
3000788 Intro to Comp Molec Biol

Lecture 8: Metagenomics

Fall 2025



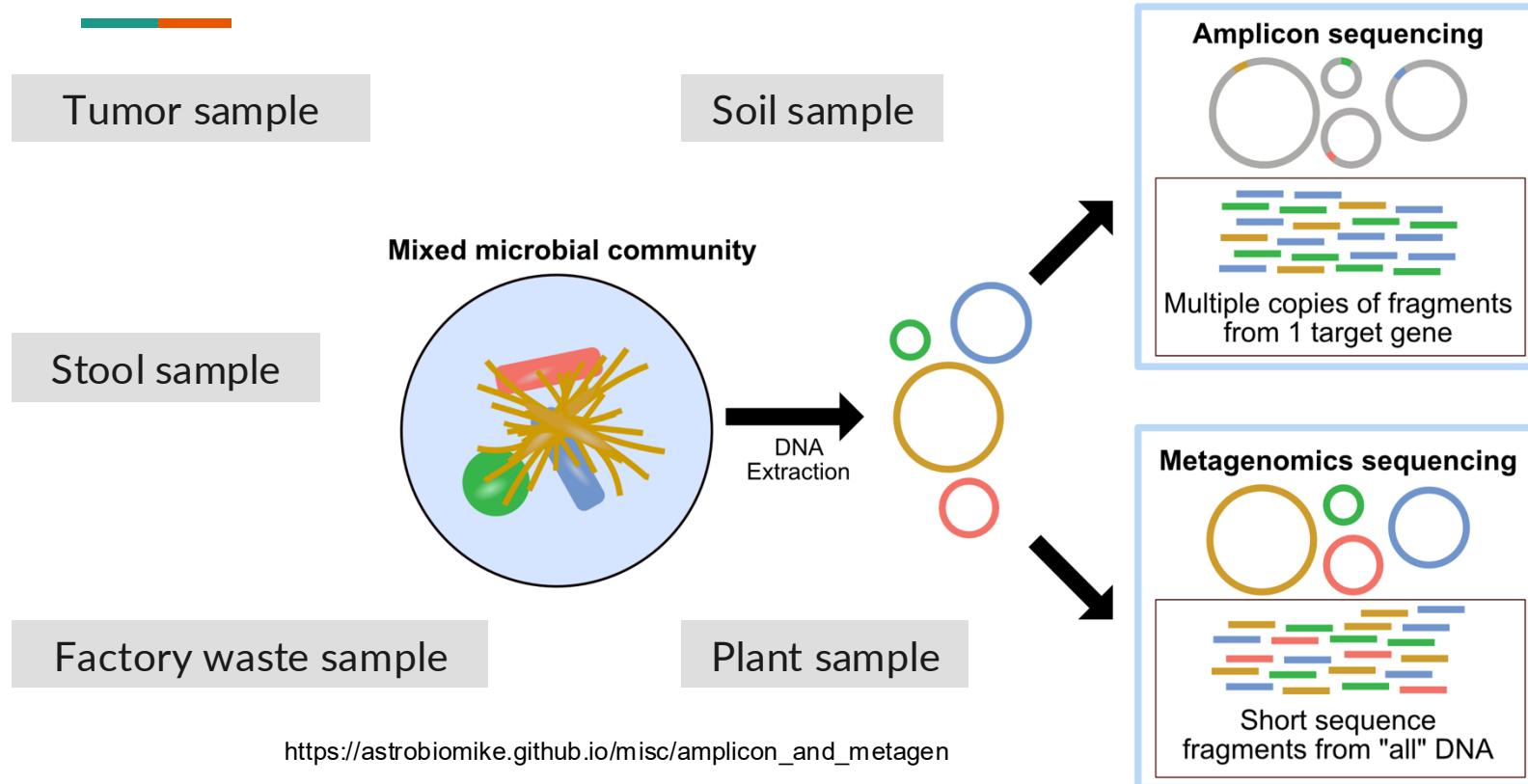
Sira Sriswasdi, PhD

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

Today's agenda

- Microbiome and metagenomics
- Relevant bioinformatics techniques

Metagenomics: analysis of a collection of organisms



Challenges in metagenomics

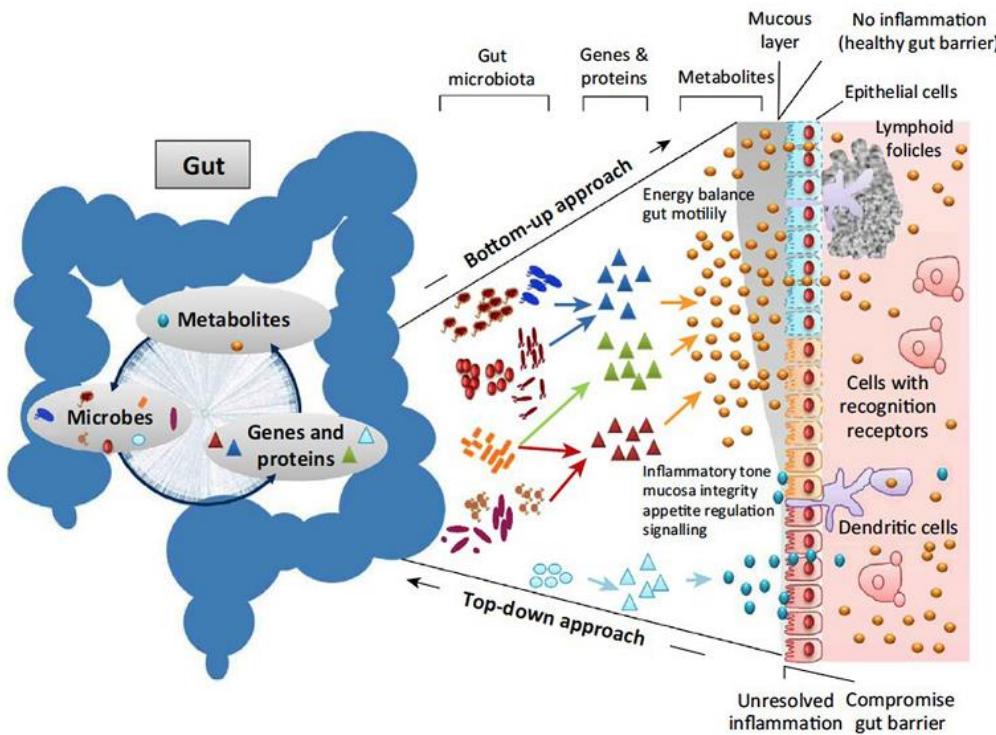


- Reads come from different genomes, how to distinguish?
 - **Clinical samples:** Host versus pathogens
 - **Environmental samples:** Mixed of animals and microorganisms
- **Genome assembly >> alignment**
 - Reference genomes are unavailable
 - Need to know gene sequences to infer functions
- Especially difficult when there are **multiple strains** of the same organisms!



Examples of microbiomes

Gut microbiomes



- Microbes digest food and synthesize nutrients and metabolites
- As a community, protect against harmful pathogens
- Reflect host health status

Human Microbiome Project

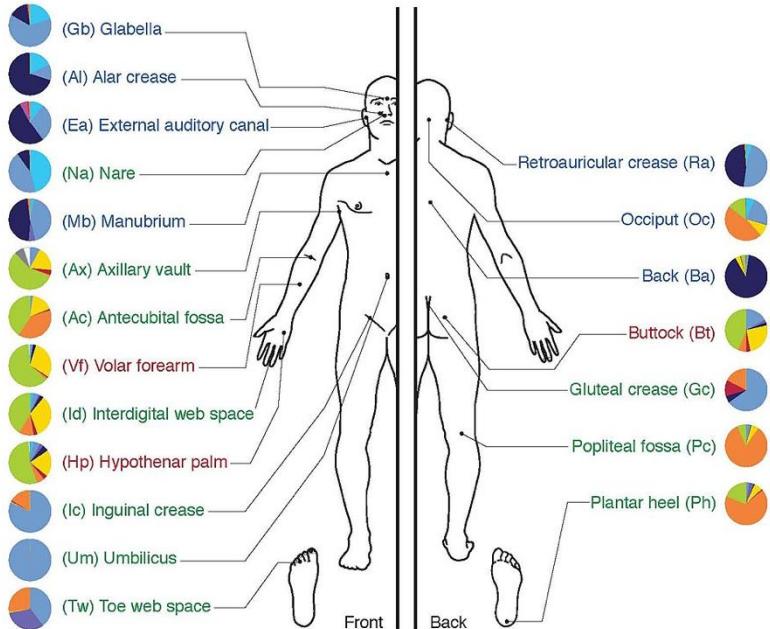
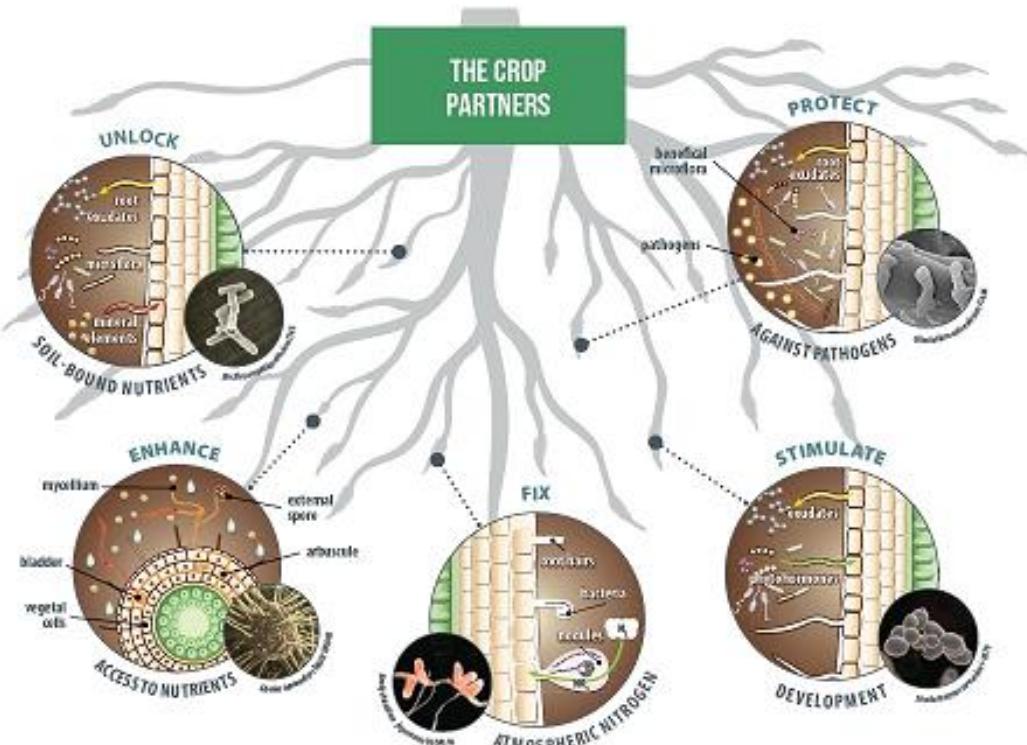


