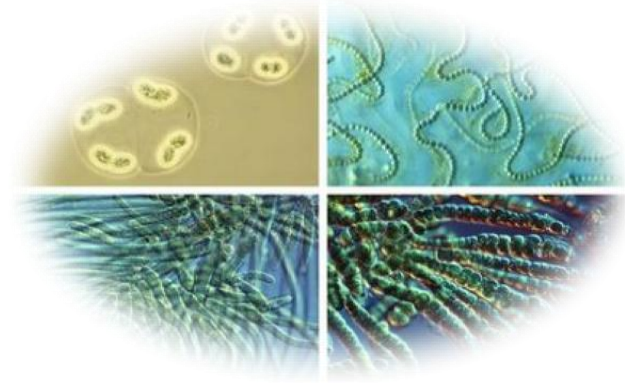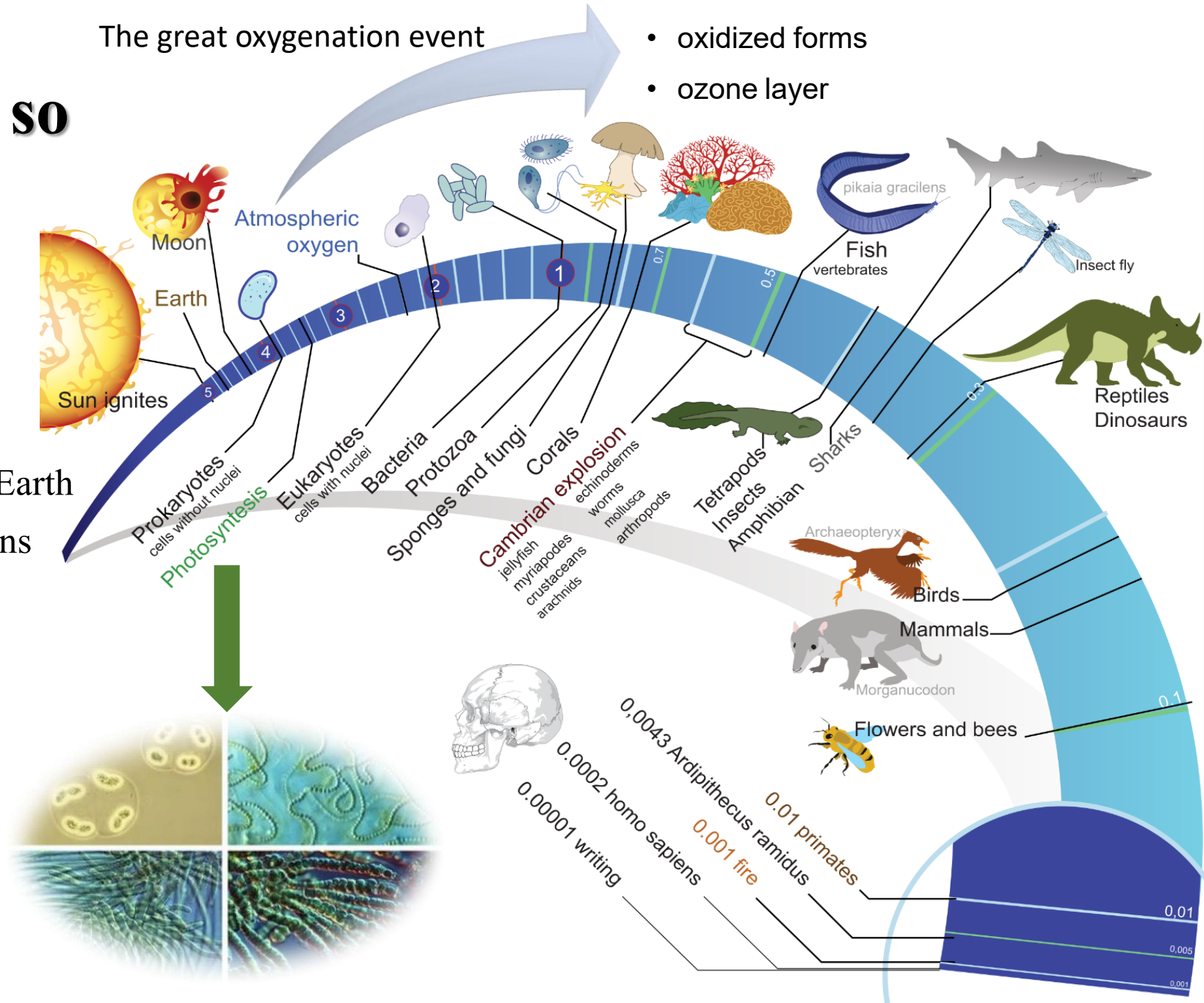
if all the $1 \times 10^{31}$ viruses on earth were laid end to end, they would stretch for 100 million light years

Astronomical numbers!



Major habitats

Total bacterial and archaeal cell numbers on Earth: $\sim 1 \times 10^{30}$

Minor habitats
• Groundwater: $5 \times 10^{27}$
• Phyllosphere: $2 \times 10^{26}$
• Cattle: $4 \times 10^{24}$
• Termites: $6 \times 10^{23}$
• Pigs: $7 \times 10^{23}$
• Humans: $4 \times 10^{23}$
• Sea surface layer: $2 \times 10^{23}$
• Atmosphere: $5 \times 10^{22}$
• Etc.

Soil: $3 \times 10^{29}$

Deep continental subsurface: $3 \times 10^{29}$

Ocean: $1 \times 10^{29}$

Upper oceanic sediment: $5 \times 10^{28}$

Deep oceanic subsurface: $4 \times 10^{29}$

(Flemming *et.al.* 2018. *Nat Rev Microbiol*)

# Why are microbes so diverse?

- Evolved early
- Initially access to all habitats on Earth
- Survided a large span of conditions
- More habitas, more niches
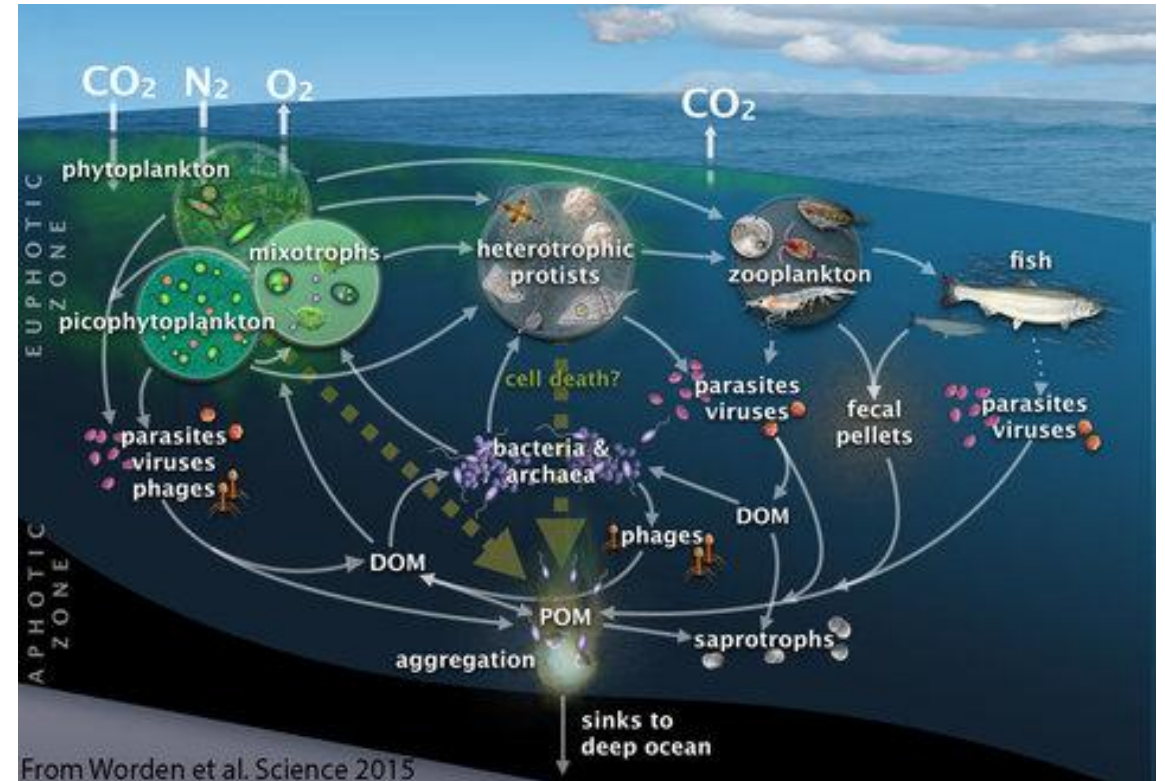- Short generation times
- Inter-species gene transfer

