

Analyzing Microbiome Data

BIOC 6102 Special Topics

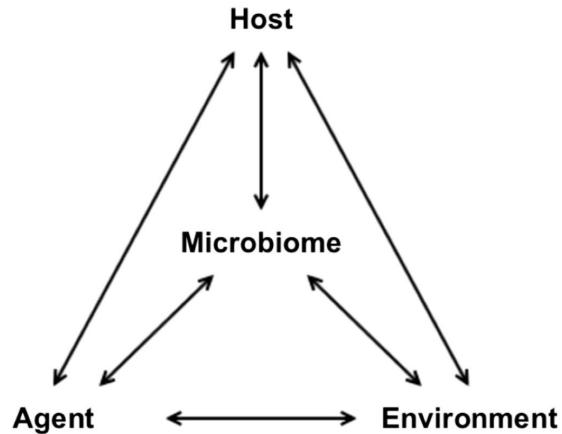
September 6, 2022



Michael S. Robeson II, Ph.D.
University of Arkansas for Medical Sciences
Department of Biomedical Informatics

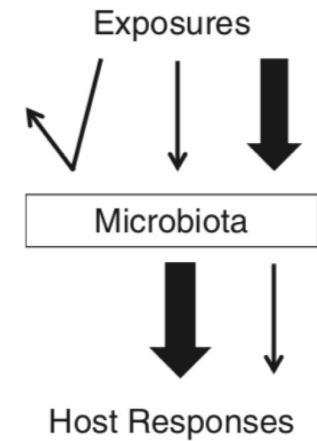


What is Microbiome (Data) Science?



*The study of all facets of the microbiome, e.g. microbial composition, diversity, and function as they interact with the biotic and abiotic features of the environment in which they live, is often referred to as the field of **microbiome science**, or **microbiomics**.*

*The concept of **microbiome data science** refers to the application of best practices from open data science to microbiome bioinformatics.*



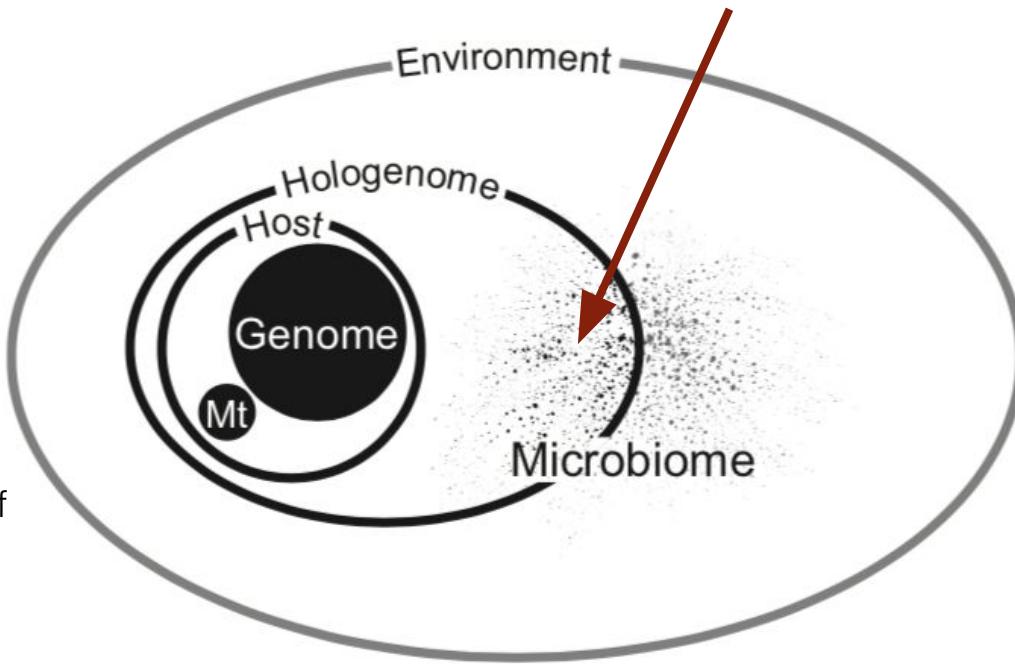
- Bokulich *et al.* 2020. "Measuring the Microbiome: Best Practices for Developing and Benchmarking Microbiomics Methods." Computational and Structural Biotechnology Journal 18 (January): 4048–62. <http://dx.doi.org/10.1016/j.csbj.2020.11.049>
- Shetty & Lahti. 2019. "Microbiome Data Science." Journal of Biosciences 44 (5). 10.1007/s12038-019-9930-2.
- Foxman, B., and Martin, E.T. 2015. Use of the Microbiome in the Practice of Epidemiology: A Primer on -Omic Technologies. American Journal of Epidemiology 182(1): 1–8. doi:10.1093/aje/kwv102.
- Hanson, B.M., and Weinstock, G.M. 2016. The importance of the microbiome in epidemiologic research. Annals of Epidemiology 26(5): 301–305. doi:10.1016/j.annepidem.2016.03.008.

Holobiont

“The 2nd Genome”
“The Extended Genome”

Holobiont:

the host and all of its associated microorganisms.



Hologenome:

the sum of the genetic information of the host and its microbiota.

Brucker & Bordenstein (2013) “The capacious hologenome.”, *Zoology* 116(5): 260–261.

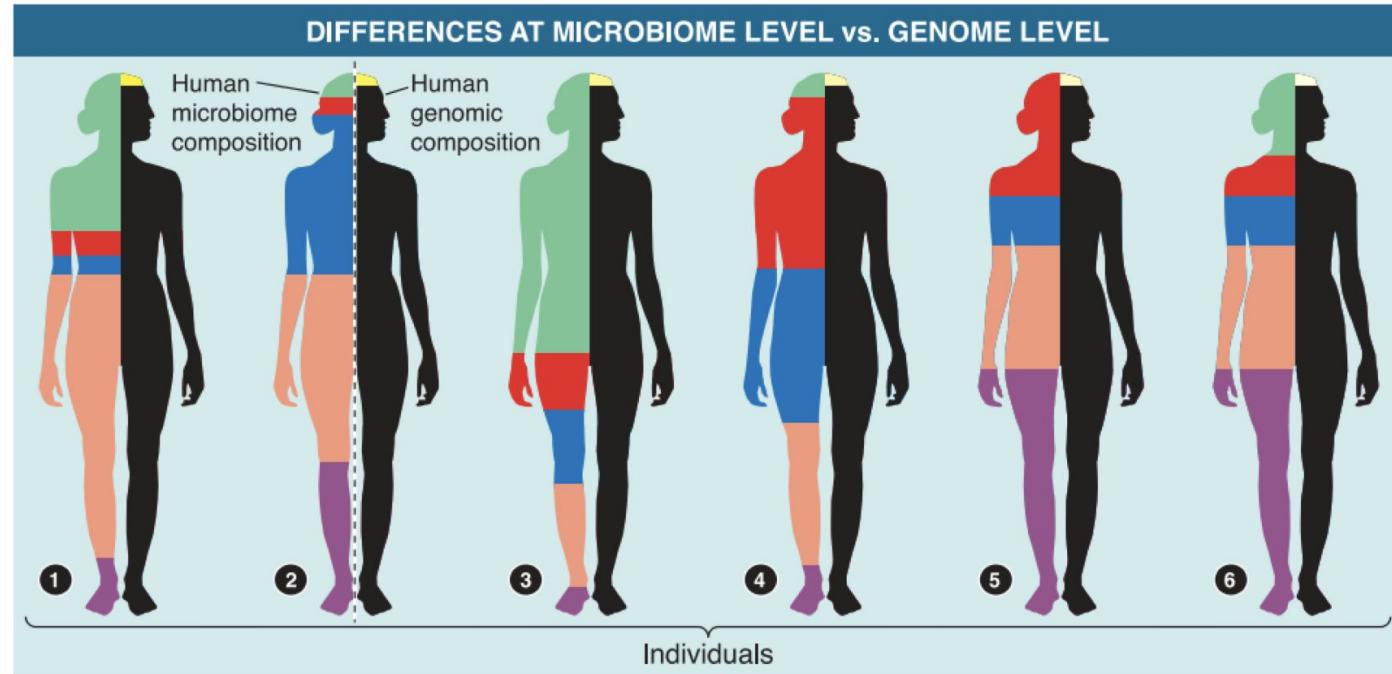
Zilber-Rosenberg & Rosenberg (2008) “Role of microorganisms in the evolution of animals and plants: the hologenome theory of evolution.”, *FEMS Microbiol. Rev.* 32(5): 723–735.

