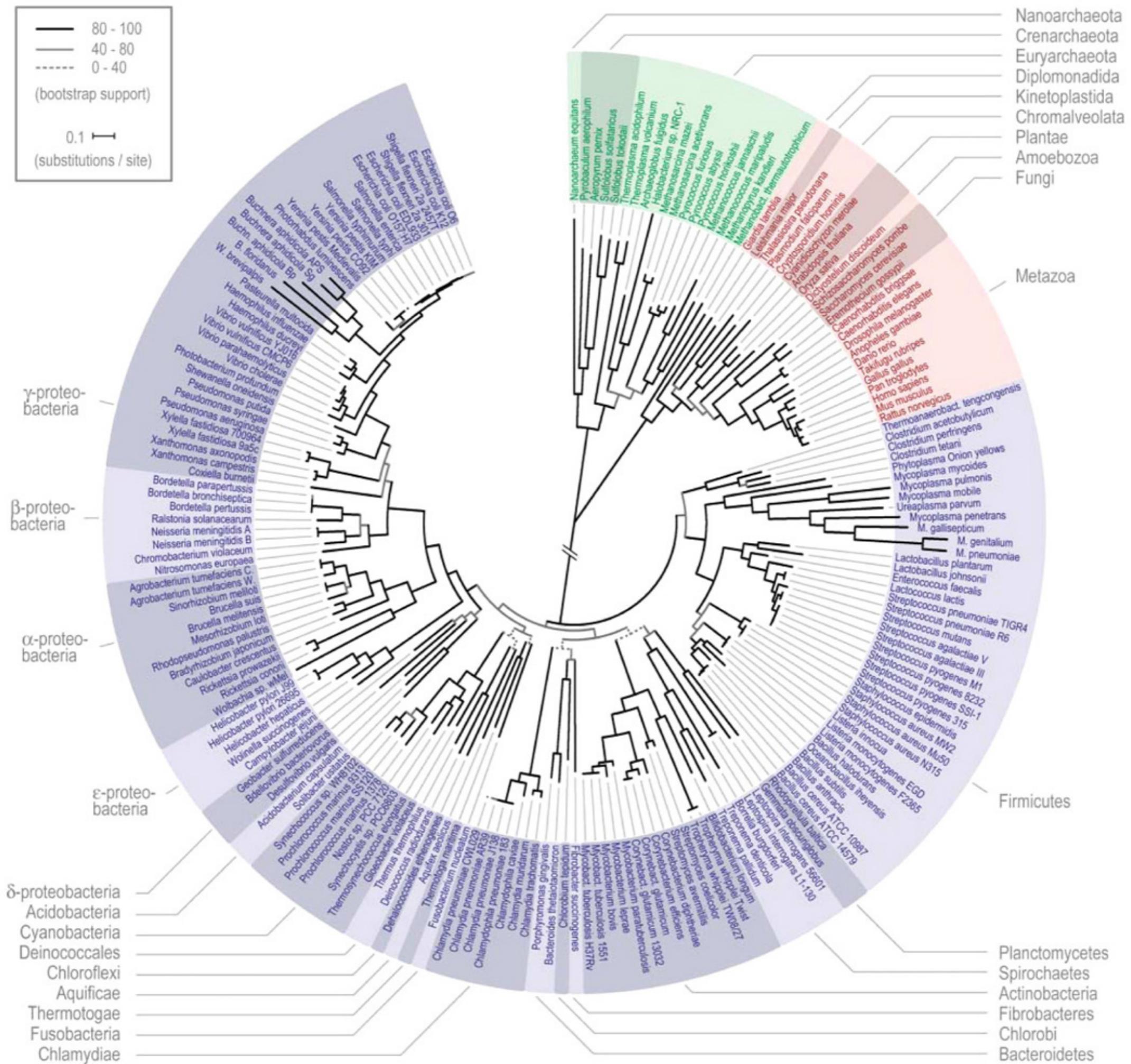


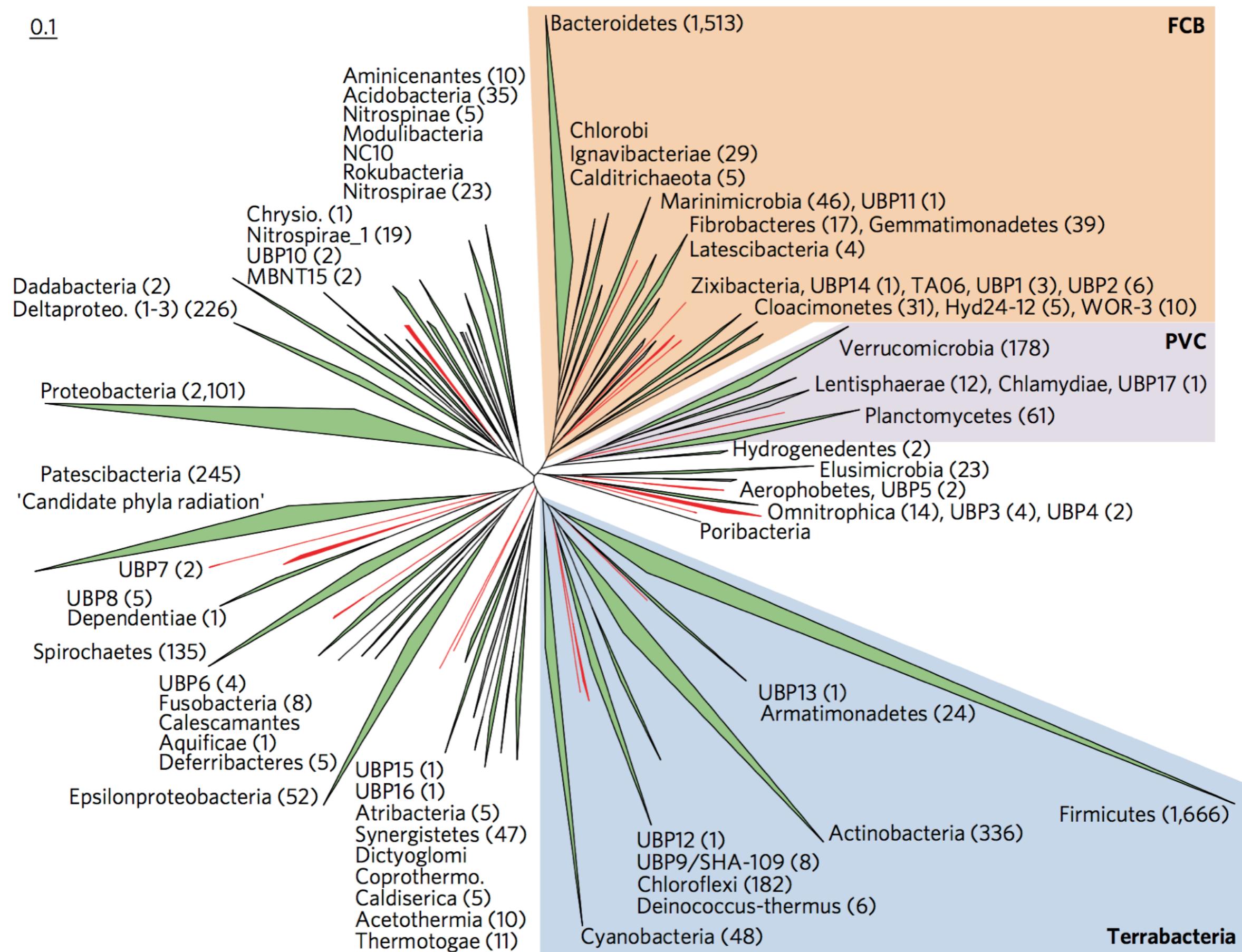
# **Microbiome Module**

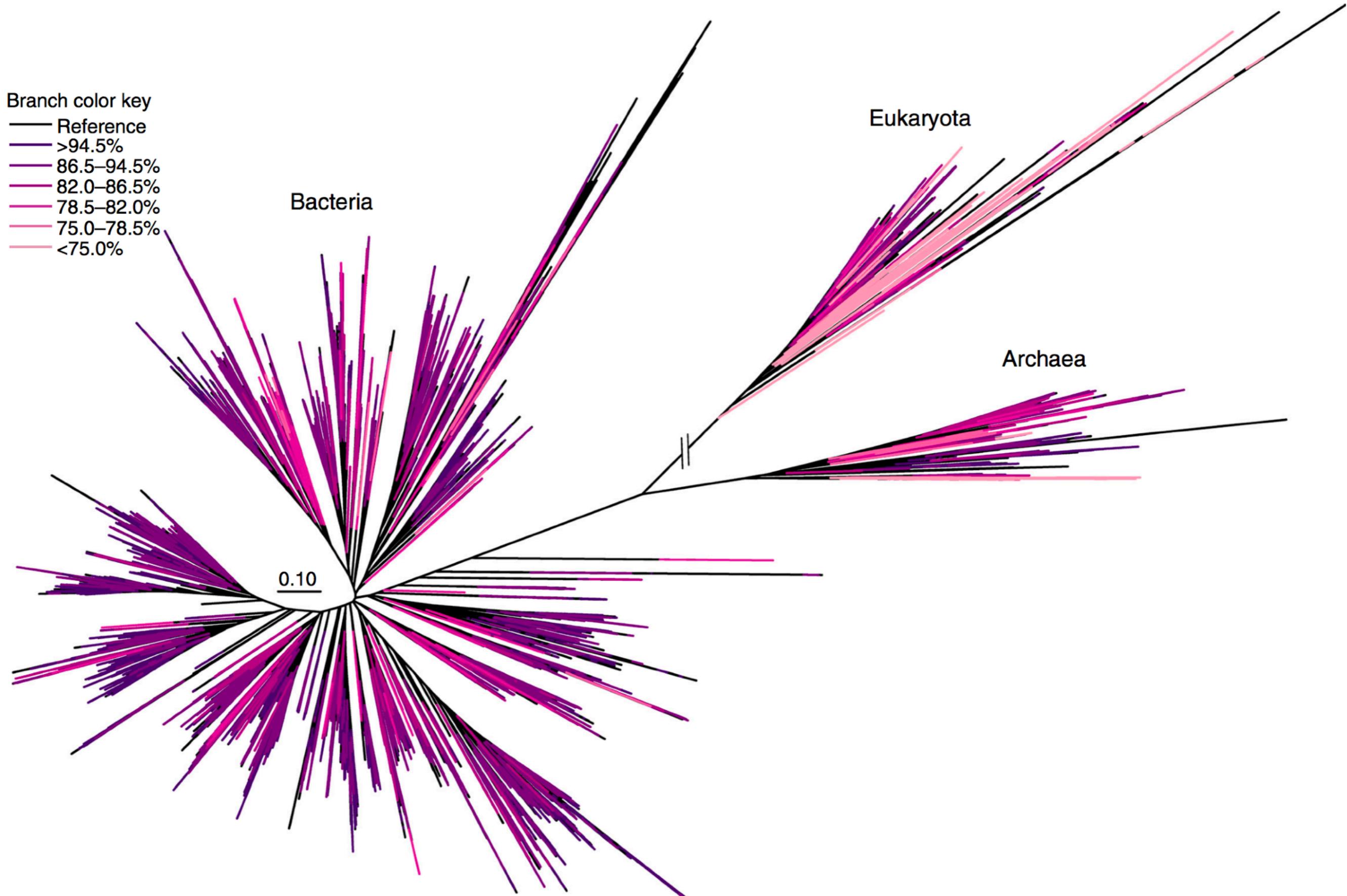
**Instructor: Erik Wright**

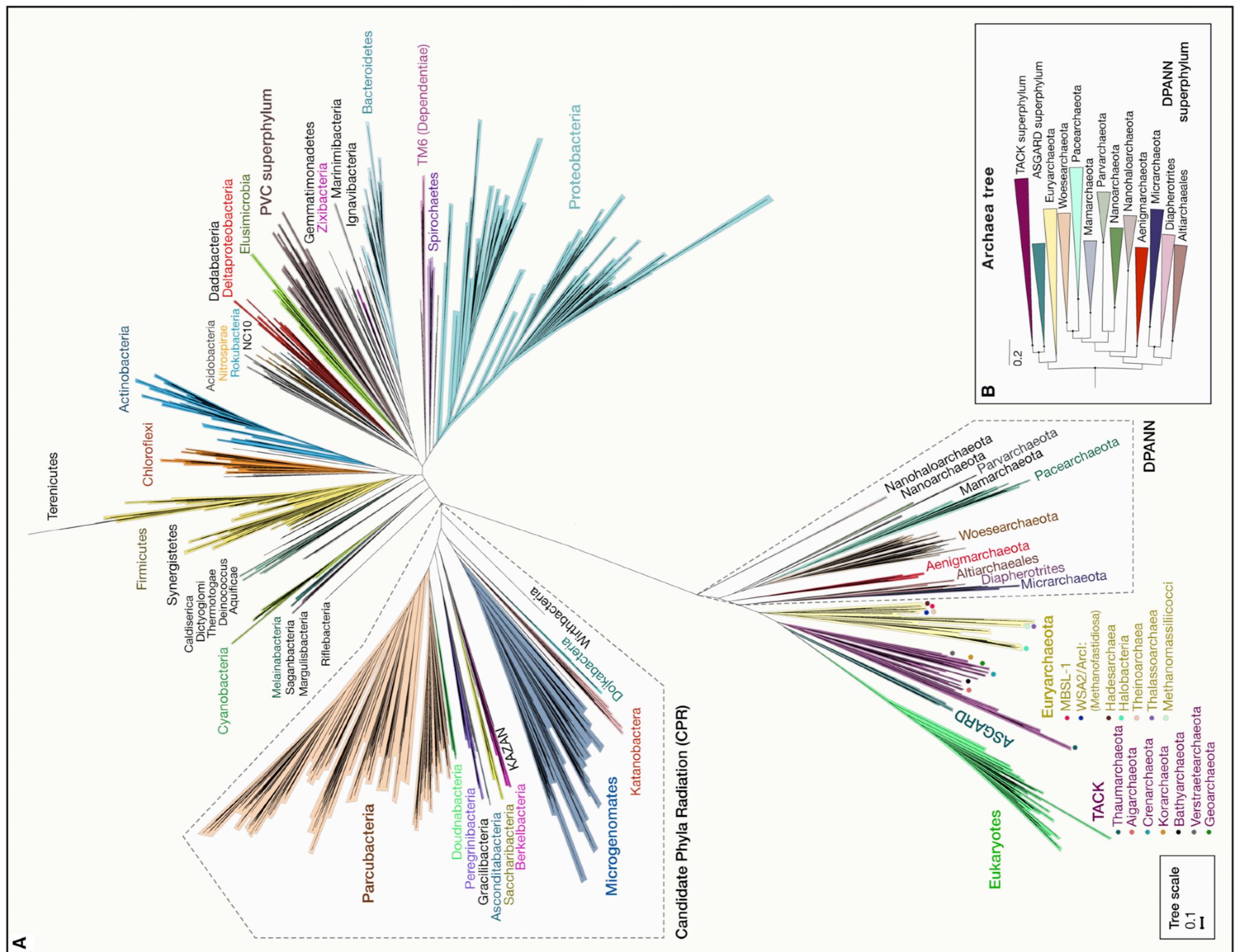