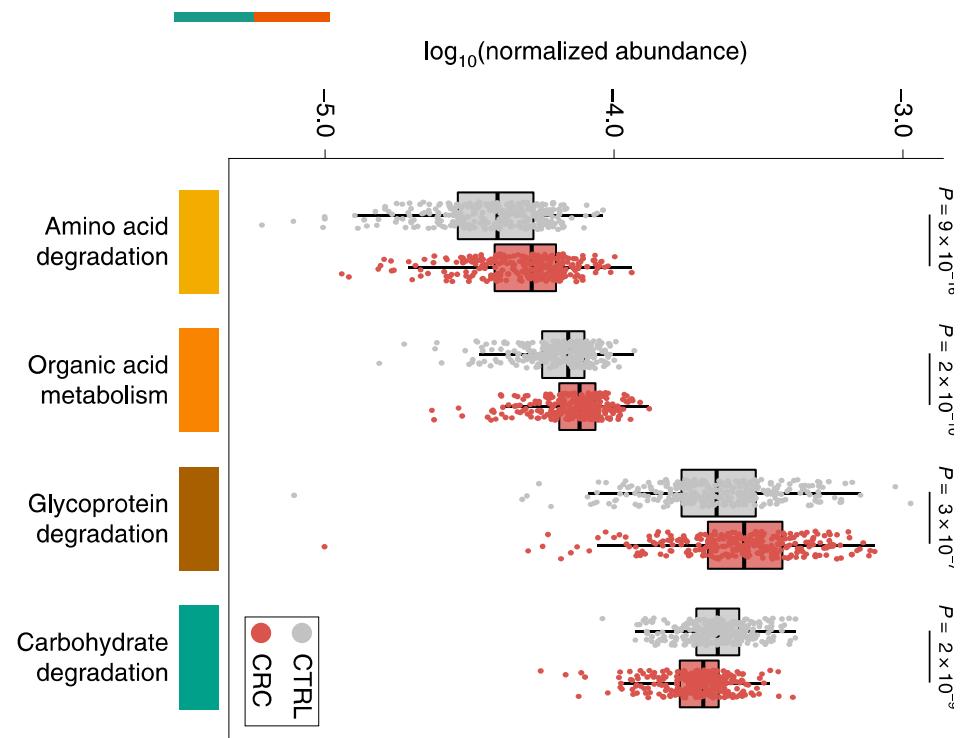
Image from Wikipedia.com

Plant and soil microbiomes



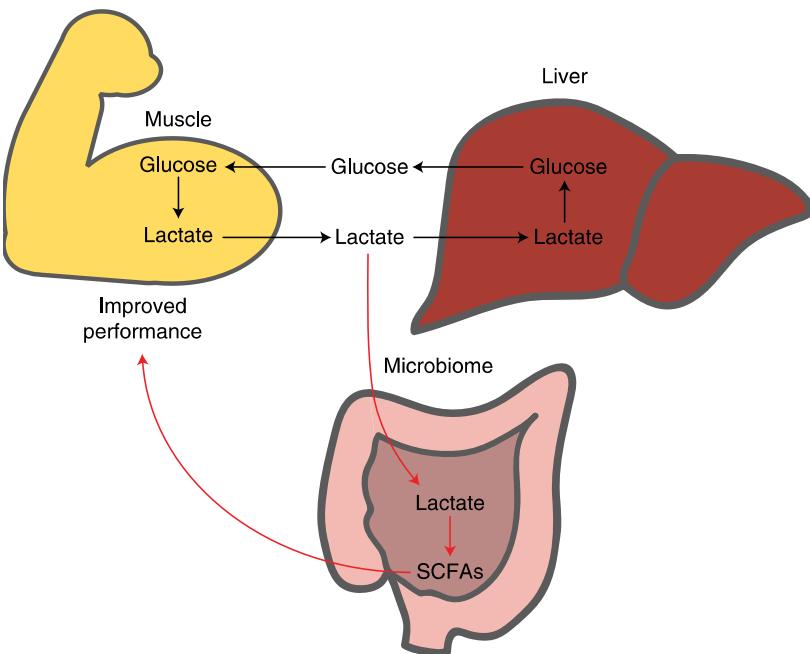
- Decompose and synthesize nutrients for plants
- As a community, protect against harmful pathogens
- Stimulate plant growth

Tumor microbiome



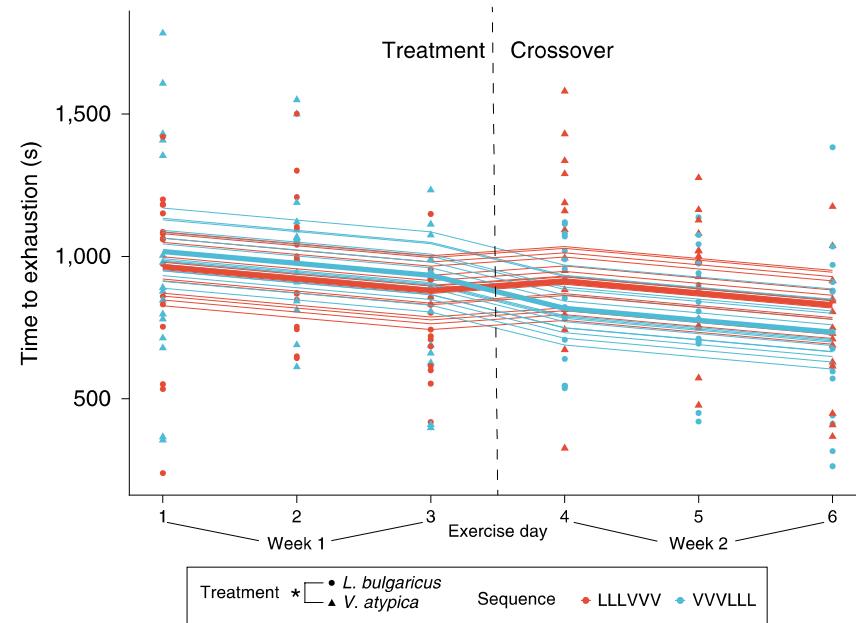
- Exchange nutrients with cancer cells in tumor microenvironment to sustain tumor growth
- Influence host immune response
- Can serve as biomarkers and targets for cancer treatment
 - Disrupt microbiome to disrupt cancer growth

Lactate-utilizing bacteria in athlete guts

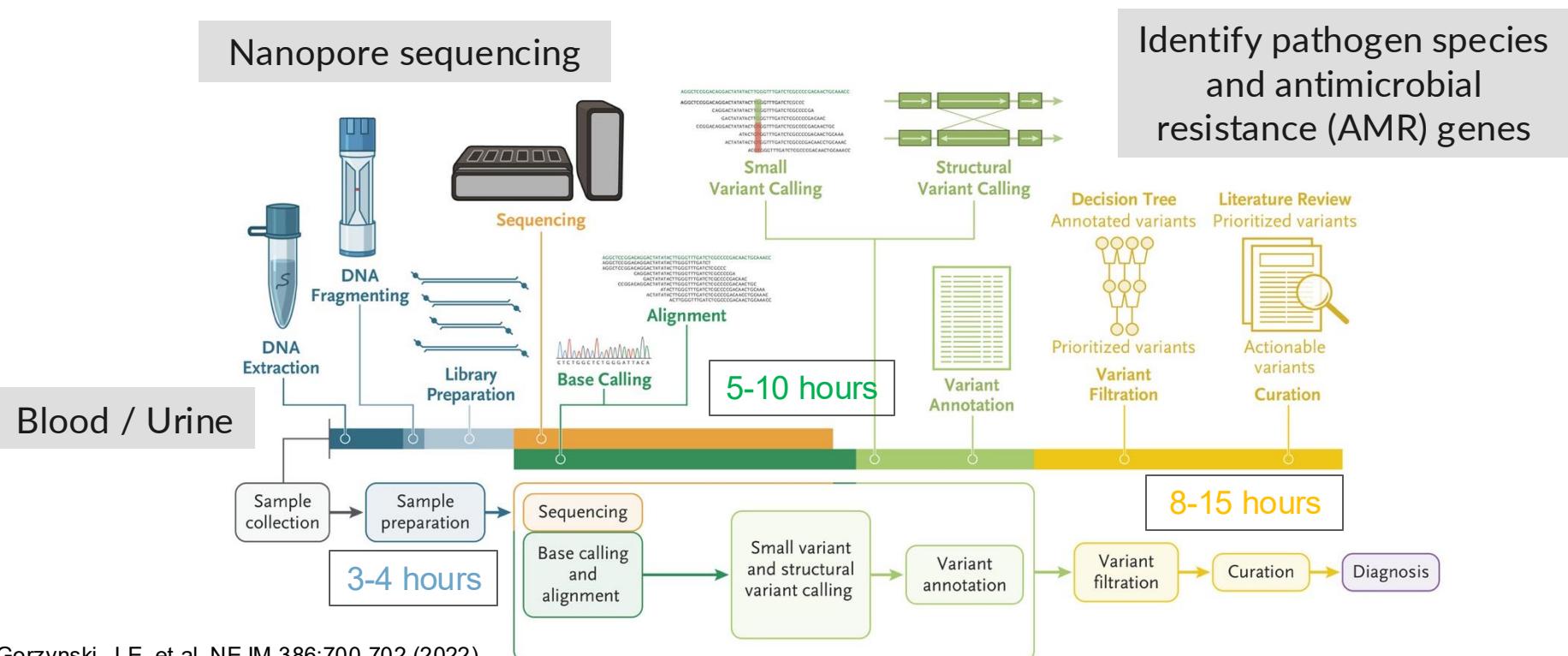


Scheiman et al. Nat Med 25:1104-1109 (2019)

Increased athletic ability in mice with transplanted microbiome



Rapid pathogen detection for clinical decision



Wildlife conservation via metagenomics monitoring



Different
dietary plants
(nutrition,
isotopic
evidences)

Père David's deer and their gut microbiome

Dissimilarity

Gut microbial
composition

*Next-generation sequencing
bioinformatics' analysis*

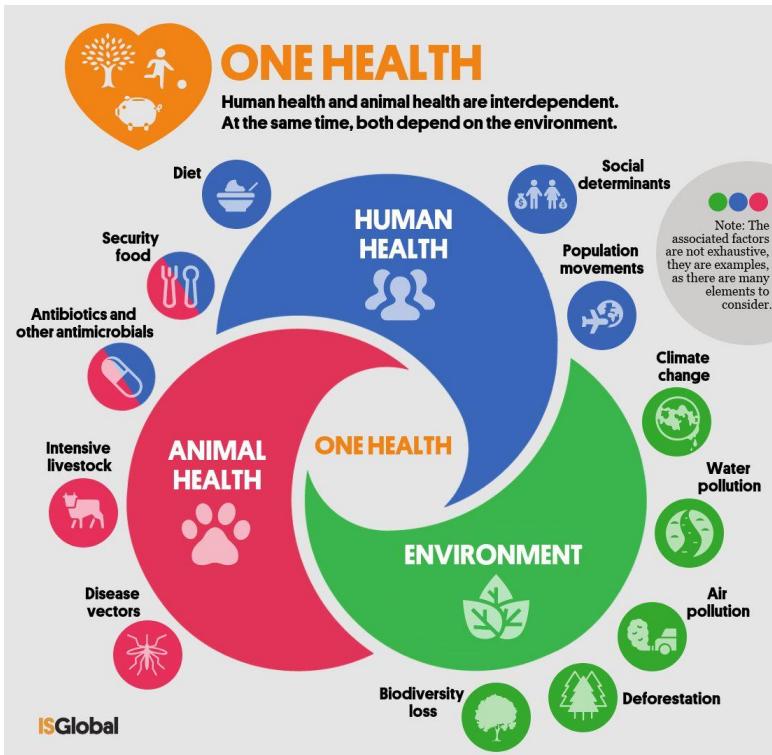
Gut microbial function
(cellulose digestion,
salt-related
metabolism)

Conservation

Reflecting increased
evolutionary potential
and resilience in response
to environment changes

Helping us select a
putative translocation
region

Metagenomics for One Health



- Monitoring of pathogens in environmental, wildlife, livestock, and agricultural microbiomes
- Monitoring of pathogen evolution during an active outbreak
- **Resistome:** Targeted sequencing of known AMR genes