O’Hara, Ann M., and Fergus Shanahan. 2006. “The Gut Flora as a Forgotten Organ.” *EMBO Reports* 7 (7): 688–93.

Relman, D. A., and S. Falkow. 2001. “The Meaning and Impact of the Human Genome Sequence for Microbiology.” *Trends in Microbiology* 9 (5): 206–8.

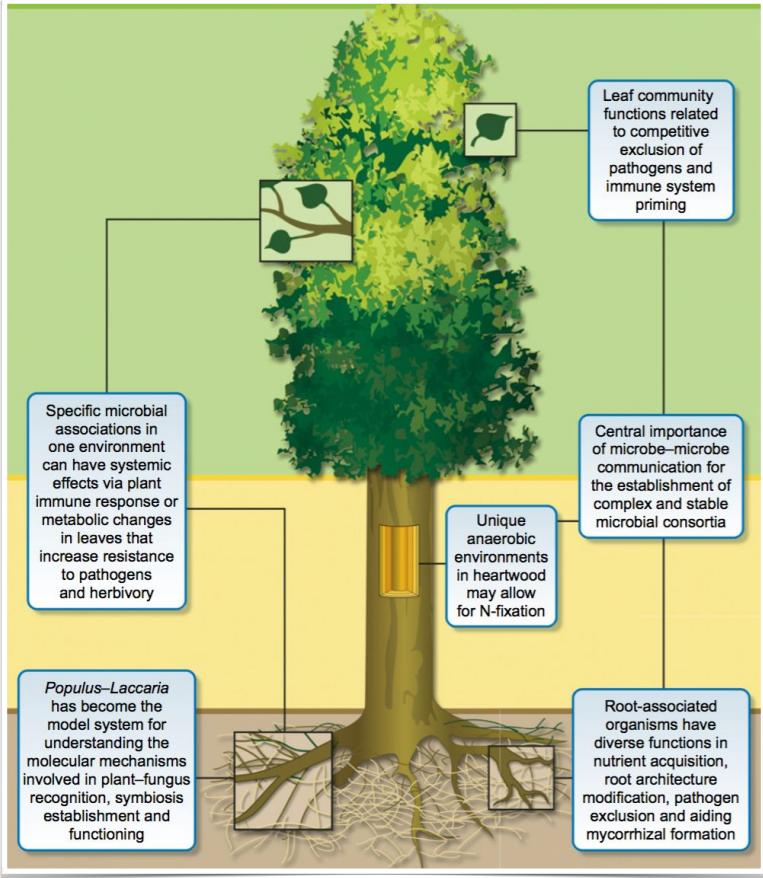
Genetic diversity.

FIGURE 1



Humans are far more different from each other in their microbiome composition than in their genomic composition. The colors in the left side of each individual represent bacterial phyla, while the colors on the right indicate host genomic similarity. For the most part, we contain similar phyla living in and on our bodies, but their relative abundances can be drastically different. On the other hand, our genomic composition is nearly identical, with only a small fraction (around 0.1%) differing across individuals.

The Microbiome



The Father of Microbiology

Antonie van Leeuwenhoek



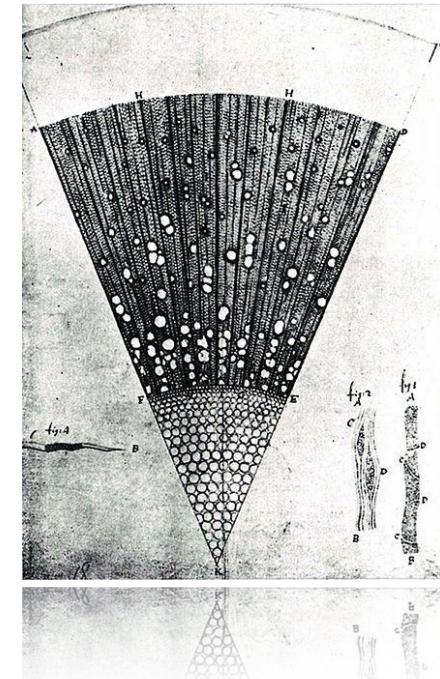
1632-1723

- Helped establish Microscopy / Microbiology as a scientific discipline.
- Used single-lens microscopes

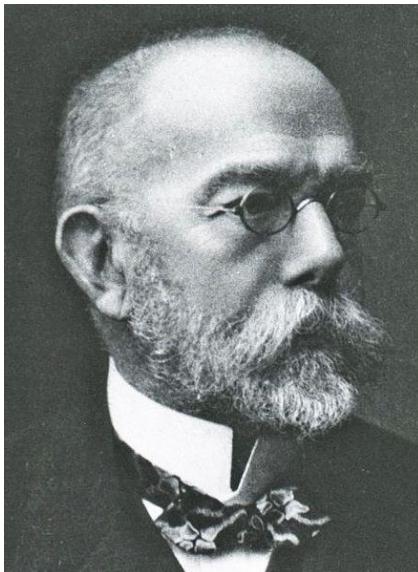


Discoveries:

- Infusoria
- Bacteria (human mouth)
- Vacuole
- Spermatozoa
- Muscle fibers



Robert Koch
(1843-1910)



Proposed the initial criteria by which to establish a causative relationship between a microbe and a disease. This had a profound impact on public health.

https://en.wikipedia.org/wiki/Koch%27s_postulates

Julius Richard Petri
(1852-1921)



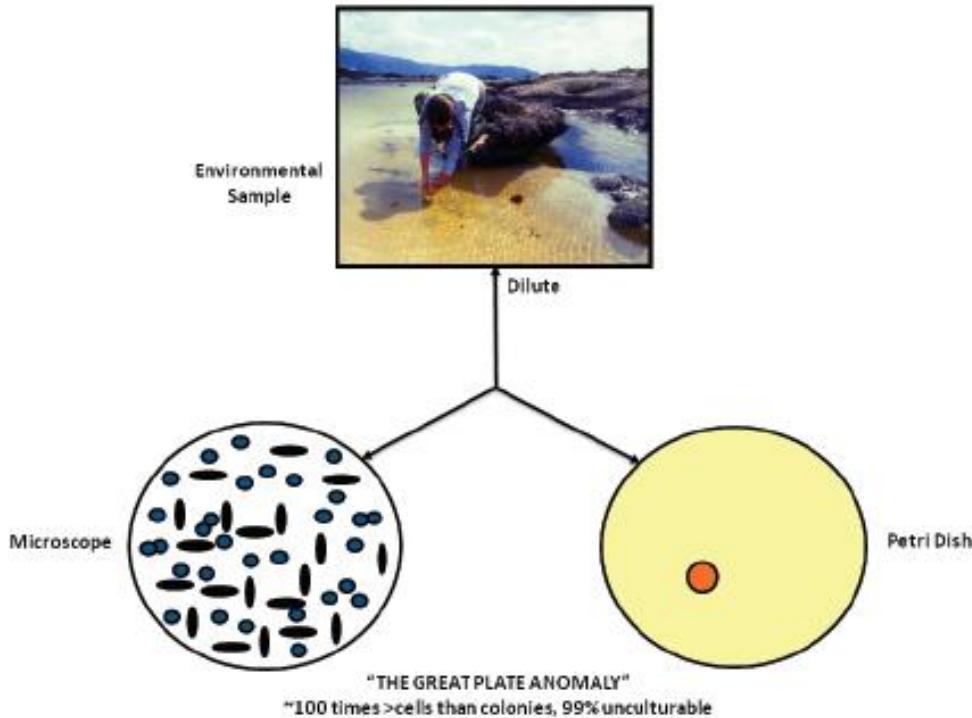
Invented petri dish while working as an assistant to Robert Koch

https://en.wikipedia.org/wiki/Julius_Richard_Petri
https://en.wikipedia.org/wiki/Growth_medium



The Great Plate Count Anomaly.