# The microbiome



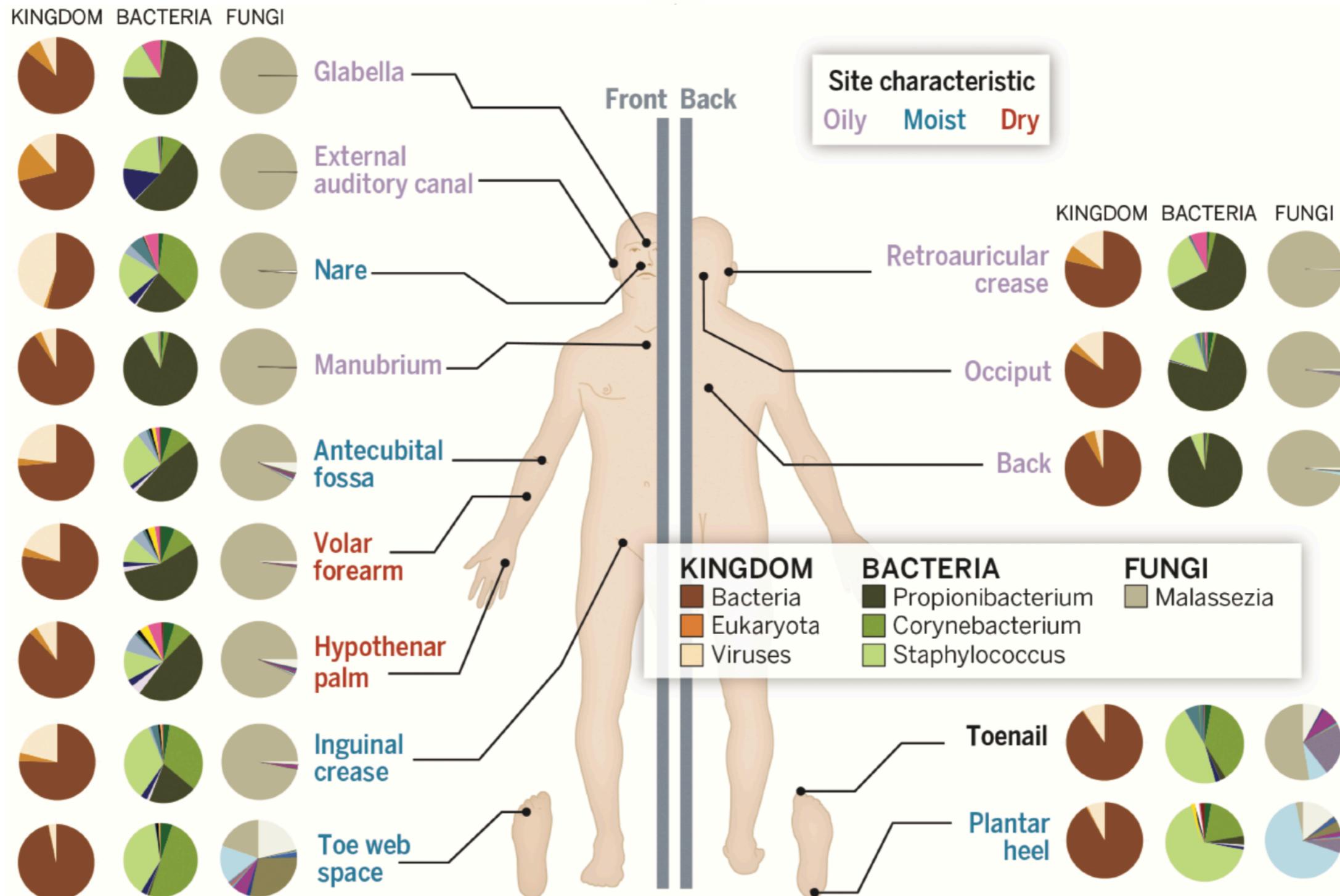
0.1



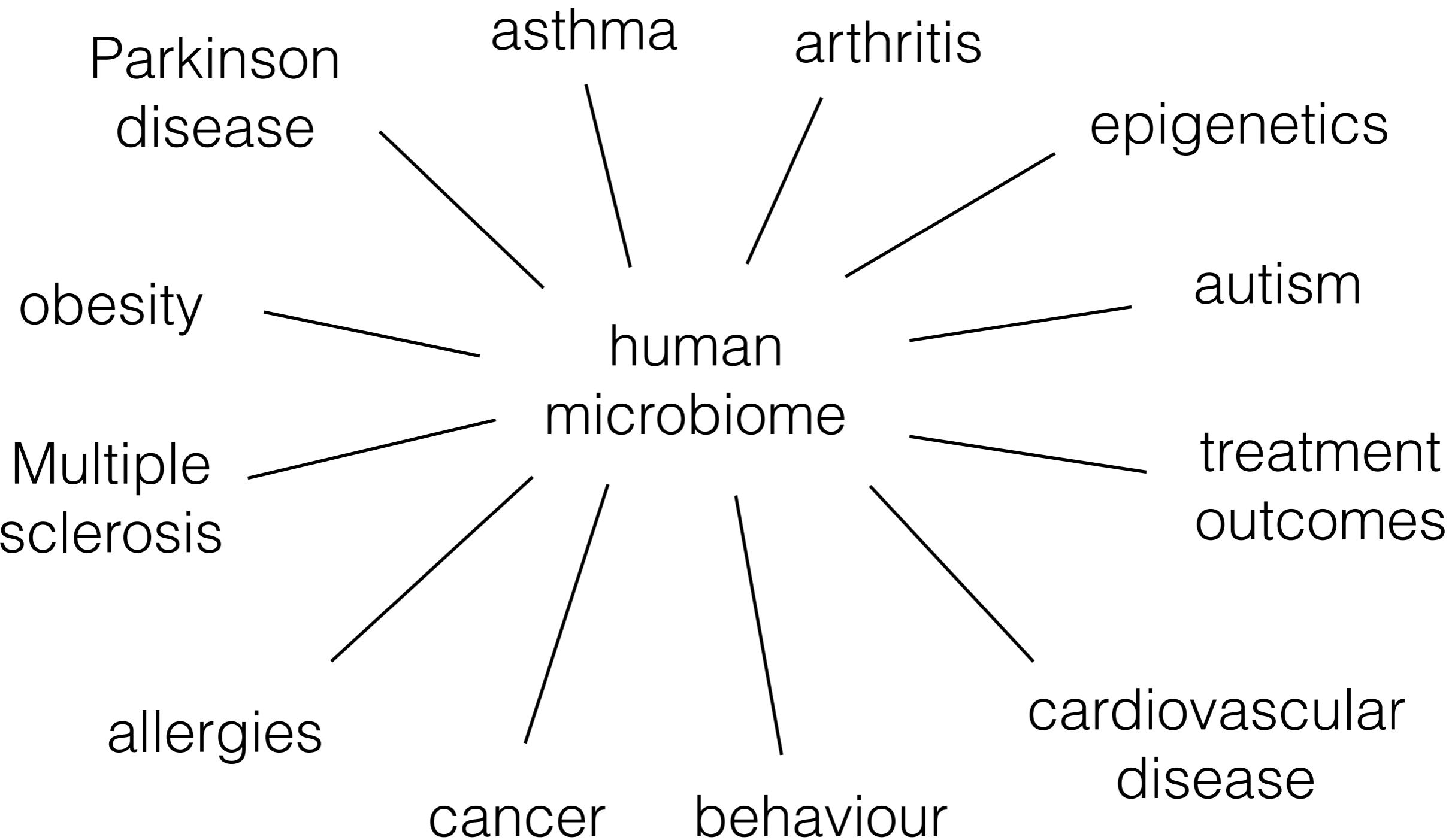




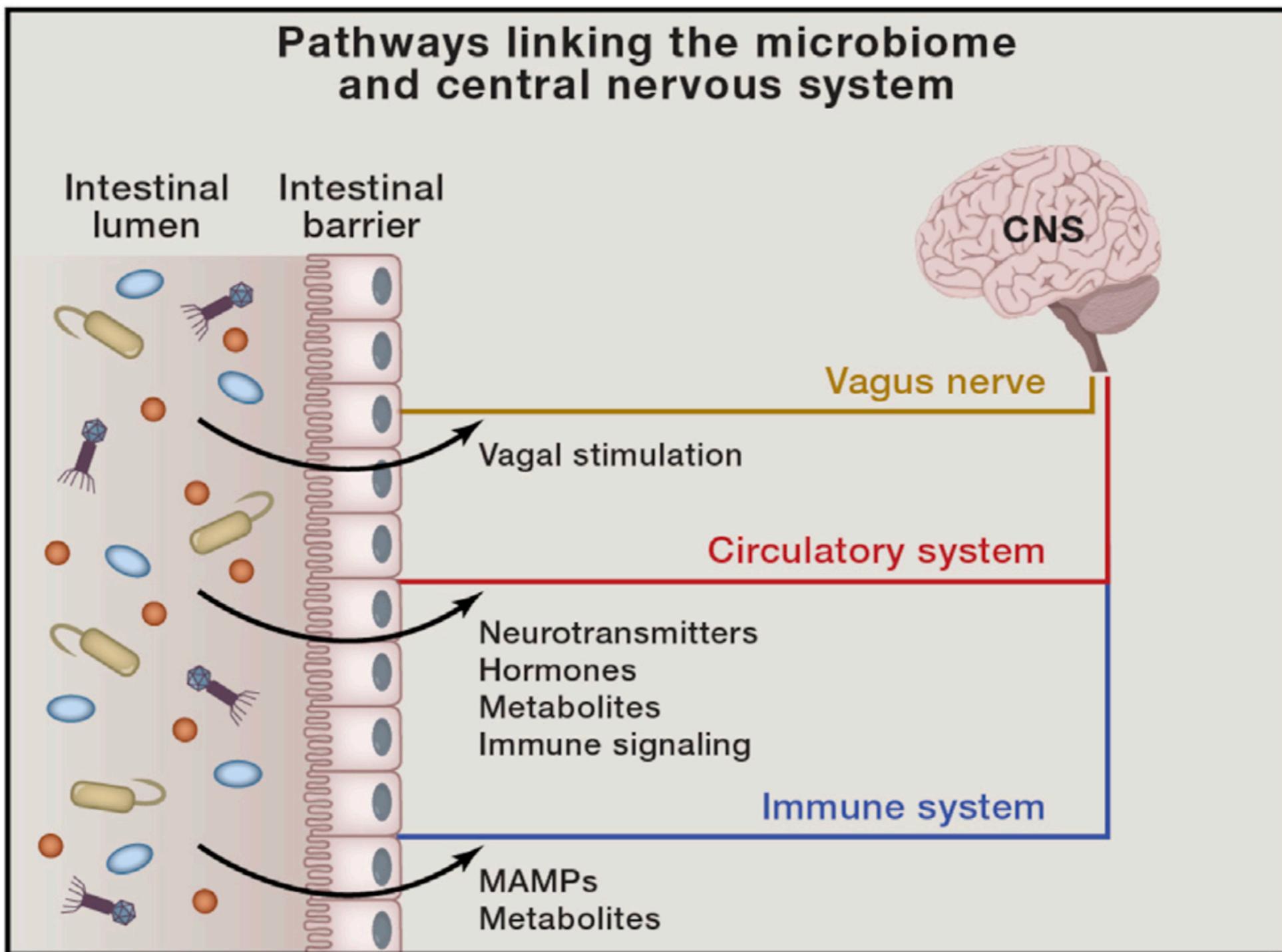
# Microbiome composition is important to health