# Microbes in carbon cycle

Microbes transfer an enormous flow of carbon through
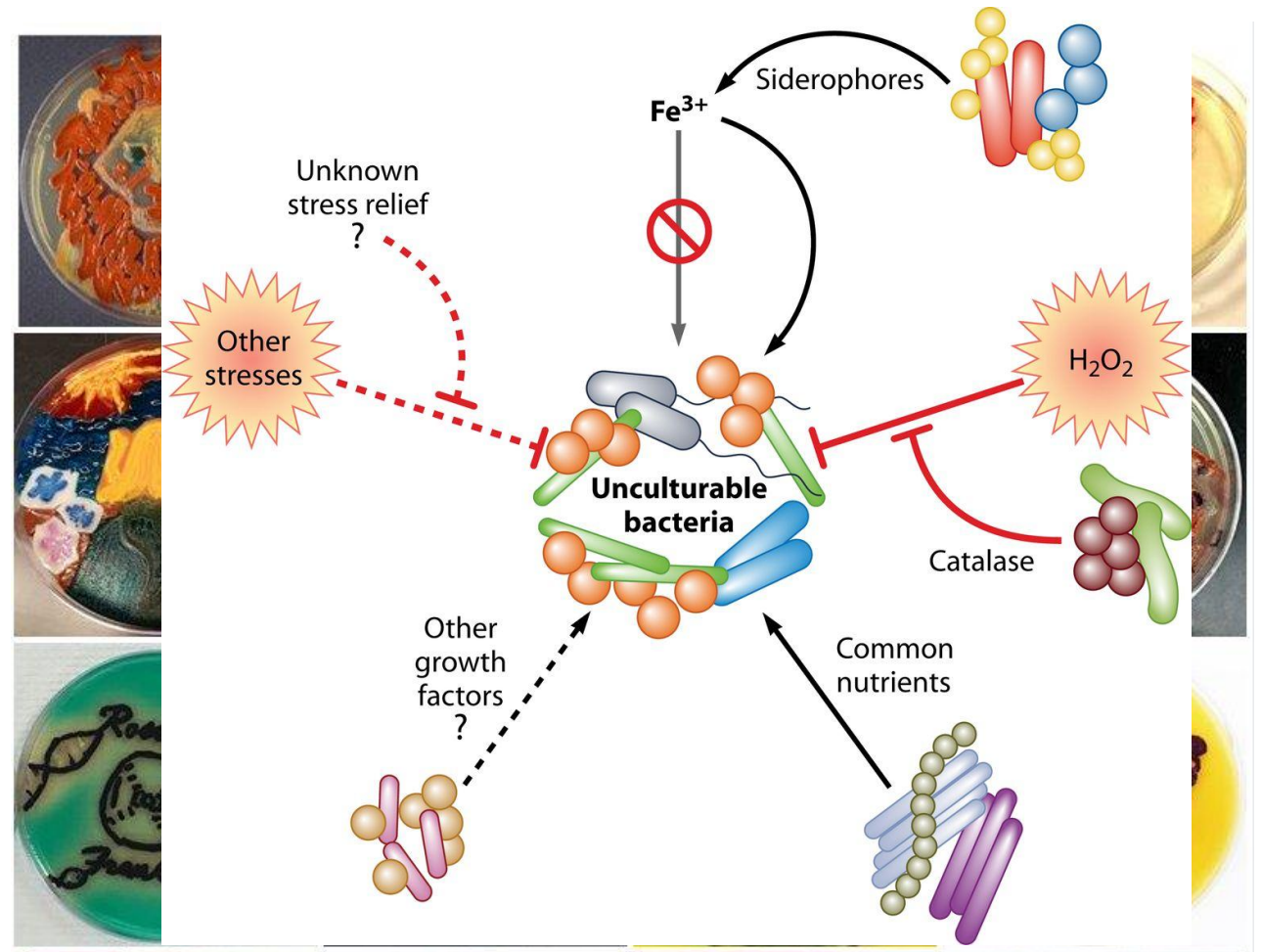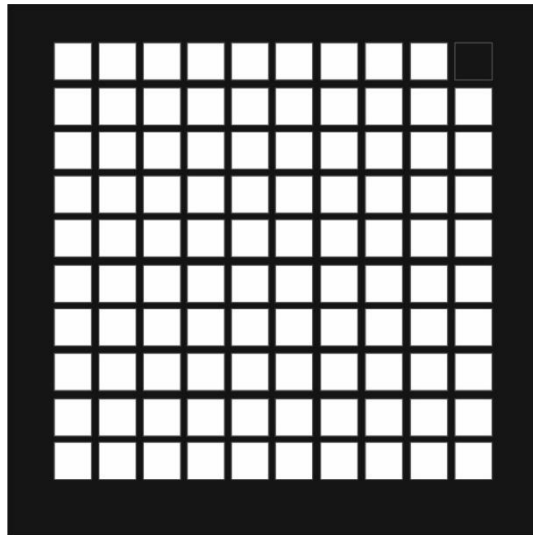
Trophic interaction
Metabolism

Greenhouse gas sink/emission control
Methane
Nitrous Oxide
CO2

**Climate impact**
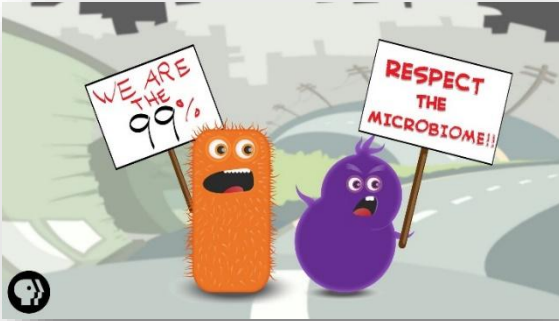


From Worden et al. Science 2015

# Challenge: Great plate count anomaly

Using conventional cultivation techniques only **0.1-1%** of prokaryotes are Culturable in laboratory conditions.
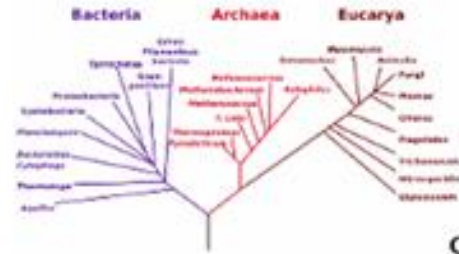
# Microbiome analysis

Desired efficiency anything more that **0.1-1%**

**Bacteria** **Archaea** **Eucarya**

Yarza et al., 2014

**Giovannoni et al., perform the first microbial community study by 16S rRNA libraries**

**Robert Koch isolates microorganisms using solid cultures**

**Carl Woese propose rRNA as marker for taxonomy**

1676    1931    1980    1998    2005    2006    2008    2011    2015

1888    1977    1990

**Leeuwenhoek reports his observations about oral microbiota**

**Winogradsky microbial ecology experiments**

**Kary Mullis develops PCR**

**Handelsman et al., propose the term 'metagenomics'**

**GA sequencer from Solexa is released**

**PacBio RS sequencer is released**

**Fred Sanger develops DNA sequencing**

**First NGS machine released by Roche**

**Human Microbiome Project publication**

**Ocean Samplig Day**

Modified from https://www.hhmi.org/biointeractive/

# Global Ocean Sampling



● Past Routes: 2003–2008     ● Europe Expedition: 2009–2010

# Tara Ocean Sampling

# Tara Ocean Sampling



Data type

- **Metagenome** (51 samples, 18 stations)
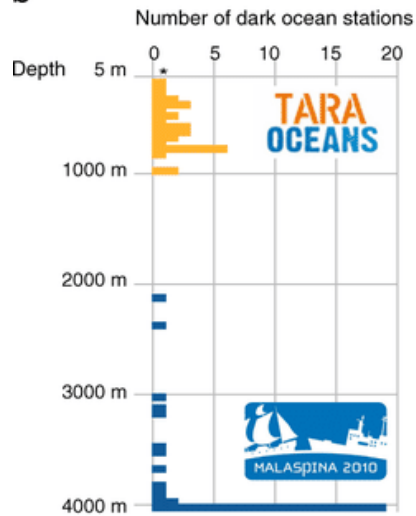- **Metatranscriptome** (58 samples, 40 stations)
- **Both** (129 samples, 68 stations)

# Malaspina Sampling

# Soil metagenomes

**a**

Soil metagenomes
( *n* = 3,304 )

- Assembly
- Binning
- Completeness ≥ 50%
- Contamination < 10%

36,398 medium quality MAGs
✓ completeness ≥ 50%
✓ contamination < 10%

**40,039 MAGs**

3,641 high quality MAGs
✓ completeness > 90%
✓ contamination < 5%
✓ tRNA number ≥ 18
✓ presence of 23S, 16S and 5/5.8S rRNA

**b**

**Ecosystem
(Number of samples)**

| | | |
|---|---|---|
| ○ Cultivated land (889) | ○ Artificial Surfaces (405) | ○ Tundra (147) |
| ○ Forest (621) | ○ Wetland (356) | ○ Shrubland (29) |
| ○ Grassland (567) | ○ Bare Land (224) | ○ Glacier (2) |

**Number of MAGs**   · 1   ○ 10   ◯ 100   ◯ 500

**c**

Completeness (%)

Contamination (%)

Genome size (bp)

N50 (bp)

Strain heterogeneity (%)

**MAG quality** ● High quality ● Medium quality

**d**

| | |
|---|---|
| Almeida et al. (2019) | 92143 |
| Nayfach et al. (2019) | 60664 |
| Marine | 52325 |
| SMAG (current study) | 40039 |
| Host associated (human) | 16441 |
| Freshwater | 7335 |
| UHGG | 4644 |
| Engineered (built environment) | 2640 |
| Engineered wastewater | 2627 |
| GEM soil | 2461 |
| Host associated (mammals) | 1955 |
| Aquatic (non marine saline and alkaline) | 1735 |
| Aquatic (thermal springs) | 1579 |
| Host associated (plants) | 1131 |
| Engineered (lab enrichment) | 800 |
| Sediment | 42 |

Number of MAGs

**MAG Quality**
High
Medium

# Deep groundwater metagenomes



a

**Alberta**
73-135 m
103Gb

**Soudan Iron Mine**
715 m
54 Gb

**McLaughlin Reserv**
76.2 m
6 Gb

**San Joaquin Valley**
100 m
20 Gb

**Illinois Basin-Decatur**
1800 m
1.4 Gb

**Crystal Geyser**
320-800 m
1.304 Gb

**Olkiluoto Island**
330-531 m
258 Gb

**Tomsk**
2600-2800 m
52 Gb

**Äspö HRL**
70-454 m
1,157 Gb

**Mont Terri**
226-562 m
417 Gb

**Horonobe URL**
140-250 m
30 Gb

**Huelva**
420-468 m
55 Gb

**TauTona mine**
3048-3136 m
13 Gb

**Driefontein mine**
1046 m
6 Gb

**Thabazimbi**
2100 m
5 Gb

**Masimong mine**
1900 m
4 Gb

**Finsch mine**
1056 m
5 Gb

**Beatrix mine**
1339-1340 m
13 Gb

**Welkom**
1339 m
161 Gb

b

Bases (Gb)

Groundwater depth (mbsl)

Continent
- Africa
- America
- Asia
- Europe

Depth categories
- (i) 70-999 mbsl
- (ii) 1,000-1,999 mbsl
- (iii) ≥ 2,000 mbsl

# How is Metagenomics done and what does it tell us?

# Capturing the unseen Majority or their footprint…



**Great Plate Count Anomaly**
Only ~1% of Bacteria is Culturable

Environmental Sample

Microscope

Agar Plate

Sequencer

Genomes

GTACATGACTAGATCAT
AGACTGGATCGATCCAG

ACGTGTACGTACGTAAG
GTACATGACTAGATCAT
AGACTGGATCGATCCAG
GGACCTAGCTAAGCTAG

Uncultivated Sequencing

Parks *et.al.* NatMicrobiol, 2017

# Amplicon vs. Metagenomics

- Less complex
  - Better coverage
  - More samples
- Extensive database
- Same fragment
  - Comparable phylogenetic info
- Qualitative
- PCR and primer bias
- Limited phylogenetic info
- Limited functional information

# Metagenomics Workflow

# What to think of when designing a experiment...



Optimized sample collection, preparation, and DNA extraction

Scientific question

Sample type

Sample should be representative

Remember all metagenomics values are RELATIVE

Process samples fast

and preserve DNA properly

& Metagenomics



What is GIGO?