1. We see millions of various microbes under the microscope, but can only cultivate a small fraction of these microbes.
2. That is, there are **many viable but non-culturable (VBNC) microorganisms**.
3. We simply do not know how to grow many of them.
4. It can be an onerous process to figure out how to cultivate many members of an entire microbial community.



But how do we identify microbial taxa?

Analytical Profile Index (API) Strips

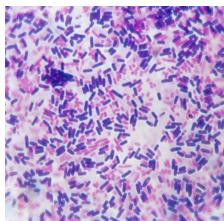


| culture no. | O | N | A | L | O | C | H | U | T | I | V | G | G | M | I | S | R | S | M | A | A | identification | |
|-------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|--------------------------|--|
| | O | N | P | D | D | C | I | 2 | R | D | N | P | E | L | A | N | O | H | A | E | M | R | |
| | P | H | C | C | C | T | S | E | A | D | L | U | N | N | R | A | C | L | Y | M | Y | A | |
| 8101 | + | - | + | + | - | - | - | - | + | - | - | + | + | - | + | + | + | + | - | + | + | Escherichia coli | |
| 5B | + | - | - | - | + | - | - | - | - | + | - | + | + | + | + | + | + | - | - | + | + | Enterobacter agglomerans | |
| 8P44 | - | - | + | + | - | + | - | - | + | - | + | + | - | - | - | - | + | - | - | + | + | Edwardsiella hoshinae | |

<https://www.jlindquist.com/generalmicro/102bactid2.html>

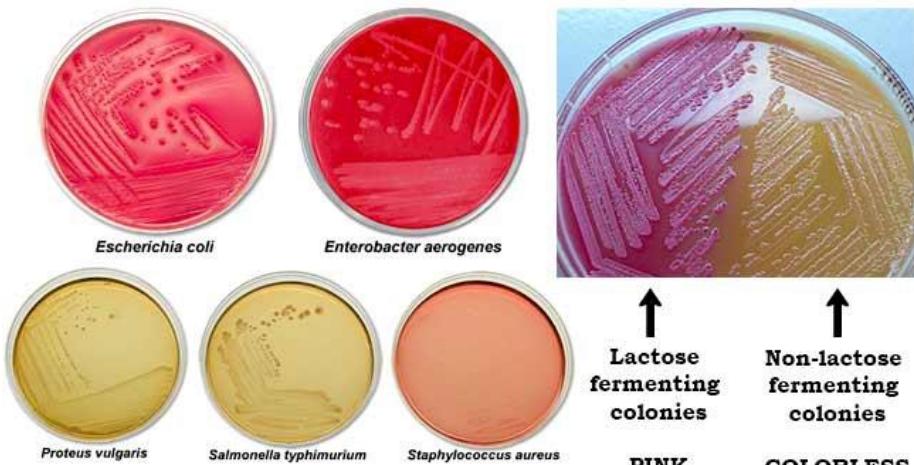
https://en.wikipedia.org/wiki/Analytical_profile_index

<https://www.thoughtco.com/gram-stain-procedure-41476>



Differential & Selective Media

MacConkey agar is
both



- Inhibits or reduces growth of gram positive bacteria
- Promotes gram negative bacterial growth.
 - Lactose fermenting colonies are pink
 - non-lactose fermenting ones are colorless / same color as the medium.

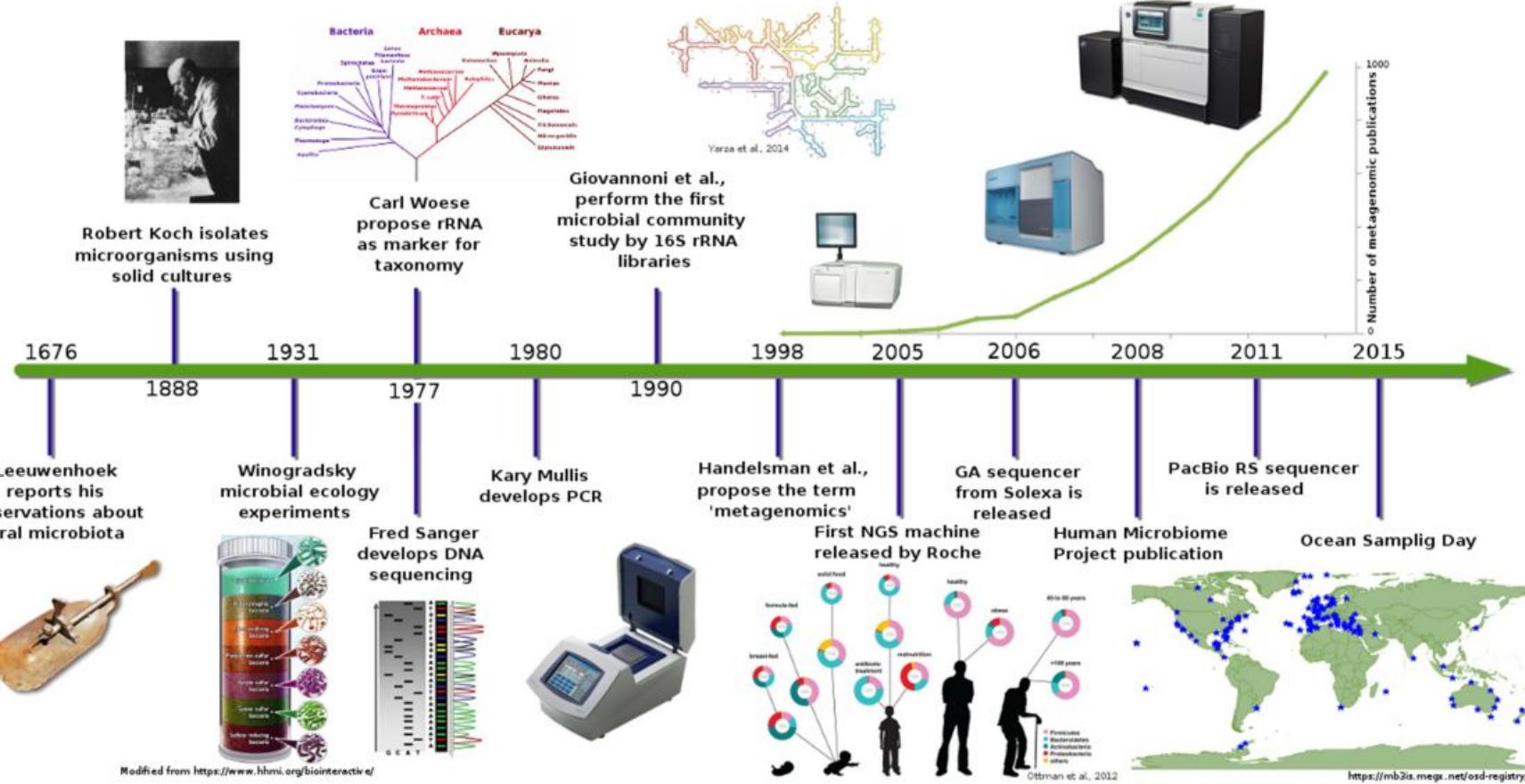
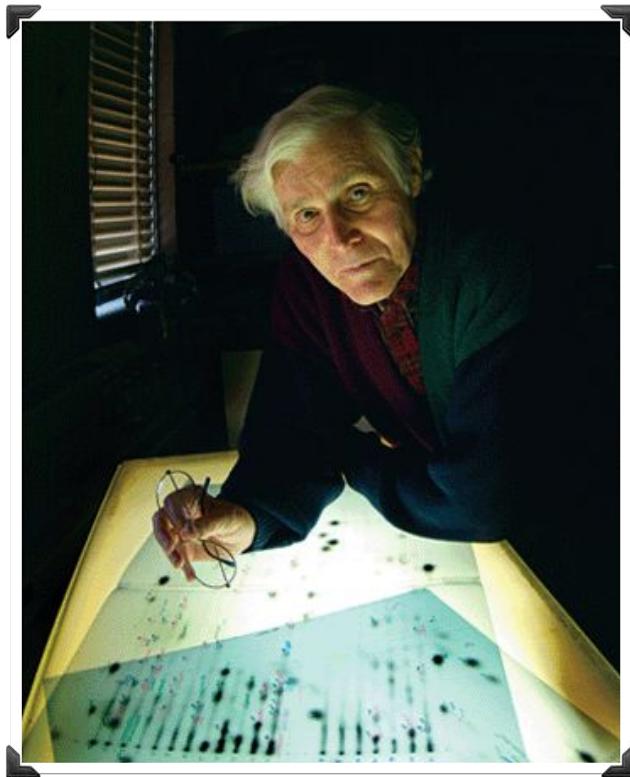