# Microbiome composition is important to health



# The gut microbiota has many connections



# Manipulations of the gut microbiome

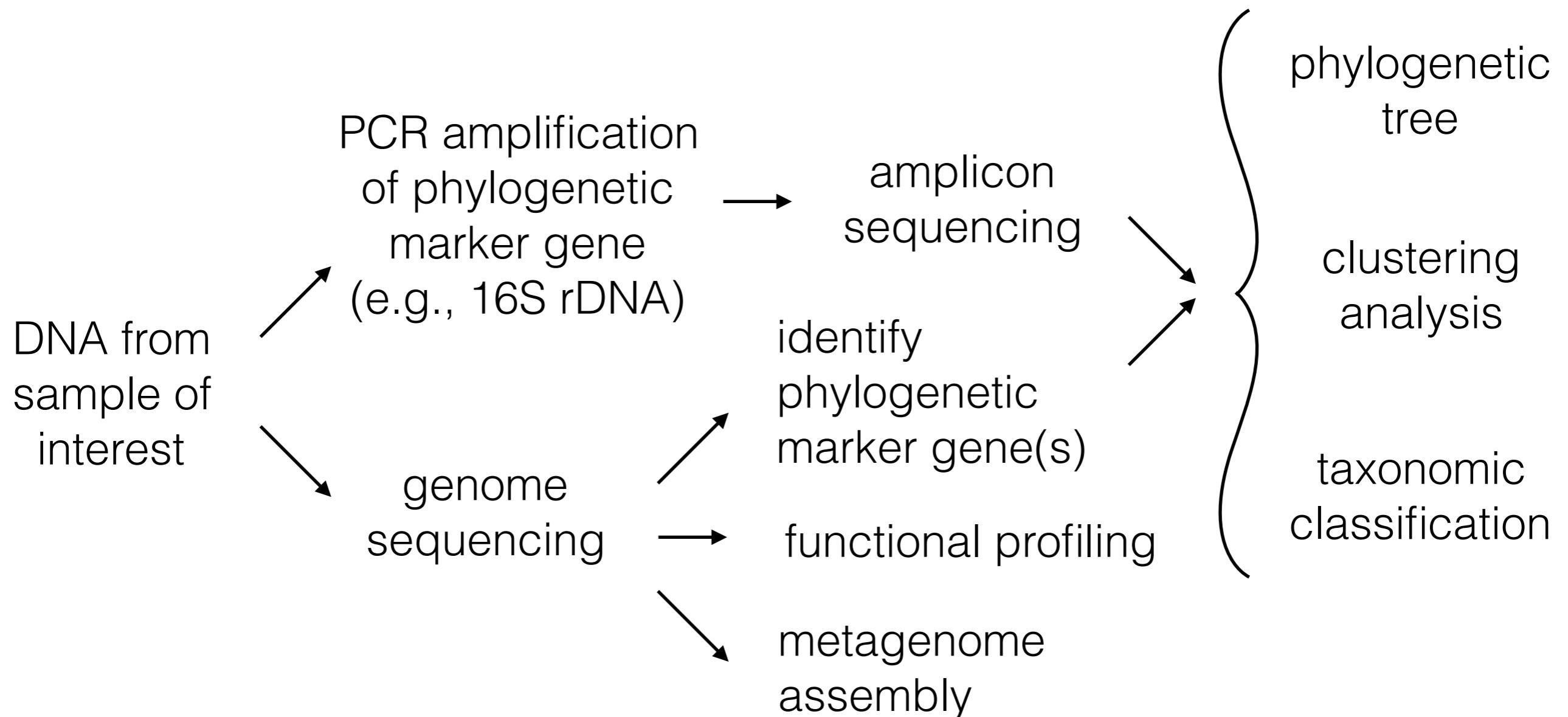


fecal  
transplants

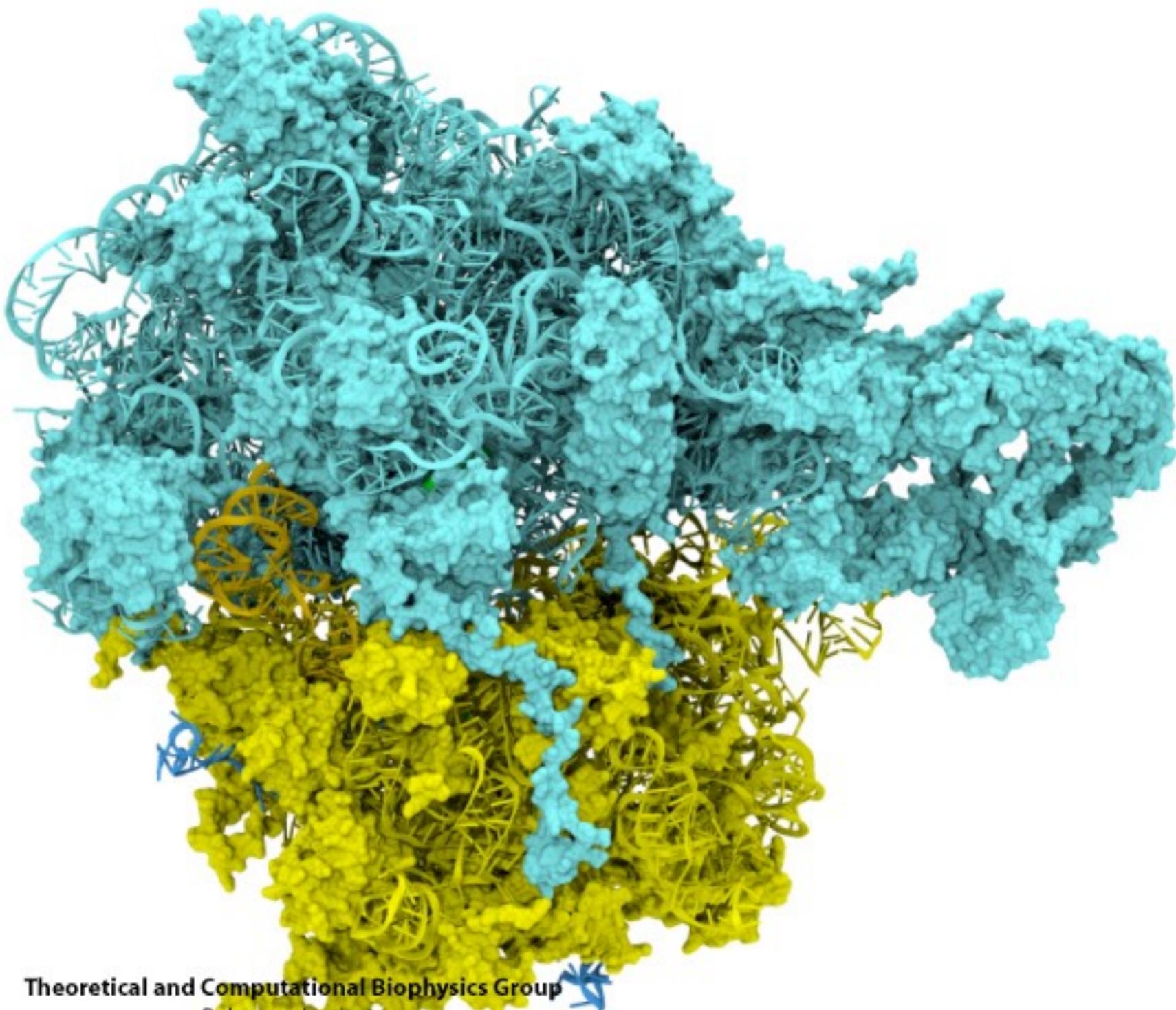


# Popular approaches to studying the microbiome

# Typical microbiome bioinformatic workflows



# Quaternary structure of the ribosome

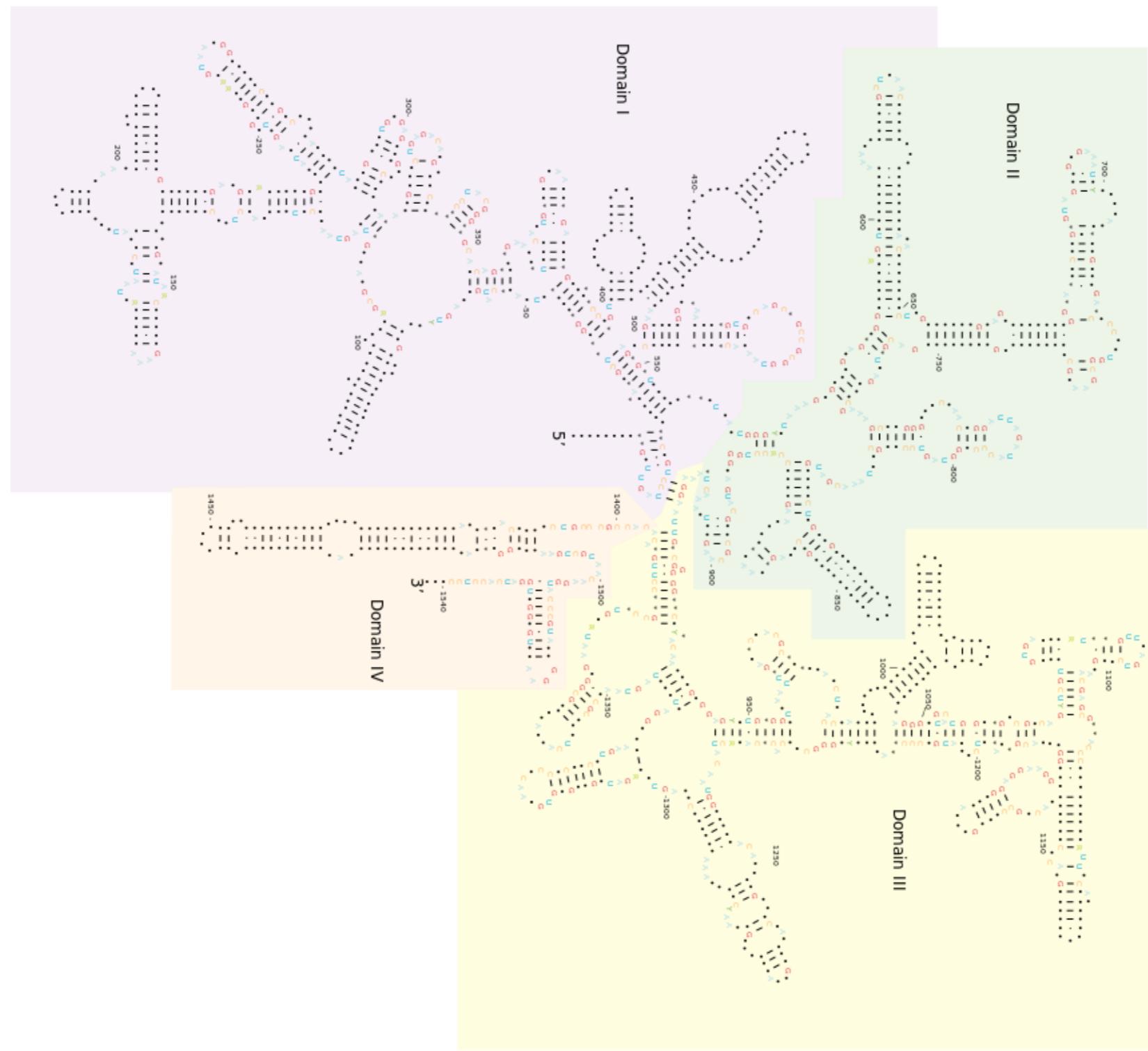


Theoretical and Computational Biophysics Group

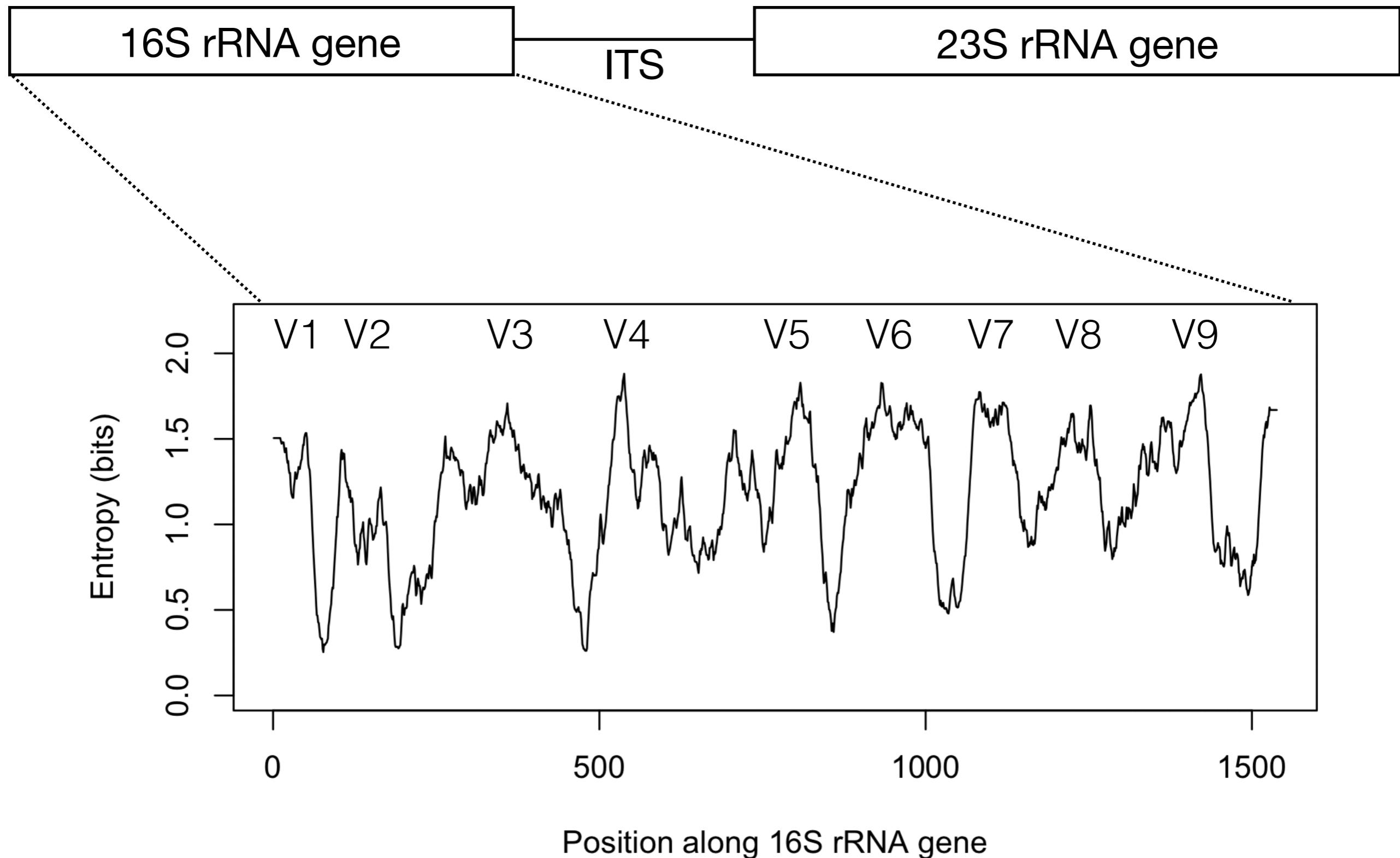
Bekman Institute

University of Illinois at Urbana-Champaign

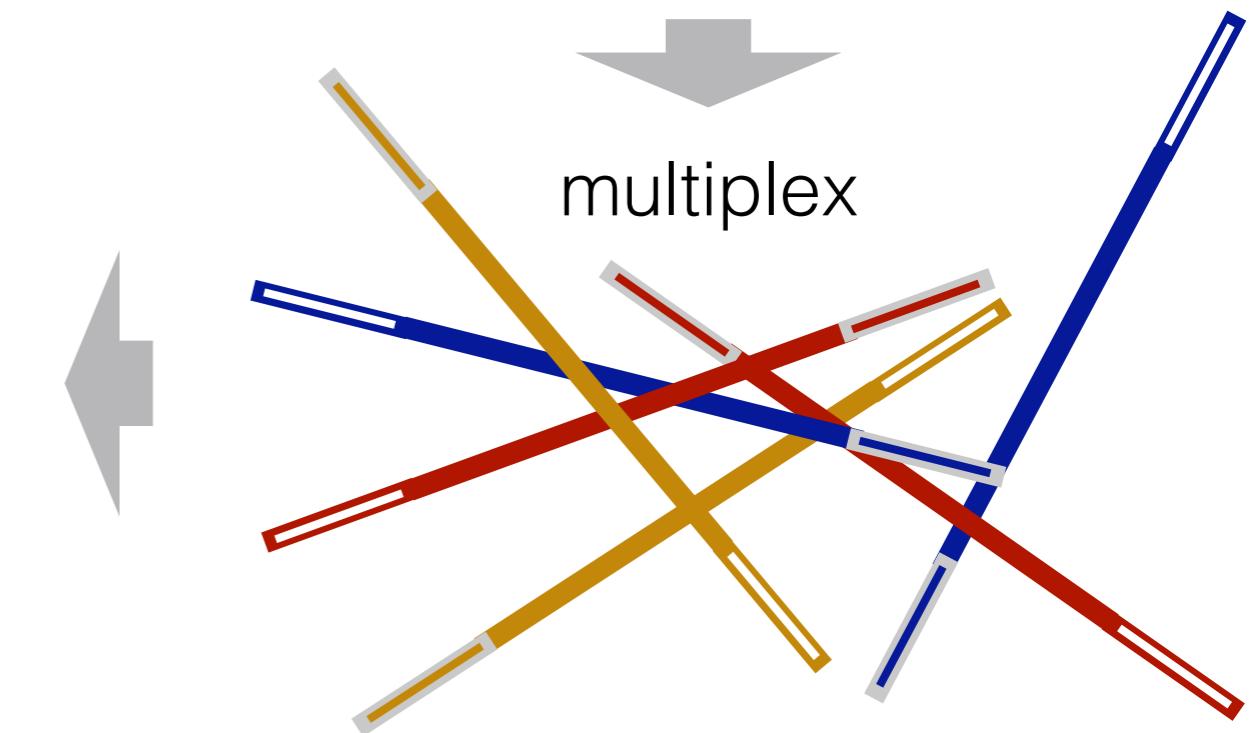
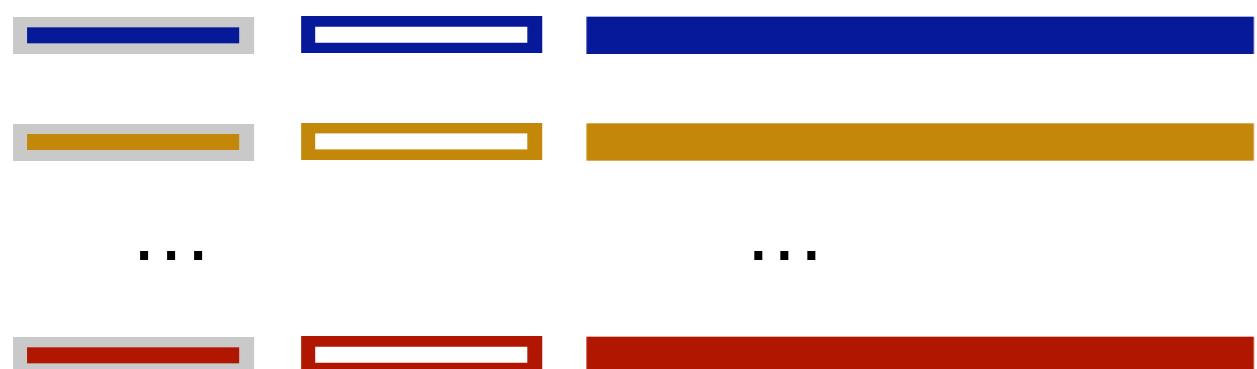
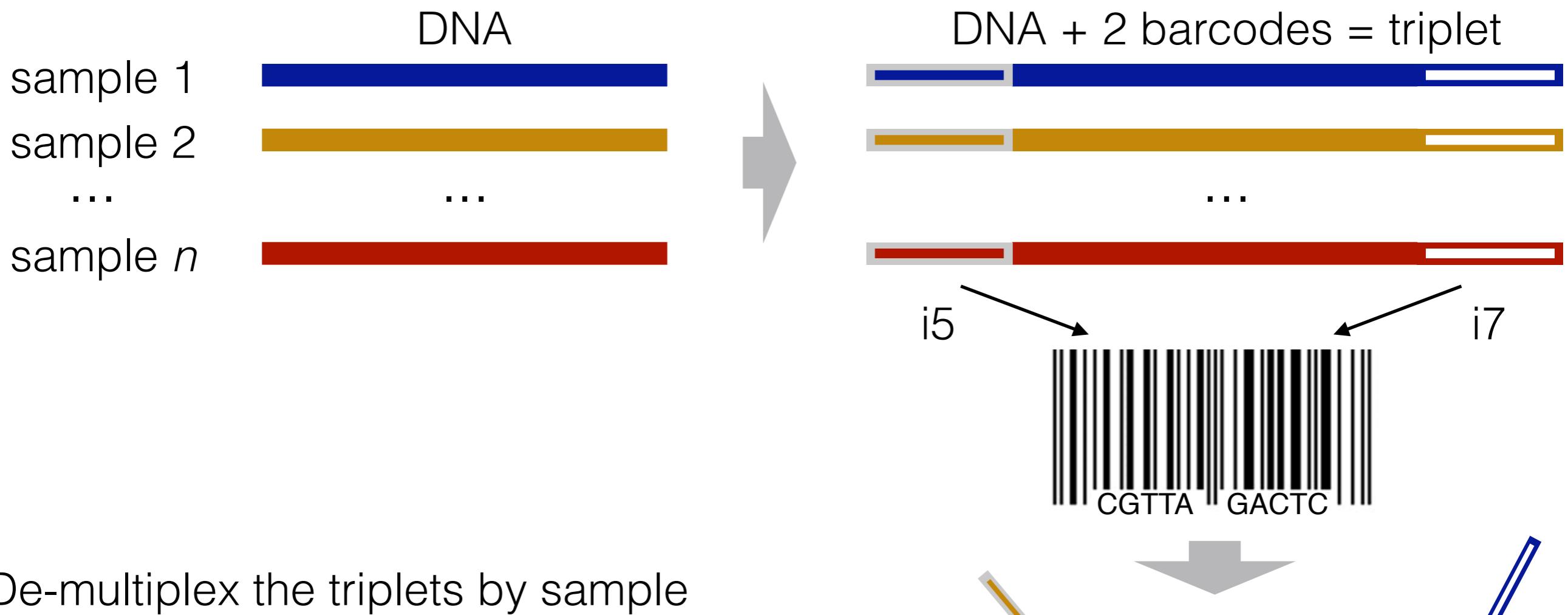
# Secondary structure of the 16S rRNA



# 16S marker gene sequencing



# Multiplexing on the Illumina platform



# Reference-free analysis approaches

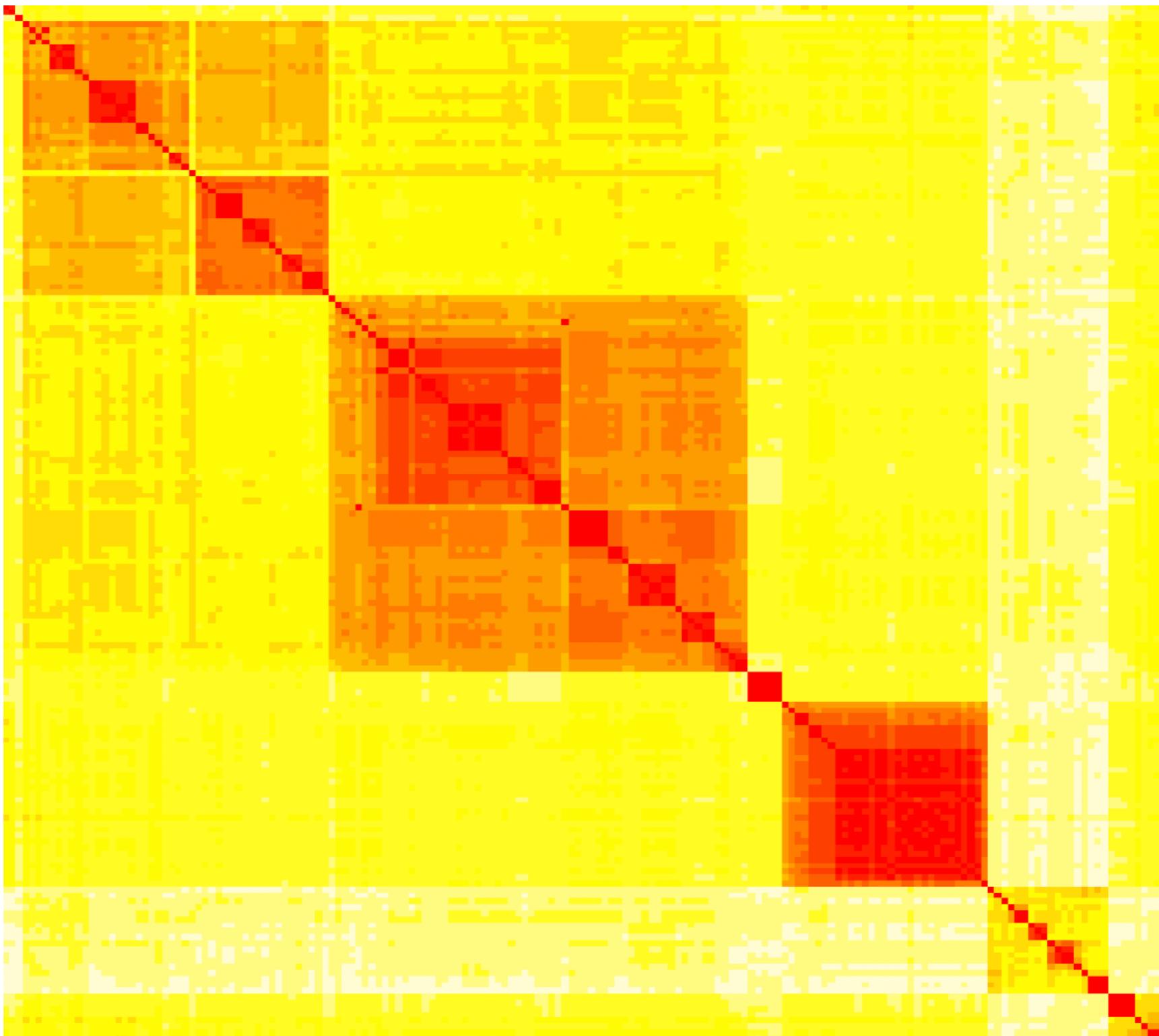
# Constructing/interpreting phylogenetic trees

The figure displays a sequence logo for a genomic region spanning from position 560 to 640. The x-axis marks positions 560, 580, 600, 620, and 640. The y-axis represents the four nucleotides: Adenine (A, green), Thymine (T, blue), Cytosine (C, orange), and Guanine (G, red). Each column of bars corresponds to a position along the sequence, showing the relative frequency of each nucleotide. High conservation is evident at positions 560, 580, 600, and 620, where specific nucleotides (e.g., A at 560, T at 580, C at 600, G at 620) reach near 100% probability. Position 640 shows lower conservation, with all four nucleotides appearing with significant frequency.

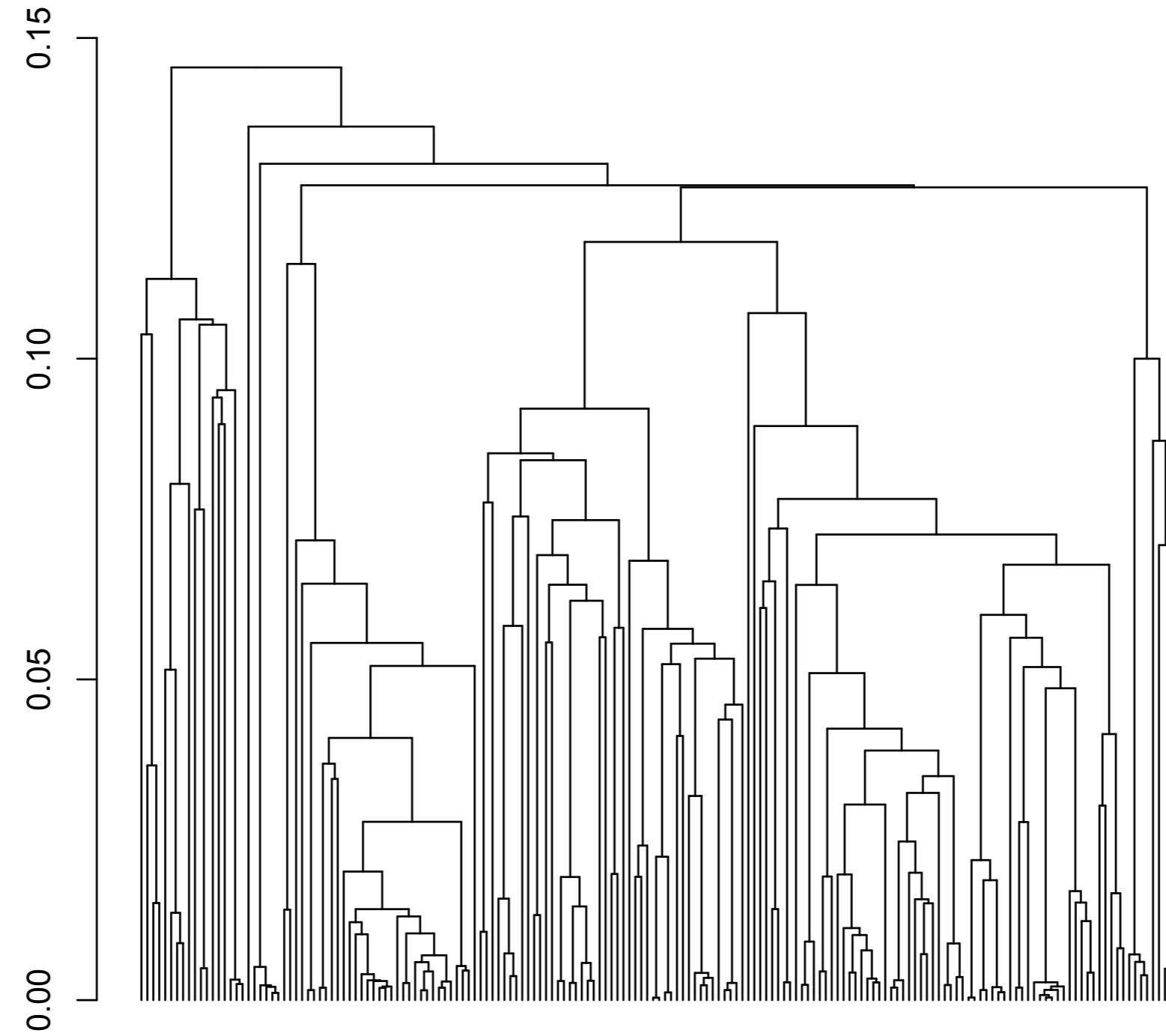
# Constructing/interpreting phylogenetic trees

The figure displays a sequence logo for a set of DNA sequences. The x-axis represents positions 560, 580, 600, 620, and 640. The y-axis shows the frequency of each nucleotide (A, T, C, G) at each position. The colors represent the probability of each nucleotide: A (green), T (red), C (blue), and G (yellow). The height of each bar indicates the relative frequency of that nucleotide at a given position. The sequence logo shows a strong bias at position 560 towards 'C' (blue), position 580 towards 'T' (red), position 600 towards 'G' (yellow), position 620 towards 'T' (red), and position 640 towards 'C' (blue).