The quality of information coming out cannot be better than the quality of information that went in.

GARBAGE IN

GARBAGE OUT

GIGO is used in IT and mathematics

Garbage In, Garbage Out

# Metagenomics Workflow

# DNA recovery method impacts the output

- Choice of DNA extraction method
  - Consistent method for all experiments we want to compare
- DNA extraction quality
  - Gel electrophoresis
    - The integrity and size of genomic DNA
  - Spectrophotometry
    - Pure DNA has an A260/A280 ratio of 1.7–1.9
    - DNA concentration has been determined using nanodrop
  - Fluorometry concentration measurements

# Metagenomics Workflow

# Sequencing

**Long read**

**Short read**

Consider the sequencing quality, read length, and price ❗

Deeper sequencing = higher resolution also = computationally intensive

Let's have a break and come back in 20 min

# Metagenomics Workflow

# Sequence analysis

**Assembly**

# Assembly tools

- **MEGAHIT:** makes use of succinct *de Bruijn* graphs (SdBG; Bowe *et al.*, 2012), which are compressed representation of *de Bruijn* graphs.

- **metaSPAdes:** first constructs the de Bruijn graph of all reads using SPAdes, transforms it into the assembly graph using various graph simplification procedures, and reconstructs paths in the assembly graph that correspond to long genomic fragments within a metagenome.



1. Fragment DNA and sequence

2. Find overlaps between reads

...AGCCTAGACCTACAGGATGCGCGACACGT
GGATGCGCGACACGTCGCATATCCGGT...

3. Assemble overlaps into contigs

4. Assemble contigs into scaffolds

# Metagenomics Workflow

# Sequence analysis

**Binning**

MAG (Metagenome assembled genome)

MAG quality check



**Preprocessing**

1. Samples from multiple sites or times

2. Metagenome libraries

3. Initial de-novo assembly using the combined library

**MetaBAT**

4. Calculate TNF for each contig

5. Calculate Abundance per library for each contig

6. Calculate the pairwise distance matrix using pre-trained probabilistic models

7. Forming genome bins iteratively

TetraNucleotides Frequency          Abundance

# Binning tools

- **MetaBAT2:** uses the same raw TNF and abundance (ABD) scores

- **CONCOCT:** does unsupervised binning of metagenomic contigs by using nucleotide composition - kmer frequencies - and coverage data for multiple samples.

- **MaxBin:** algorithm utilizes two different genomic features: tetranucleotide frequencies and scaffold coverage levels to populate the genomic bins using single-copy maker genes and an expectation-maximization algorithm.

# Genome/MAG quality check

- CheckM provides robust estimates of genome completeness and contamination by using collocated sets of genes that are ubiquitous and single-copy within a phylogenetic lineage.

- Assessment of genome quality using plots depicting key genomic characteristics (e.g., GC, coding density) which highlight sequences outside the expected distributions of a typical genome.

- CheckM also identifies genome bins that are likely candidates for merging based on marker set compatibility, similarity in genomic characteristics, and proximity within a reference genome tree.

- https://ecogenomics.github.io/CheckM/

# Genome completeness standards

## Table 1 Genome reporting standards for SAGs and MAGs

| Criterion | Description |
| --- | --- |
| **Finished (SAG/MAG)** | |
| Assembly quality[a] | Single contiguous sequence without gaps or ambiguities with a consensus error rate equivalent to Q50 or better |
| **High-quality draft (SAG/MAG)** | |
| Assembly quality[a] | Multiple fragments where gaps span repetitive regions. Presence of the 23S, 16S, and 5S rRNA genes and at least 18 tRNAs. |
| Completion[b] | >90% |
| Contamination[c] | <5% |
| **Medium-quality draft (SAG/MAG)** | |
| Assembly quality[a] | Many fragments with little to no review of assembly other than reporting of standard assembly statistics. |
| Completion[b] | ≥50% |
| Contamination[c] | <10% |
| **Low-quality draft (SAG/MAG)** | |
| Assembly quality[a] | Many fragments with little to no review of assembly other than reporting of standard assembly statistics. |
| Completion[b] | <50% |
| Contamination[c] | <10% |

This is a compressed set of genome reporting standards for SAGs and MAGs. For a complete list of mandatory and optional standards, see **Supplementary Table 1**.

[a]Assembly statistics include but are not limited to: N50, L50, largest contig, number of contigs, assembly size, percentage of reads that map back to the assembly, and number of predicted genes per genome. [b]Completion: ratio of observed single-copy marker genes to total single-copy marker genes in chosen marker gene set. [c]Contamination: ratio of observed single-copy marker genes in ≥2 copies to total single-copy marker genes in chosen marker gene set.

# Genome taxonomy

- Taxonomy and nomenclature

- https://gtdb.ecogenomic.org/

- https://ncbiinsights.ncbi.nlm.nih.gov/2021/12/10/ncbi-taxonomy-prokaryote-phyla-added/

# Sequence analysis

## Annotation

Gene prediction => ORF finding

    Prodigal



Function assignment

    BLAST

    HMM models

# Lists of IDs lists of names

- **BioCyc**
- **KEGG**
- **Ensembl Bacteria**
- **Kbase**
- **IMG**
- **PATRIC**

# KEGG

- The Kyoto Encyclopedia of Genes and Genomes is a resource for understanding high-level functions of a biological system from molecular-level information.

- 

- Tools for analysis of large-scale molecular datasets generated by high-throughput experimental technologies.
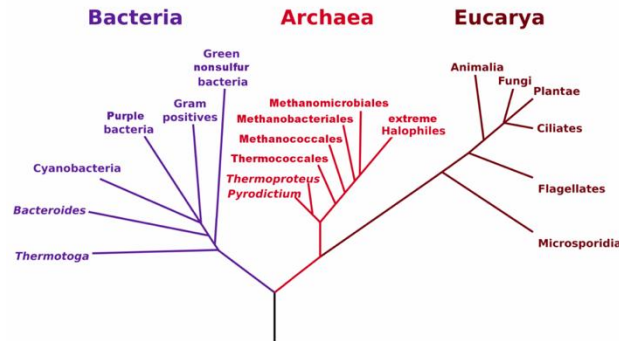

- Home page: https://www.kegg.jp/

# From Protein to metabolism …

# The chase takes time …
## The case of SAR202



*Proc. Natl. Acad. Sci. USA*
Vol. 93, pp. 7979–7984, July 1996
Microbiology

## 16S rRNA genes reveal stratified open ocean bacterioplankton populations related to the Green Non-Sulfur bacteria

(molecular ecology/phylogeny/thermophily)

STEPHEN J. GIOVANNONI*, MICHAEL S. RAPPÉ, KEVIN L. VERGIN, AND NANCI L. ADAIR

**1996**

Lim et. al. *Nat. Commun,*

**2023**

Mehrshad et. al. *ISME J,*

**2018**

# Tree of Life has evolved



1837 **Charles Darwin**'s ideas on evolution, species descend from common ancestors and evolve over time.

1990 **Carl Woese** tree with LUCA and three domains. Based on rRNA gene. Later elaborated by **Norman Pace**.

2016 **Jillian Banfield** and **Laura Hug** included genomes derived from metagenomes and based on 16 ribosomal proteins for the "Hug tree".
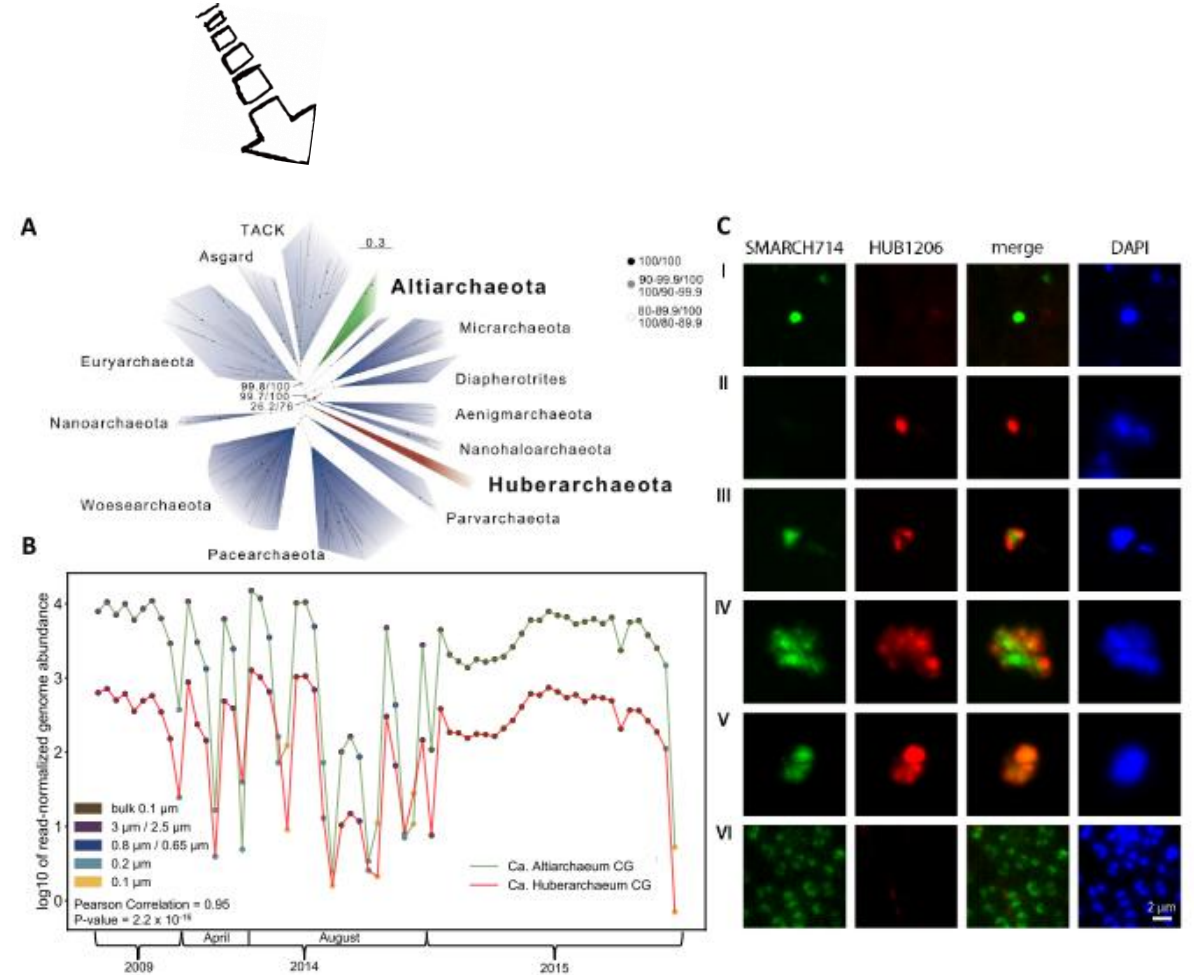
(Hug et al., 2016. Nature Microbiology)