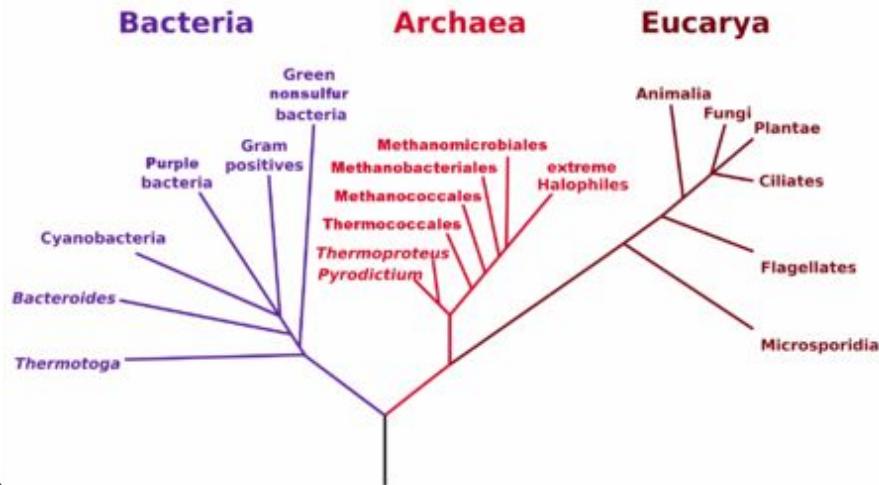
FIGURE 1 | Metagenomics timeline and milestones. Timeline showing advances in microbial communities studies from Leeuwenhoek to NGS (Ottman et al., 2012; Yarza et al., 2014).

Carl Woese

July 15, 1928 – December 30, 2012



Phylogenetic Tree of Life



[http://science.sciencemag.org/content/339/6120/661.full](https://science.sciencemag.org/content/339/6120/661.full)

https://en.wikipedia.org/wiki/Carl_Woese

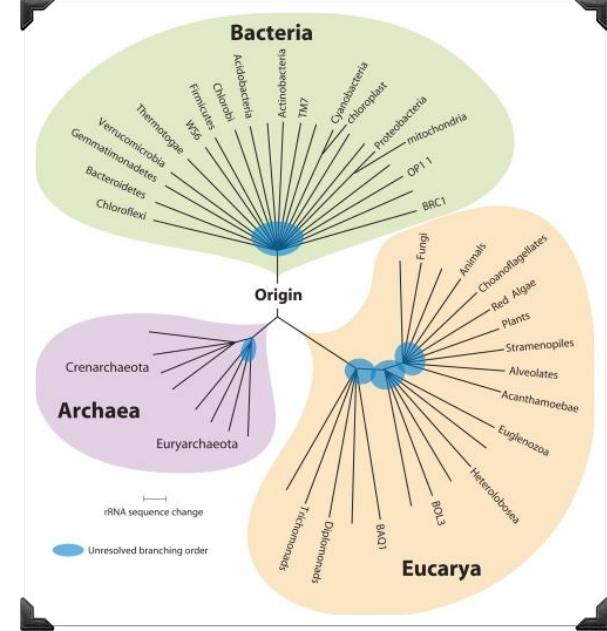
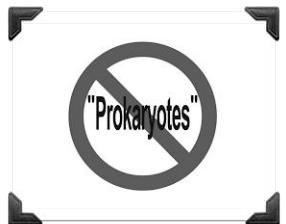
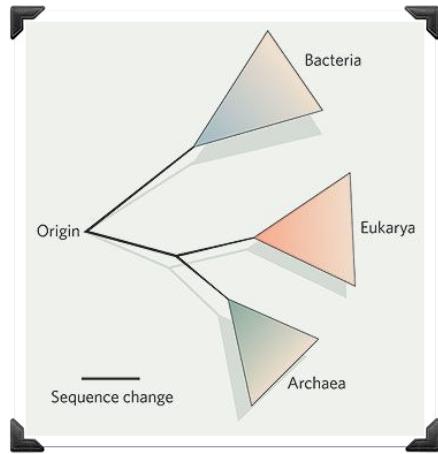
Woese C, Fox G (1977). "Phylogenetic structure of the bacteria domain: the primary kingdoms." Proc Natl Acad Sci USA. 74 (11): 5088–90

Woese CR, Kandler O, and Wheelis ML. (1990). "Towards a Natural System of Organisms: Proposal for the Domains Archaea, Bacteria, and Eucarya." Proc Natl Acad Sci USA 87 (12): 4576–79.

Norm Pace

(1942 -)

"The Man Who Blew The Door Off The Microbial World"

