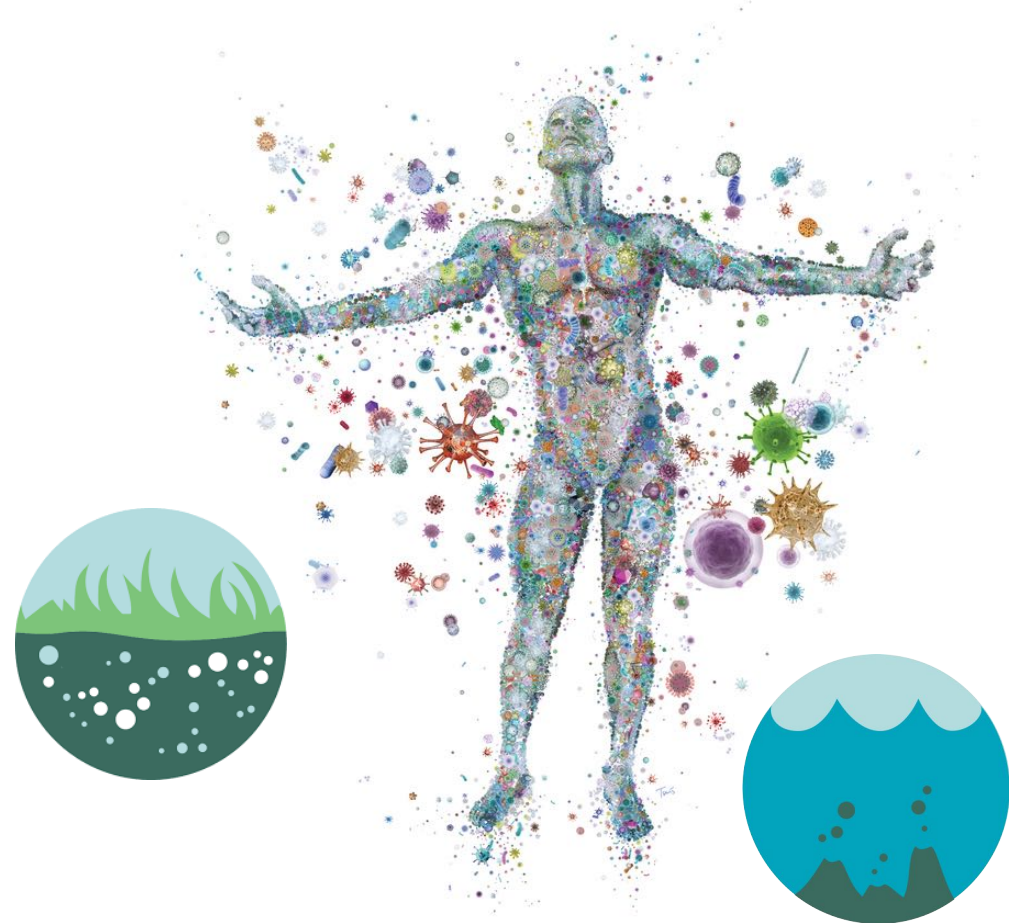National Institute of Allergy and Infectious Diseases

# METAGENOMICS OVERVIEW

## Metataxonomics

MSB7105

March, 2021

NIAID

**Angelina Angelova, PhD**

**Bioinformatics and Computational Biosciences Branch (BCBB)**

**OCICB/OSMO/OD/NIAID/NIH**

National Institute of Allergy and Infectious Diseases

# Today's instructor

**Angelina Angelova**, PhD
Metagenomics Analysis Specialist

Bioinformatics and Computational Biosciences Branch (BCBB)
National Institute of Allergies and Infectious Diseases (NIAID)
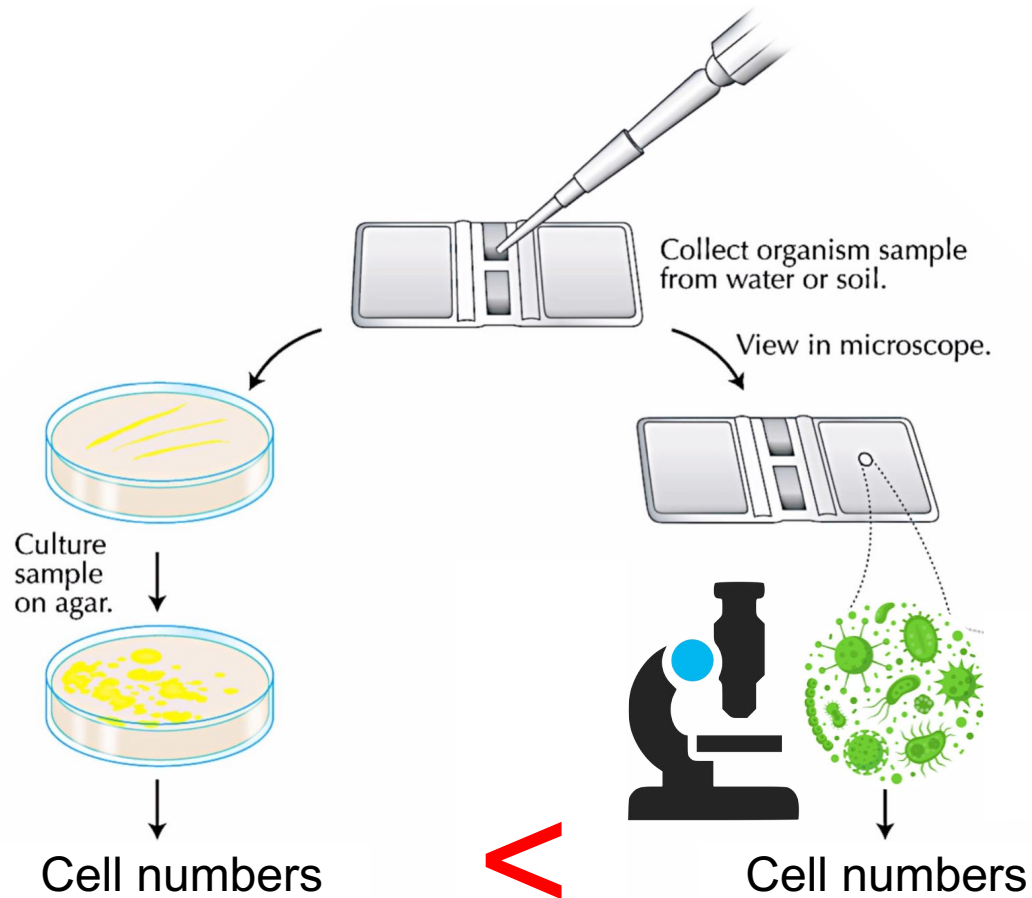National Institute of Health (NIH)
Bethesda, MD, USA

Contact instructor:
angelina.angelova@nih.gov
Contact our team:
bioinformatics@niaid.nih.gov

# The great plate count anomaly


Collect organism sample from water or soil.

View in microscope.

Culture sample on agar.

Cell numbers < Cell numbers

A term coined by Staley & Konopka in 1985 to describe the difference (in orders of magnitude) between the number of cells from natural environments countable by microscopy and those observed after culturing on common agar media.

**"Simple" genomics is not enough**
- The microbial world is **extremely diverse** (construct ~1/2 of Earth's biomass) and **largely unknown**
- Less than 1% of organisms are culturable due not only to lack of proper growth conditions for them in the lab, but also due to proper social interaction
- Microbes exist in complex communities & have complex relations between each other and larger organisms!

# Define "Metagenomics"

- NGS made the field of metagenomics possible
- Metagenomics: Refers to the idea that the collection of genes (the metagenome), obtained directly from a community in its natural habitat (the microbiome), can provide an understanding of the function and characteristics of the whole community, in a similar way as the collection of genes from a single organism can provide an understanding of the function and identity of that organism.
- Metagenomics bypasses the need for isolation or cultivation of individual microbes.
- Allows for exploration of the structure (abundance & identities), interactions, strategies (communication, survival, etc.), functionality and dynamics of a community

Example microbiomes:

Human | Digestive system | Aquatic | Marine

Plants | Soil | Skin | Wastewater
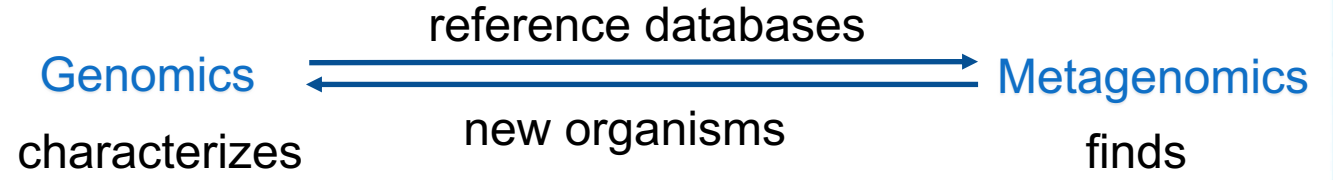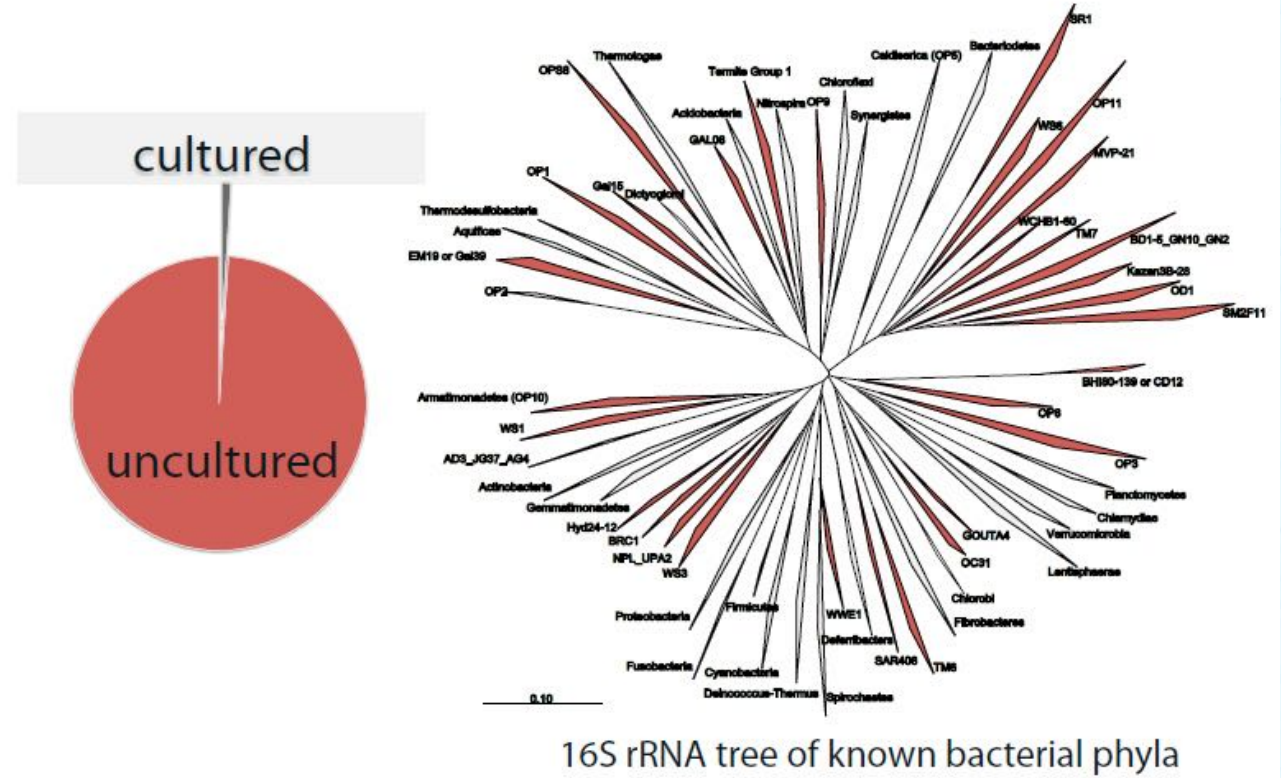
NIH > National Institute of Allergy and Infectious Diseases

# 'Shining a light' on microbial 'dark matter'

- Culture-based techniques are very limited to only what is cultivatable and produce a strong bias towards exploration of only cultivatable organisms, excluding > 99% of microorganisms from exploration
- Metagenomics enables scientists to explore this microbial 'dark matter'
- Vast applications:
  - Biotechnology & Medicine
  - Environmental preservation & recovery
- The more organisms can recognize, the more we expand our capacity to 'see' new species.



Our skewed view of the microbial world

cultured

uncultured

16S rRNA tree of known bacterial phyla

National Institute of Allergy and Infectious Diseases

reference databases

Genomics ← → Metagenomics

new organisms

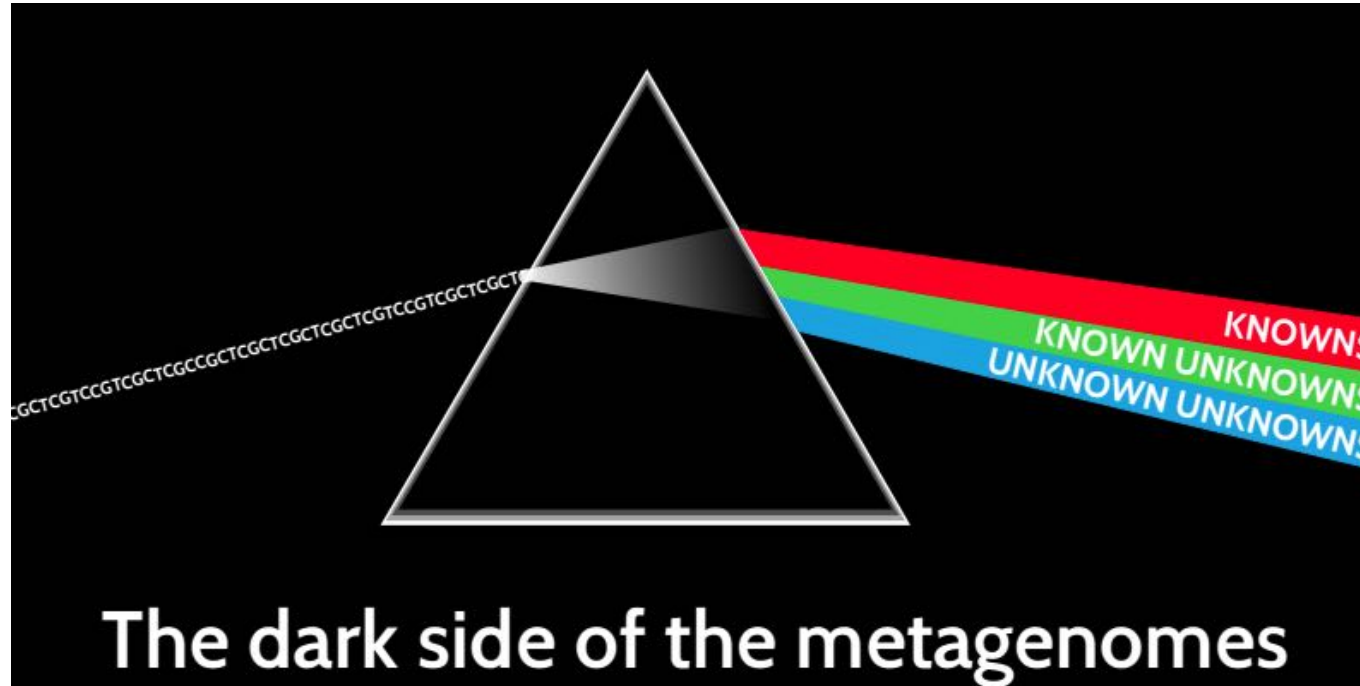characterizes                finds

# Reference genomic databases

A reference genomic databases are a collection of DNA sequences that are idealistic genomic representations of recognized organisms. These sequences are sourced either from individual cultivated organisms (a type strain representing that lineage) or in case of more complex organisms – from multiple organisms from the same species (e.g. human).
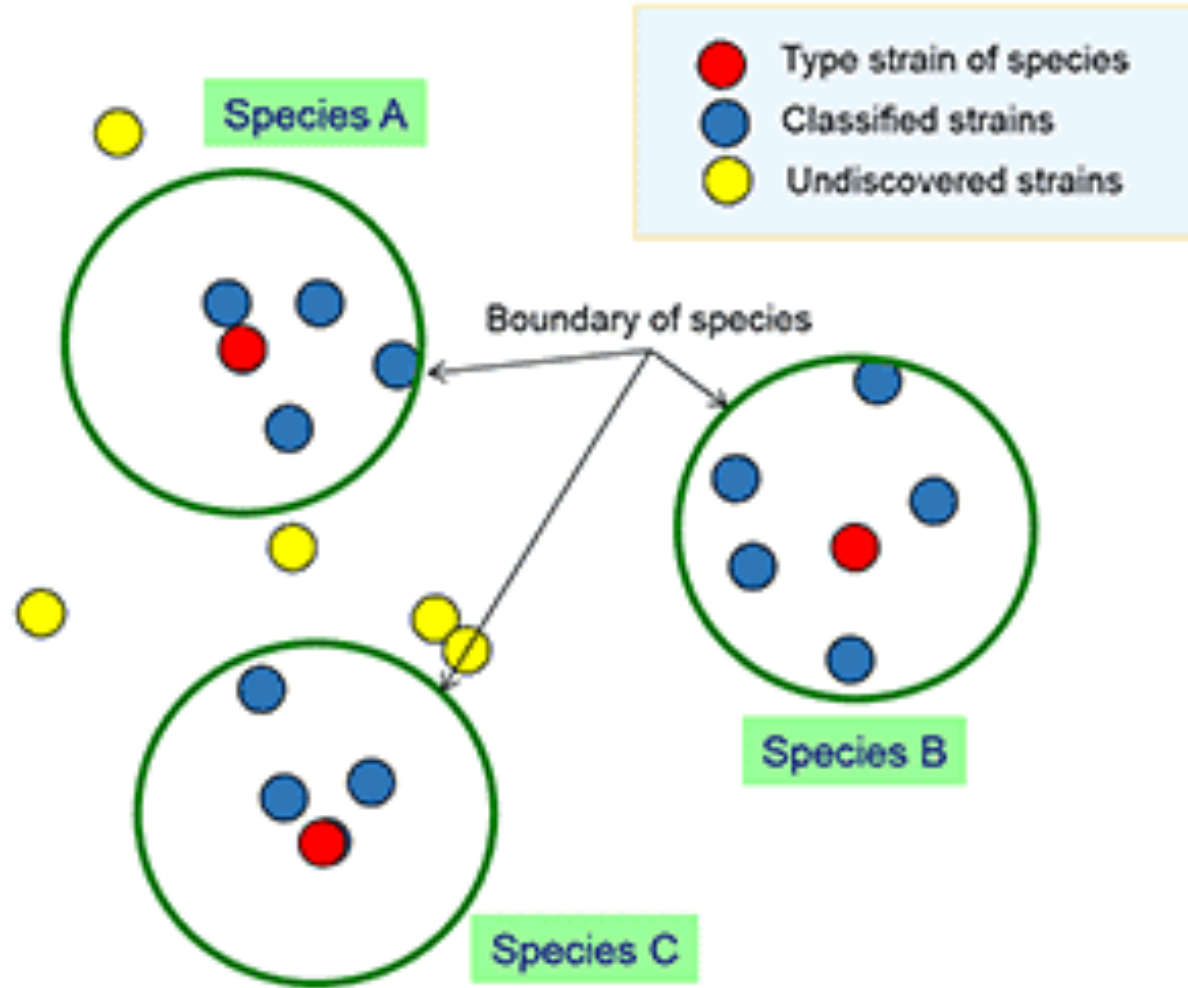