# smallest genomes with epi-symbiotic lifestyle
# CPR & DPANN



Burstein et. al. *Nat Commun*,2019

Schwank et. al. ISMEJ,2019

# Remarkable aspect of the tree of life

- Candidate phyla radiation (CPR)

- DPANN (an acronym of the names of the first included phyla, 'Candidatus Diapherotrites', 'Candidatus Parvarchaeota', 'Candidatus Aenigmarchaeota', Nanoarchaeota and 'Candidatus Nanohaloarchaeota')

  - Small genomes
  - Small cell sizes
  - Notable gaps in core metabolic potential
  - Mostly symbiotic lifestyle
  - Their ecological role is not yet well understood

# Asgard archaea illuminate the origin of eukaryotic cellular comple



Spang et. al. *Nature*,2015
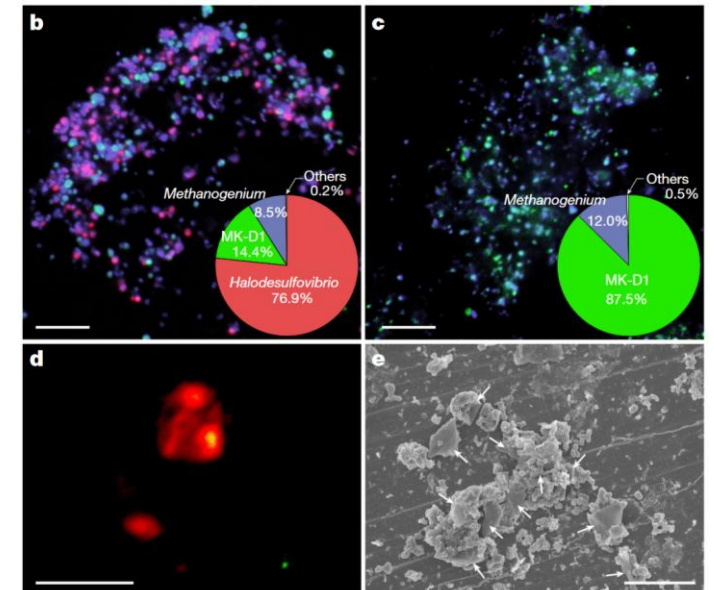
López-García et. al. *NatMicrobiol*,2019
https://www.nature.com/articles/s41564-020-0710-4

Imachi & Nobu et. al. *nature*,2020

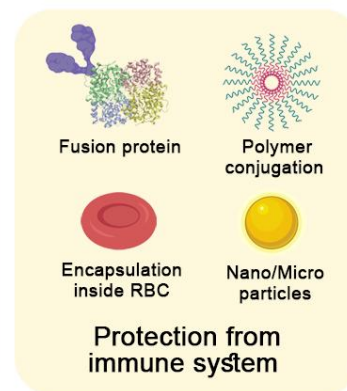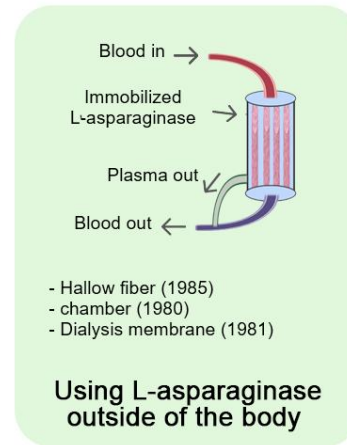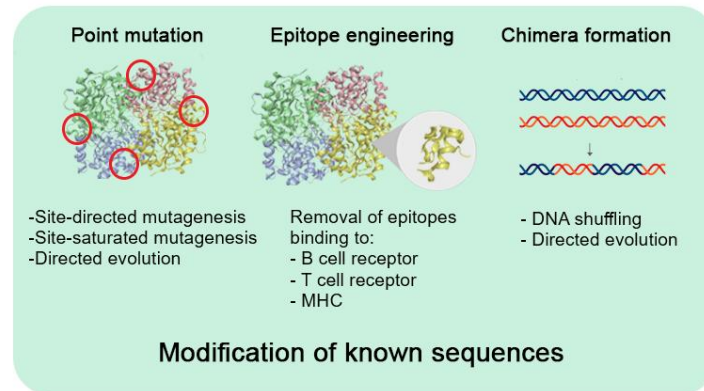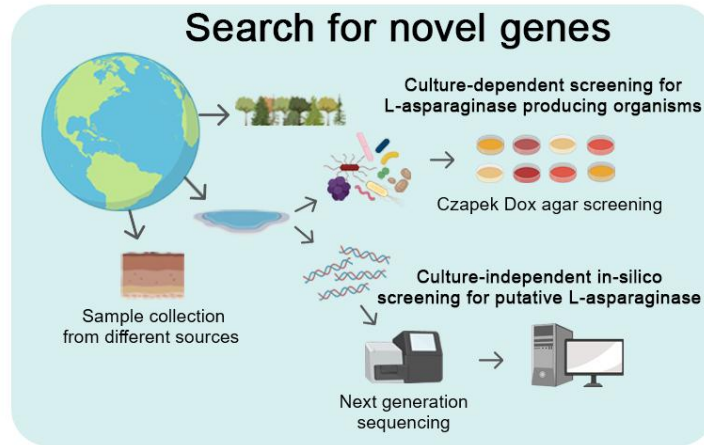# 'CandidatusPrometheoarchaeum syntrophicum'

- Pure co-culture of the target archaeon MK-D1 and Methanogenium after a 12-year study

- From bioreactor-based pre-enrichment of deep-sea sediments to a final 7 years of in vitro enrichment.

- Extremely slow growth rate and low cell yield.

- The culture consistently had a **30–60-day lag phase** and required more than 3 months to reach full growth: around $10^5$ 16S rRNA gene copies $ml^{-1}$