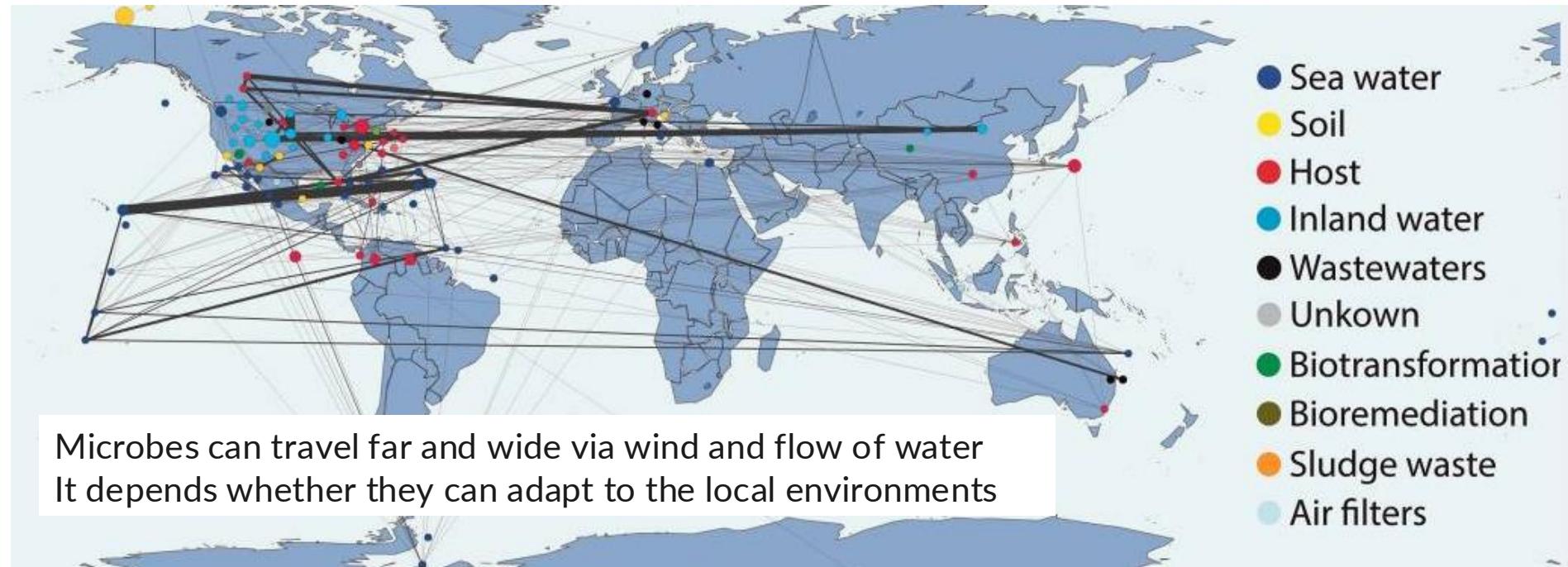
Some research questions in metagenomics

- **Health:** microbiome functional markers for diseases and health statuses
- **Ecology:**
 - Change in microbiome due to human actions
 - Factory and hospital wastes
 - Global warming
 - Microbiome of extreme conditions → new biotechnology tools
- **Agriculture:** soil-plant microbiome interactions
- Pathogen surveillance



Dynamics of microbiomes

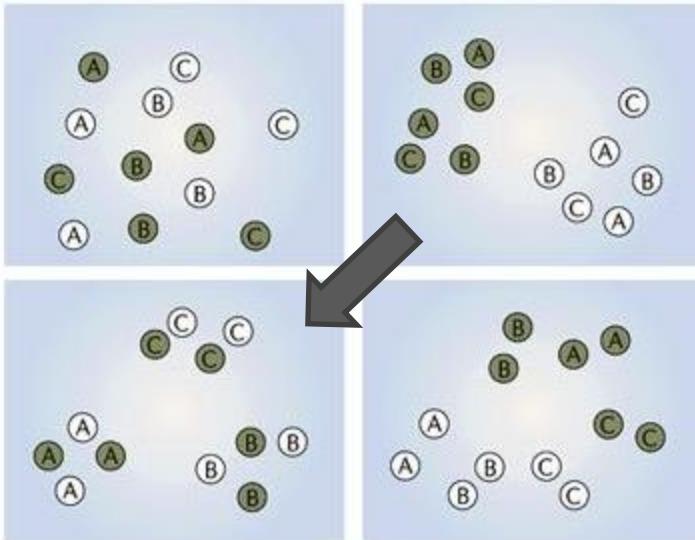
Everything is everywhere, but the nature selects



Environment defines microbiome more than geography

Environmental impact

Geographical impact

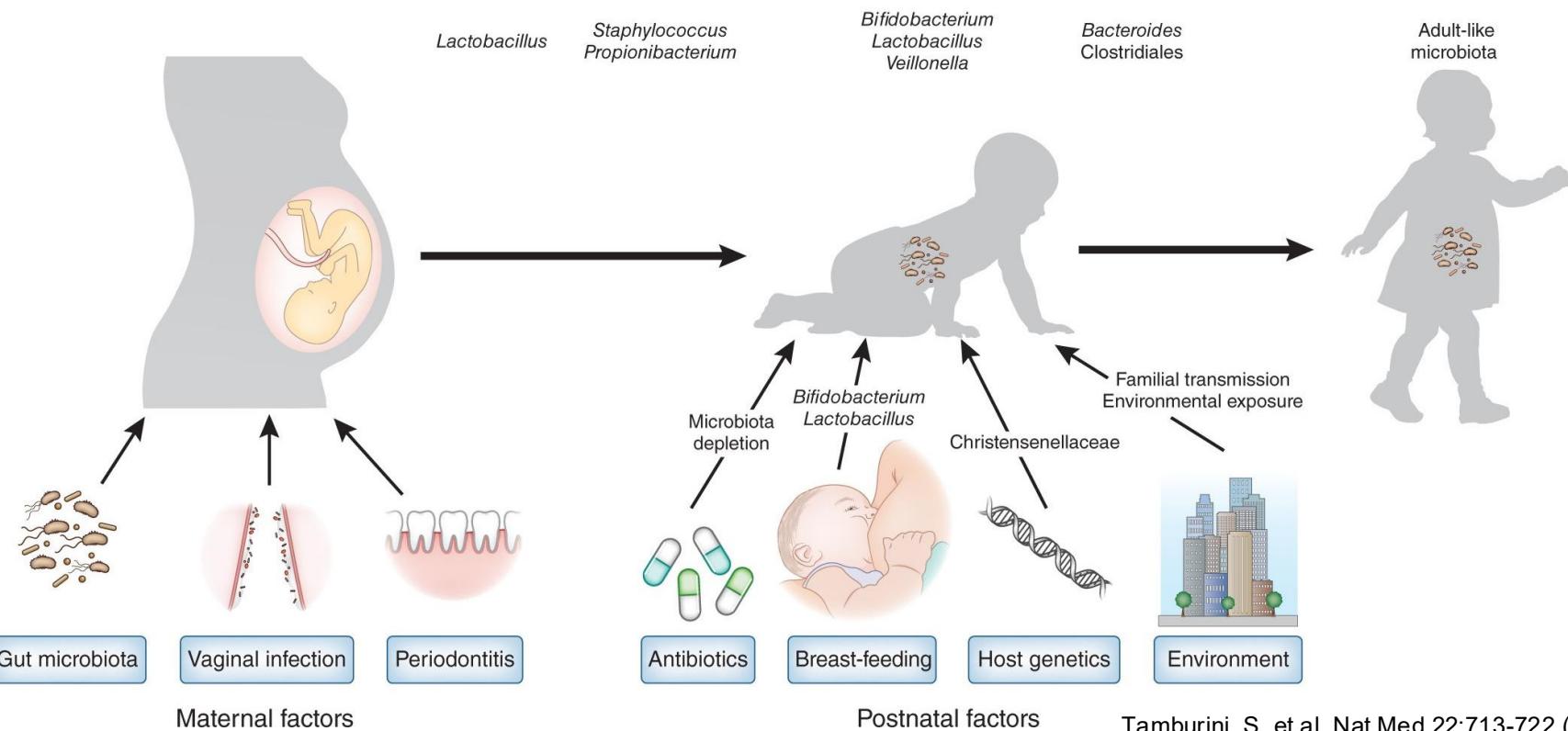


● ○ indicate same geography

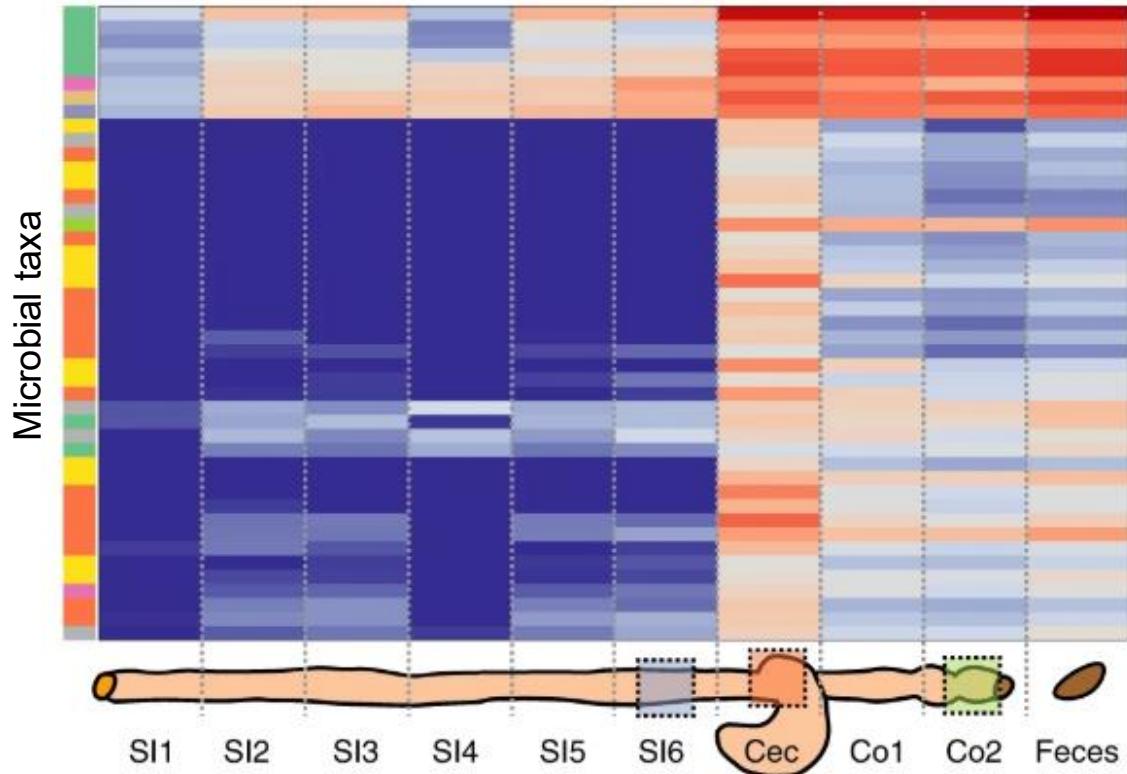
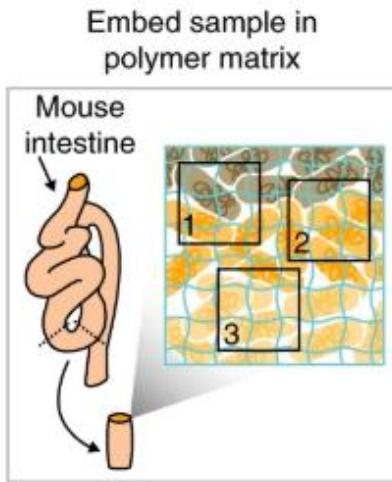
A, B, C indicate same environment

- **Large animals** can withstand changing environments but cannot travel world-wide
- **Microbes** can travel world-wide but cannot withstand changing environments