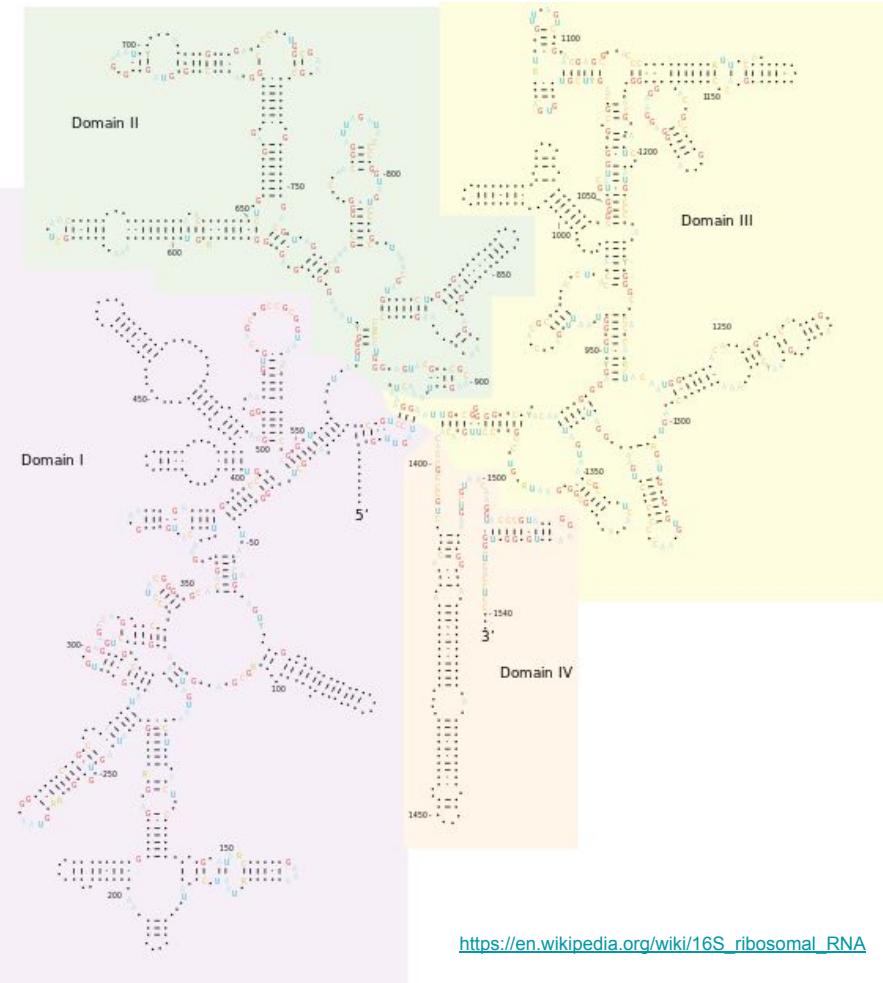


Pace (2006) Nature 441 (7091): 289.

Pace (1986) Adv. Microbial Ecol. 9:1-55.

16S rRNA gene / amplicon markers

- low cost enables sequencing of hundreds to thousands of samples simultaneously
 - useful for tracking microbial compositions in large surveys
- low complexity of amplicon sequence data make it much more computationally tractable for a wide range of research applications
- short reads are less reliable for species-level identification and
 - though sequence variant information can be used to differentiate “phylotypes”



Different variable regions vary in their ability to discriminate among taxa.

V1-V3 : skin

V3-V5

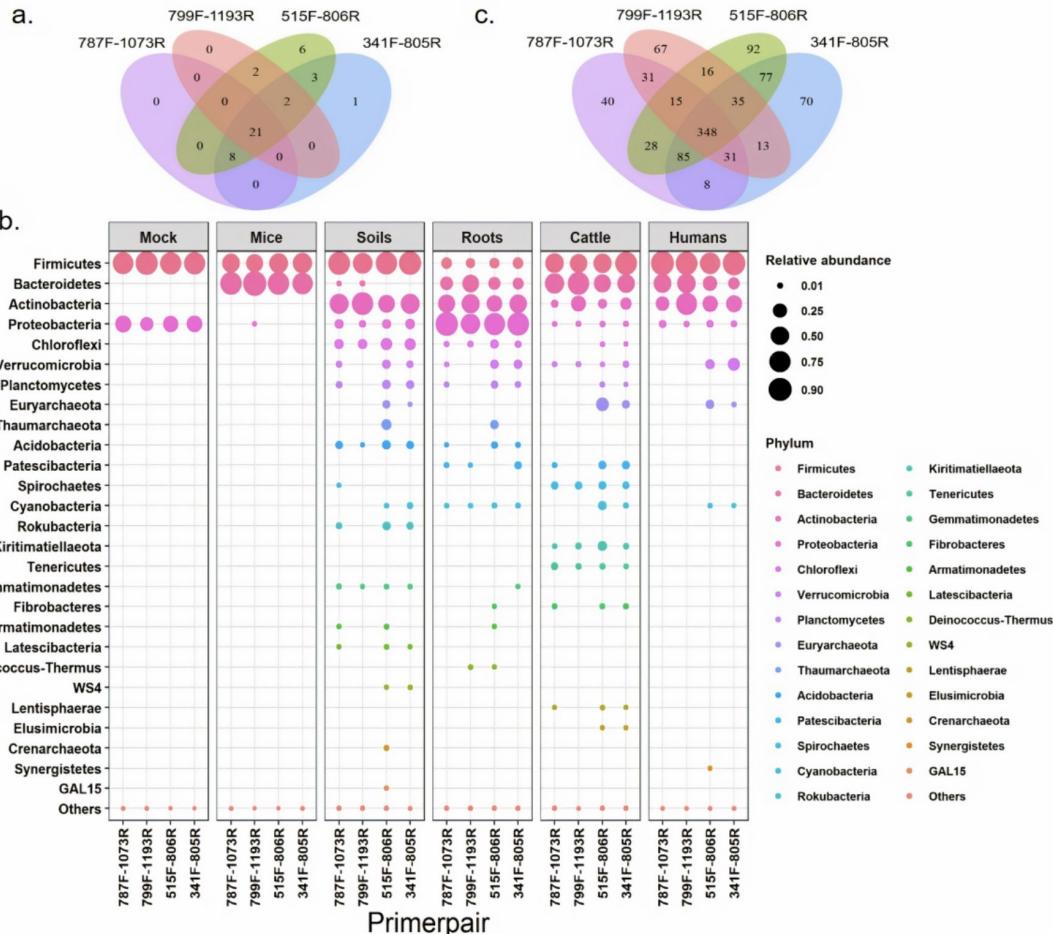
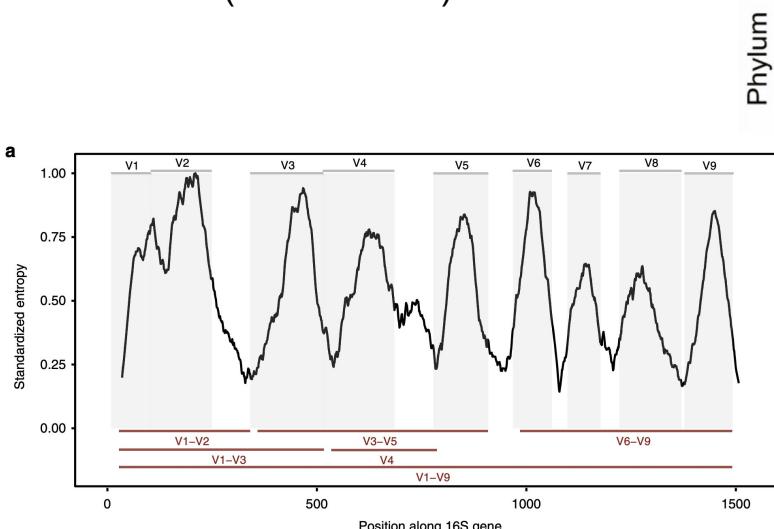
V4 (515F-806R) : *EMP*

V3V4 (341-805R)

V4V5 : *EMP*

V5V6 (787F-1073R)

V5V6V7 (799F-1193R)



How do you define:

- 1) Microbiota?
- 2) Metagenome?
- 3) Microbiome?

Microbiota

“The assemblage of microorganisms present in a defined environment.”

A “microbial census” is predominantly performed by 16S rRNA gene sequencing.

Amplicon sequencing surveys are often, and incorrectly, referred as metagenomics surveys!

Metagenome

“The collection of genomes and genes from the members of a microbiota.”

Often obtained through shotgun sequencing of DNA extracted from a sample.

Microbiome

“This term refers to the entire habitat, including the micro-organisms (bacteria, archaea, lower and higher eukaryotes, and viruses), their genomes (i.e., genes), and the surrounding environmental conditions.”

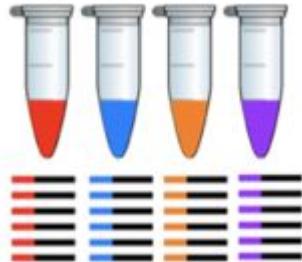
Based on the term “biome”. That is, the biotic and abiotic factors of given environments.