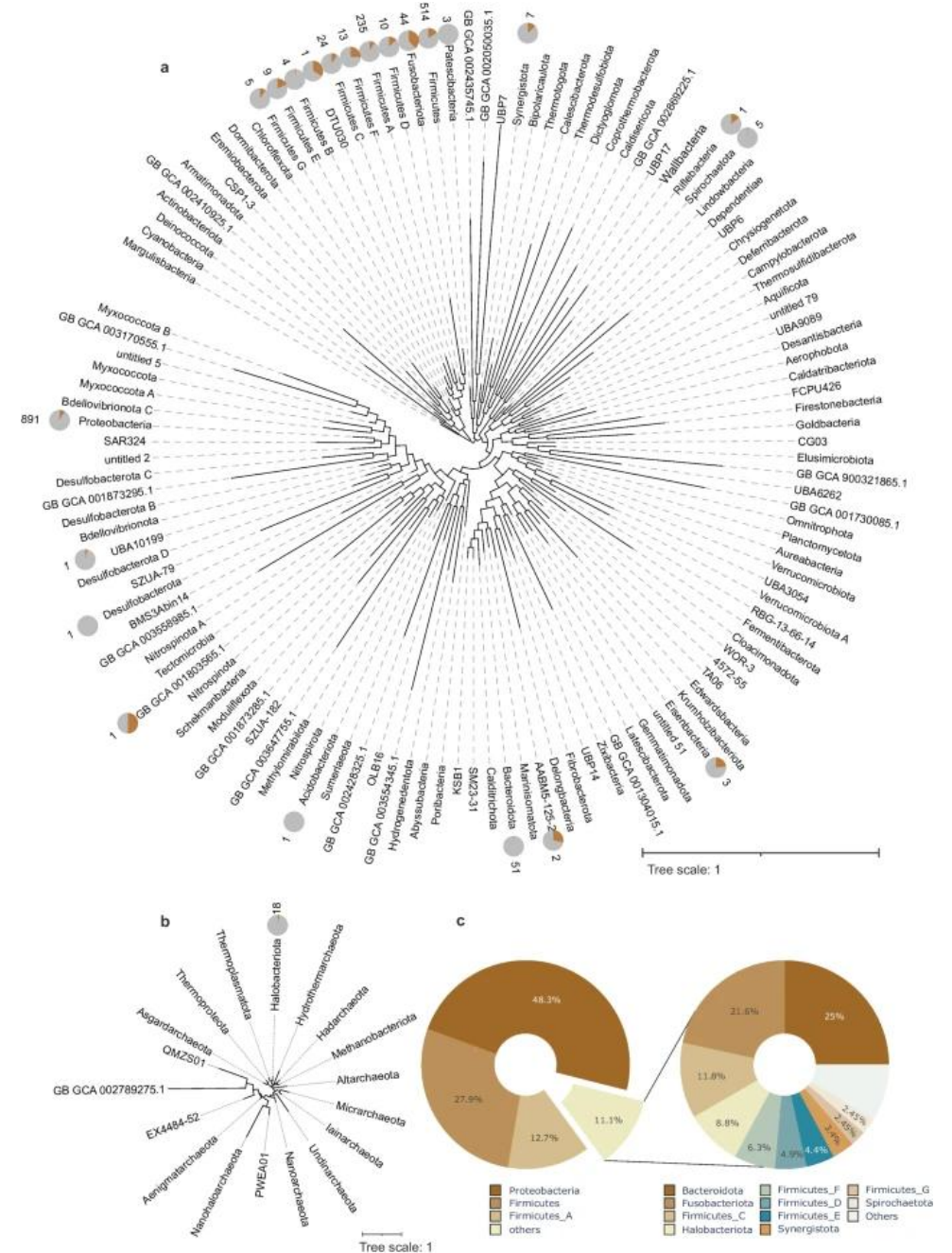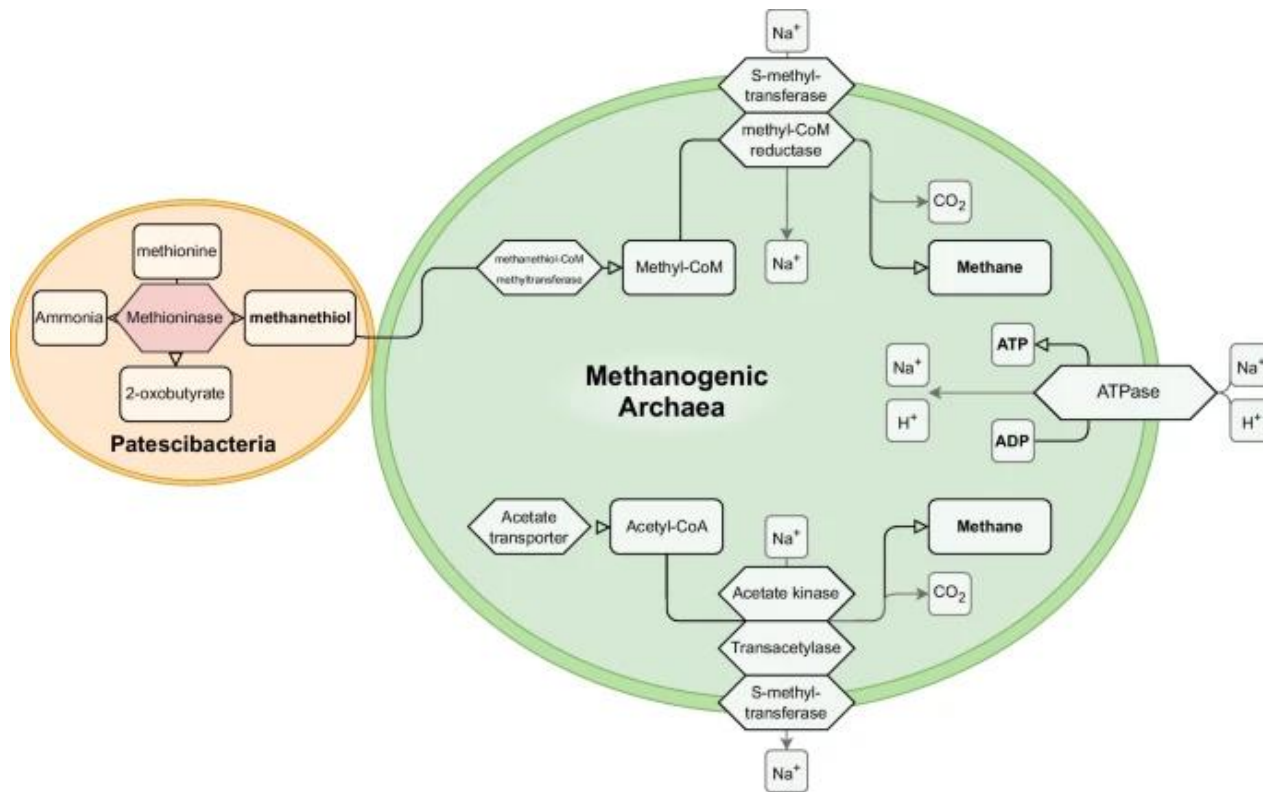- The doubling time was estimated to be approximately **14–25 days**.
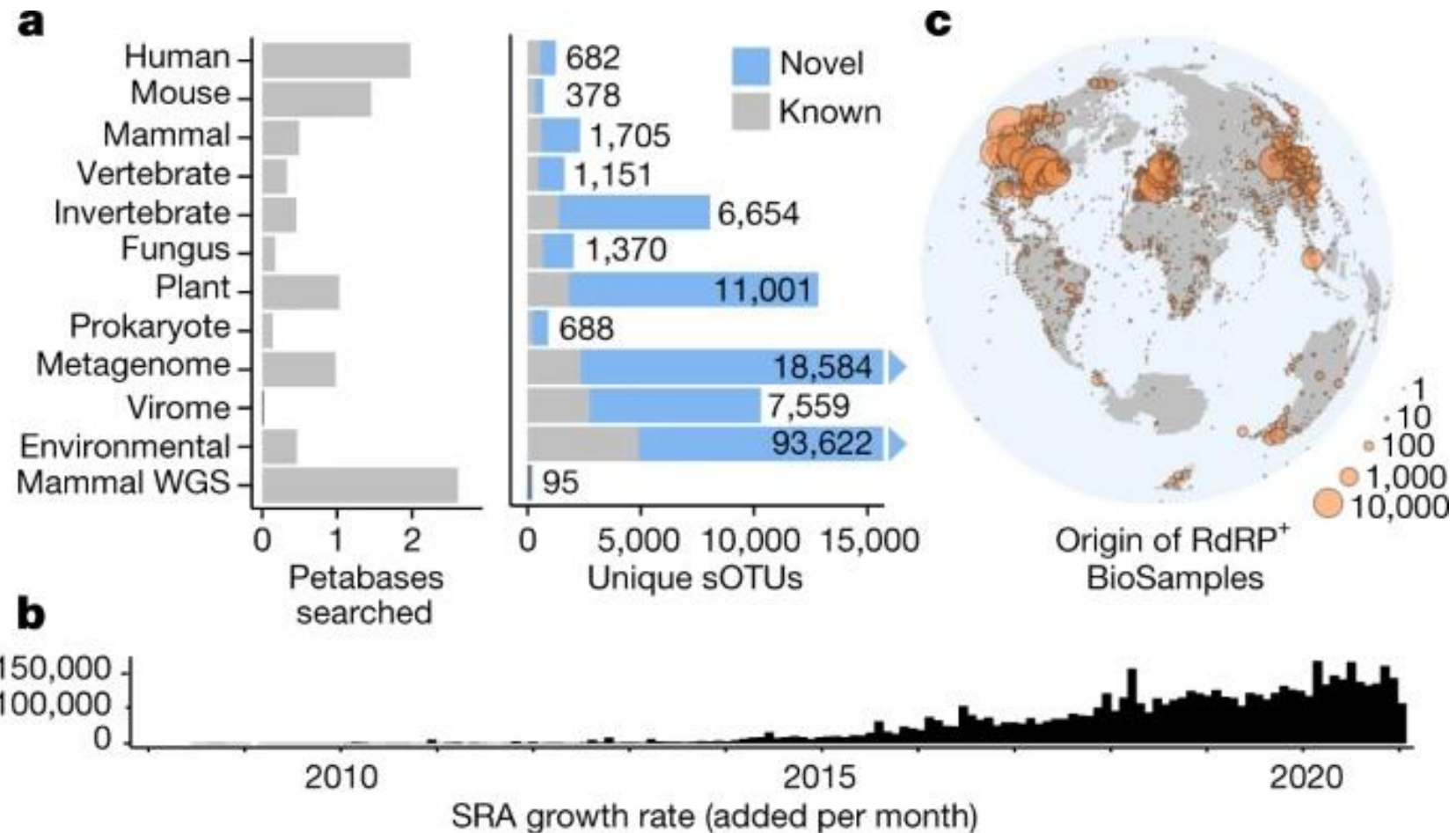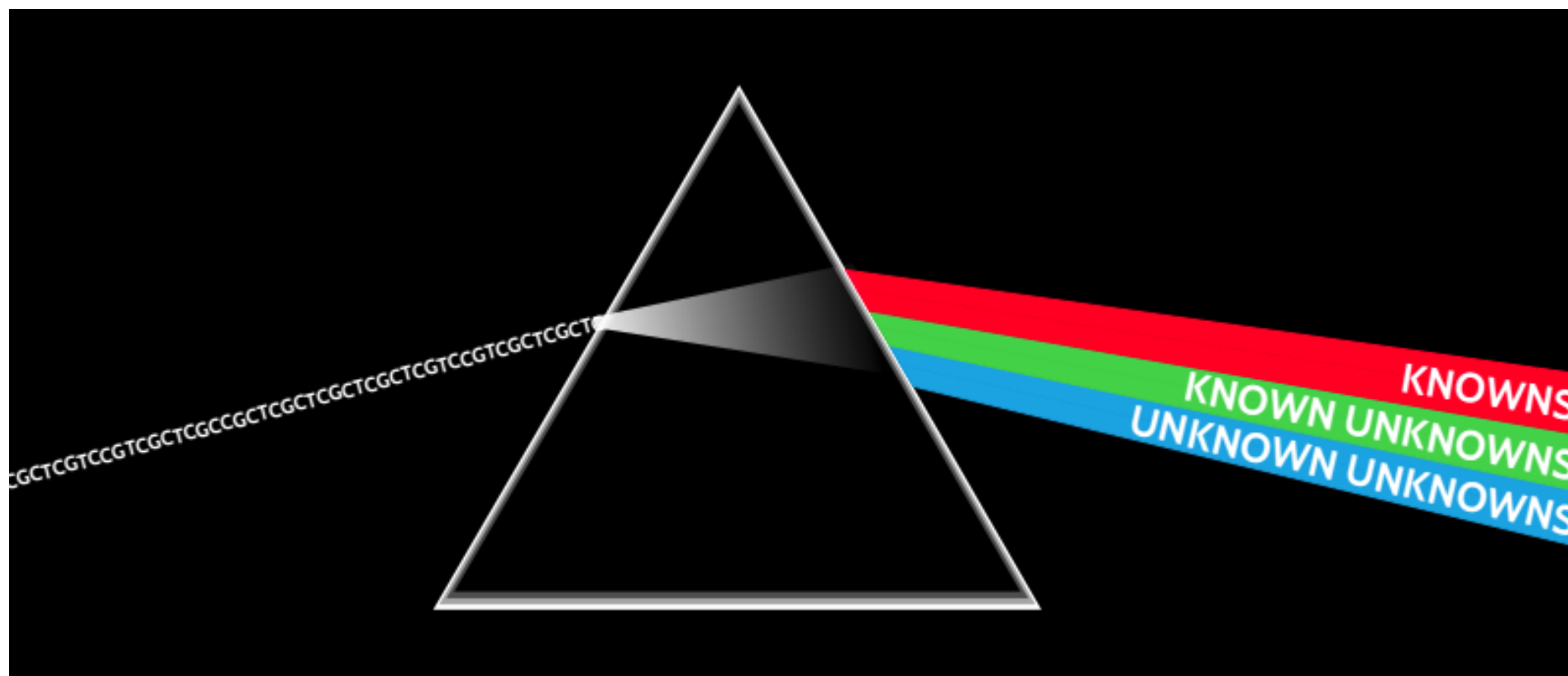
- Imachi, H., *et al. Nature,* 2020