Microbiomes are extremely dynamics



Spatial variability of gut microbiome



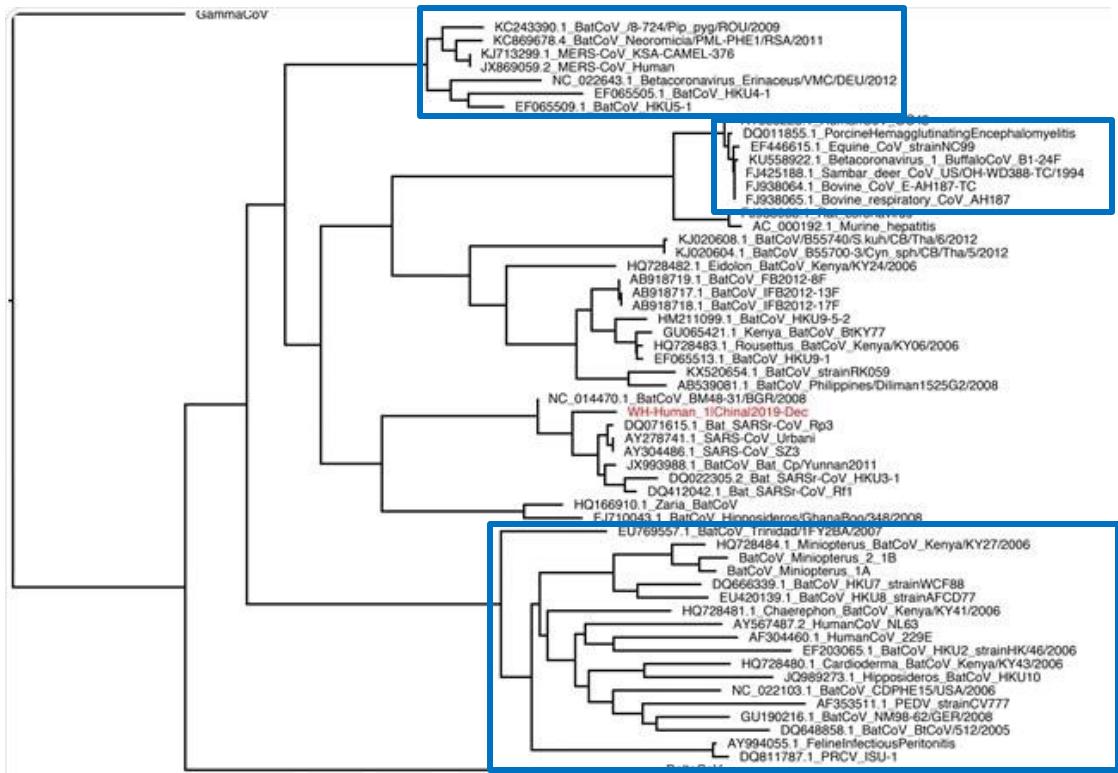
Sheth et al., Nature Biotech 37: 877-883 (2019)



Operational taxonomic unit (OTU)

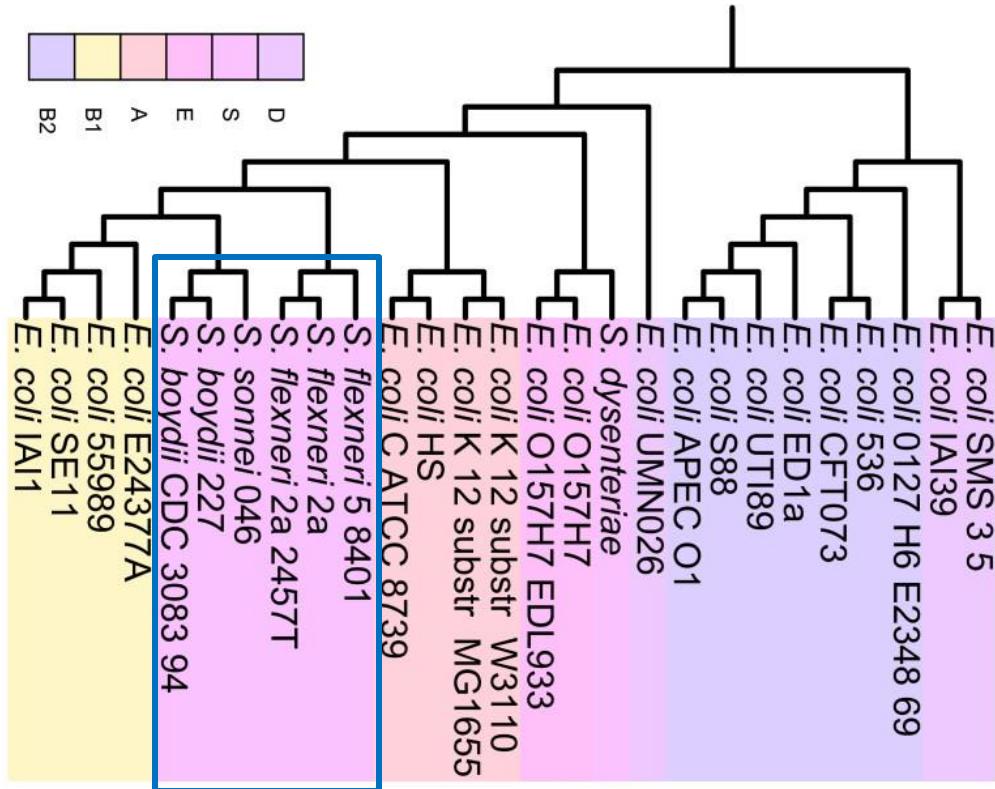
Traditional phylogeny: Cluster of similar sequences

- Is sequence similarity alone sensitive enough to distinguish **closely-related microbes**?
- If sequences are all very similar, where to draw the cutoff?

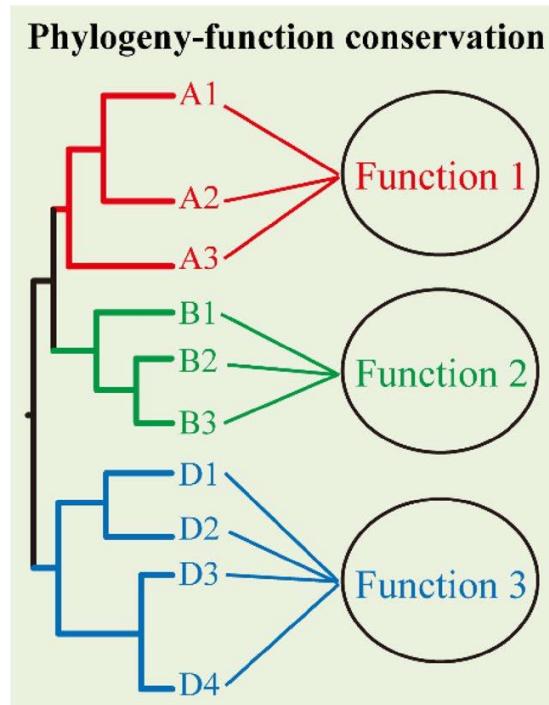


When sequence similarity alone is not enough

- Genus *Shigella* are pathogens that evolved from an *E. coli* ancestor
- 80-90% similarity to some *E. coli* clades
- Definition of taxonomy may require both genotypes and phenotypes



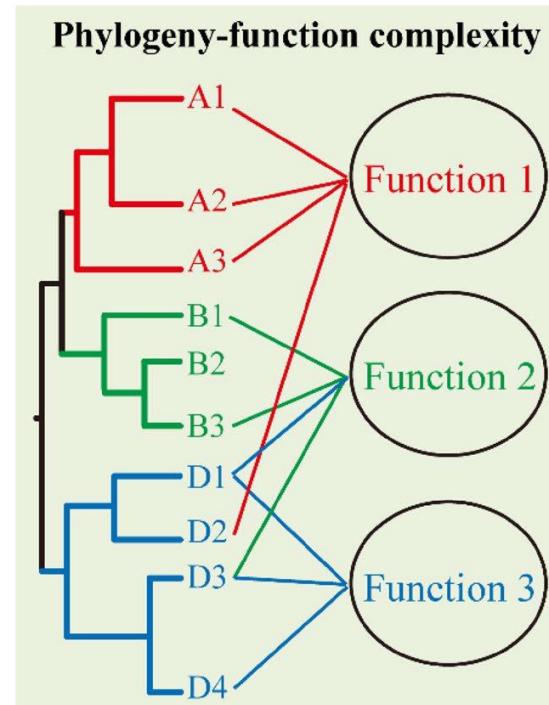
Phylogeny based on function (inferred from genes)



→

Phylogeny-function decoupling

- Horizontal gene transfer
- Gene gain and loss
- Functional type
- Functional rate



A rough OTU similarity cutoff for microbes

nature reviews microbiology

Published: 14 August 2014

Uniting the classification of cultured and uncultured bacteria and archaea using 16S rRNA gene sequences

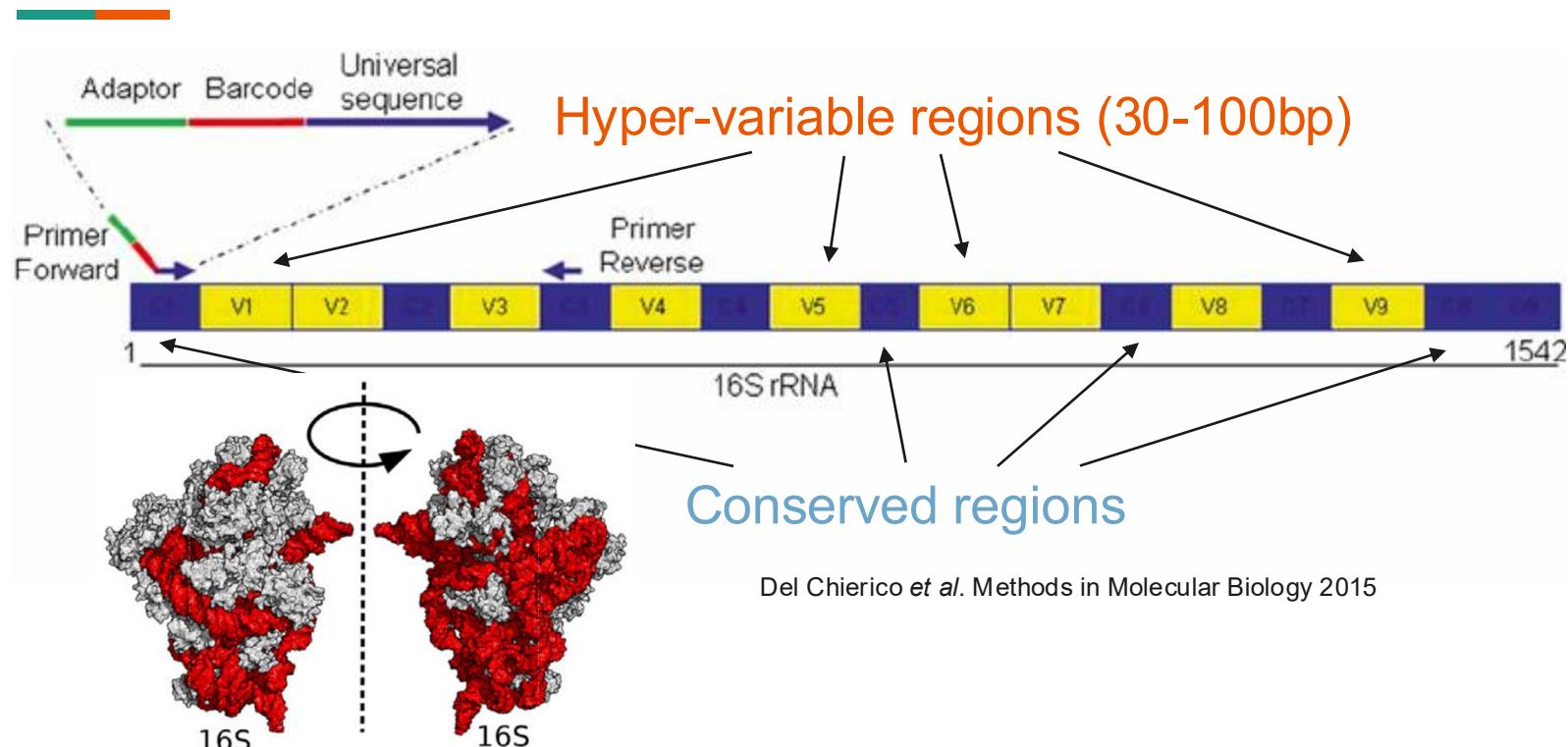
Pablo Yarza , Pelin Yilmaz, Elmar Pruesse, Frank Oliver Glöckner, Wolfgang Ludwig, Karl-Heinz Schleifer,
William B. Whitman, Jean Euzéby, Rudolf Amann & Ramon Rosselló-Móra 

- Based on annotated **16S rRNA genes**
 - 94.5% similarity for genus
 - 86.5% similarity for family
 - 75% similarity for phylum
- Near-complete sequences are required to get accurate classification