# Constructing/interpreting phylogenetic trees



# Constructing/interpreting phylogenetic trees



# Constructing/interpreting phylogenetic trees

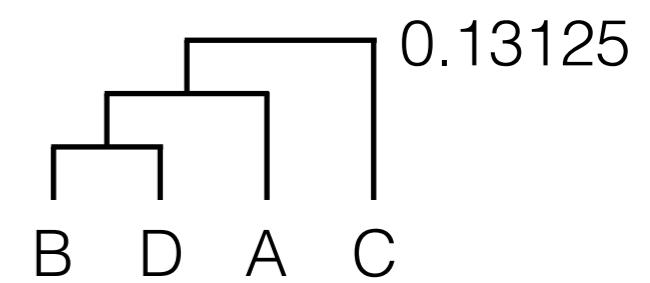
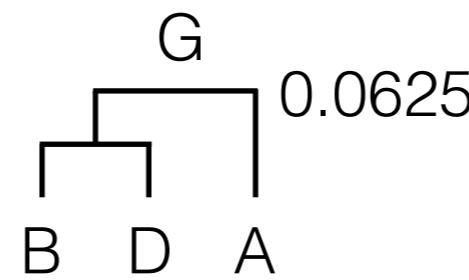
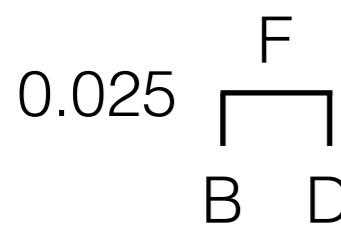
	A	B	C	D
A	0			
B	0.1	0		
C	0.3	0.2	0	
D	0.15	0.05	0.25	0

# Constructing/interpreting phylogenetic trees

	A	B	C	D
A	0			
B	0.1	0		
C	0.3	0.2	0	
D	0.15	<b>0.05</b>	0.25	0

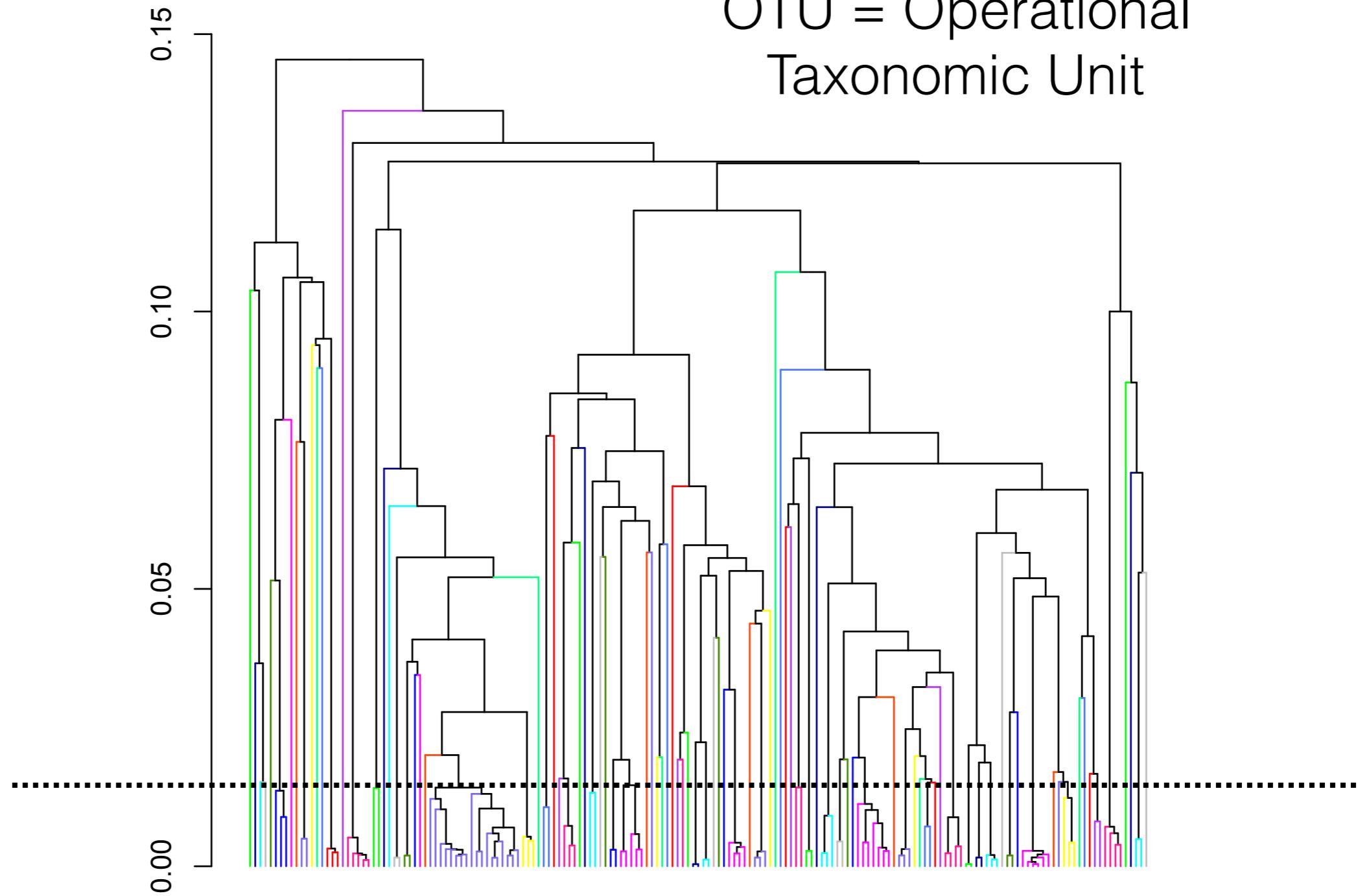
	A	F	C
A	0		
F	<b>0.12</b>	0	
C	0.3	0.23	0

	G	C
G	0	
C	<b>0.26</b>	0



# Clustering into OTUs at $\geq 97\%$ similarity

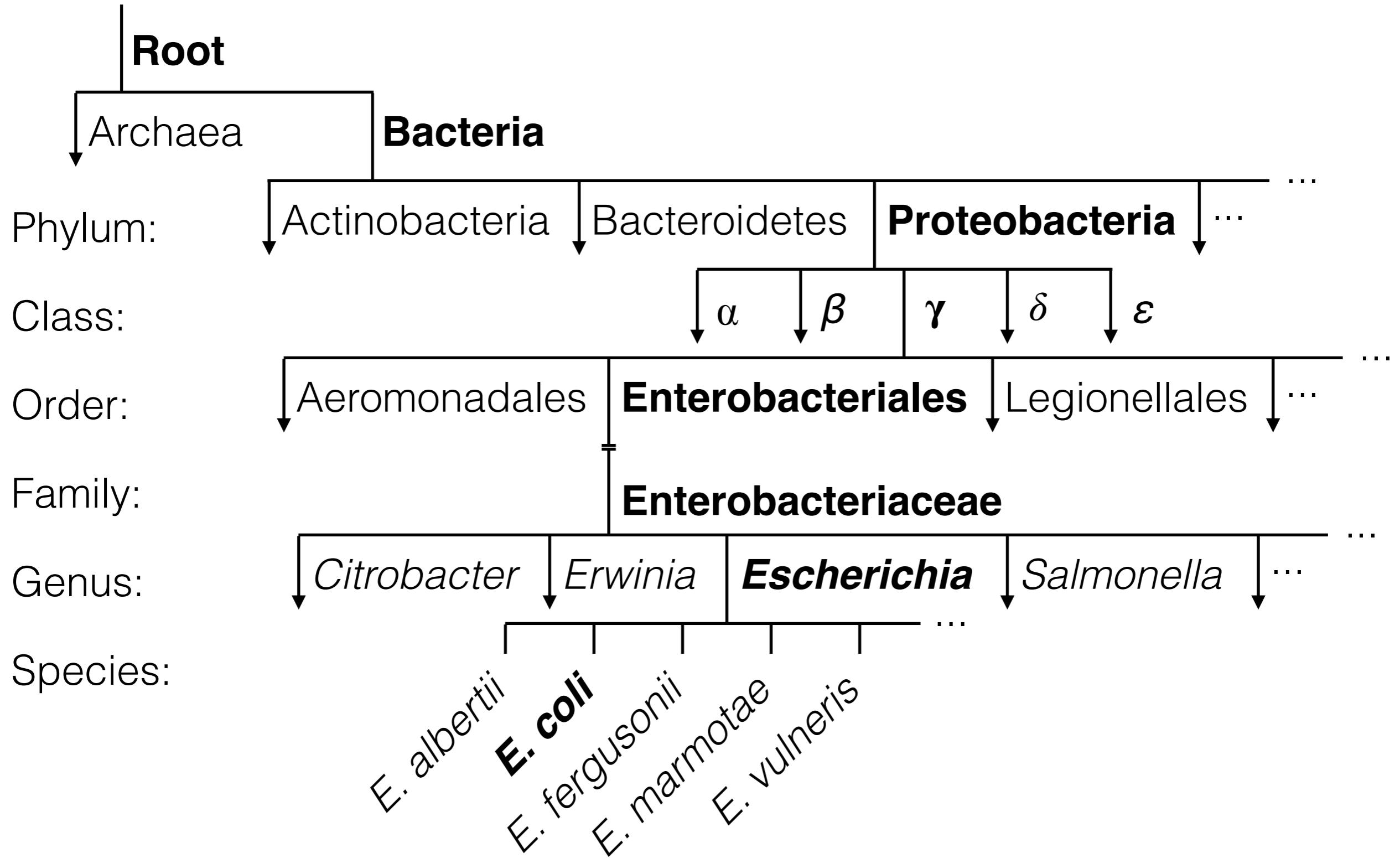
OTU = Operational  
Taxonomic Unit





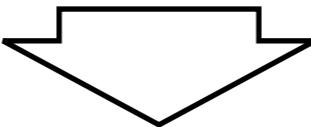
# Reference-based analysis

# Classifying sequences into a taxonomy



# Classifying sequences into a taxonomy

AGCGGGCAGCACAGAGGAACCTTGTTCCTTGGGTGGCGAGCG ...

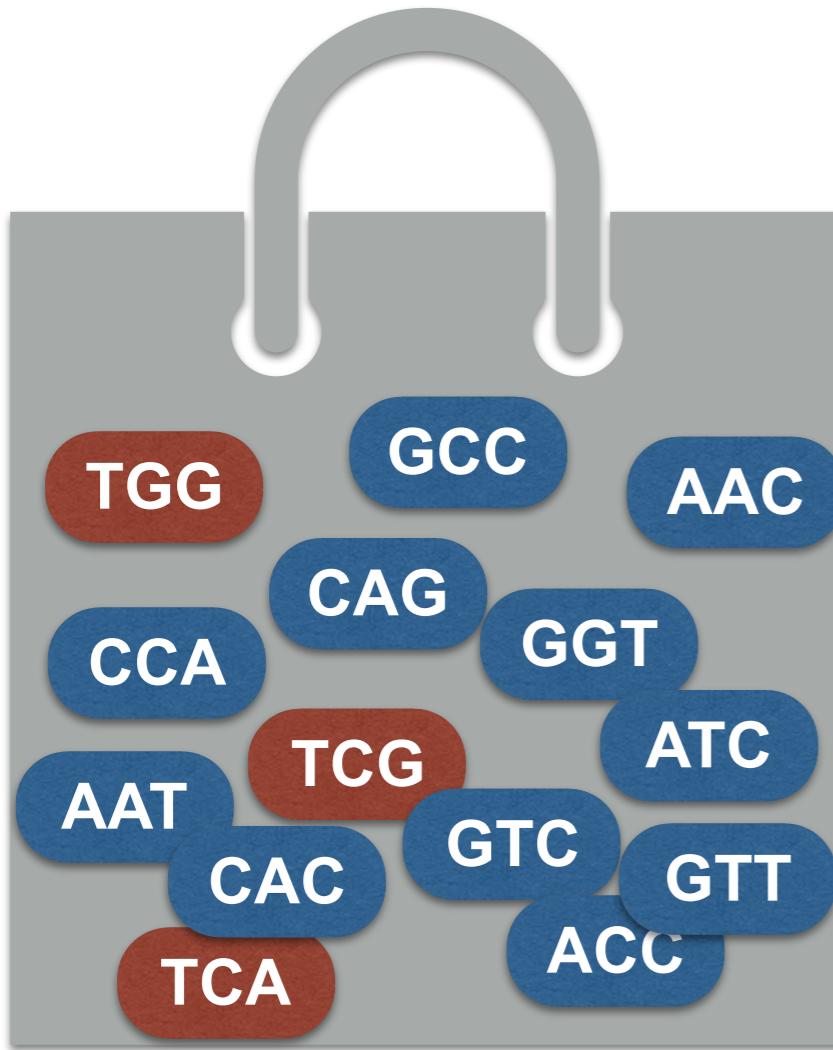


Root (97.8%);  
Bacteria (97.8%);  
Proteobacteria (97.8%);  
Gammaproteobacteria (97.8%);  
Enterobacteriales (97.8%);  
Enterobacteriaceae (97.8%);  
**Escherichia (95%)**

# Classifying via machine learning

CACCGGGTTCAGTCG...