- RefDBs allow for the characterization of
  - a specific species through identification of conserved genes within that organism's genome (genetic markers).
  - a specific function through identification of known genes/proteins observed within the organism

- RefDB are important for
  - Propper phylogenic & functional assignments of unknown sequences
  - Understanding genomic structures, functional capacity & survival strategies of organisms
  - Guiding assembly software & genome mining tools

National Institute of Allergy and Infectious Diseases

# Can you think of any genetic markers commonly used in microbial identification?

Genetic markers are used for identification of organisms or function. These are ideally a single-copy genes with universal presence but internal variability in all organisms



The dark side of the metagenomes

# Can you recall a name of a reference sequence database?



Legend:
- Type strain of species (red)
- Classified strains (blue)
- Undiscovered strains (yellow)

Species A, Species B, Species C

Boundary of species

# Shotgun **VS** Amplicon

## Shotgun Strategy

ALL the DNA from ALL the genomes within the ENTIRE community, is fragmented to the "bite-size" capacity of a sequencing platform. ALL DNA is sequenced. The sequences are used to explore taxonomic composition *and* functional capacity of the entire community

## Long-read Strategy

ALL the DNA from ALL the genomes within the ENTIRE community, is sequenced in "large bites". The sequences are used to explore taxonomic composition *and* functional capacity of the entire community

## Common platforms

For long reads:
- ○ PacBio, Nanopore (MinIon)

## Amplicon Strategy

One gene (a marker gene or a fraction of it) from ALL the genes from within ALL the genomes of ALL the organisms in a community, is targeted for amplification. Its sequence is used to explore the taxonomic composition of the entire community.

## Common marker genes:

- ➤ For Bacterial & Archaeal organisms:
  - ○ 16S rRNA gene
- ➤ For Eukaryotic organisms:
  - ○ 18S rRNA *gene* (less conserved)
  - ○ ITS: internal transcribed spacer region

## Common platforms

For Short reads:
- ○ Illumina HiSeq, NextSeq, NovoSeq

## Common platforms

For Short reads:
- ○ Illumina MiSeq, NextSeq
- ○ TermoFisher IonTorrent



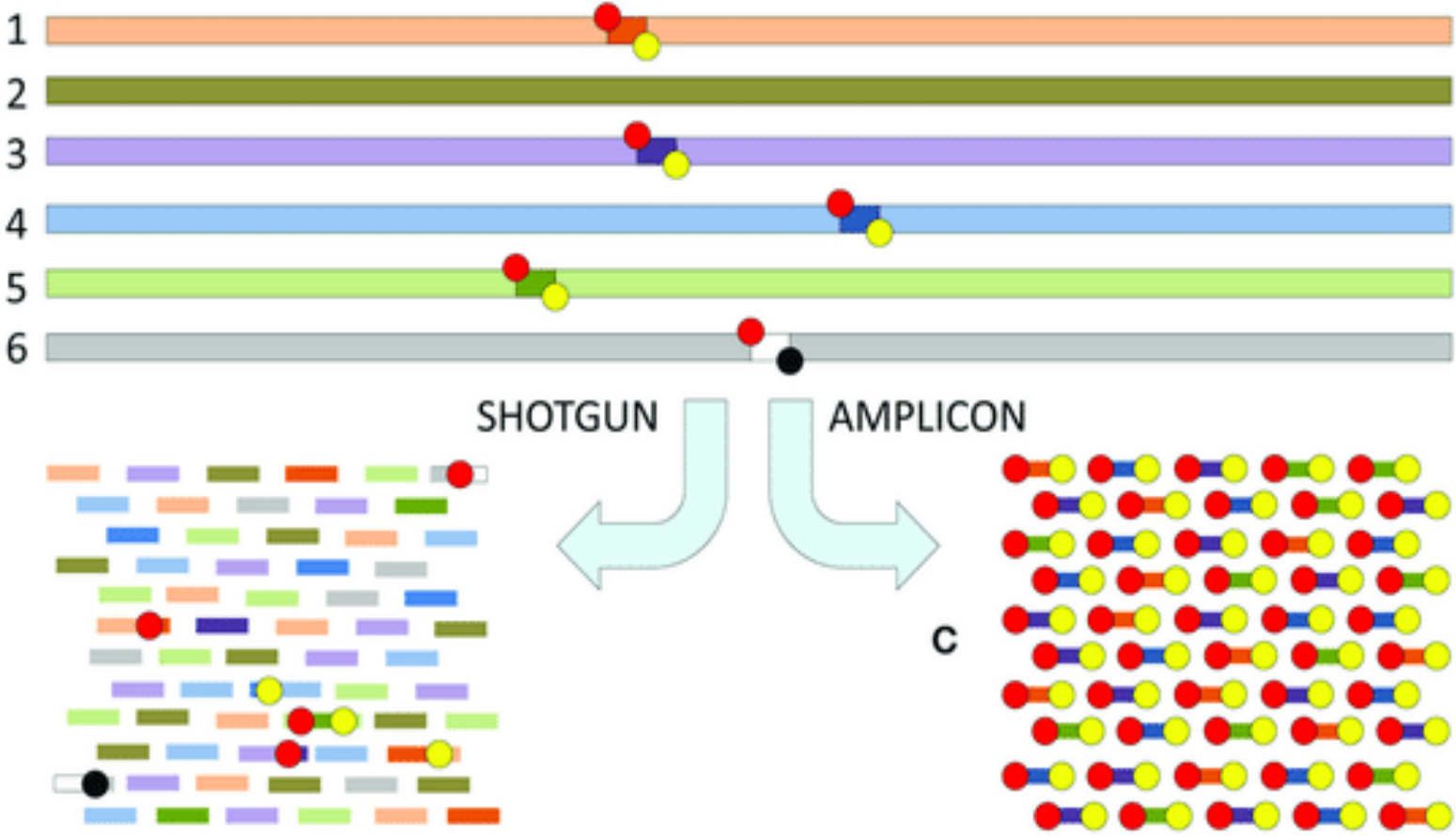National Institute of Allergy and Infectious Diseases

# Why is long-read sequencing not commonly used?

o Because DNA gets commonly shredded during extraction, anyway.

o Because **long-read platforms have a lower resolution for the rarer organisms or organisms with smaller genomes**, so such will often be unrepresented

o Because bacterial genomes are often not "that big" and long-read strategy is often overkill.

o Because long-read strategy still creates sequence errors which inflates the diversity of a community

Long-read sequencing can still be used & be quite useful depending on what one is looking to explore (e.g. micro-eukaryotic communities!)