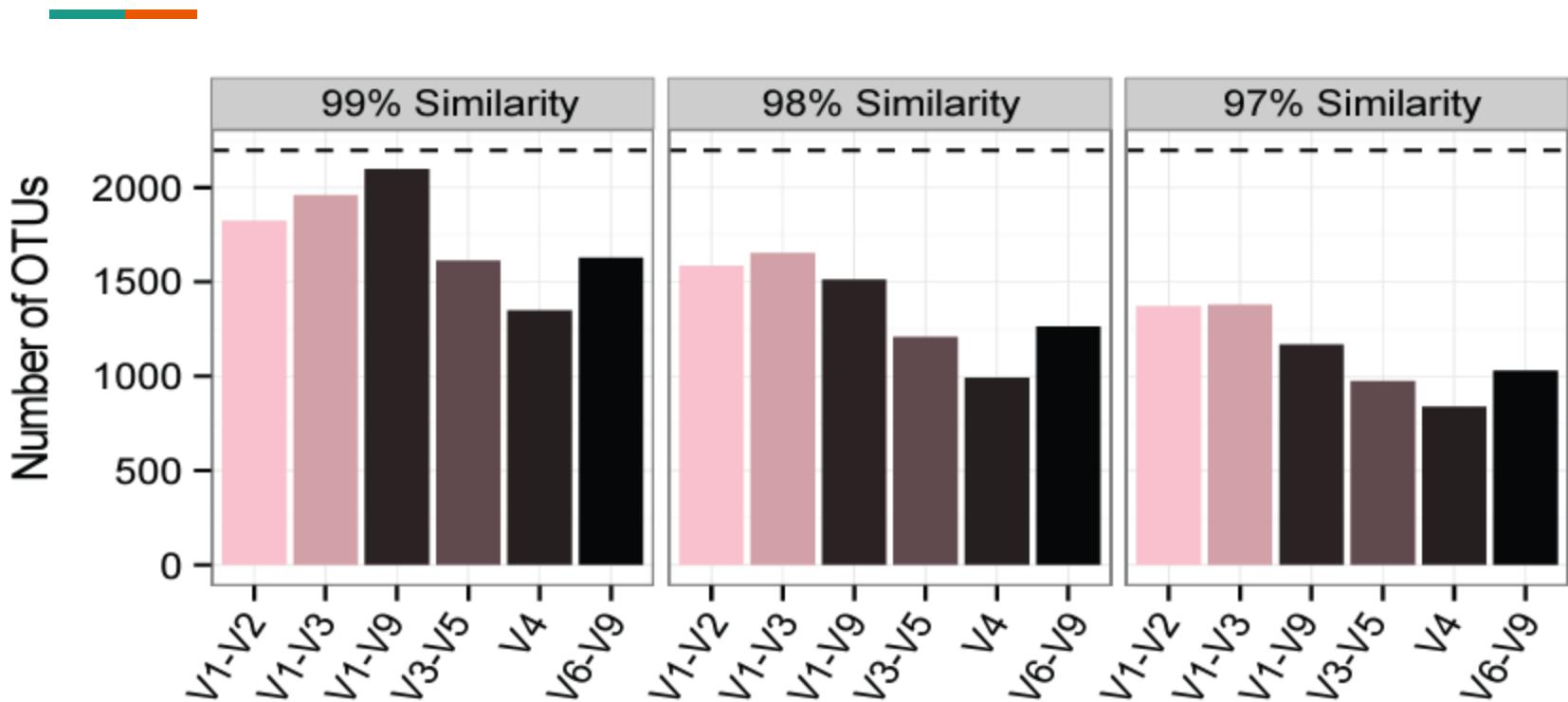
Ribosomal RNA analysis

16S rRNA in prokaryotes



Del Chierico et al. Methods in Molecular Biology 2015

Different variations across hypervariable regions



rRNA sequence alignment

The screenshot shows a search interface titled "Choose Search Set". Under the "Database" section, the "rRNA/ITS databases" option is selected. A dropdown menu is open, showing several options under "16S ribosomal RNA sequences (Bacteria and Archaea)": "16S ribosomal RNA sequences (Bacteria and Archaea)", "18S ribosomal RNA sequences (SSU) from Fungi type and reference material", "28S ribosomal RNA sequences (LSU) from Fungi type and reference material", and "Internal transcribed spacer region (ITS) from Fungi type and reference material". The first option is highlighted. On the right side of the interface, there is a "Targeted Loci Project Information" section and a "Add organism" button.

- Endpoint of rRNA sequence analysis is taxonomic assignment
- **Short read:** accurate only up to family and genus
- **Long read:** species and sub-species level classification is possible

rRNA databases



high quality ribosomal RNA databases
Home SILVangs Browser Search ACT Download Doc

SILVA

Welcome to the SILVA rRNA database project

A comprehensive on-line resource for quality checked and aligned ribosomal RNA sequence data.

SILVA provides comprehensive, quality checked and regularly updated datasets of aligned small (16S/18S, SSU) and large subunit (23S/28S, LSU) ribosomal RNA (rRNA) sequences for all three domains of life (*Bacteria*, *Archaea* and *Eukarya*).

SILVA are the official databases of the software package ARB.

ANNOUNCEMENTS

RDP News

12/12/2018 RDP and Fungene Pipelines are back

RDP Release 11, Update 5 :: September 30, 2016

3,356,809 16S rRNAs :: 125,525 Fungal 28S rRNAs
Find out what's new in RDP Release 11.5 [here](#).

GREENGENES

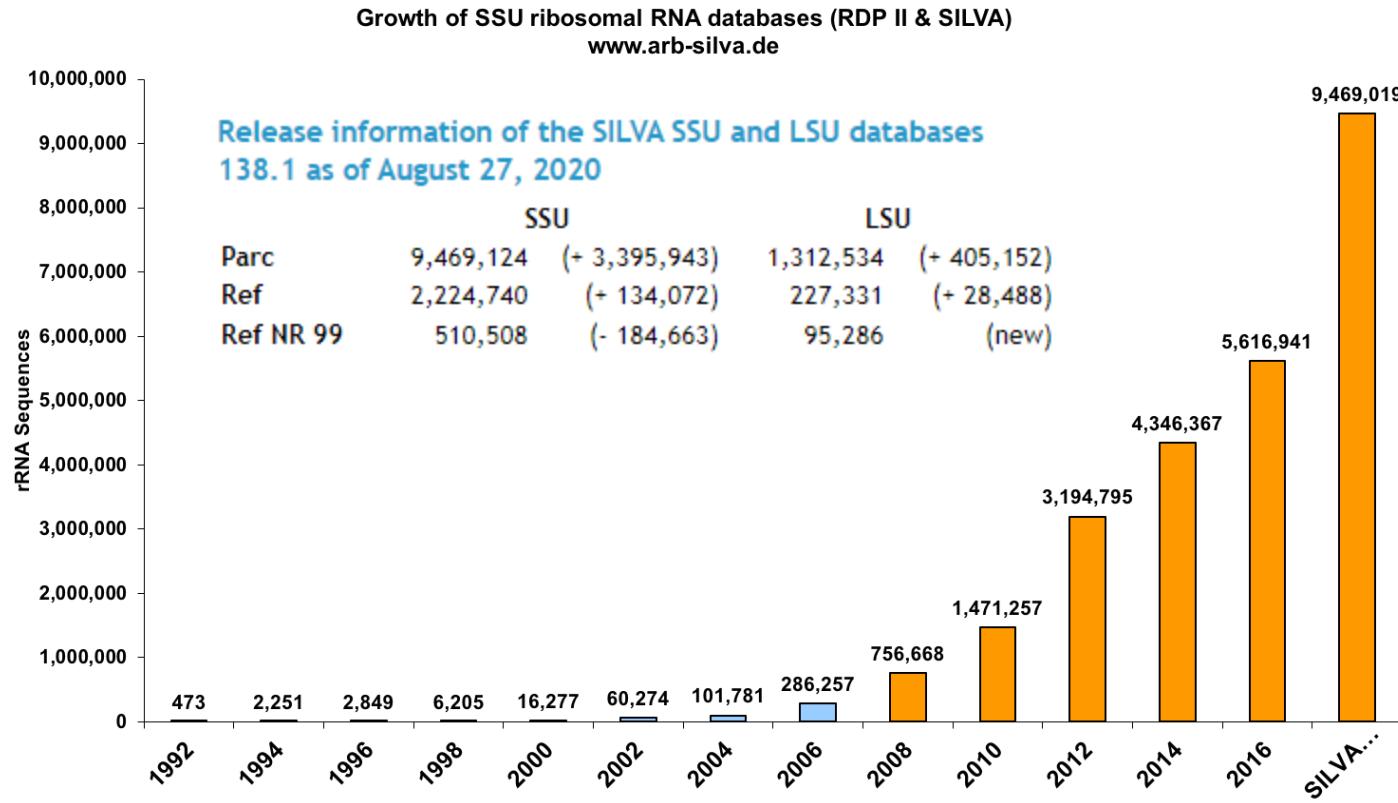
The Greengenes Database

Browse links below to download versions of the Greengenes 16S rRNA gene database or experimental datasets created with the PhyloChip 16S rRNA microarray. Beware that these publicly available versions of the Greengenes database utilize taxonomic terms proposed from phylogenetic methods applied years ago between 2012 and 2013. Since then, a variety of novel phylogenetic methods have been proposed for Archaea and Bacteria. For a recent example see Yokono 2018.

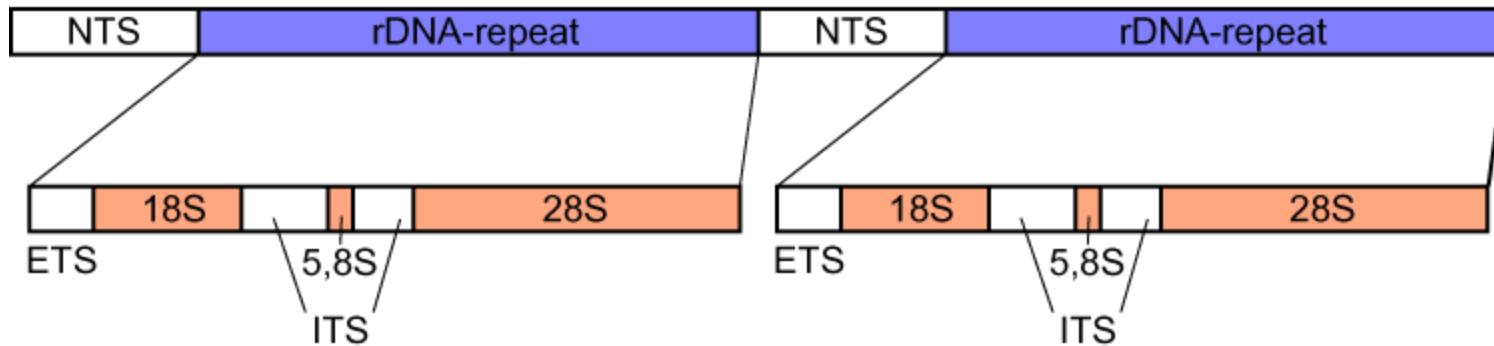
rDNA ITS based identifica

Current version: 8.2; Last updated: 2020-01-15 ([read more](#))
Number of ITS sequences (UNITE+INSD): 2 480 043; Number of UNITE fungal S

Explosion of known rRNA sequences



Internal transcribed spacer (ITS) in fungi and algae



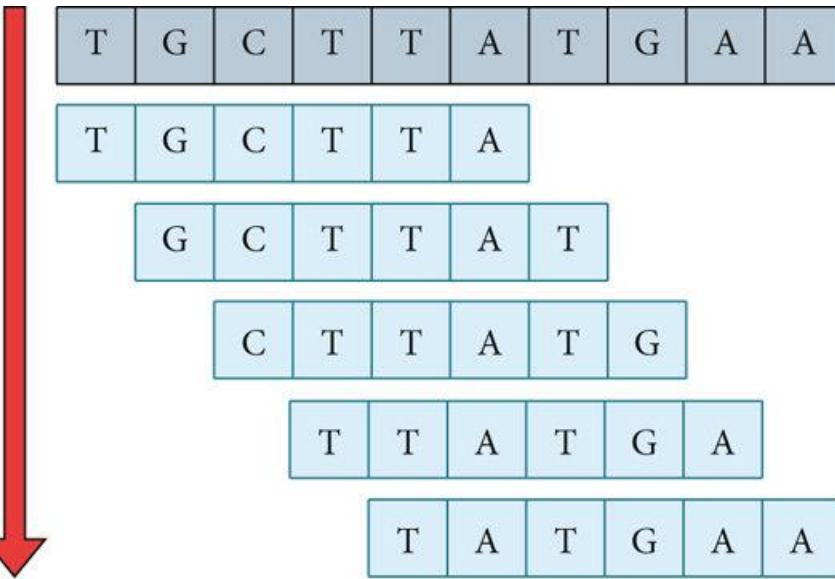
wikipedia.com

- Located between rRNA repeats
- Typically study ITS1 and ITS2 regions
- 400-1000 bp