↓  
k-mers ( $K=3$ )

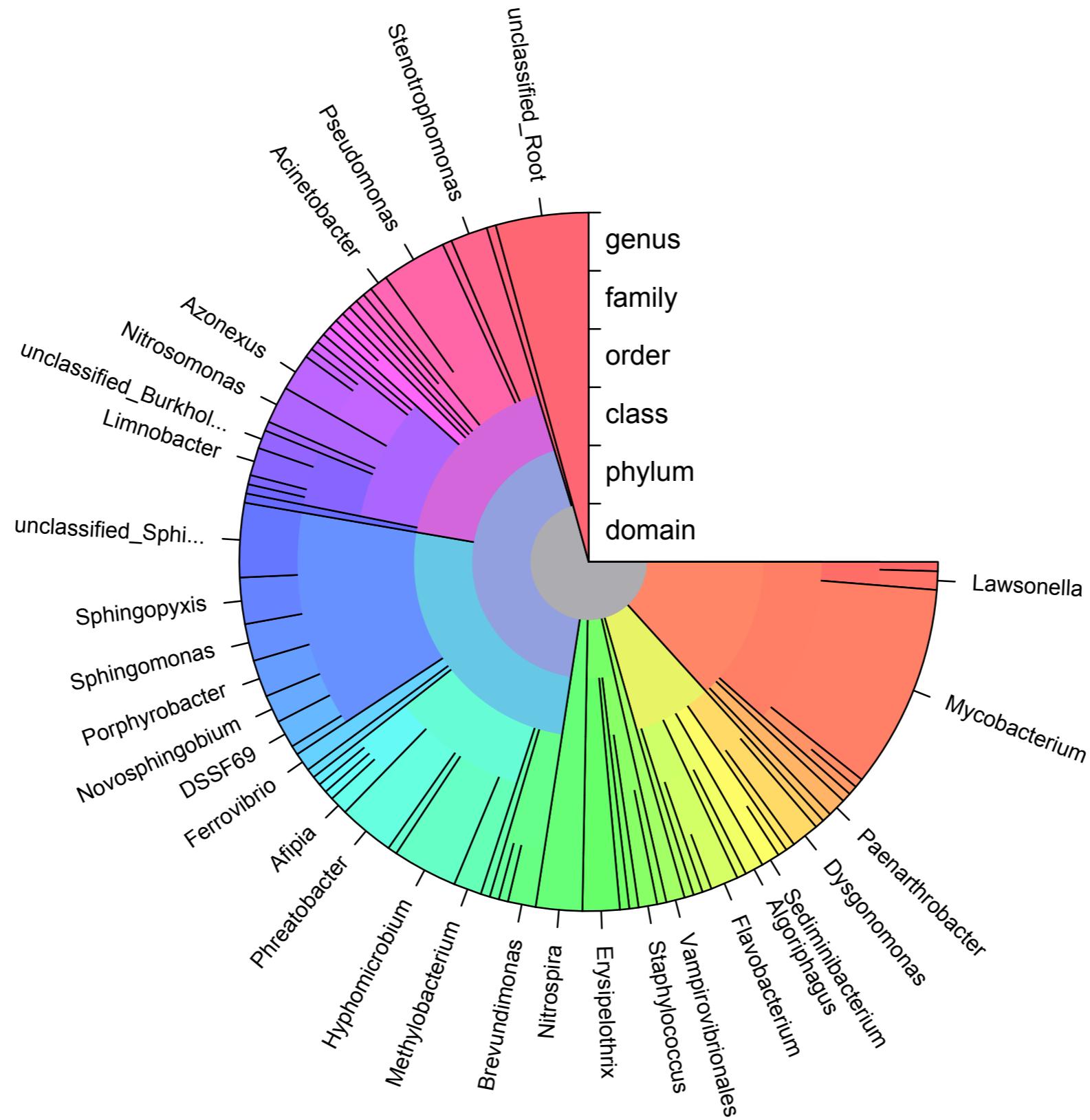


Bootstrap replicates ( $B=100$ )

Subsamples ( $S=3$ )

1	GCC	CAG	GGT	Escherichia
2	CAG	TCG	AAT	Escherichia
3	GCC	GCC	TCA	Escherichia
4	TTT	TTT	TGG	Salmonella
5	CCA	GCC	GTC	Escherichia
6	AAT	TCG	TTT	Salmonella
7	GCC	CAG	AAC	Escherichia
8	AAT	GTC	GGT	Escherichia
...	...	...	...	...
100	CAC	TTT	GGT	Escherichia

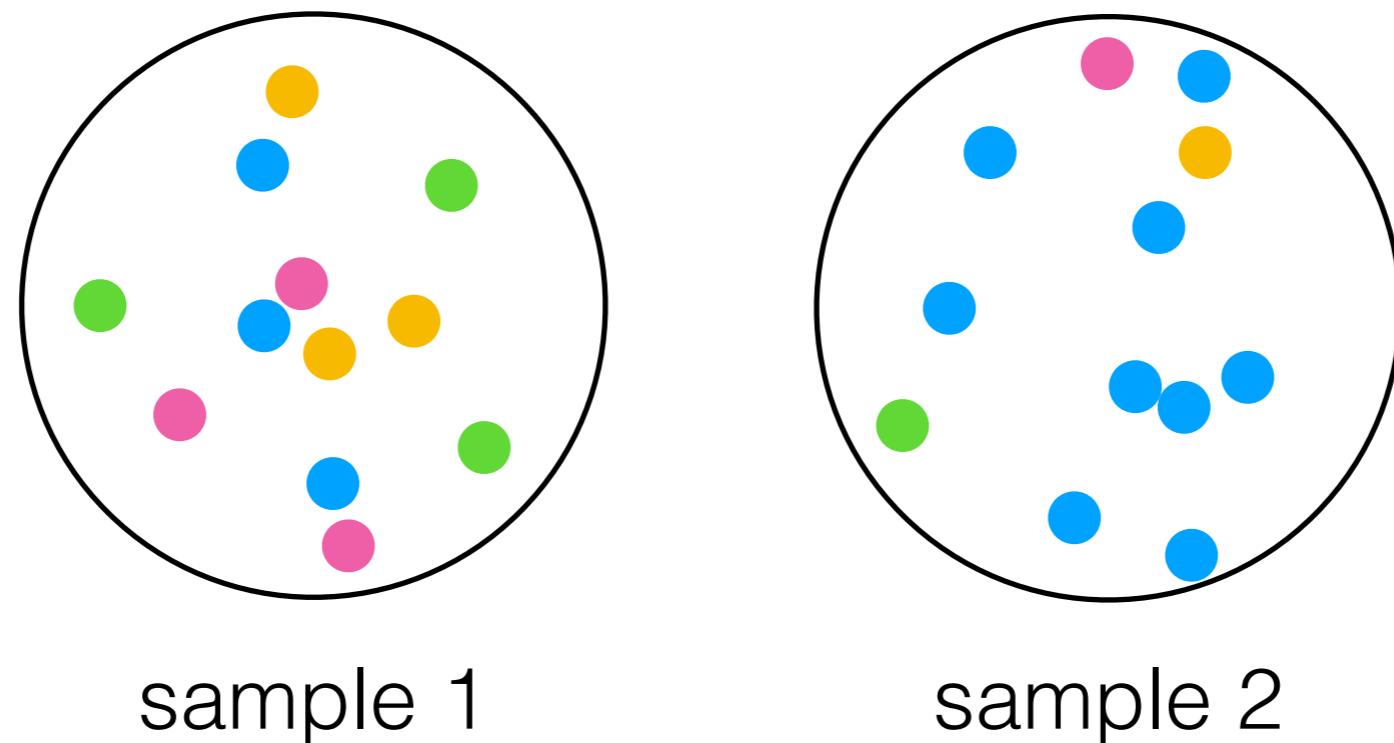
# Visualizing the results of classification



# Analyzing multiple samples

# Measures of diversity within samples

- $\alpha$ -Diversity: within-sample taxonomic diversity
  - Simply count the number of OTUs
  - only measures richness
- Shannon Index =  $-\sum_i(p_i * \ln(p_i))$ ,  $p_i$  = proportion of each OTU
  - measures richness and evenness
- Simpson Index =  $\sum_i(p_i^2)$ ,  $p_i$  = proportion of each OTU
  - measures richness and evenness



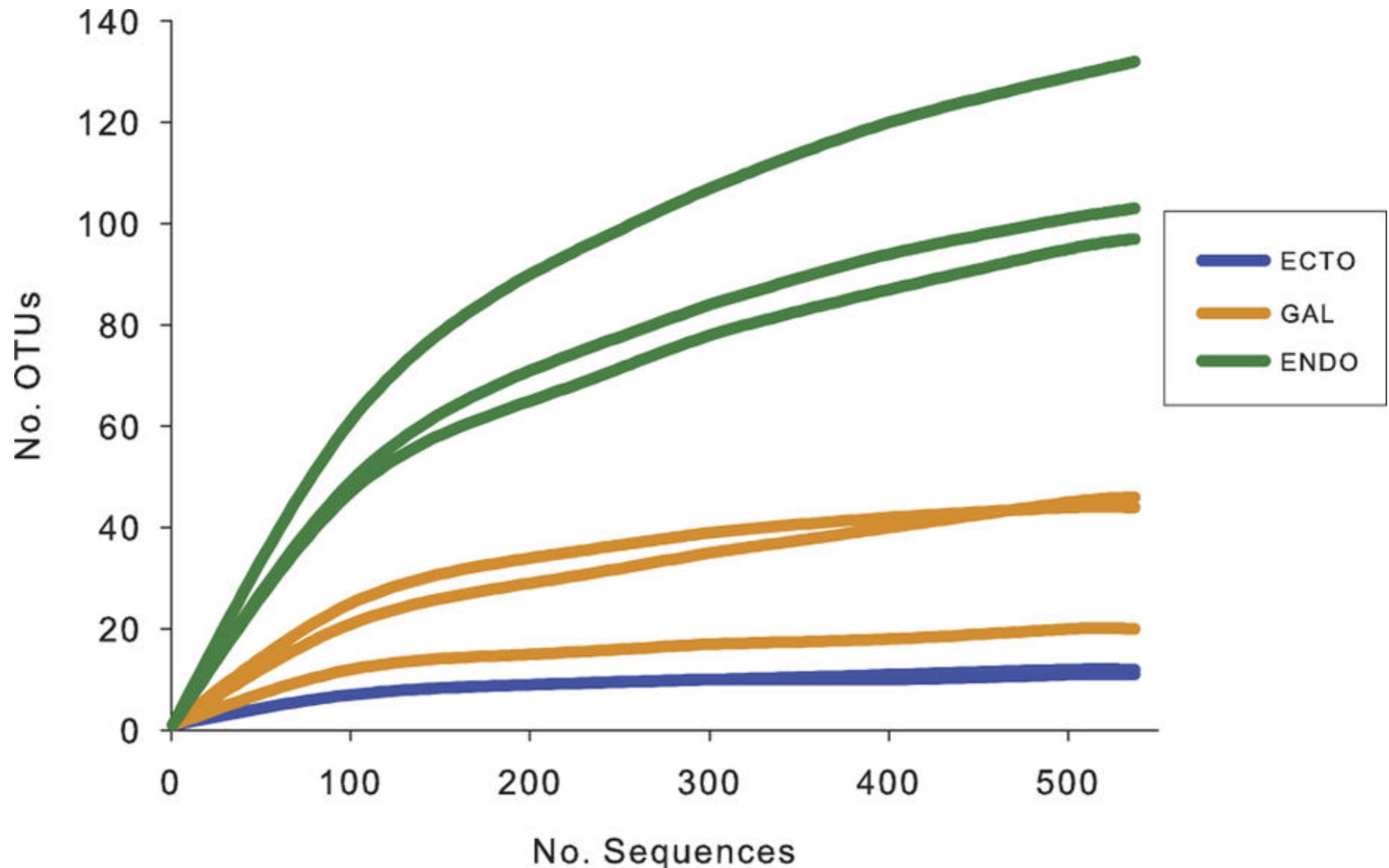
# Measures of diversity within samples

	Sample 1	Sample 2	Sample 3
OTU1	0.0	0.1	0.1
OTU2	0.4	0.2	0.0
OTU3	0.2	0.1	0.8
OTU4	0.0	0.4	0.0
OTU5	0.2	0.2	0.1
$\alpha$ -Diversity			
Shannon Index			
Simpson Index			

# Measures of diversity within samples

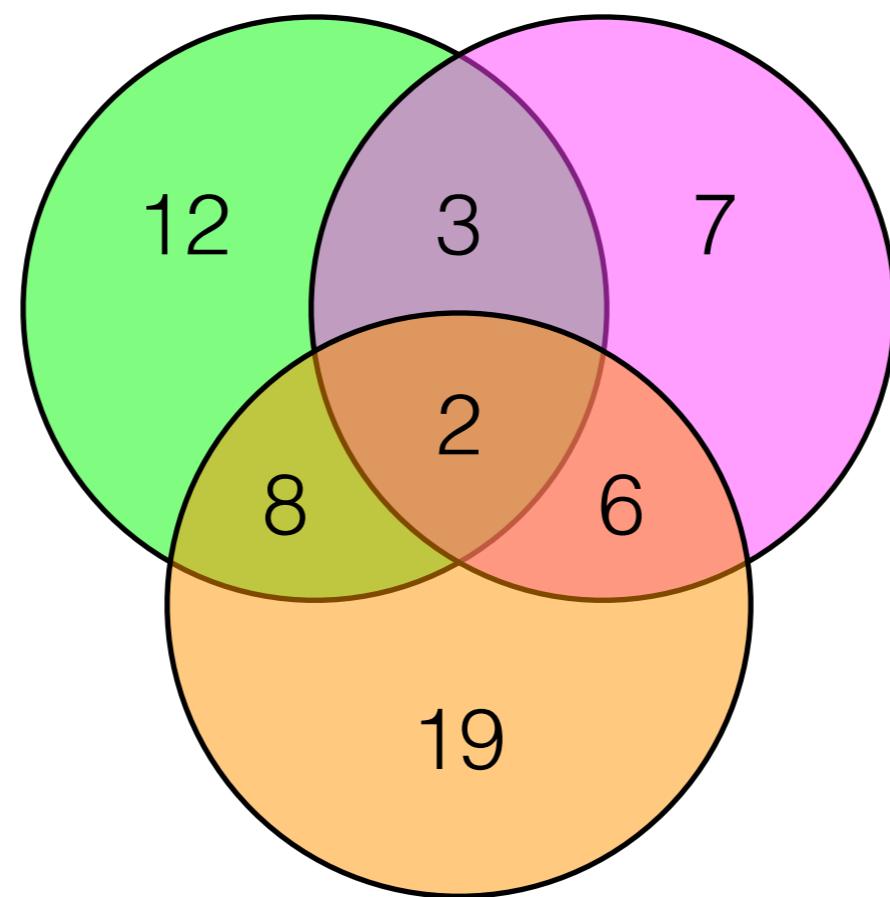
	Sample 1	Sample 2	Sample 3
OTU1	0.0	0.1	0.1
OTU2	0.4	0.2	0.0
OTU3	0.2	0.1	0.8
OTU4	0.0	0.4	0.0
OTU5	0.2	0.2	0.1
$\alpha$ -Diversity	3	5	3
Shannon Index	1.01	1.47	0.64
Simpson Index	0.24	0.26	0.66

# Rarefaction curves to estimate total diversity



# Measures of diversity among samples

- $\beta$ -Diversity: between-sample taxonomic diversity
  - Simply look at OTU overlap between samples
  - Bray-Curtis Dissimilarity =  $1 - 2*C/(S_i + S_j)$ 
    - $S_i$  &  $S_j$  are the sum of OTUs at sites  $i$  and  $j$
    - $C$  is the number of OTUs in common between  $S_i$  &  $S_j$
  - UniFrac = similar idea but take into account phylogenetic breadth



# Measures of diversity among samples

	Sample 1	Sample 2	Sample 3
OTU1	0.0	0.1	0.1
OTU2	0.4	0.2	0.0
OTU3	0.2	0.1	0.8
OTU4	0.0	0.4	0.0
OTU5	0.2	0.2	0.1