National Institute of Allergy and Infectious Diseases

# Microbiome exploration strategies



Metagenomics — SHOTGUN — AMPLICON — Metataxonomics

# Amplicon-based community exploration



Metataxonomics

DNA indexing
DNA barcoding
Marker gene sequencing
Amplicon sequencing

Down rabbit hole #1

# 16S rRNA gene characteristics
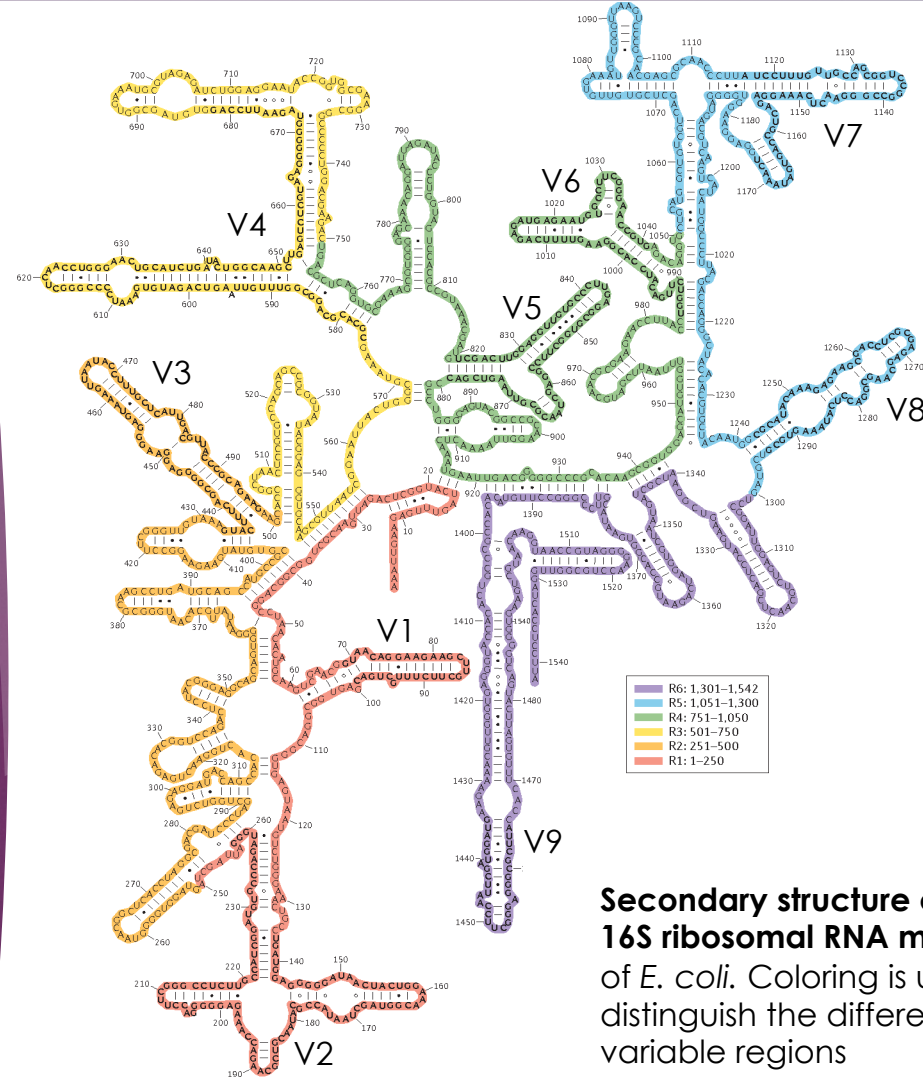
Ubiquitous gene, found in all prokaryotes

Most sequenced gene!

16S rRNA gene-specific primers are not only numerous, but also well characterized.
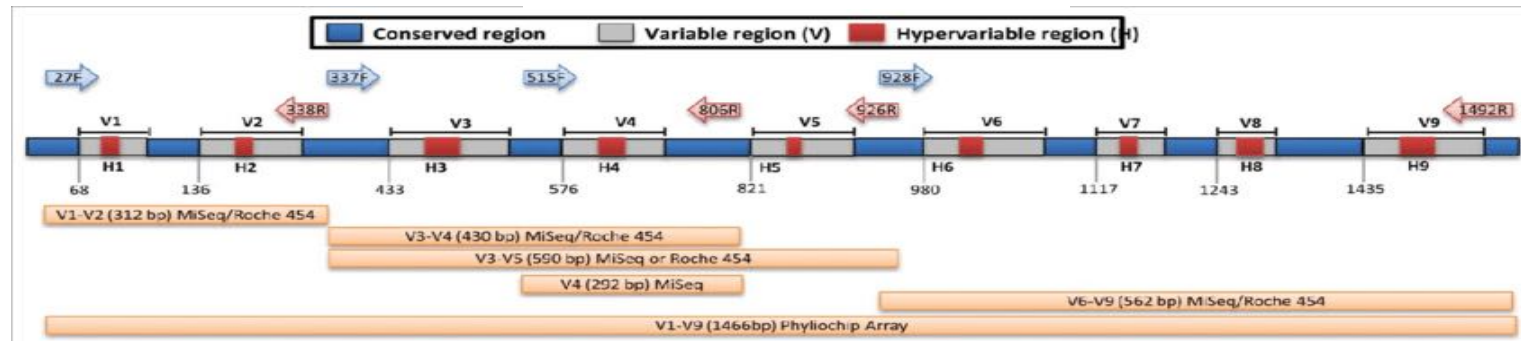
Current gold standard as evolutionary marker

Cost effective approach for exploring a community

Availability of large 16S rRNA gene databases (e.g. RDP, SILVA).



**Length:** ~1,600bp

**Structure:** contains
- highly <u>conserved regions</u>: enabling gene targeting across species; and
- highly <u>variable regions</u>: enabling taxonomic characterization

R6: 1,301–1,542
R5: 1,051–1,300
R4: 751–1,050
R3: 501–750
R2: 251–500
R1: 1–250

**Secondary structure of the 16S ribosomal RNA molecule** of *E. coli*. Coloring is used to distinguish the different variable regions

Int J Syst Bacteriol. 1992. doi: 10.1099/00207713-42-1-166
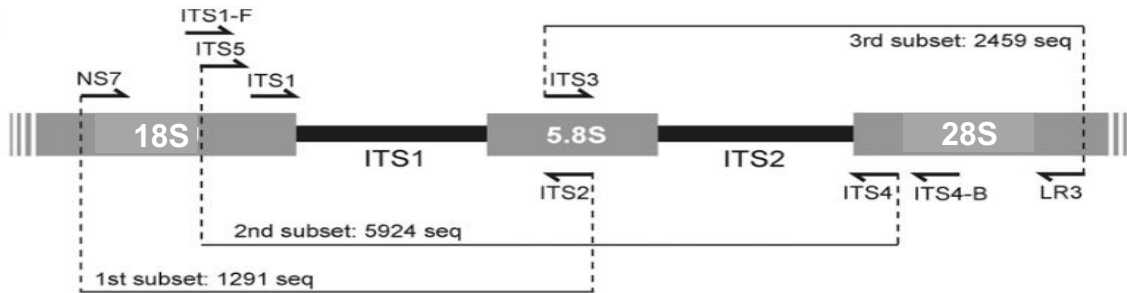Nat Rev Microboil, 2014.  https://doi.org/10.1038/nrmicro3330

# Common universal 16S rRNA Primers



- F27: AGAGTTTGATCCTGGCTCAG
- F343: CCTACGGGNGGCWGCAG
- F515: GTGYCAGCMGCCGCGGTAA
- F968: ACGCGHNRAACCTTACC
- F1177: ACGTCATCCCCACCTTCC

- R1492: CGGNTACCTTGTTACGACTT
- R926: CCGYCAATTCMTTTRAGTTT
- R785: GACTACHVGGGTATCTAATCC
- R534: GTWTTACCGCGGCTGCTGG
- R338: GCTGCCTCCCGTAGGAGT

# Common universal ITS & 18S primers



Fungal equivalent of 16S rRNA marker gene is the **ITS region** (Internal transcribed spacer region)
Unlike 16S or 18S rRNA amplicons, the ITS regions can drastically vary in length and sequence between species, complicating sequence processing

### Common ITS primers

| Primer | Author | Primer sequence | Position |
|--------|--------|-----------------|----------|
| *Forward primers* | | | |
| NS7 | [19] | GAGGCAATAACAGGTCTGTGATGC | 1403-1426 |
| ITS1-F | [18] | CTTGGTCATTTAGAGGAAGTAA | 1723-1744 |
| ITS5 | [19] | GGAAGTAAAAGTCGTAACAAGG | 1737-1758 |
| ITS1 | [19] | TCCGTAGGTGAACCTGCGG | 1761-1779 |
| ITS3 | [19] | GCATCGATGAAGAACGCAGC | 2024-2045 |
| *Reverse primers* | | | |
| ITS2 | [19] | GCTGCGTTCTTCATCGATGC | 2024-2043 |
| ITS4 | [19] | TCCTCCGCTTATTGATATGC | 2390-2409 |
| ITS4-B | [18] | CAGGAGACTTGTACACGGTCCAG | 2526-2548 |
| LR3 | [13] | CCGTGTTTCAAGACGGG | 3029-3045 |

**Figure 1 Commonly used primers for amplifying parts or the entirety of the ITS region.** a) Relative position of the primers, design of the subsets and number of sequences in each subset. b) Primer sequences, references and position of the primer sequence according to a reference sequence of *Serpula himantioides* (AM946630) stretching the entire nrDNA repeat.