PSA:

Sometimes the best way to study microbial communities may require neither DNA sequencing nor “-omics”!

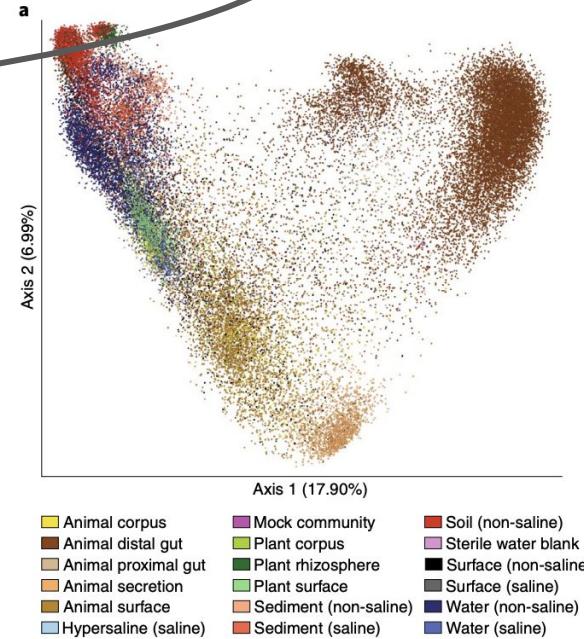
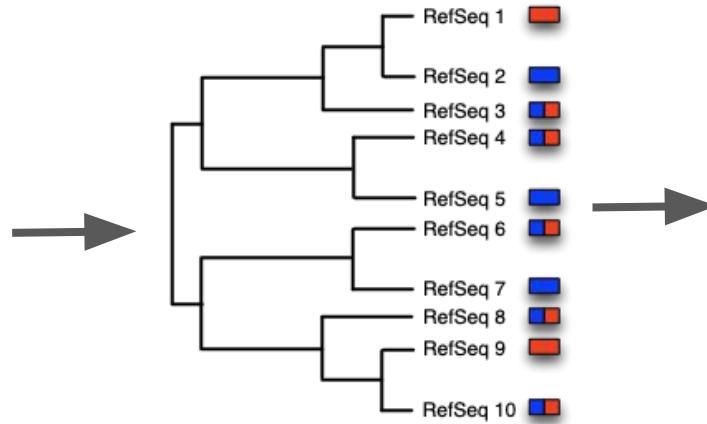
FISH, qPCR, etc., may be all you need!

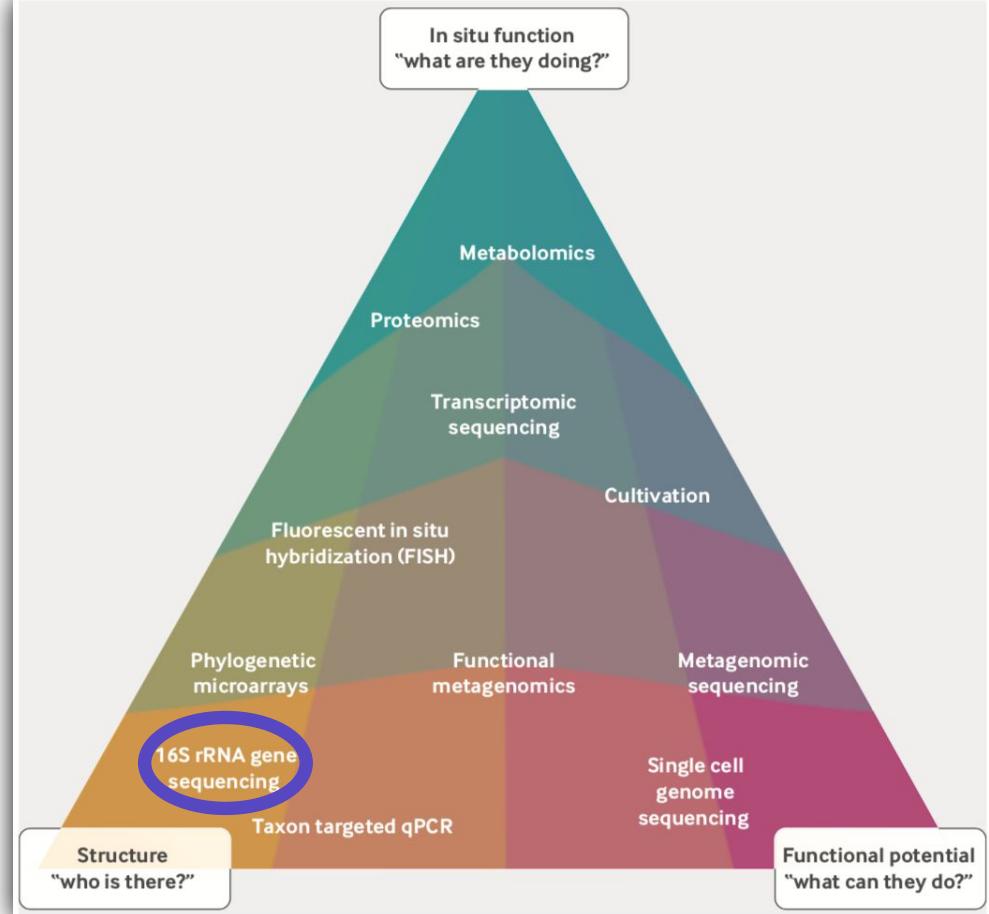
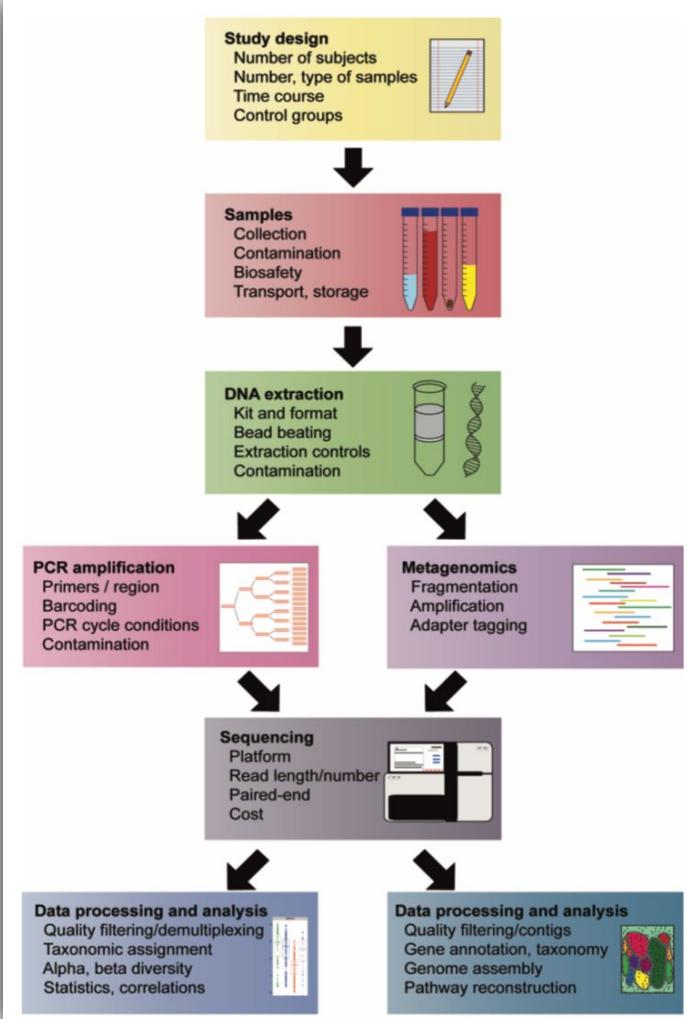
Use the right tool for the question or hypothesis at hand!