Bray-Curtis  
Dissimilarity

	S1	S2	S3
S1	0		
S2		0	
S3			0

# Measures of diversity among samples

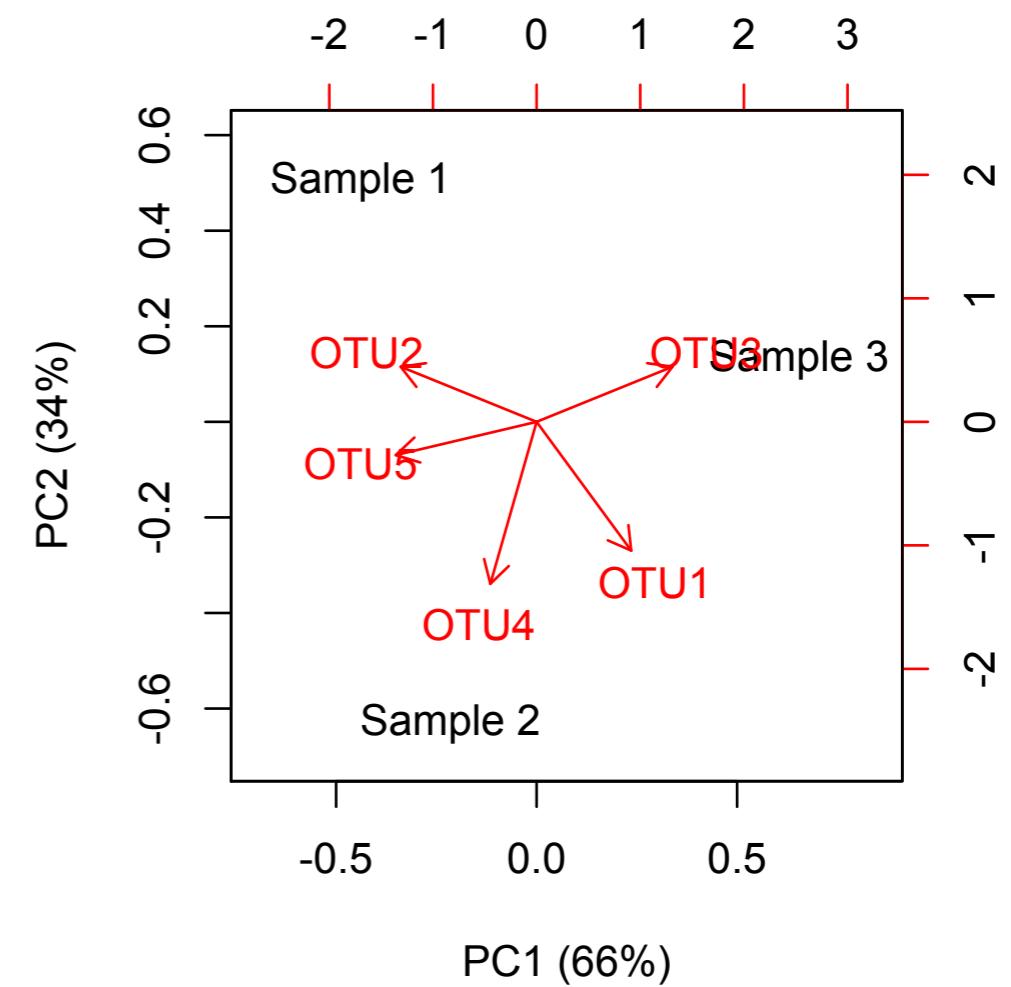
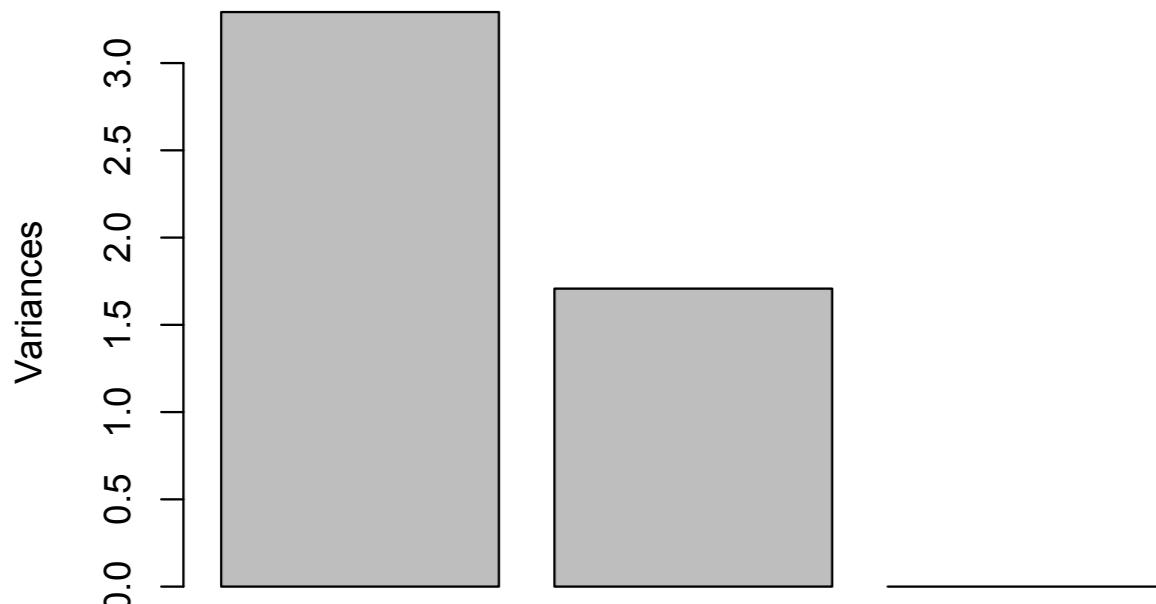
	Sample 1	Sample 2	Sample 3
OTU1	0.0	0.1	0.1
OTU2	0.4	0.2	0.0
OTU3	0.2	0.1	0.8
OTU4	0.0	0.4	0.0
OTU5	0.2	0.2	0.1

Bray-Curtis  
Dissimilarity

	S1	S2	S3
S1	0		
S2	0.25	0	
S3	0.33	0.25	0

# Principle components analysis

	<b>PC1</b>	<b>PC2</b>	<b>PC3</b>
<b>OTU1</b>	0.3630540	-0.5757433	0.19782508
<b>OTU2</b>	-0.5217308	0.2465473	0.72185580
<b>OTU3</b>	0.5217308	0.2465473	0.01140291
<b>OTU4</b>	-0.1775562	-0.7244540	0.22067749
<b>OTU5</b>	-0.5406103	-0.1487108	-0.62526863



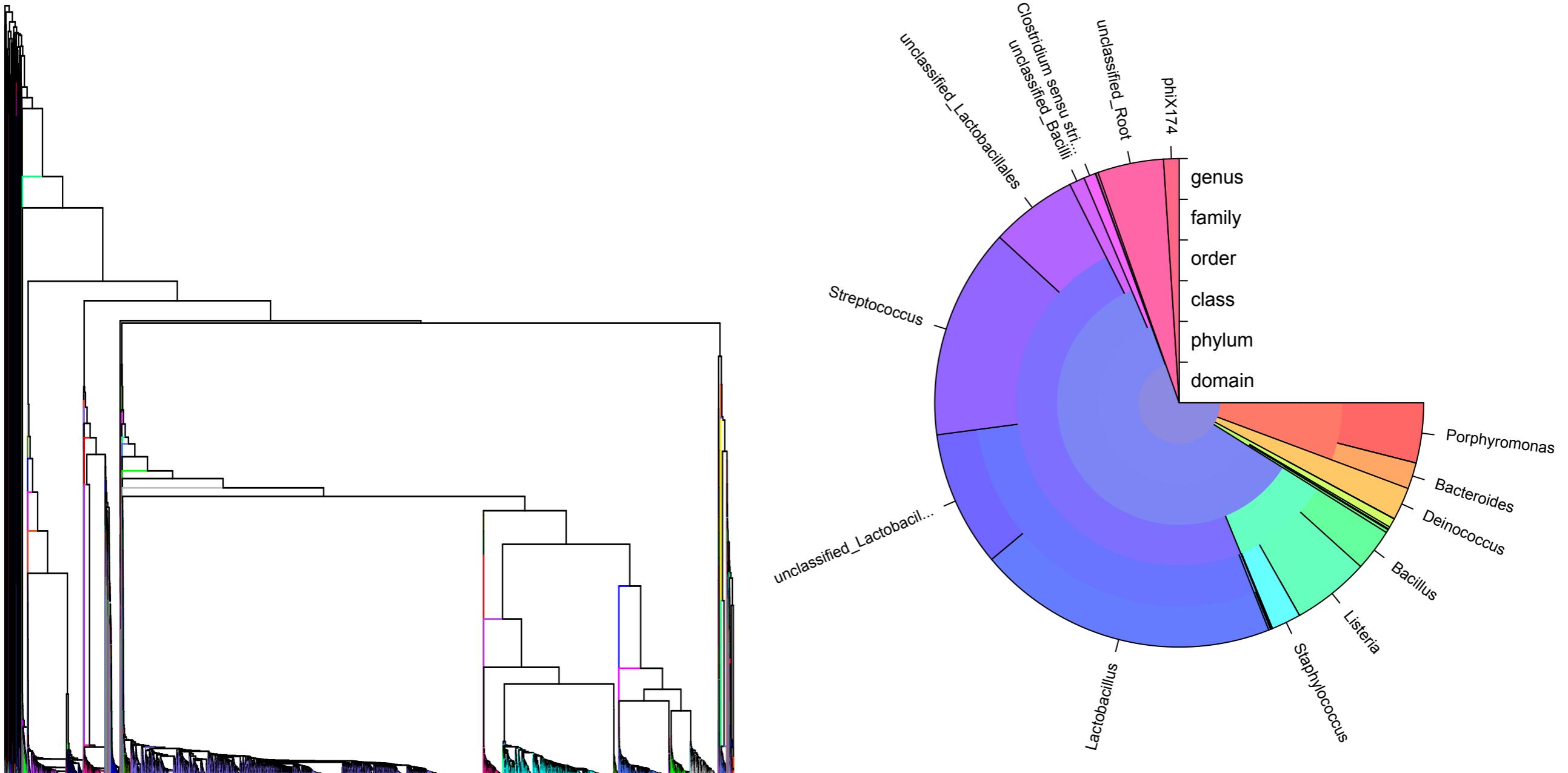
# Pitfalls and caveats in microbiome analyses

**What do you think might be some  
limitations to 16S sequence  
surveys of the microbiome?**

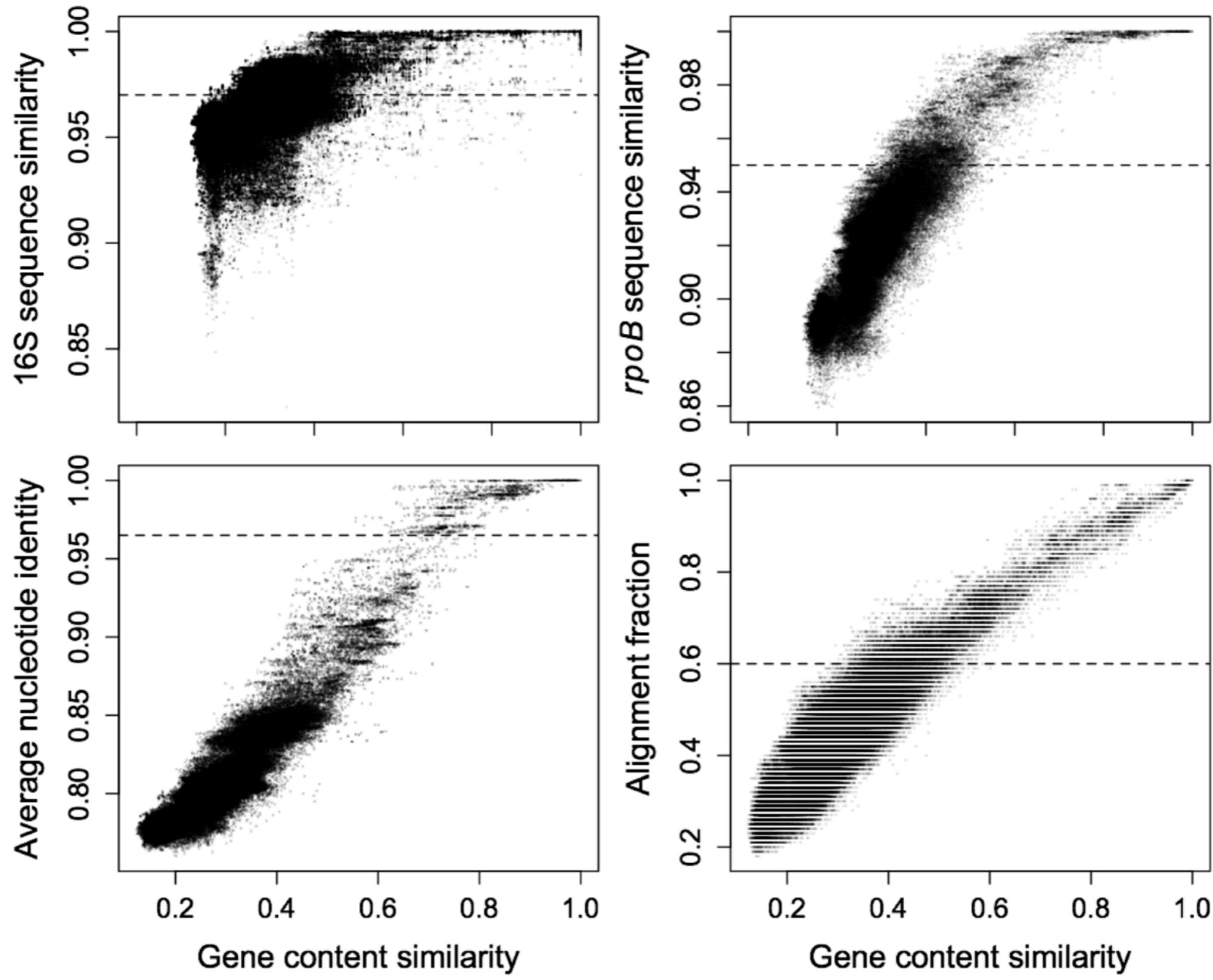
# Potential experimental challenges in 16S studies

- "Kitome"
- PCR artifacts
  - Primer-dimer
  - Chimeras
- Biases:
  - DNA extraction
  - Amplification
  - Primers
  - Sequencing
- Multiple copies of 16S gene
- Sequencing error
- Demultiplexing error
- PhiX spike-in
- Limited length of the 16S gene
  - Short amplicon length
- Slow rate of evolution of the 16S gene
- Misses a lot of organisms
  - viruses
  - eukaryotes