## Common 18S rRNA gene primers

➤ F563: TGC CAG CAG CCG CGG TAA TTC C

➤ R1150: CCG TCA ATT CCT TTA AGT TT

➤ F1267: GGT GGT GCA TGG CCG TTC TTA G

➤ R1644: GAC GGG CGG TGT GTA CAA AGG

National Institute of Allergy and Infectious Diseases

# Selecting your primers (TestProbe 3.0)

Mystery primers

➢ Df514     TCC AGC TCC AAT AGC GTA

➢ Dr1069    TCT TTA AGT TTC AGC CTT GC

Primers you select will drastically affect the resolution and visibility of different phylogenies. You can test ANY primer at https://www.arb-silva.de/search/testprobe/ to determine *in silico* a primer's ability to "recognize" specific lineages of organisms, prior to sequencing

# PCR primer scaffolds for high-throughput sequencing of 16S rRNA gene

Primers scaffold:

Adaptors help amplicon for sequencing attach to flow cell

Primer pads help avoid hairpin and primer dimers within the primer scaffolds

Locus-specific sequence is targets & anneals to the region of interest

**FWD scaffold**

5'- AATGATACGGCGACCACCGAGATCTACACGCT XXXXXXXXXXXX TATGGTAATT GT GTGYCAGCMGCCGCGGTAA -3'

| 5' Illumina adapter P5 | Barcode | Primer pad | Primer linker | Locus-specific sequence | 515F primer sequence |

**REV scaffold**

5'- CAAGCAGAAGACGGCATACGAGAT XXXXXXXXXXXX AGTCAGCCAG CC GGACTACNVGGGTWTCAAT -3'

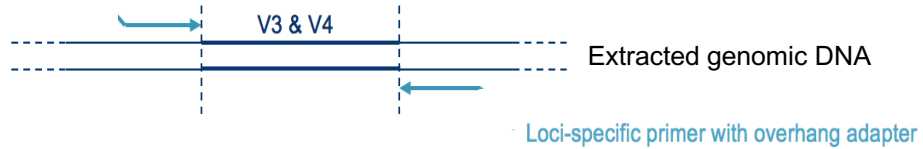| 3' Illumina adapter P7 | Barcode | Primer pad | Primer linker | Locus-specific sequence | 806R primer sequence |

Barcodes are attached to the primers (indexing) to differentiate amplicons from each sample. Afterwards libraries can be been pooled together for sequencing. Distinguishing the produced reads based on the inserted barcodes after sequencing, is called demultiplexing

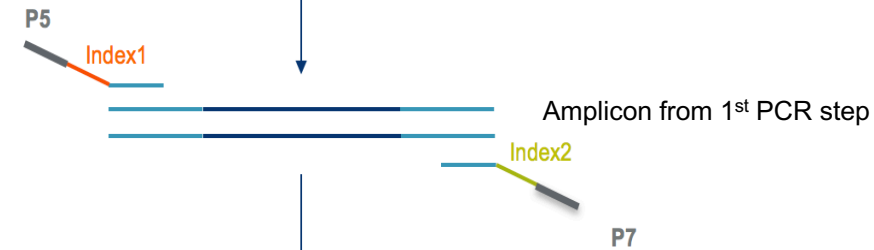Primer linkers help introduce a kink / bend in the molecular shape of the primer & provide spacing between the locus-specific primer & the rest of the scaffold

National Institute of Allergy and Infectious Diseases

# Library preparation for high-throughput sequencing of 16S rRNA gene



https://www.idtdna.com/pages/education/decoded/article/16s-rrna-indexed-primers-amplify-phylogenic-markers-for-microbiome-sequencing-analysis

# Problems with PCR: inherent biases

- Inherent biases:
  - PCR primer bias: bias introduced from the different primer annealing efficiency to different templates /phylogenic groups
  - PCR efficiency bias: bias introduced from the different PCR rates among different templates/ phylogenic groups
- Resulting in inaccurate representation of template distributions
- PCR bias cannot be avoided in amplification-based studies, but it can (and should!) be:
  - constrained with molecular biology techniques (e.g. 2 step PCR, degenerate bases); or
  - assessed in the downstream analysis, acknowledged and presented in each study!



ZymoBIOMICS™ Microbial Community Standards

Legend:
- g:Bacillus
- g:Enterococcus
- g:Escherichia
- g:Lactobacillus
- g:Listeria
- g:Pseudomonas
- g:Staphylococcus
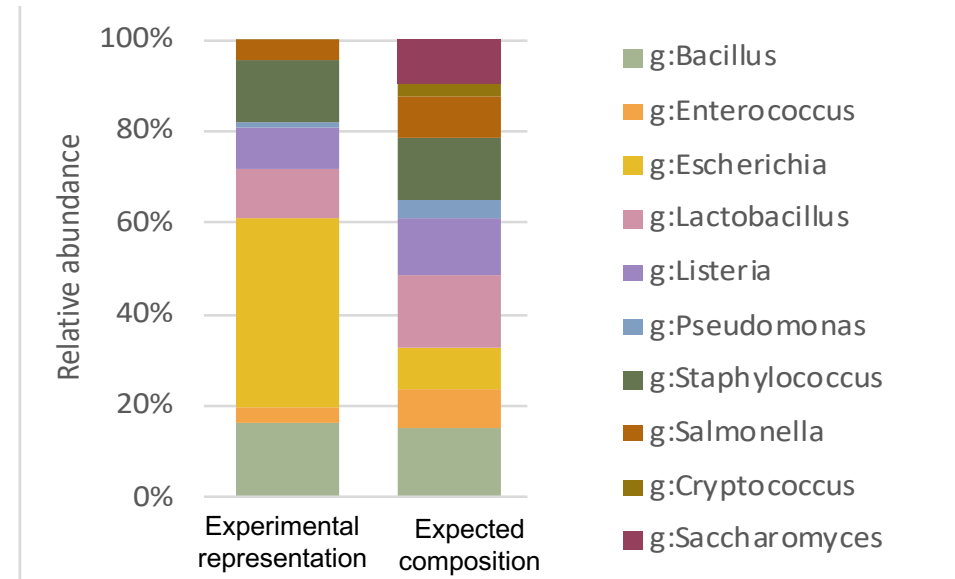- g:Salmonella
- g:Cryptococcus
- g:Saccharomyces

**It is pivotal to include microbial community standards to your experimental samples to assess PCR bias**

e.g. ATCC microbiome standards, ZymoBiomics Microbial Community standards

National Institute of Allergy and Infectious Diseases

Acinas *et al*. 2005. Appl Environ Microbiol 7(12). doi: 10.1128/AEM.71.12.8966-8969.2005

# Problems with PCR: Chimera formation

- Chimeras:
  - Artifacts of PCR amplification, sequencing or read merging
  - Hybrid products between multiple parent sequences
- Inflate apparent community organismal diversity, by suggesting presence of non-existent organisms
- Chimeric formation is more common between more closely related organisms

PCR



**Figure 1.** Formation of chimeric sequences during PCR. An aborted extension product from an earlier cycle of PCR can function as a primer in a subsequent PCR cycle. If this aborted extension product anneals to and primes DNA synthesis from an improper template, a chimeric molecule is formed.



**Figure 4.** Alignment of sequences corresponding to chimeras between *Streptococcus* and *Staphylococcus* 16S rRNA genes. Only columns from the NAST multiple alignment containing nonidentical nucleotides between the reference sequences (*top* and *bottom*) are shown. Nucleotides matching *Streptococcus* sequences are colored red. Sequence prefixes correspond to the four experimental replicates A–D.

Computational tools are used to identify and remove chimeric sequences

# Raw sequences: FastQ file structure

Output files provided from sequencer (normally): Fastq files in archived format
Sample1…_**R1**_...fastq.gz (Forward Reads)
Sample1…_**R2**_...fastq.gz (Reverse Reads)
Sample1…_**R0**_...fastq.gz (Undetermined)

What you will mostly need

A text-based format for both sequence and associated *phred quality scores*, developed by Sanger Institute

| | |
|---|---|
| Identifier | @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50 |
| Sequence | TTGCCTGCCTATCATTTTAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT |
| '+' sign | + |
| Quality scores | hhhhhhhhhghhghhhhhfhhhhhfffffe'ee['X]b[d[ed'[Y[^Y |
| Identifier | @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50 |
| Sequence | GATTTGTATGAAAGTATACAACTAAAACTGCAGGTGGATCAGAGTAAGTC |
| '+' sign | + |
| Quality scores | hhhhgfhhcghghggfcffdhfehhhhcehdchhdhahehffffde'bVd |

Data lines

Description lines

National Institute of Allergy and Infectious Diseases

# General Bioinformatics workflow
## (metataxonomics)



**PRIMARY**

**SECONDARY**

**TERTIARY**

Library Prep → Sequencer → De-multiplex → QC Check (No / Yes)

F & R raw reads → QC check, filter & trim (No / Yes) → F & R read merging + Dereplication + Clustering or Denoising

Abundance scoring → Taxonomic annotation → Abundance & tax table → Community analysis

Reference DBs
NCBI 16S Microbial
SILVA, RDP, GG, Custom

# Processing strategies

- Cluster-based strategy: OTU formation
  - OTU: Operational taxonomic unit
  - A sequence (usually the longest) representing a cluster of similar sequences
  - recently depreciated for carrying over too many errors
- Denoising strategy: ASV formation
  - ASV: Amplicon Sequence Variants
  - A sequence (error-corrected) representing an amplicon variant
  - common practice / more correct representation

# Clustering (OTU formation)

**OTUs = Operational Taxonomic Units**
**Cluster representative sequence**

- Created by clustering of all sequences, based on a fixed similarity threshold (97% seq similarity). A **representative read** is chosen (usually the longest one within the cluster) to represent the whole cluster and its abundance is determined by the number of reads in its cluster.