# L-asparaginase

# Methioninase involved in Metabolic Syntrophy

# Petabase-scale sequence alignment catalyses viral discovery



Edgar. *et al. Nature* **2022**.

**"multi-omics" approach to answer eco-evolutionary questions**

Single-cell amplified genomes
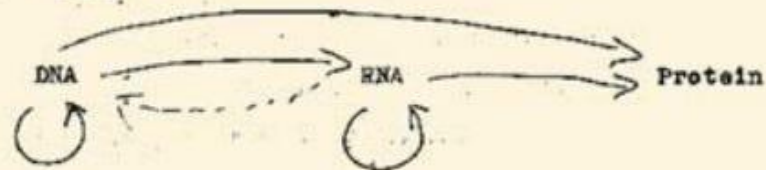
Metagenomics

Metatranscriptome

# The transfer of information
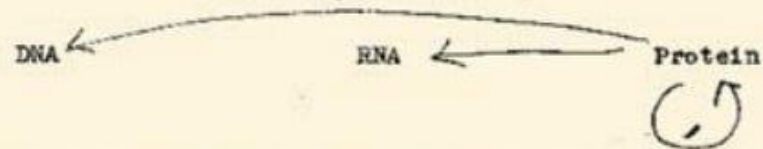


Ideas on Protein Synthesis (Oct. 1956)

The Doctrine of the Triad.

The Central Dogma: "Once information has got into a protein it can't get out again". Information here means the sequence of the amino acid residues, or other sequences related to it. That is, we may be able to have

DNA → RNA → Protein

but never

DNA ← RNA ← Protein

where the arrows show the transfer of information.

# Metatranscriptome

## What to consider?

- Sampling

- Sample processing

- RNA extraction

- Replicates

- Library preparation (rRNA depletion)

# Metatranscriptome

## What to consider?

- Sampling

- Sample processing

- RNA extraction

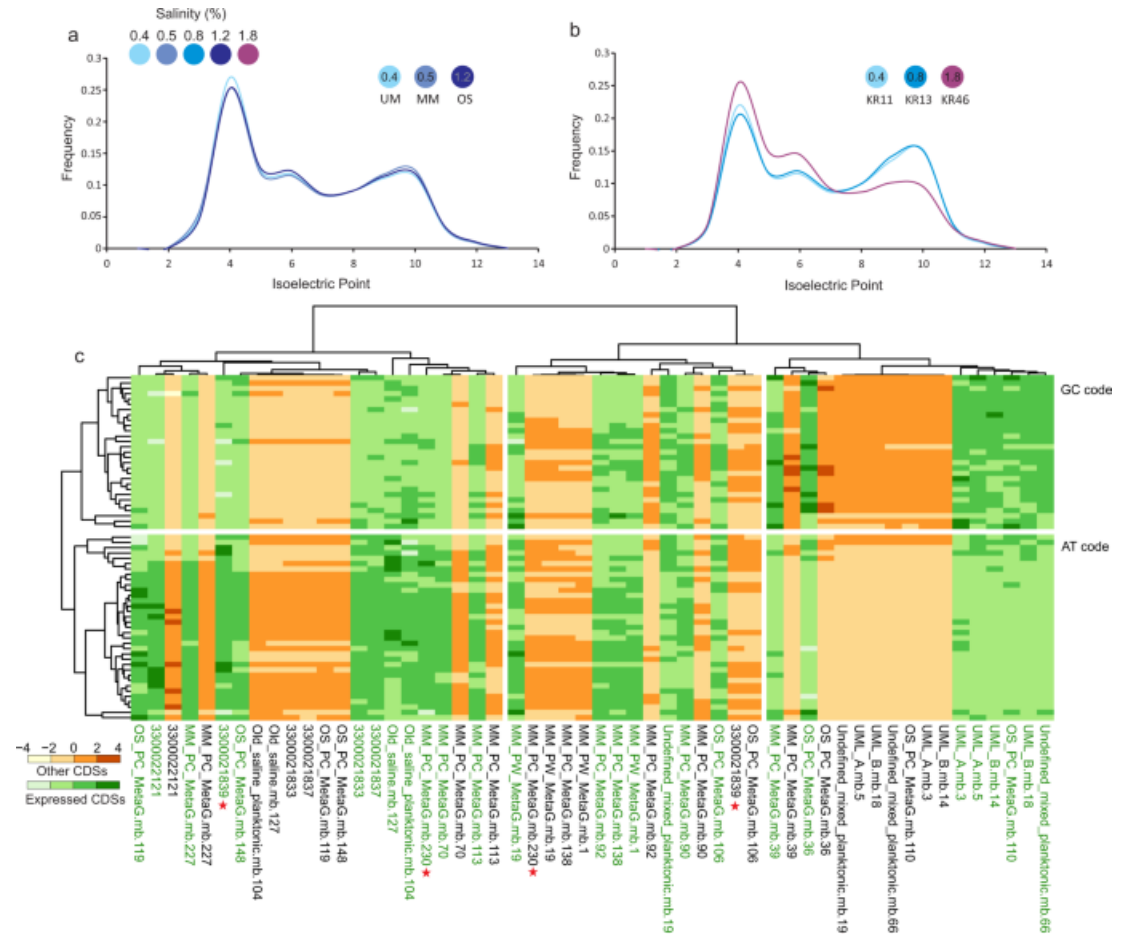- Replicates

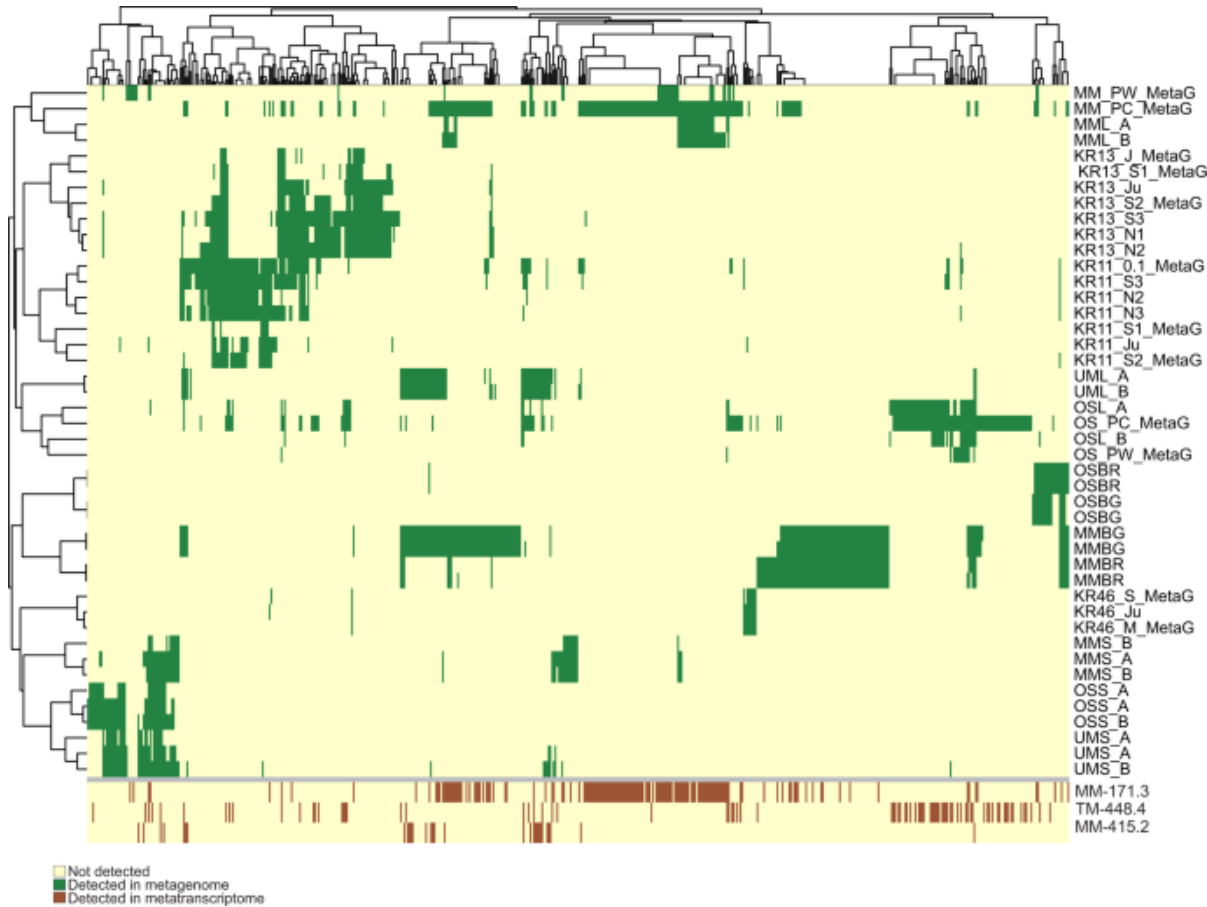- Library preparation (rRNA depletion)

## How to analyze?

- Remove rRNA or not?