# Sequencing a mock community of 20 strains



# 16S is weakly correlated with gene content



# Potential experimental challenges in 16S studies

- "Kitome"
- PCR artifacts
  - Primer-dimer
  - Chimeras
- Biases:
  - DNA extraction
  - Amplification
  - Primers
  - Sequencing
- Multiple copies of 16S gene
- Sequencing error
- Demultiplexing error
- PhiX spike-in
- Limited length of the 16S gene
  - Short amplicon length
- Slow rate of evolution of the 16S gene
- Misses a lot of organisms
  - viruses
  - eukaryotes

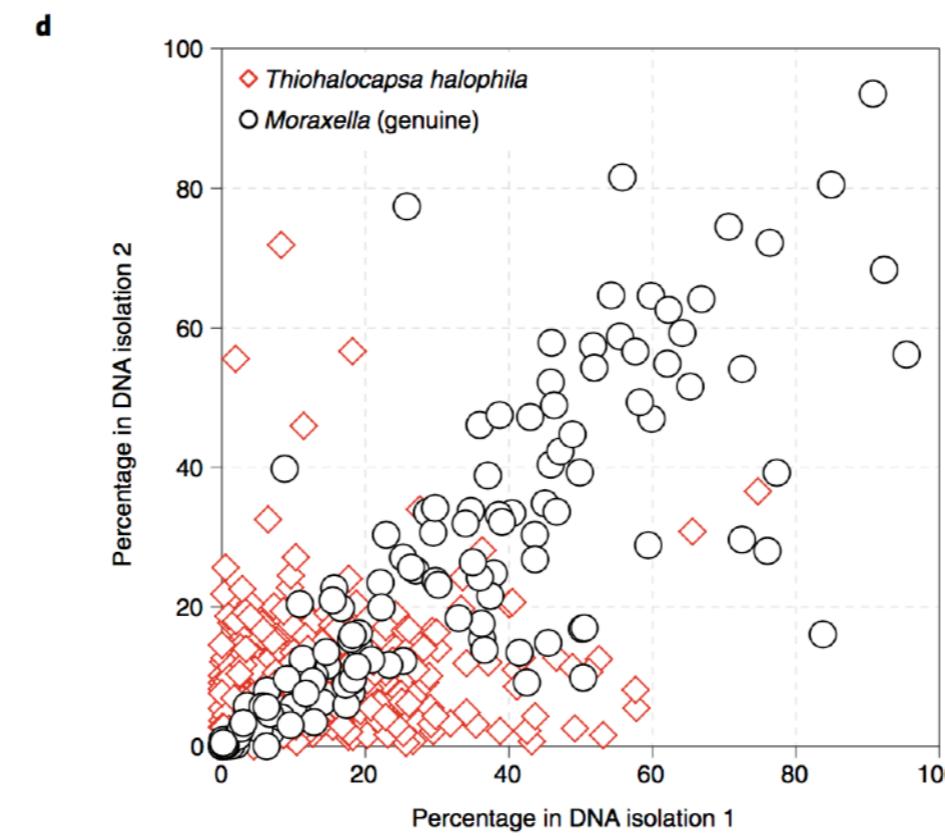
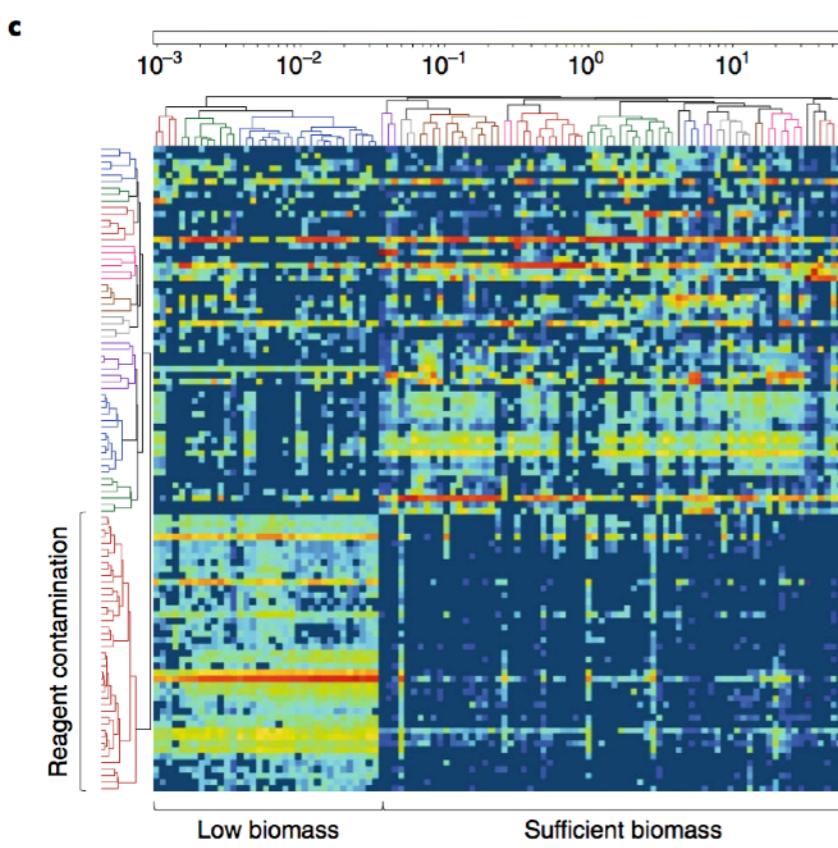
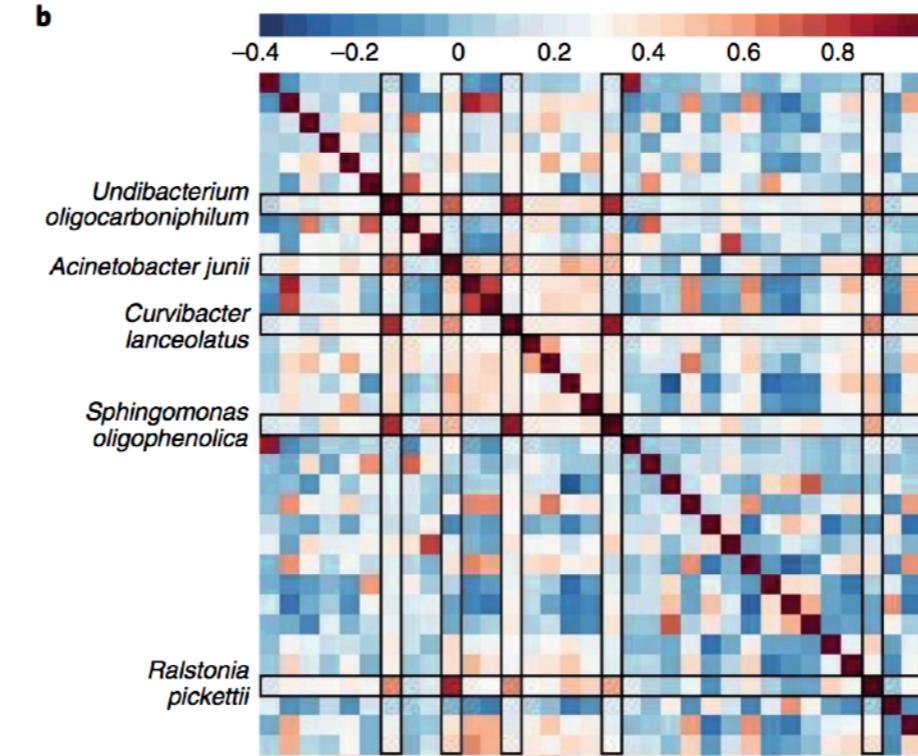
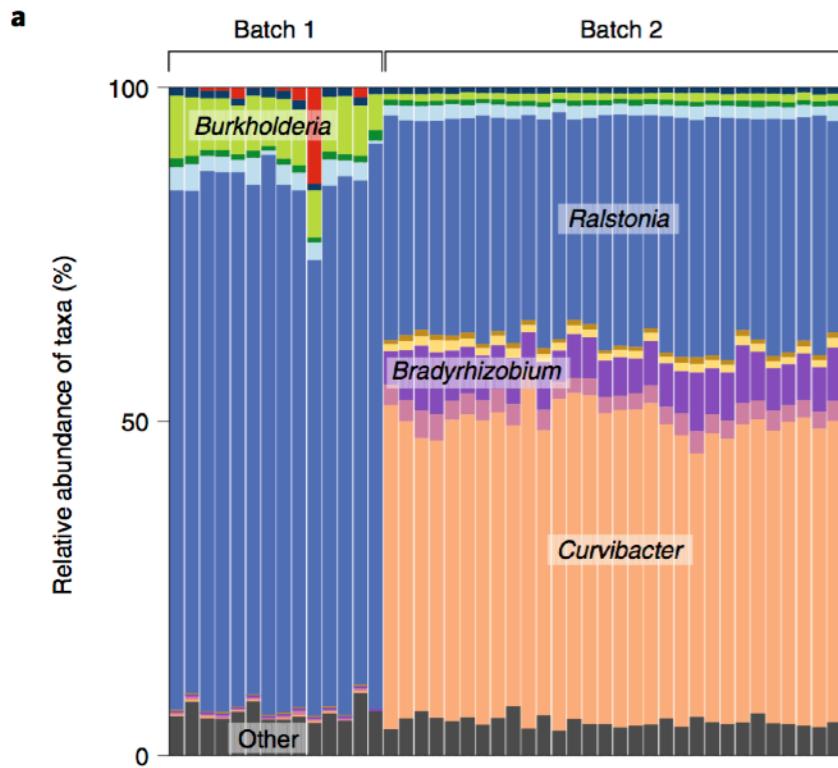
**How would you  
control for each of  
these error modes?**

# Addressing experimental challenges

---

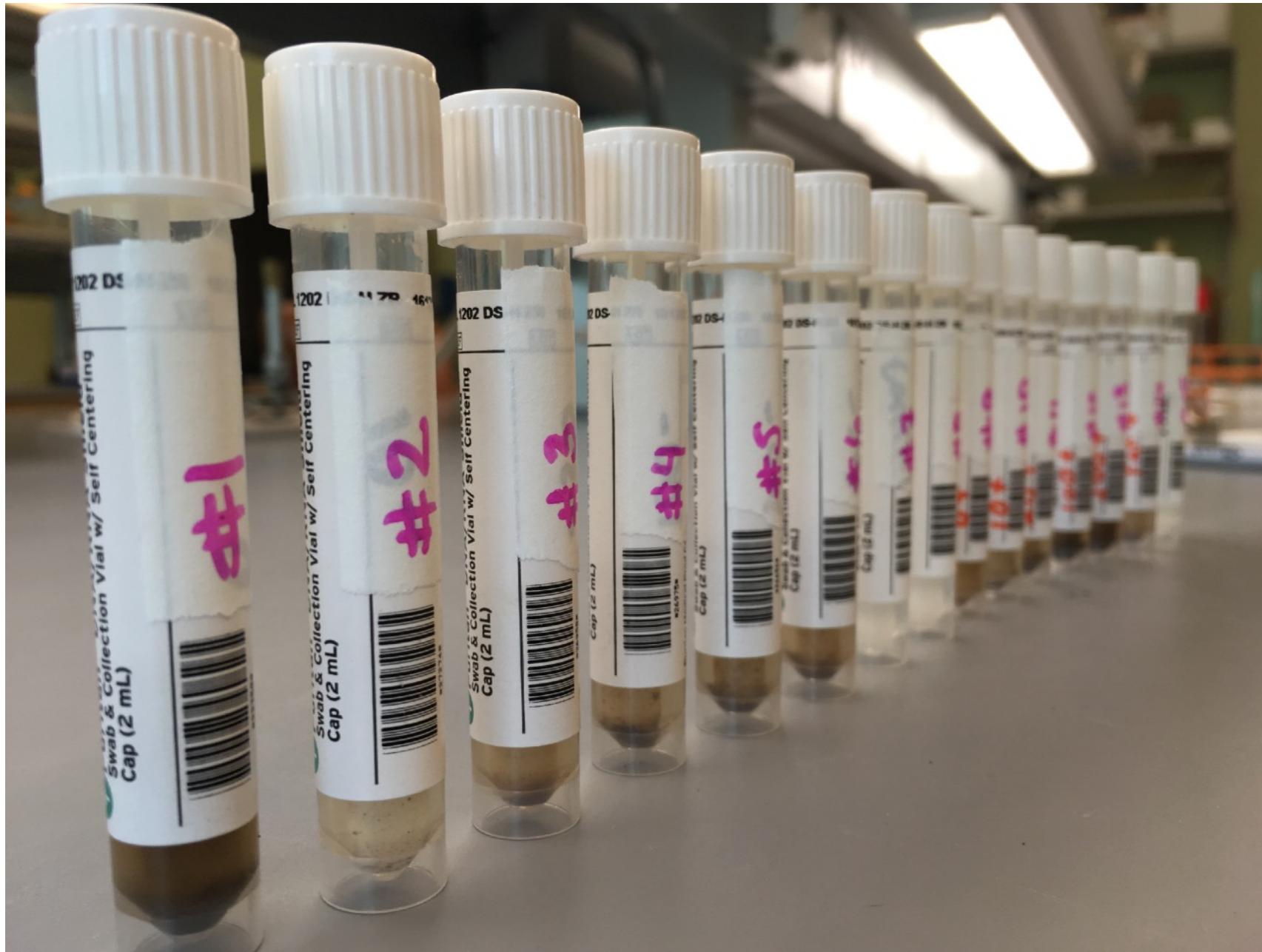
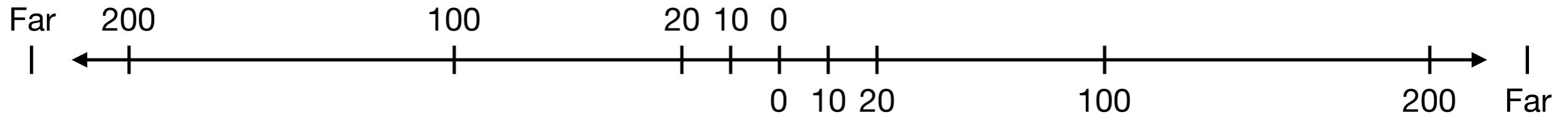
- Negative controls
- Positive controls
  - Mock microbial communities
  - Single strain controls
- Programs for correcting artifacts

# Methods for detecting contamination



# Case study: sequencing the soil microbiome

# Soil collected at varying distances



**Also:**

- + Control: *Klebsiella*
- Control: Nothing
- 20 spp. mock community