DNA k -mers

k -mer: substring of length k



K-mer of size = 6

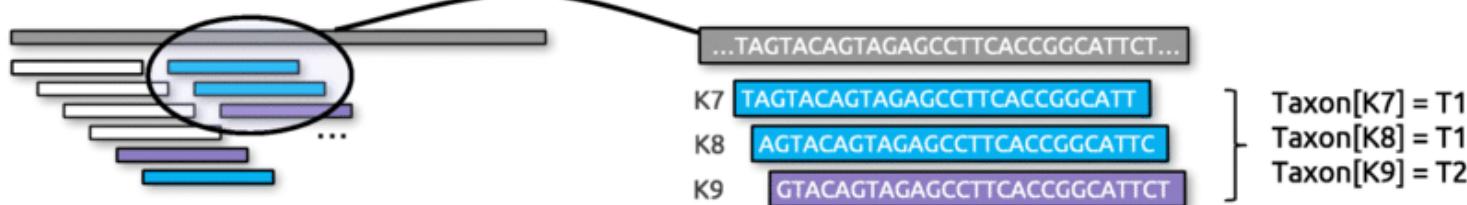
- Similar to the concept of word in BLAST
- A long alignment must also contain shorter matches
- **How about searching for short matches instead?**

***k*-mer-based taxonomic classification**

- We want to estimate $P(\text{genus} | \text{observed } w_1, w_2, \dots, w_n)$, where w_i 's are *k*-mer
- But $P(\text{observed } w_1, w_2, \dots, w_n | \text{genus})$ is much easier to calculate since we know the rRNA sequences
- **Bayes' rule:** $P(\text{genus} | \text{words}) \propto P(\text{words} | \text{genus}) \times P(\text{genus})$, where $P(\text{genus})$ is the prior probability of observing a genus

k-mer matching result

A Read *k*-mers are looked-up in the database and assigned to taxa:



C K-mer count and coverage in taxonomic report show evidence behind classifications:

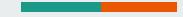
reads	kmers	dup	cov	taxID	rank	name
122	112	144	0.0004	11855	species	<i>Clostridioides difficile</i>
9650	7129	74.5	0.192	10632	species	Human polyomavirus 2
15	1570	1	0.0002	7643	species group	<i>Mycobacterium tb</i> complex

Bad classification with few k-mers
Good classification, reads cover genome

Number of distinct *k*-mers for taxon, and coverage of the taxon's *k*-mers

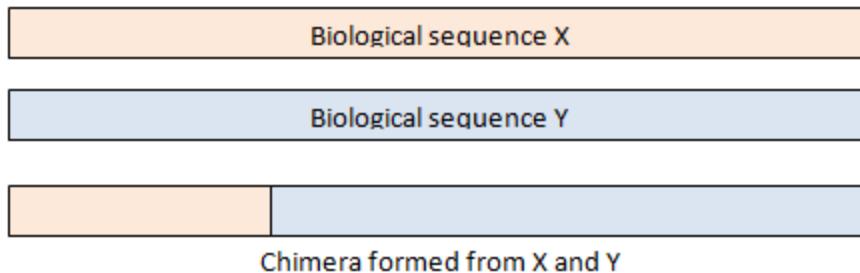
Benefits of using k -mer matching

- Much faster than regular alignment
 - Find short perfect matches instead of dynamic programming
- Filter reads for downstream analysis
 - Group reads for assembly
 - Remove host contamination
 - Remove reads mapped to abundant bacteria

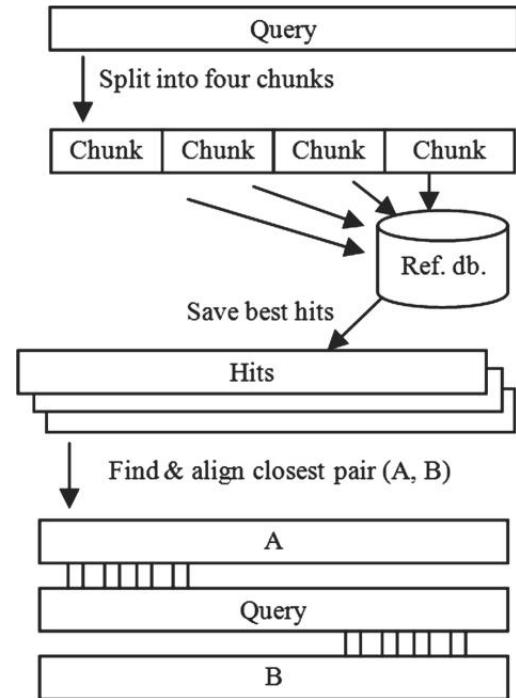


Metagenomics read processing

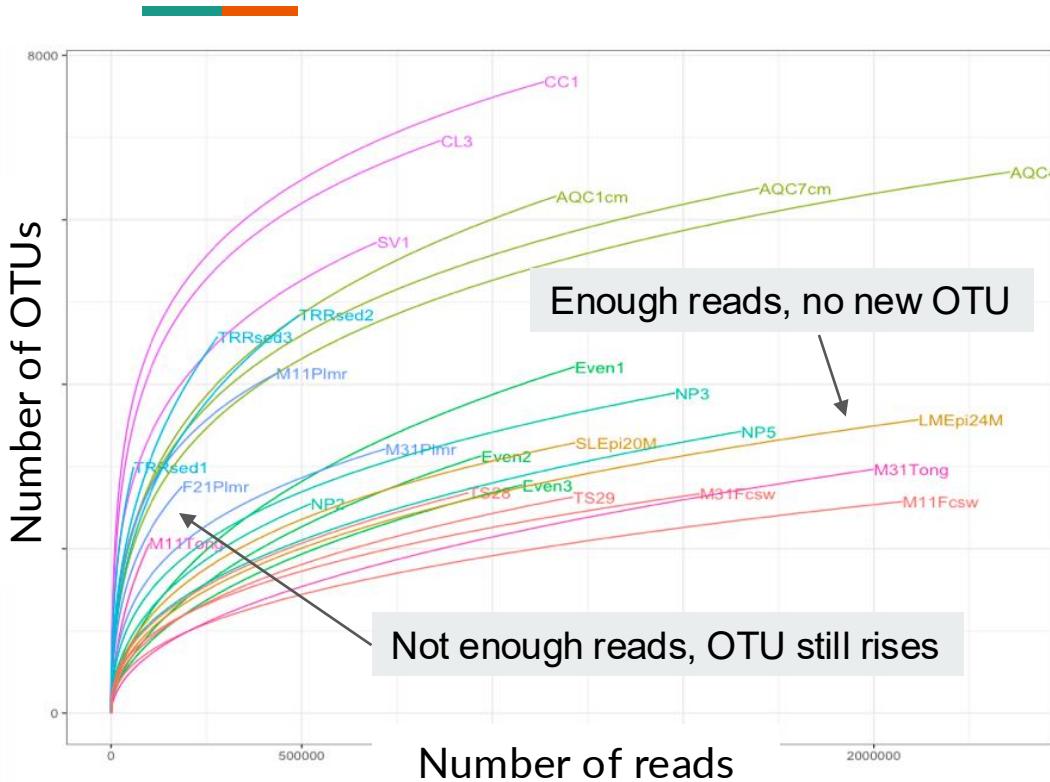
Chimeric reads in amplicon sequencing



- Produced during PCR amplification
- Detected by alignment different portion of the reads to rRNA databases
 - Same idea as gapped alignment



Rarefaction curve



- One curve = one sample
- Show the gain in unique OTU as more reads are obtained
- **Plateau:** additional reads do not detect more OTU = enough reads

Complexity of microbiome composition

Diversity indices/ Parameters	Description	Formula
Shannon diversity index (H)	Estimator of species richness and species evenness: more weight on species richness	$H = -\sum_{i=1}^s (p_i \ln p_i)$ <p>where s is the number of OTUs and p_i is the proportion of the community represented by OTU i.</p>
Simpson's index (D)	Estimator of species richness and species evenness: more weight on species evenness	$D = \frac{1}{\sum_{i=1}^s p_i^2}$ <p>S = number of taxa</p> <p>where s is the total number of species in the community and p_i is the proportion of community represented by OTU i.</p>

Kim, B.-R. et al. J Microbiol Biotechnol 27:2089-2093 (2017)

- **Richness** = number of distinct species
- **Evenness** = lack of dominant species

Behavior of Shannon entropy

Diversity indices/ Parameters	Description	Formula
Shannon diversity index (H)	Estimator of species richness and species evenness: more weight on species richness	$H = -\sum_{i=1}^s (p_i \ln p_i)$
Simpson's index (D)	Estimator of species richness and species evenness: more weight on species evenness	$D = \frac{1}{\sum_{i=1}^s p_i^2}$ S = number of taxa where s is the total number of species in the community and p_i is the proportion of community represented by OTU i .

Kim, B.-R. et al. J Microbiol Biotechnol 27:2089-2093 (2017)

- Entropy is maximized when $p_i = \frac{1}{S} \rightarrow H_{max} = \ln(S)$
- Entropy is minimized when there is one dominant species $H_{min} = 0$

Behavior of Simpson's index

Diversity indices/ Parameters	Description	Formula
Shannon diversity index (H)	Estimator of species richness and species evenness: more weight on species richness	$H = - \sum_{i=1}^s (p_i \ln p_i)$ <p>where s is the number of OTUs and p_i is the proportion of the community represented by OTU i.</p>
Simpson's index (D)	Estimator of species richness and species evenness: more weight on species evenness	$D = \frac{1}{\sum_{i=1}^s p_i^2}$ <p>S = number of taxa</p> <p>where s is the total number of species in the community and p_i is the proportion of community represented by OTU i.</p>

Kim, B.-R. et al. J Microbiol Biotechnol 27:2089-2093 (2017)

- D is maximized when $p_i = \frac{1}{S} \rightarrow D_{max} = S$
- D is minimized when there is only dominant species $D_{min} = 1$



Similarity between microbial compositions

Bray-Curtis dissimilarity (intersection over union)

Sample 1

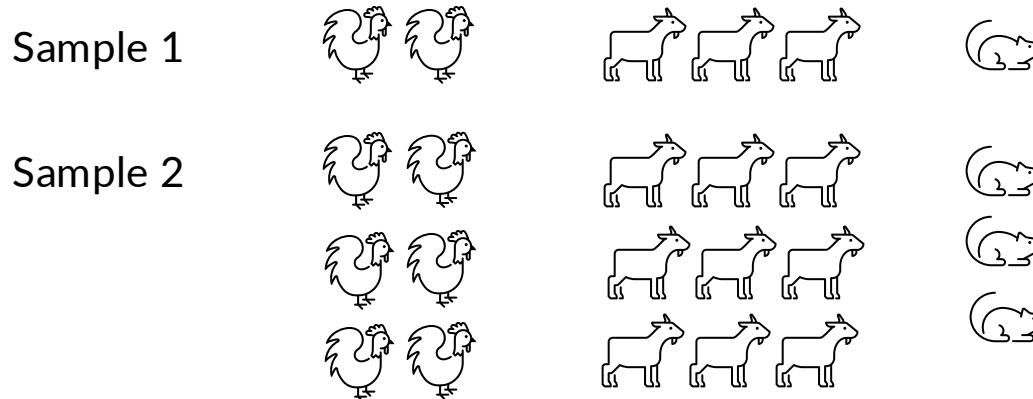


Sample 2



- $S_1 = \# \text{ of individuals in sample 1} = 6$
- $S_2 = \# \text{ of individuals in sample 2} = 5$
- Overlap = 1 chicken + 2 goat + 1 mouse = 4
- Bray-Curtis dissimilarity = $1 - \frac{2 \times \text{Overlap}}{S_1 + S_2} = 1 - \frac{8}{11} = \frac{3}{11}$