- This is a depreciated method because: Lingering sequencing errors artificially inflate community diversity

- More stringent downstream filtering is required, loosing information about rarer species

# Denoising (ASVs)

**ASVs = Amplicon Sequence Variant**

- Using nucleotide identity and quality, an error model is created individually for F & R reads, to identify artificial variations. These models are then used to correct the sequence errors, prior to read merging and abundance scoring.

- The amplicon sequence variants (ASVs), provide a more accurate and fine-scale resolution into the *real* diversity of the amplicons, than any clustering algorithm.

- The error model may vary from one sequencing run to another (batch effect)

- Algorithm is reference-free and works on any genetic locus, highly similar sequences (amplicons of the same locus).

# Denoising (ASVs) VS Clustering (OTUs)

```
s: ATTAACGAGATTATAACCAGAGTACGAATA...
       |                 |
r: ATCAACGAGATTATAACAAGAGTACGAATA...
```

$$p(r|s) = \prod_{i=1}^{L} p(r(i)|s(i), q_r(i), Z)$$

**Error rates depend on....**

- Substitution (eg. A->C)
- Quality score (eg. Q=30)
- Batch effect (eg. run)          *Using more data!*

In a good model, the observed error rates (black line) will decrease with increase of quality score (the x axis), keeping the trend of the expected error rates (red line)

### Denoising error model

# Bioinformatics DADA2 workflow

# Denoising (ASVs) VS Clustering (OTUs)



- Denoising algorithms
  - DADA2 (DADA2,R, QIIME2)
  - UNOISE3 (USEARCH, VSEARCH)
  - Deblur (Deblur, QIIME)
- Clustering algorithms
  - UCLUST (QIIME)
  - UPARSE (USEARCH, VSEARCH)
  - Mothur (mothur)

Divisive Amplicon Denoising Algorithm

# Databases for taxonomic assignment

## 16S / Bacterial & Archaeal Databases

SILVA

▶ 2,225,272 full 16S & 18S rRNA gene sequences + guide tree

▶ Latest release: v138.1 from August 2020

RDP

▶ 3,356,809 full 16S rRNA gene sequences + 125,525 fungal 28S rRNA gene sequences

▶ latest release: v18 from August 2020

NCBI's 16S Microbial

▶ 20,845 full 16S rRNA gene sequences from type strains

▶ Latest release: Oct 202 (regularly updated)

GreenGenes

▶ Outdated, latest release 2013

# Databases for taxonomic assignment

## ITS / Fungal Databases

### UNITE Community

► 35,077 ITS gene region sequences

► Latest release: v8.2 from Feb 2020

### GlobalFungi

► 145,873,740 ITS sequence variants

► Latest release: v0.9.8 from Jan 2020

### FungiDB

► 6,632 ITS gene regions

► Latest release: 50 beta from Dec 2020

### R-Syst DBs

► A collection of custom databases specific for different phylogenies across kingdoms

## 18S protozoan databases

### SILVA

► 2,225,272 full 16S & 18S rRNA gene sequences + guide tree

► Latest release: v138.1 from August 2020

### PR2 database

► 184,000 18S rRNA gene region sequences

► Latest release: v14.12.0, from August 2019

► manually curated & metadata available

### PhytoRef

► Uses *plastidal* 16S rRNA gene to identify photosynthetic micro**eukaryotes**

► 6,490 plastidial 16S rDNA reference sequences

► Latest release: 2015

PR²

FungiDB
Fungal and Oomycetes Genomics Resources
http://FungiDB.org

R-SYST

National Institute of Allergy and Infectious Diseases

silva
high quality ribosomal RNA databases

PhytoRef
A reference database of the plastidial 16S rRNA gene of photosynthetic eukaryotes with curated taxonomy

NIAID

# Structure of ASV abundance table

Samples | Taxonomic Assignments

ASVs

Raw sequence counts (reads assigned per ASV in each sample)

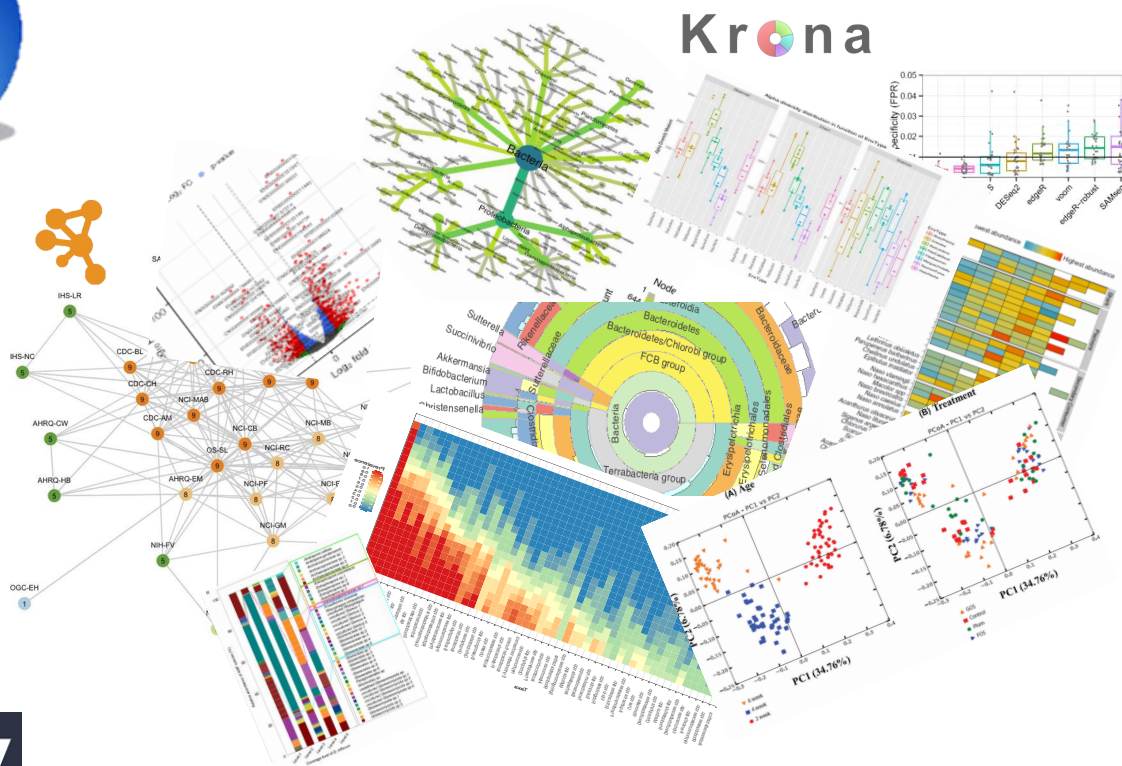| ASVID | MNAW0514 | MNAW0515 | NAW0514 | NAW0515 | NSAIW0514 | NSAIW0515 | NSDW0514 | NSDW0515 | Kingdom | Phylum | Class | Order | Family | Genus | Species |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ASV1 | 0 | 0 | 49 | 55 | 786 | 913 | 620 | 1071 | Bacteria | Proteobacteria | Gammaproteobacteria | Oceanospirillales | Alcanivoracaceae | Alcanivorax | borkumensis |
| ASV10 | 270 | 457 | 129 | 430 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | SAR11 | Pelagibacter | NA | NA |
| ASV100 | 13 | 0 | 59 | 0 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | Alteromonadales | Pseudoalteromonadaceae | Pseudoalteromonas | denitrificans |
| ASV101 | 0 | 0 | 28 | 87 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | SAR11 | Pelagibacter | NA | NA |
| ASV102 | 15 | 0 | 12 | 41 | 0 | 0 | 0 | 0 | Archaea | Euryarchaeota | Thermoplasmata | Methanomassiliicoccales | Methanomassiliicoccaceae | Methanomassiliicoccus | NA |
| ASV103 | 0 | 0 | 71 | 0 | 0 | 0 | 0 | 0 | Bacteria | Bacteroidetes | Flavobacteriia | Flavobacteriales | Flavobacteriaceae | NA | NA |
| ASV104 | 0 | 0 | 31 | 0 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | Rhodospirillales | Rhodospirillaceae | Nisaea | NA |
| ASV105 | 0 | 0 | 30 | 0 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | SAR11 | Pelagibacter | NA | NA |
| ASV106 | 22 | 0 | 0 | 0 | 0 | 0 | 0 | 24 | Bacteria | Proteobacteria | Betaproteobacteria | Burkholderiales | Comamonadaceae | Acidovorax | temperans |
| ASV107 | 130 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | Sphingomonadales | Sphingomonadaceae | Sphingobium | xenophagum |
| ASV108 | 0 | 0 | 30 | 29 | 0 | 0 | 0 | 0 | Bacteria | Actinobacteria | Actinobacteria | Acidimicrobiales | NA | NA | NA |
| ASV109 | 0 | 0 | 0 | 0 | 46 | 0 | 69 | 18 | Bacteria | Proteobacteria | Gammaproteobacteria | Oceanospirillales | Oceanospirillaceae | Oleispira | lenta |
| ASV11 | 0 | 0 | 0 | 0 | 224 | 312 | 235 | 170 | Bacteria | Marinimicrobia | NA | NA | NA | NA | NA |
| ASV110 | 0 | 0 | 0 | 0 | 0 | 0 | 69 | 30 | Bacteria | Proteobacteria | Betaproteobacteria | Burkholderiales | Comamonadaceae | Acidovorax | temperans |
| ASV111 | 0 | 0 | 72 | 29 | 0 | 0 | 0 | 0 | Archaea | Euryarchaeota | Thermoplasmata | NA | NA | NA | NA |
| ASV112 | 53 | 0 | 0 | 23 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | NA | NA | NA | NA |
| ASV113 | 89 | 19 | 15 | 0 | 0 | 0 | 0 | 0 | Bacteria | Bacteroidetes | Cytophagia | Cytophagales | Cytophagaceae | Pseudarcicella | NA |
| ASV114 | 0 | 0 | 0 | 95 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | Alteromonadales | Pseudoalteromonadaceae | Pseudoalteromonas | denitrificans |
| ASV115 | 0 | 0 | 0 | 85 | 0 | 0 | 0 | 0 | Bacteria | Cyanobacteria | Chloroplast | Chloroplast | Bacillariophyta | NA | NA |
| ASV116 | 25 | 0 | 29 | 41 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | NA | NA | NA | NA |
| ASV117 | 0 | 0 | 0 | 0 | 26 | 44 | 34 | 29 | Bacteria | Proteobacteria | Alphaproteobacteria | Rhodospirillales | Rhodospirillaceae | Magnetospira | NA |
| ASV118 | 92 | 0 | 36 | 0 | 0 | 0 | 0 | 0 | Bacteria | Bacteroidetes | Flavobacteriia | Flavobacteriales | Flavobacteriaceae | NA | NA |
| ASV119 | 0 | 0 | 86 | 0 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | NA | NA | NA | NA |
| ASV12 | 0 | 0 | 0 | 0 | 159 | 154 | 79 | 323 | Archaea | Euryarchaeota | Thermoplasmata | Methanomassiliicoccales | Methanomassiliicoccaceae | Methanomassilii | coccus |
| ASV120 | 0 | 0 | 0 | 38 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | Sphingomonadales | Sphingomonadaceae | Sphingobium | xenophagum |
| ASV121 | 0 | 0 | 0 | 0 | 30 | 0 | 34 | 37 | Bacteria | Proteobacteria | Gammaproteobacteria | Oceanospirillales | Oceanospirillaceae | Oleispira | lenta |
| ASV122 | 0 | 0 | 0 | 54 | 0 | 0 | 0 | 0 | Archaea | Euryarchaeota | Thermoplasmata | Methanomassiliicoccales | Methanoma |  |  |
| ASV123 | 0 | 0 | 0 | 0 | 53 | 26 | 36 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | Alteromonadales | Pseudoalter |  |  |
| ASV124 | 40 | 0 | 59 | 0 | 0 | 0 | 0 | 0 | Bacteria | Bacteroidetes | Flavobacteriia | Flavobacteriales | NA |  |  |
| ASV125 | 0 | 0 | 0 | 0 | 64 | 0 | 0 | 93 | Bacteria | Proteobacteria | Alphaproteobacteria | SAR11 | Pelagibacter |  |  |
| ASV126 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | SAR11 | Pelagibacter |  |  |
| ASV127 | 0 | 0 | 0 | 0 | 0 | 42 | 64 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | Oceanospirillales | Alcanivorac |  |  |
| ASV128 | 0 | 0 | 0 | 0 | 33 | 0 | 29 | 44 | Bacteria | Bacteroidetes | Bacteroidia | Bacteroidales | Prevotellace |  |  |
| ASV129 | 52 | 0 | 0 | 0 | 0 | 21 | 0 | 0 | Bacteria | Proteobacteria | Gammaproteobacteria | Oceanospirillales | Oceanospiril |  |  |
| ASV13 | 0 | 0 | 0 | 0 | 136 | 300 | 43 | 107 | Bacteria | Actinobacteria | Actinobacteria | Acidimicrobiales | Acidimicrob |  |  |
| ASV130 | 0 | 0 | 0 | 0 | 50 | 64 | 0 | 0 | Bacteria | Bacteroidetes | Bacteroidia | Bacteroidales | Prevotellace |  |  |
| ASV131 | 0 | 0 | 0 | 0 | 0 | 96 | 0 | 0 | Bacteria | Proteobacteria | Alphaproteobacteria | Caulobacterales | Hyphomona |  |  |