qiime2

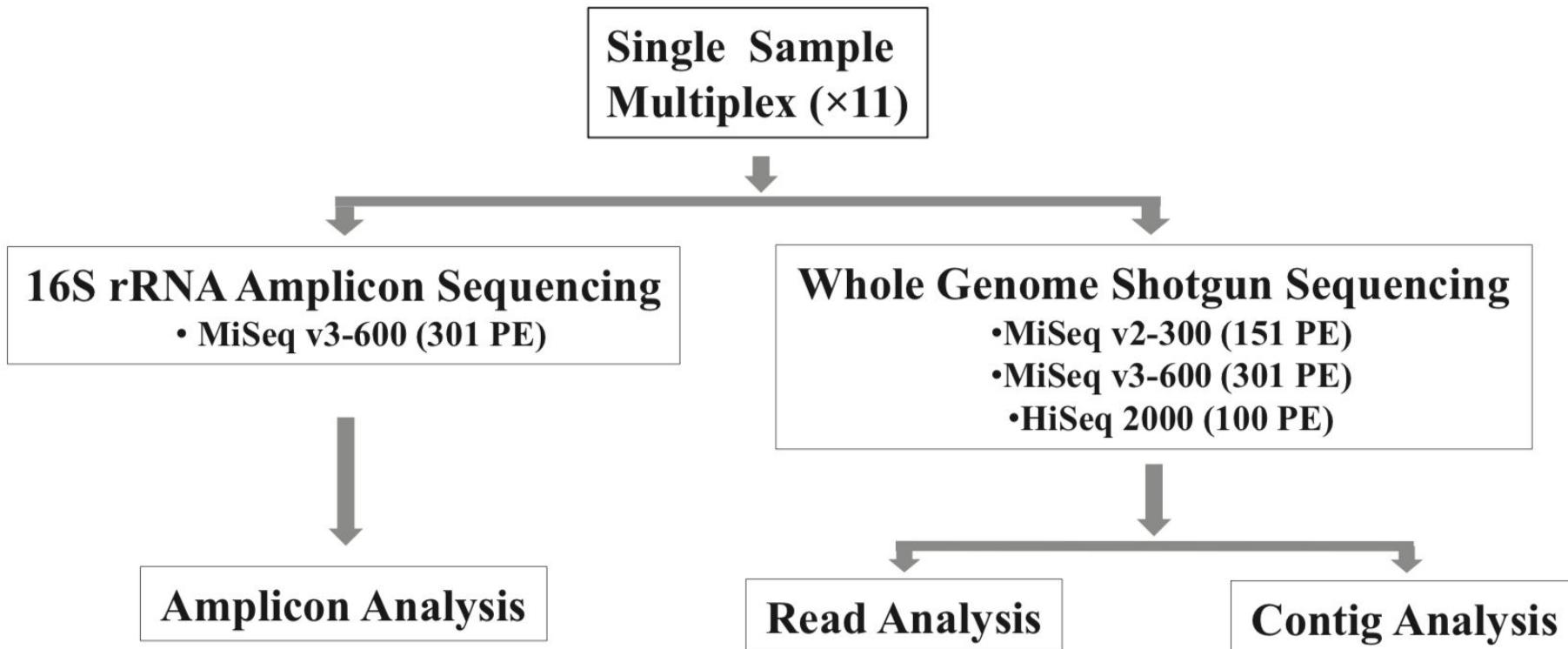
>GCACCTGAGGACAGGCATGAGGAA...
>GCACCTGAGGACAGGGGAGGAGGA...
>TCACATGAACCTAGGCAGGACGAA...
>CTACCGGAGGACAGGCATGAGGAT...
>TCACATGAACCTAGGCAGGAGGAA...
>GCACCTGAGGACACGCAGGACGAC...
>CTACCGGAGGACAGGCAGGAGGAA...
>CTACCGGAGGACACACAGGAGGAA...
>GAACCTTCACATAGGCAGGAGGAT...
>TCACATGAACCTAGGGCAAGGAA...
>GCACCTGAGGACAGGCAGGAGGAA...





- Bik, E.M. 2016. The Hoops, Hopes, and Hypes of Human Microbiome Research. *The Yale Journal of Biology and Medicine* 89(3): 363–373. Yale Journal of Biology and Medicine.
- Young, V.B. 2017. The role of the microbiome in human health and disease: an introduction for clinicians. *BMJ* 356: j831. British Medical Journal Publishing Group. doi:10.1136/bmj.j831.

Experimental Strategy



Ranjan *et al.* (2016) Analysis of the microbiome: Advantages of whole genome shotgun versus 16S amplicon sequencing. Biochem. Biophys. Res. Commun. 469(4): 967–977



Quantitative Insights Into Microbial Ecology 2

Popularity primarily due to:

- Supportive and interactive community.
- Cutting edge functionality quickly added via plugins.

Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2

7134 *

E Bolyen, JR Rideout, MR Dillon, NA Bokulich, CC Abnet, GA Al-Ghalith, ...
Nature biotechnology 37 (8), 852-857

Reproducible, interactive, scalable and extensible microbiome data science using QIIME 2

6413

E Bolyen, JR Rideout, MR Dillon, NA Bokulich, CC Abnet, GA Al-Ghalith, ...
Nature biotechnology 37 (8), 852-857, 2019

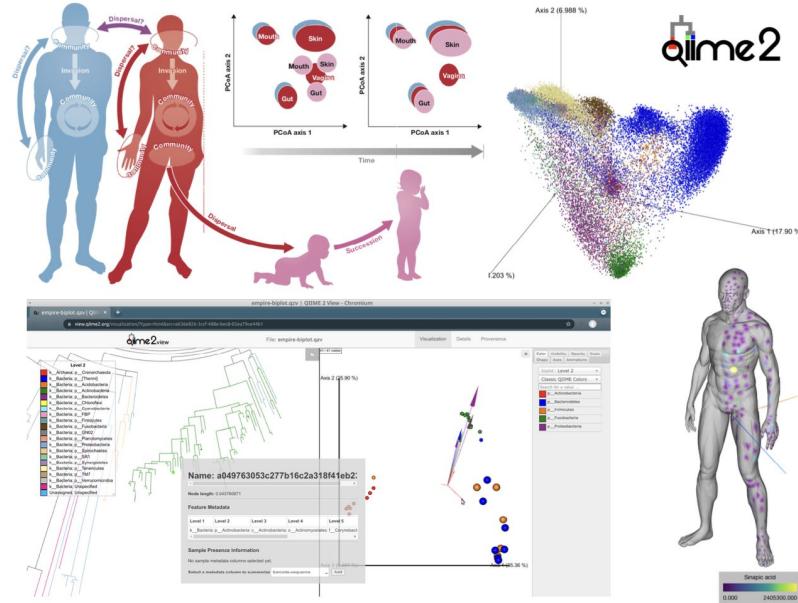
QIIME 2: Reproducible, interactive, scalable, and extensible microbiome data science

794 *

E Bolyen, JR Rideout, MR Dillon, NA Bokulich, C Abnet, GA Al-Ghalith, ...
PeerJ Preprints, 2018

Currently uploading QIIME 2 workshop videos each week to YouTube:
<https://www.youtube.com/c/QIIME2>

Foundations of the Human Microbiome BMIG 6202 - Spring 2022



Course Description:

This graduate course will provide students with foundational knowledge and practical analytical skills required for analyzing microbiome data.

Students will explore the microbial inhabitants of the human body, with an emphasis on how microbial communities affect human health and disease progression.

Students will learn how to perform current leading-edge microbiome analysis directly from a QIIME 2 developer!

- ✓ No course prerequisites.
- ✓ Review foundational research.
- ✓ Discuss host-microbiome interactions.
- ✓ Learn how to plan a microbiome study.
- ✓ Learn how to analyze microbiome data.