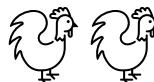
Impact of sequencing depth on Bray-Curtis



- Bray-Curtis is only suitable between samples with similar depths
- Both samples have 2:3:1 composition
- $\text{Bray-Curtis dissimilarity} = 1 - \frac{2 \times \text{Overlap}}{S_1 + S_2} = 1 - \frac{12}{24} = \frac{1}{2}$

Bray-Curtis does not consider taxonomic similarity

Sample 1



Sample 2



Sample 3



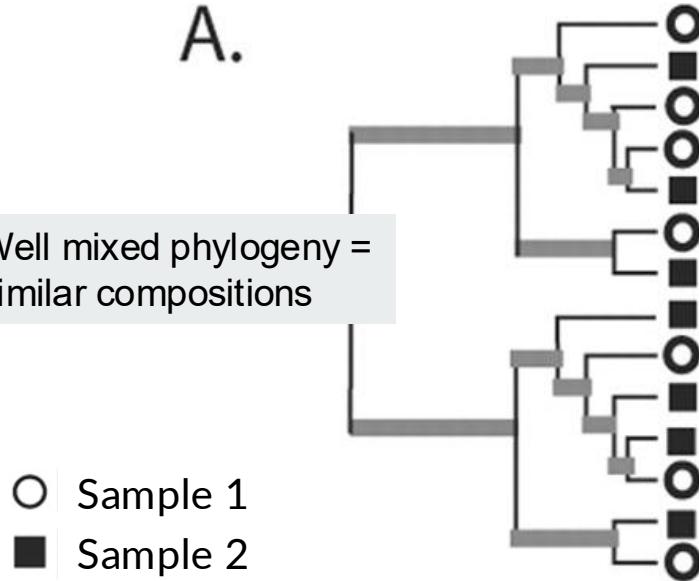
- Difference between white goats and black goats is treated the same as the difference between goats and birds
- Discordant with the view of functional interpretation of metagenomics

UniFrac distance

Lozupone, C. and Knight, R. Applied and Environmental Microbiology 71 (2005)

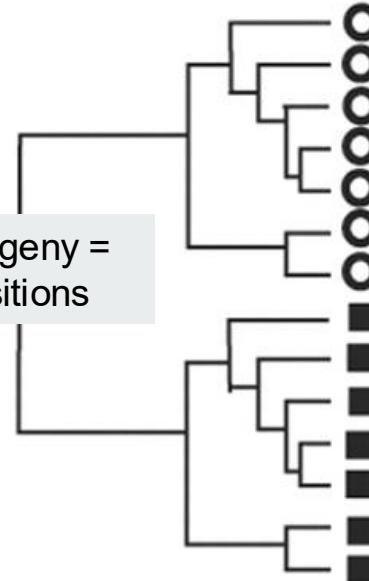
A.

Well mixed phylogeny =
similar compositions



B.

Segregated phylogeny =
dissimilar compositions

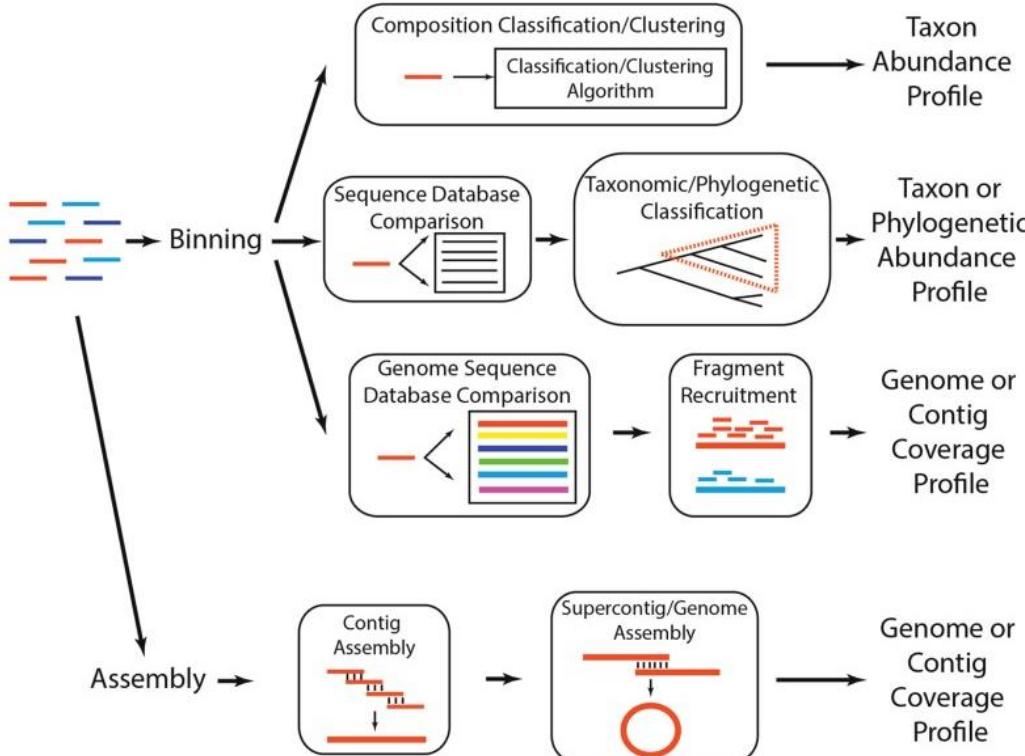


- UniFrac = fraction of shared phylogenetic branches between samples
- Can be weighted or unweighted by taxa abundances



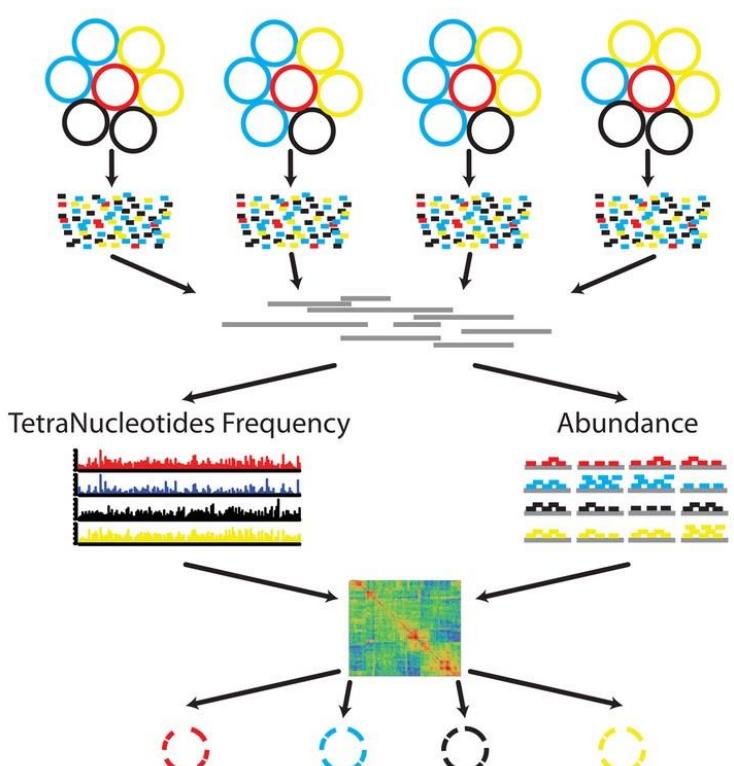
Shotgun metagenomics

Shotgun metagenomics



- **Binning:** grouping DNA/RNA from the same host organisms together
- Direct assembly is possible for abundant species and long reads

Read binning strategies

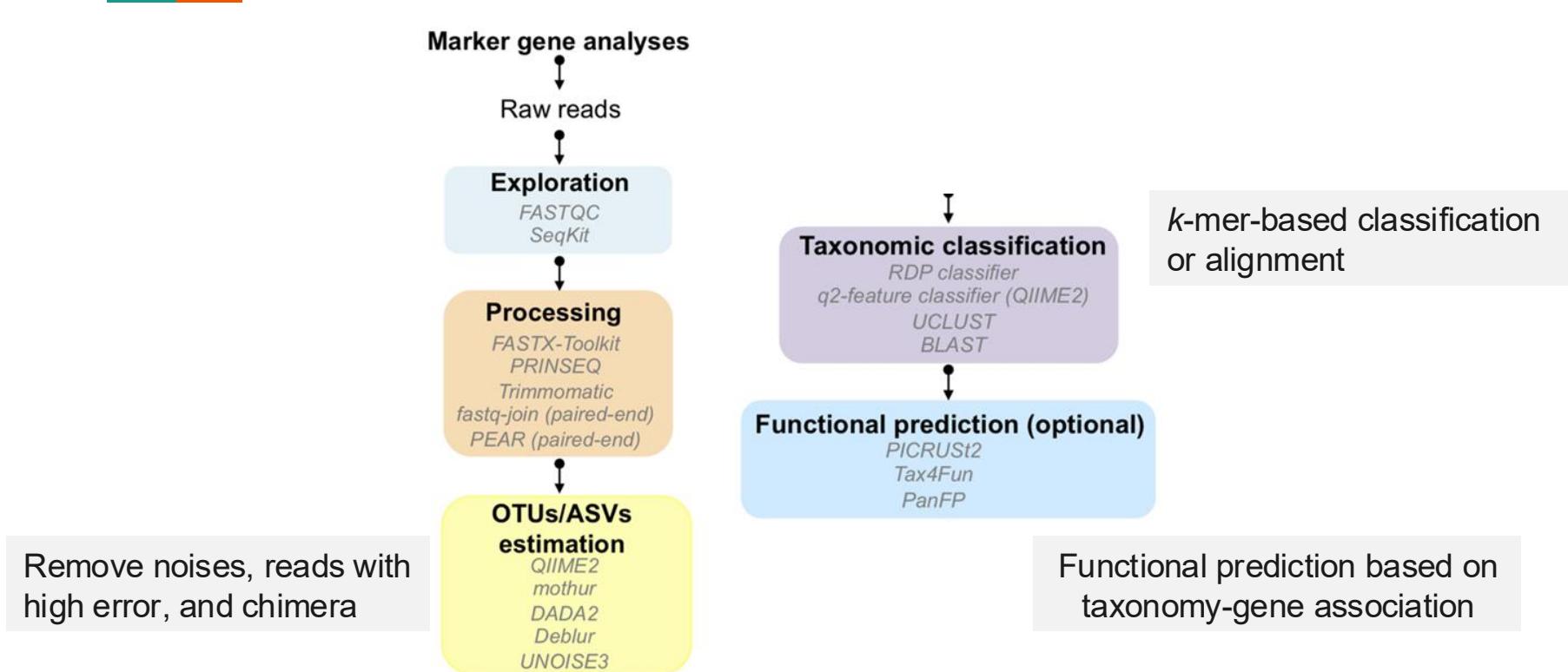


- **k-mer:** reads originated from the same species have similar k-mer profiles (like GC content)
- **Abundance:** reads originated from the same organism have correlated abundance profiles across samples