ASV sequences

```
>ASV1
GGAATATTGCACAATGGAGGAAACTCTGATGCAGCAATGTCGCGTGAGTGAAGAAGGCCCT
ATGATGACGGTACCCCAAGAATAAGCACCGGCTAACTATGTGCCAGCAGCCGCGGTAATACATAGGGTGCGAGCGTTGTTCGGAATTACT
GGGCGTAAAGGGCGCGCAGGCGGAATAGTAAGTCGGAGGTGAAAGCCCGGGGCTCAACCCCGGAGGGTCTTTCGAAACTACTAATCTAGA
GAGGGTCAGGGGCCGGCAGAATTCCTGGTGTAGAGGTGAAATTCGTAGATATCAGGAGGAATACCGGTGGCGAAGGCGGCCGGCTGGGGC
CACTCTGACGCTGAGGCGCGAAAGCGTGGGGAGCAAACAG
>ASV2
GGAATATTGCACAATGGGGGAAACCCTGATGCAGCCATGCCGCGTGTGTGAAGAAGGCCTTCGGGTTGTAAAGCACTTTCAGTTGTGAGG
AAAAGTTAGTAGTTAATACCTGCTAGCCGTGACGTTAACAACAGAAGAAGCACCGGCTAACTCCGTGCCAGCAGCCGCGGTAATACGGAG
GGTGCGAGCGTTAATCGGAATTACTGGGCGTAAAGCGCACGCAGGCGGTTTGTTAAGTCAGATGTGAAAGCCCCGGGCTCAACCTGGGAC
GGTCATTTAGAACTGGCAGACTAGAGTCTTGGAGAGGGGAGTGGAATTCCAGGTGTAGCGGTGAAATGCGTAGATATCTGGAGGAACATC
AGTGGCGAAGGCGACTCCCTGGCCAAAGACTGACGCTCATGTGCGAAAGTGTGGGTAGCGAACAG
>ASV3
GGAATATTGCACAATGGAGGAAACTCTGATGCAGCAATGTCGCGTGAGTGAAGAAGGCCCTTGGGTCGTAAAGCTCTTTTATGGGGGAAG
ATGATGACGGTACCCCAAGAATAAGCACCGGCTAACTATGTGCCAGCAGCCGCGGTAATACATAGGGTGCGAGCGTTGTTCGGAATTACT
GGGCGTAAAGGGCGCGCAGGCGGAATAGTAAGTCGGAGGTGAAAGCCCGGGGCTCAACCCCGGAGGGTCTTTCGAAACTGCTAATCTAGA
CACTCTGACGCTGAGGCGCGAAAGCGTGGGGAGCAAACAG
>ASV4
GGAATATTGGACAATGGGGGCAACCCTGATCCAGCAATACCGCGTGTGTGAAGAAGGCCTTGGGGTTGTAAAGCACTTTCAATTGTGAAG
AAAAGTTAACGGTTAATAACCGTTAGCCTTGACGTTAACTTTAGAAGAAGCACCGGCTAACTCCGTGCCAGCAGCCGCGGTAATACGGAG
GGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGTGCGTAGGCGGTTTATTAAGTCAGATGTGAAAGCCCCGGGCTTAACCTGGGAA
CTGCATTTGAAACTGGTCAACTAGAGTATGGTAGAGGGAAAGTGGAATTTCTGGTGTAGCGGTGAAATGCGTAGATATCAGAAGGAACATC
AATGGCGAAGGCAGCTTTCTGGACCAATACTGACGCTGAGGTACGAAAGCGTGGGGAGCAAACAG
>ASV5
GGAATCTTGGACGAGGAAGGCCCAGCCGTAGTCAGCCATGCCGCGTGAGTGATGAAGGCCTTAGGGTCGTAAAGCTCTTTCGCCAGAGTG
ATAATGACGATATCTGGTAAAGAACCCCGGCTAACTCCGTGCCAGCAGCCGCGGTAATACGGAGGGGGGTTAGCGTTGTTCGGAATTACT
GGGCGTAAAGCGTACGTAGGCGGATTAGTCAGTCAGAGGTGAAATCCCAGGGCTCAACCCTGGAACTGCCTTTGATACTGCTAGTCTTGA
GTTCGAGAGAGGTGAGTGGAATTCCAAGTGTAGAGGTGAAATTCGTAGATATTTGGAGGAACACCAGTGGCGAAGGCGGCTCACTGGCTC
GATACTGACGCTGAGGTACGAAAGTGTGGGGAGCAAACAG
```

OTUs: Operational Taxonomic Units: created via clustering of reads (old)
ASVs: Amplicon Sequence Variants: created via denoising of reads (new)

# Community Analysis

# Questions addressable by metataxonomics

- What organisms are present in each microbiome and in what proportion? (community structure)
- What is the natural variation of organisms within each microbiome (diversity of organisms)?
- How different is community 1 from community 2 in its composition?
- Which organisms are different and which the same between 2 microbiomes?
- Which microbiome is our organism of interest, more abundant in?
- What is the natural non-variable faction of organisms the microbiome (core microbiome)?
- How is the core microbiome different in community 1 vs community 2?
- How does the diversity of community change depending on factors (e.g. treatment, time)?
- Which organisms responds to factors 1, factor 2, etc.?