- Assemble or not?

- Gene-resolved metatranscriptomics
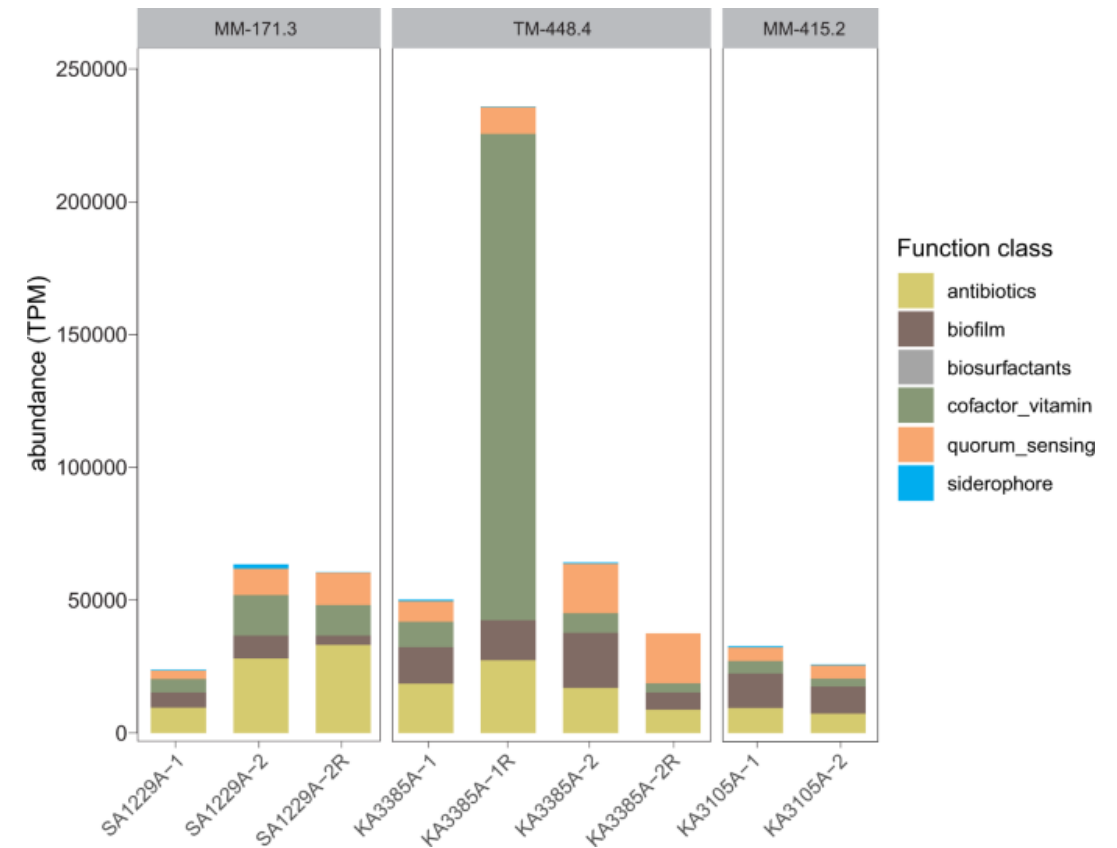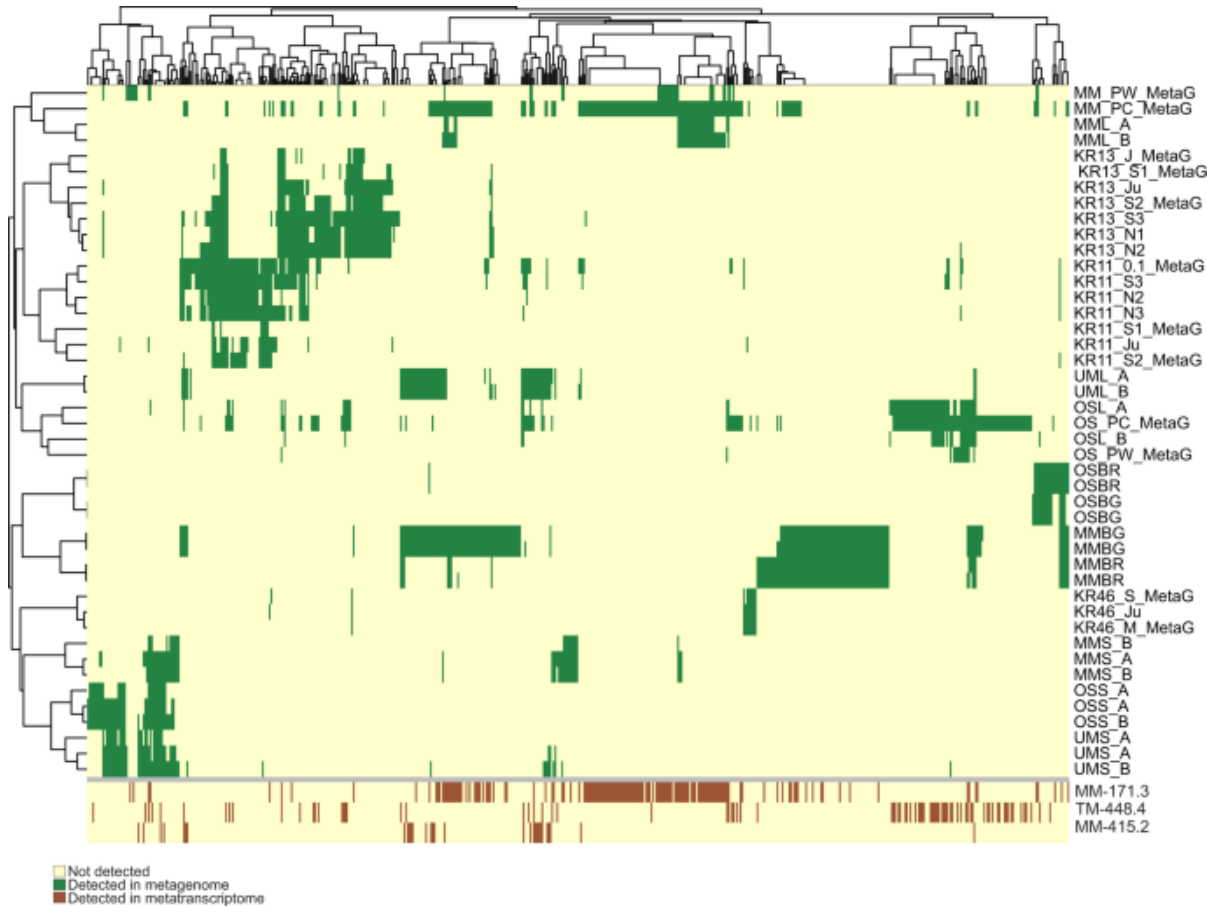
- How to normalize?

# Some normalization methods

| method | description |
|---|---|
| relative_abundance | (default) Percentage relative abundance of each genome, and the unmapped read percentage |
| mean | Average number of aligned reads overlapping each position on the genome |
| trimmed_mean | Average number of aligned reads overlapping each position after removing the most deeply and shallow-ly covered positions. See --trim-min/--trim-max to adjust. |
| coverage_histogram | Histogram of coverage depths |
| covered_bases | Number of bases covered by 1 or more reads |
| variance | Variance of coverage depths |
| length | Length of each genome in base pairs |
| count | Number of reads aligned to each genome. Note that supplementary alignments are not counted. |
| reads_per_base | Number of reads aligned divided by the length of the genome |
| anir | Average BLAST-like identity of mapped reads |
| rpkm | Reads mapped per kilobase of genome, per million mapped reads |
| tpm | Transcripts Per Million as described in Li et al 2010 https://doi.org/10.1093/bioinformatics/btp692 |