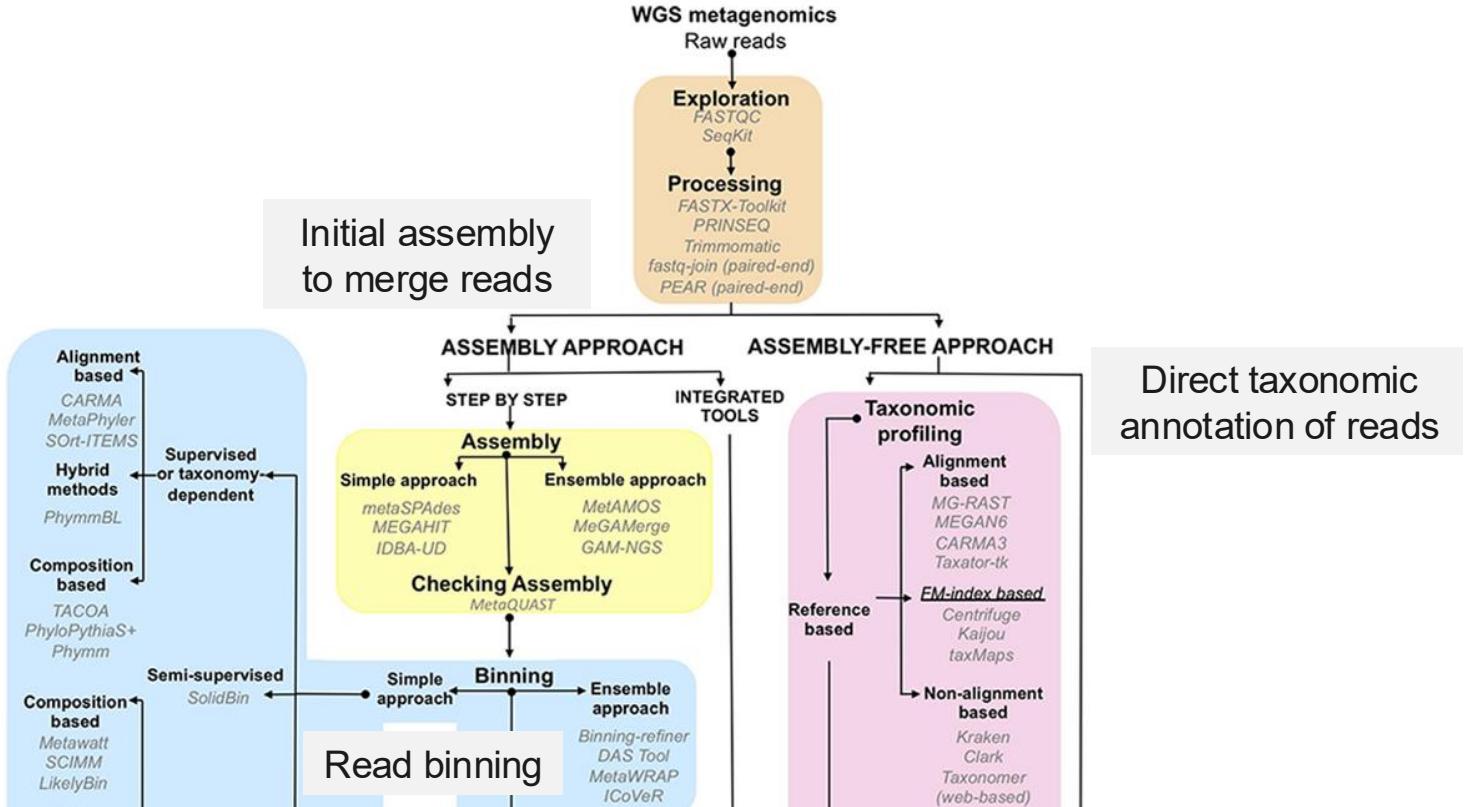


Resources for metagenomics

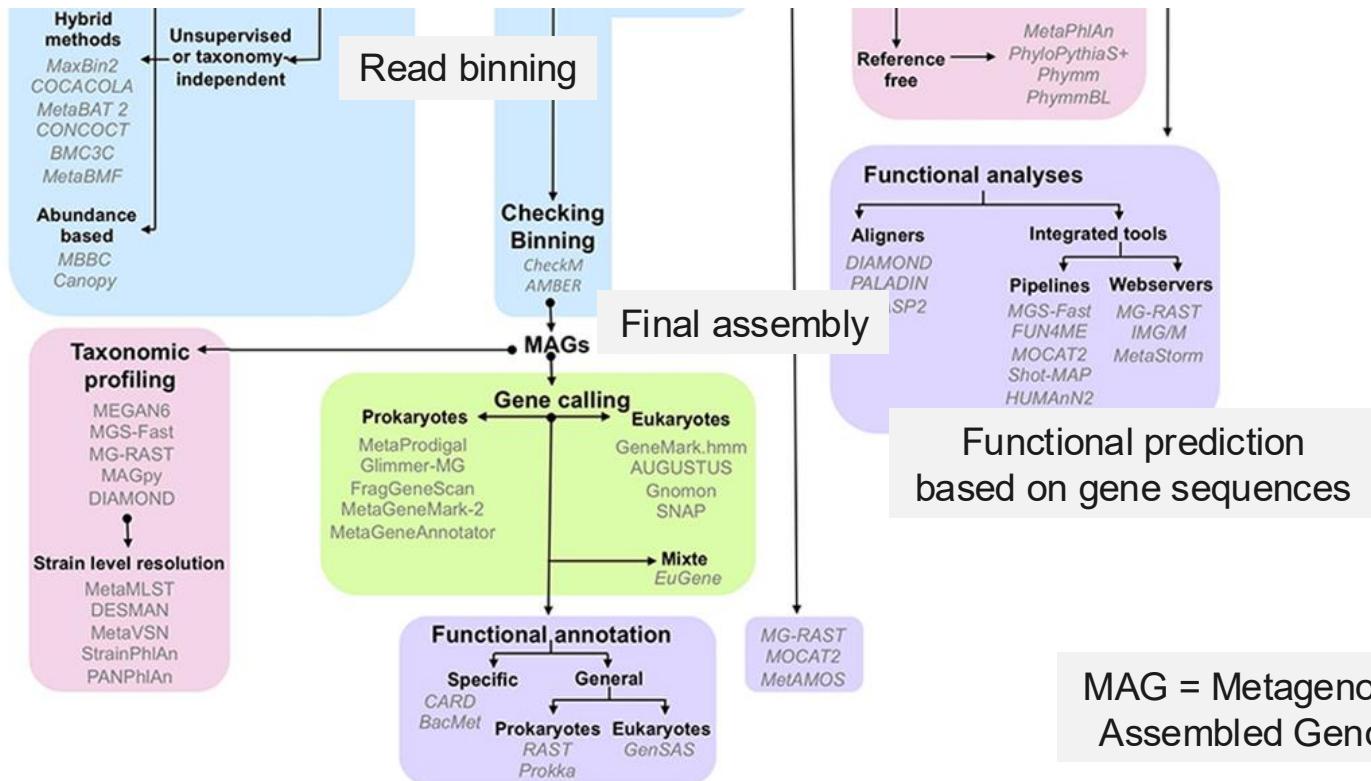
Amplicon analysis pipeline (e.g., 16S rRNA)



Shotgun analysis pipeline



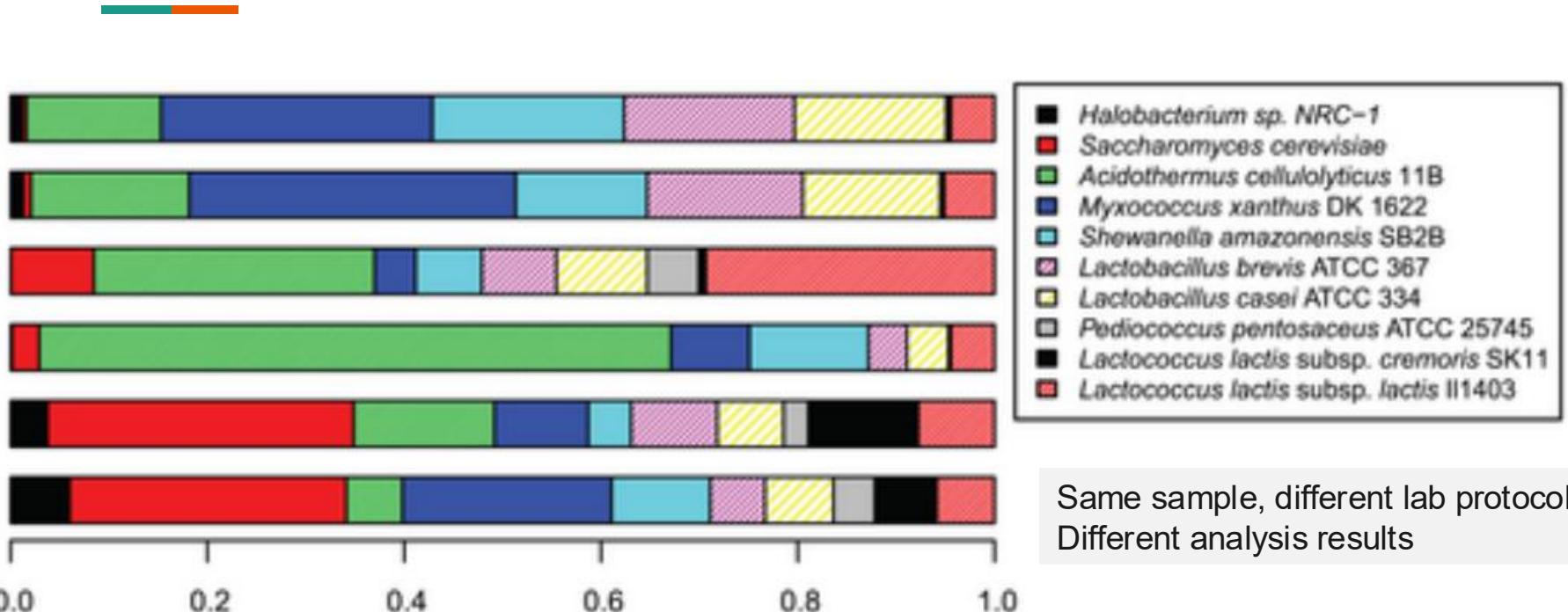
Shotgun analysis pipeline (continued)



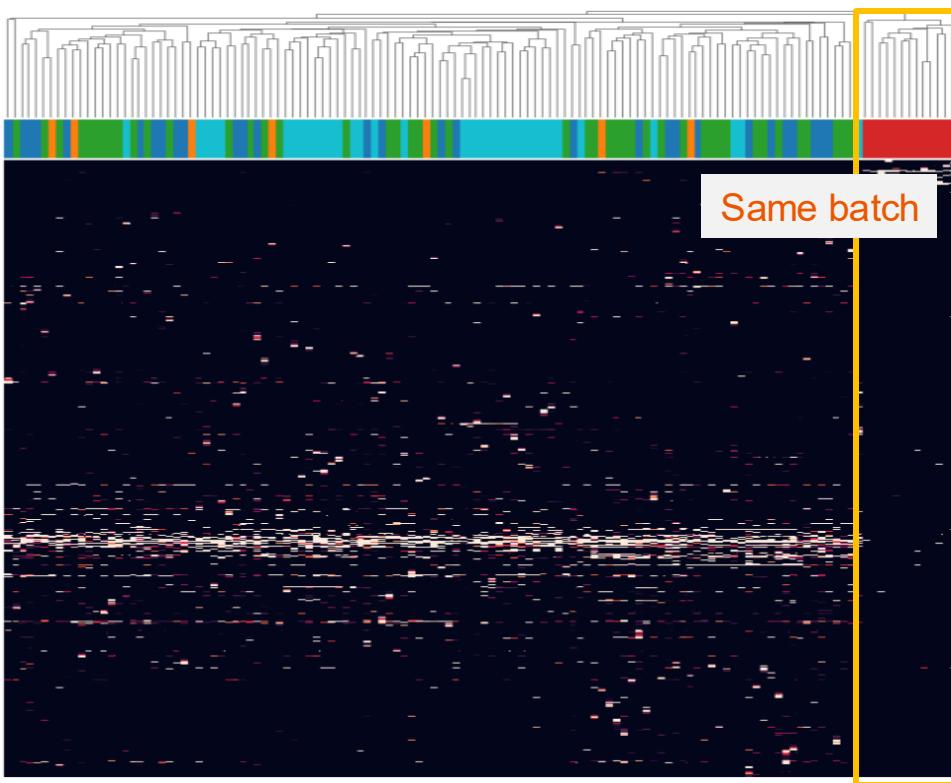


Caution: Variability of microbiome

Microbiome profile influenced by lab protocol



Batch effects



- A common problem
- **Good batch design:** distribute samples of the same class across batches
- Don't group them in the same batch

Any question?

- See you next time