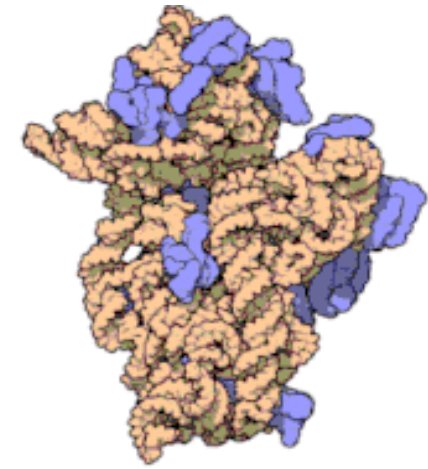Questions **NOT addressable** by 16S amplicon sequencing
- What is the functional capacity of the microbiome as a whole or for each individual organism?
- What is the finer scale (strains-level) and larger scale (cross-kingdom) diversity of a microbiome?
- How does a community adjust to factor 1 in terms of its functionally?

National Institute of Allergy and Infectious Diseases

# Limitations of 16S rRNA gene

The 16S rRNA gene has become the most sequenced taxonomic marker and the cornerstone for current systematic classification of bacteria and archaea.

**Restrictions, caveats, limitations**:

➢ No single variable region captures the variability of the full gene
➢ Different variable regions have different capacity to differentiate taxonomies
➢ Due to inherent biases of this method, abundance estimates can be askew
➢ 16S rRNA gene copy number issues can distort true abundances
  ➢ Gene number can vary from species to species, creating distorted profiles
  ➢ Gene sequence can vary between copies even within the same organism
➢ Inferring *true* phylogenic relations from a single gene can be risky!
  ➢ Even full length 16S gene cannot absolutely resolve the diversification of closely related organisms (species or strains)
  ➢ Some evolutionarily distant organisms have similar 16S rRNA gene copies, which can cluster closely in a phylogenic tree
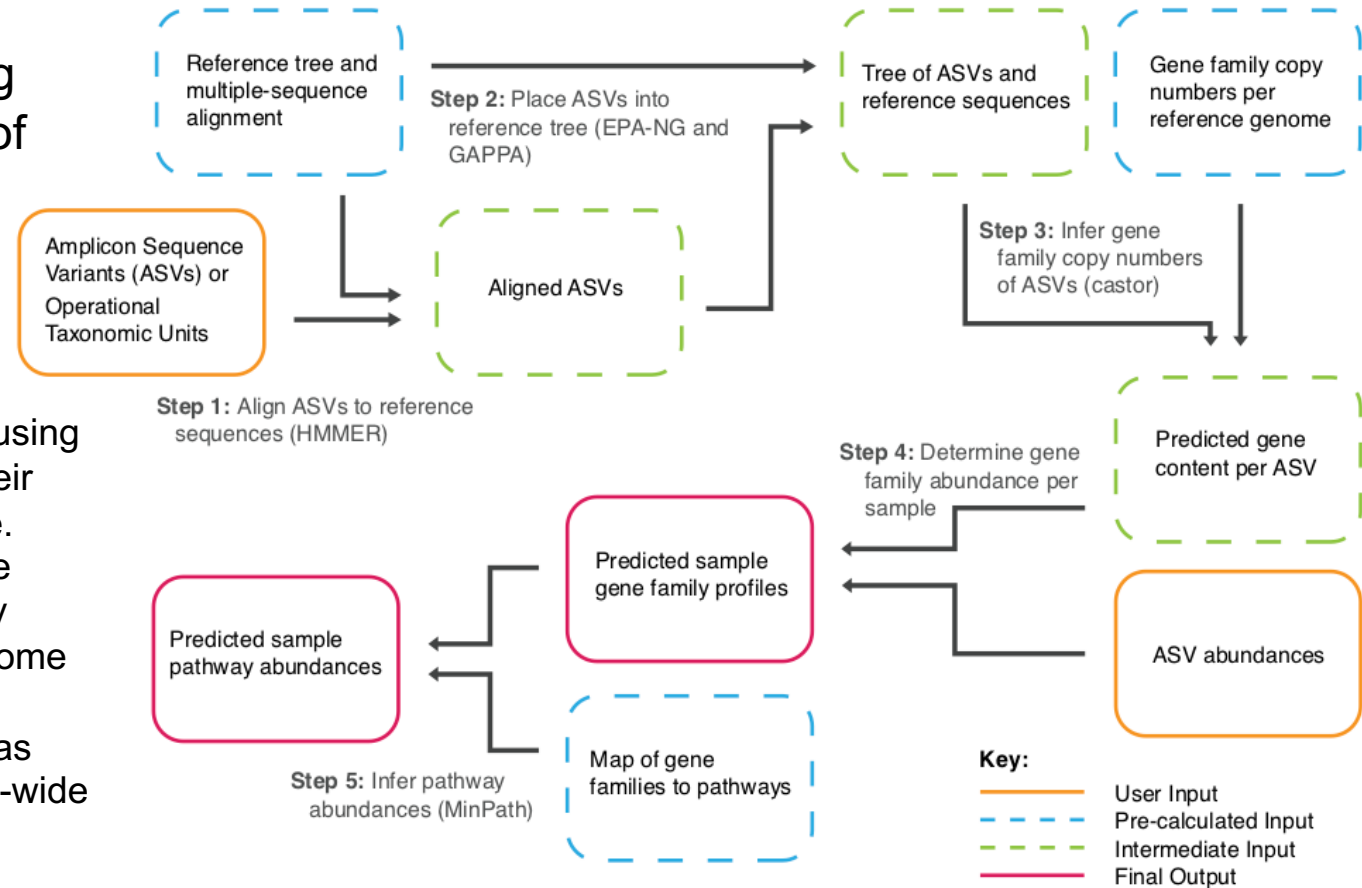
# PICRUSt2

**Phylogenic Investigation of Communities by Reconstruction of Unobserved States … 2**

An *in silico* approach to **predicting** the functional **potential** of metagenome using marker gene data (16S) and a database of annotated reference genome.

- The genetic content of each sample is *estimated* using the taxonomic identities of the input ASVs, and their most closely related genomes from IMG database.
- Gene abundances are approximated based on the ASV abundances and are corrected for gene copy numbers, again estimates from the reference genome for each taxa.
- The genetic content and abundances (presented as KO & EC numbers) are translated into community-wide pathways abundances, derived from the MetaCyc database for metabolic pathways.

# Data resource explosion

NGS has produced data explosion, causing new opportunities for exploration, but also challenges in the scalability and compatibility of data analysis.

- Opportunities: obtaining insight into the microbial "dark matter" problem without culturing

- Challenges: Vast quantities of available data pose problems in its analytical reproducibility, compatibility and comparability



Credit: Amanda Montañez; Source: "Phylogenetically Novel Uncultured Microbial Cells Dominate Earth Microbiomes," by Karen G. Lloyd et al., in *mSystems*, Vol. 3, No. 5; September/October 2018

# Data availability & quality control

- The wide variation in biological questions, coupled with technical variability during sample processing and bioinformatic tools and pipelines for analysis, makes it impossible create a single best protocol for all studies.

- Huge variations are observed between the results of analyses from studies using various DNA extraction methods, 16S primer selection, or bioinformatics pipelines, even when utilized on the same samples.

- the Microbiome Quality Control Project (MBQC) attempts to evaluate and standardize technologies and computational methods for assessing (at least the) human-associated microbial communities.

National Institute of
Allergy and
Infectious Diseases

# Metadata is just as important as the data itself!

- Metadata considered **_critical_** *to data interpretation & reproducibility.*
- Needs to be recorded & provided as accurate and concise as possible

∴        **Community Driven Metadata Standards** are being implemented (e.g. NCBI BioSample db)!

…to promote international standardization of (meta)genome quality **and accompanying metadata** (e.g. vocabulary/ontology, informational fields)

…to promote data discoverability, comparability and reproducibility within and across studies.

Checklists of **Minimum Information about any (x) Sequence (MIxS)** available to implement informational requirements (required metadata) for different types of studies (e.g. host-associated vs soil vs water-associated samples)

- Checklists for metaxtaxonomic (marker gene) studies: **MIMARKS** (**Minimal Information about a Marker Sequence)**
- Checklists for (meta)genomic studies: MIGS & MIMS (**Minimal Information about a (Meta)Genomic Sequence)**

These and other standardization checklists available at: https://gensc.org/mixs/

NIH National Institute of Allergy and Infectious Diseases

NIAID

MIxS

# Thank you

[bioinformatics@niaid.nih.gov](mailto:bioinformatics@niaid.nih.gov)

Bioinformatics and Computational Biosciences Branch

Office of Cyber Infrastructure and Computational Biology