# The transfer of information

# The transfer of information

"multi-omics" approach to answer eco-evolutionary questions
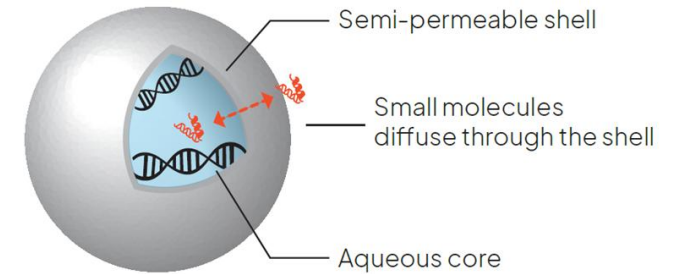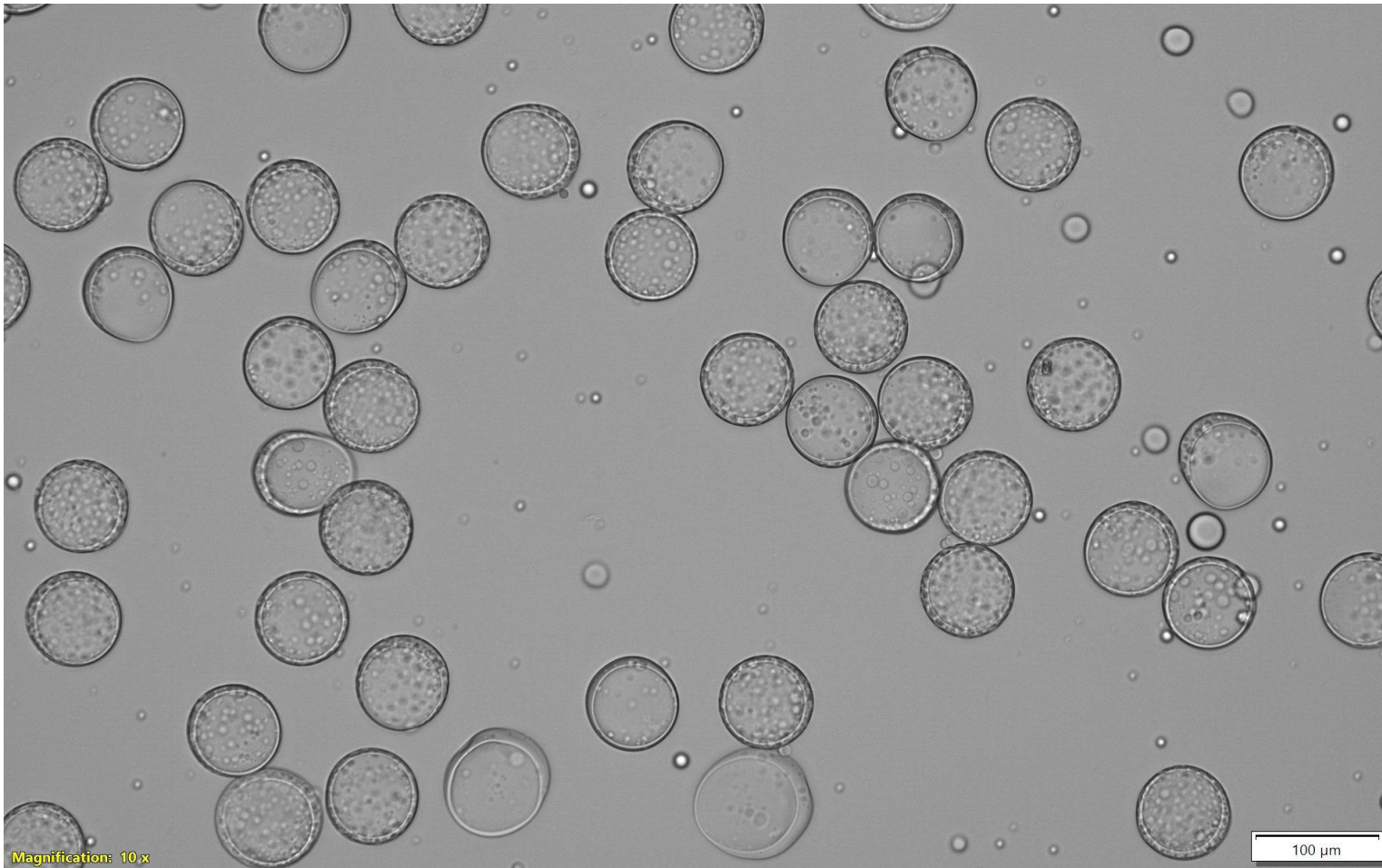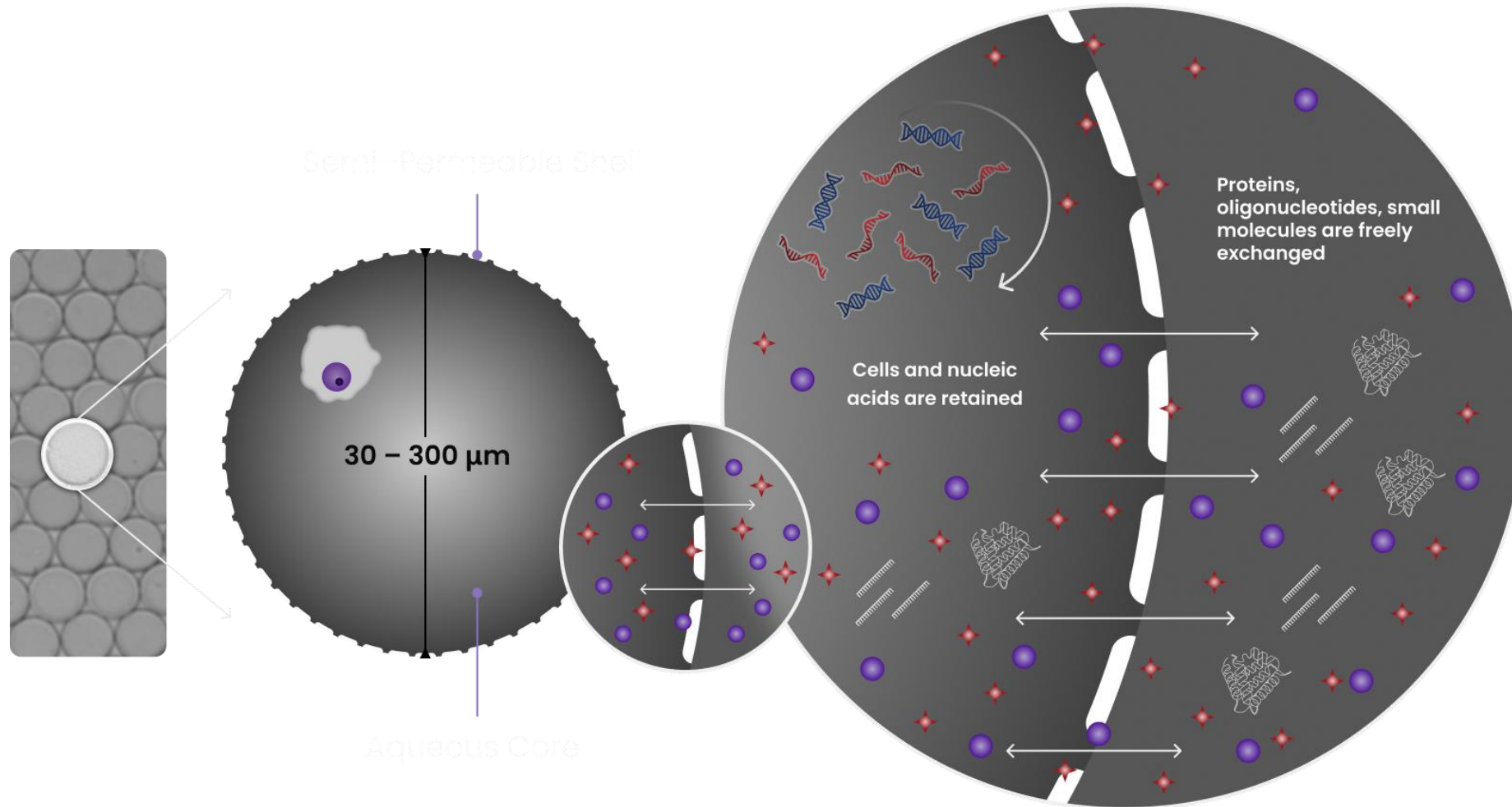
Single-cell amplified genomes
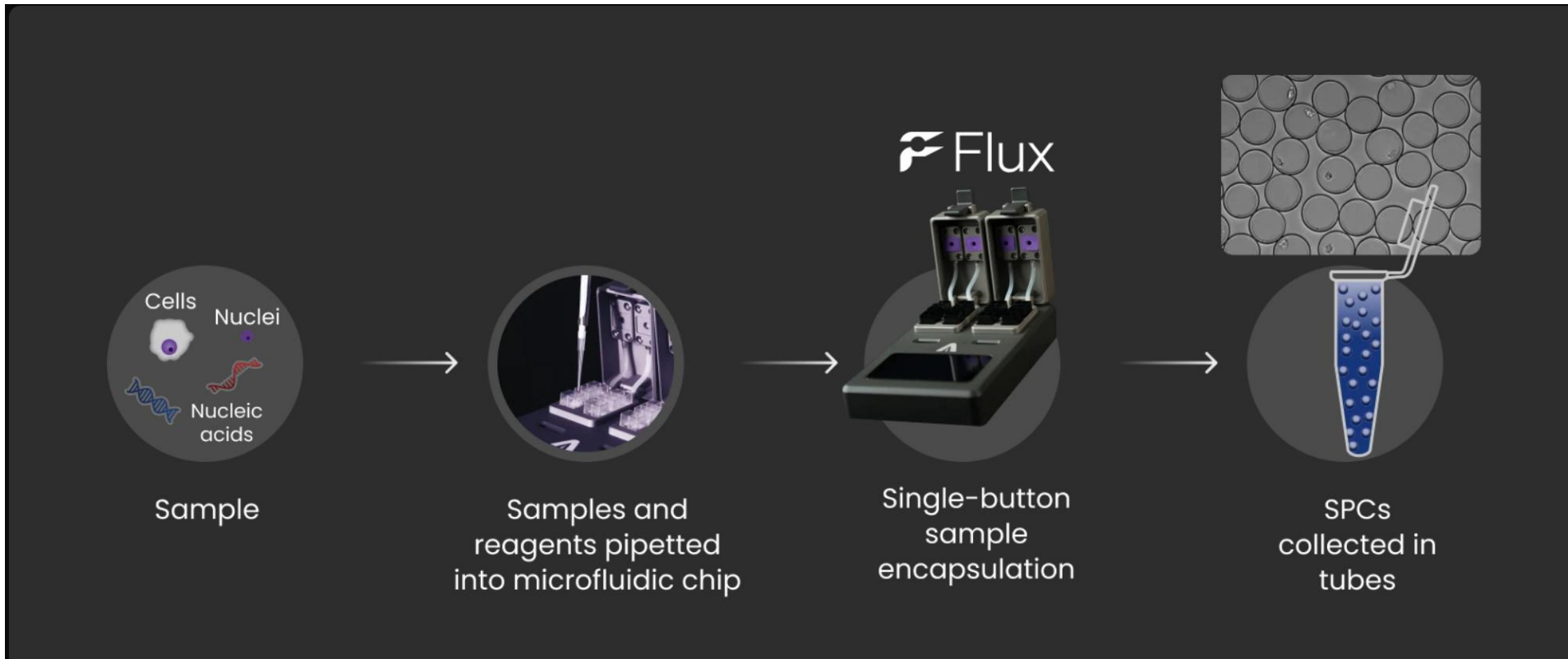
Metagenomics

Metatranscriptome

# Compartmentalized metagenomics



Magnification: 10 x

100 µm

Semi-permeable shell

Small molecules
diffuse through the shell

Aqueous core

Atrandi
BIOSCIENCES

# Semi permeable capsules



Semi-Permeable Shell

30 – 300 μm

Aqueous Core

Cells and nucleic acids are retained

Proteins, oligonucleotides, small molecules are freely exchanged

Atrandi
BIOSCIENCES

# Semi permeable capsules

# Removed unpublished results

# Metagenomics is it a big deal!?

**Maliheh Mehrshad**

**Department of Aquatic Sciences and Assessment, SLU**