Location: TDB

Times:

Tu 9:00 - 10:00 AM (Lecture)
Thu 9:00 - 11:00 AM (Computer Lab)

Contact: MRobeson@uams.edu
SJun@uams.edu



[NSF Award: 1565100](#)

Low-level features

- **Decentralized provenance tracking** automates bioinformatics record keeping facilitating reproducibility.
- **Multiple user interfaces.** The same functionality is accessible through graphical interface, command line interface, and API, which target different types of users.
- **Plugin architecture** allows the software to keep pace with the field. Any developer can create and distribute a QIIME 2 plugin.

"What did you do 5 months ago?"

notes.txt

```
echo "core_diversity_analyses.py -i  
/home/caporaso/analysis/atacama-7may2013/slout_r1/or_otus/otu_table_mc2_w_tax_no_pynast_failures.  
biom -o /home/caporaso/analysis/atacama-7may2013/slout_r1/or_otus/cd_16662/ -e 16662 -m  
/home/caporaso/analysis/atacama-7may2013/map.txt -ao 26 -t  
/home/caporaso/analysis/atacama-7may2013/slout_r1/or_otus/rep_set.tre -c  
SiteName,Depth,ExtractGroupNo, TransectID, Vegetation" | qsub -keo -N ata-cd
```

The above failed during OTU category significance (see the log file for the error - maybe one of these categories has a value that is observed only once?), so re-running without that step for now...

```
echo "core_diversity_analyses.py -i  
/home/caporaso/analysis/atacama-7may2013/slout_r1/or_otus/otu_table_mc2_w_tax_no_pynast_failures.  
biom -o /home/caporaso/analysis/atacama-7may2013/slout_r1/or_otus/cd_16662/ -e 16662 -m  
/home/caporaso/analysis/atacama-7may2013/map.txt -ao 26 -t  
/home/caporaso/analysis/atacama-7may2013/slout_r1/or_otus/rep_set.tre -c  
SiteName,Depth,ExtractGroupNo, TransectID, Vegetation --suppress_otu_category_significance  
--recover_from_failure" | qsub -keo -N ata-cd
```

```
echo "pick_open_reference_otus.py -i /home/caporaso/analysis/atacama-7may2013/slout_r2/seqs.fna  
-r /data/gg_13_5_otus/rep_set/97_otus.fasta -o  
/home/caporaso/analysis/atacama-7may2013/slout_r2/or_otus/ -ao 28 -p  
/home/caporaso/analysis/atacama-7may2013/uc_fast_params.txt" | qsub -keo -N or-3
```

Katy's analysis on first and second sequencing runs

You have read access to all files in /home/caporaso/analysis/atacama-~~7may2013~~.

The split_libraries_fastq.py input files you need are:

sequences:

```
/home/caporaso/analysis/atacama-7may2013/2014.04.30/Undetermined_S0_L0  
01_R1_001.fastq.gz
```

The screenshot shows a file browser interface with a list of files and folders. The files are listed in the following order:

- delete-me.txt
- delete-me.txt.bak
- emp_shenzen_map_collapsed_lat_long.txt
- emrakul
- mapping_master_jh_corrected_3May.xls
- master_otu_table_20jul.txt.gz
- notes.txt
- old-results
- sample-metadata-fixed-missing-controls.tsv
- sample-metadata-old.1.tsv
- sample-metadata.1.xlsx
- sample-metadata.xlsx
- sampling_notes.txt
- table.biom
- table.without-controls.biom.tsv
- table1.biom
- table2.biom
- tax-0
- taxOn0my
- taxa_plots.tgz
- taxonomy-plots.png.pdf
- unweighted_unifrac_pc1_pc...olor_by_host_associated.pdf
- unweighted_unifrac_pc1_pc...3_color_by_sample_type.pdf

QIIME 2 integrated data provenance ensures reproducibility

qiime2view

File: taxa-bar-plots.qzv

Download
SVG (bars) SVG (legend) CSV

Hover over the plot to learn more

Relative Frequency

Sample

Taxonomic Level
Level 1
Level 2
Level 3
Level 4
Level 5
Level 6
Level 7

Color Palette schemeAccent
Sort Samples By k_Bacteria:p_Firmicutes Ascending

Provenance Graph

```
graph TD; seqs[seqs] --> sequences[sequences]; sequences --> reference_taxonomy[reference_taxonomy]; demultiplexed_seqs[demultiplexed_seqs] --> reference_reads[reference_reads]; reference_reads --> reads_classifier[reads classifier]; reads_classifier --> tab[tab']; taxonomy[taxonomy] --> tab';
```

Action Details

- execution:
 - uuid: 3897fb5c-55ed-46b1-a48d-ae0651d2b597
- runtime:
 - start: 2017-09-28T21:14:34.374Z
 - end: 2017-09-28T21:23:05.935Z
 - duration: "8 minutes, 31 seconds, and 561708 microsecond s"
- action:
 - type: "method"
 - plugin: "environment:plugins:dada2"
 - action: "denoise_single"
- inputs:
 - demultiplexed_seqs: "ce7e102e-4b8c-455c-b2af-a7fb342b7fa1"
- parameters:
 - trunc_len: 120
 - trim_left: 0
 - max_ee: 2
 - trunc_q: 2
 - chimera_method: "consensus"
 - min_fold_parent_over_abundance: 1
 - n_threads: 1

Visualization Peek Provenance

File: taxa-bar-plots.qzv

Visualization Peek Provenance

Action Details

- execution:
 - uuid: 3897fb5c-55ed-46b1-a48d-ae0651d2b597
- runtime:
 - start: 2017-09-28T21:14:34.374Z
 - end: 2017-09-28T21:23:05.935Z
 - duration: "8 minutes, 31 seconds, and 561708 microsecond s"
- action:
 - type: "method"
 - plugin: "environment:plugins:dada2"
 - action: "denoise_single"
- inputs:
 - demultiplexed_seqs: "ce7e102e-4b8c-455c-b2af-a7fb342b7fa1"
- parameters:
 - trunc_len: 120
 - trim_left: 0
 - max_ee: 2
 - trunc_q: 2
 - chimera_method: "consensus"
 - min_fold_parent_over_abundance: 1
 - n_threads: 1

Use an interface based on your needs.

qiime2.view

This interface can view .qza and .qzv files directly in your browser without uploading to a server. [Click here to learn more.](#)

Drag and drop or click here to view a QIIME 2 Artifact or Visualization (.qza/.qzv) from your computer.

You can also provide a link to a file on Dropbox or a file from the web.

Gallery

Don't have a QIIME 2 result of your own to view? Try one of these!

Taxonomic Bar Plots

Explore the taxonomy of samples in the Moving Pictures Tutorial. Try selecting different taxonomic levels and metadata-based sample sorting.

Volatility Control Chart

Explore interactive line plots to assess how volatile a dependent variable is over a continuous, independent variable in one or more groups.

QIIME 2 Studio

Active Jobs 1 Finished Jobs Failed Jobs

| Action | Started | Elapsed |
|--|-------------------|----------|
| Denoise and derePLICATE paired-end sequences | 17-07-07 01:57:27 | 00:00:05 |

Artifacts 2 Visualizations Metadata 1

| Name | UUID | Type |
|--------------------------|---------------------------------------|---|
| demux | 043bdccdf-9f32-48ce-8c8d-4403bf550a59 | SampleData[PairedEndSequencesWithQuality] |
| emp-paired-end-sequences | 9c7033e-82d6-4f4a-9cf7-baebae8b642f | EMPPairedEndSequences |

+

2. ~ (zsh)

```
$ qiime info
System versions
Python version: 3.5.3
QIIME 2 release: 2017.6
QIIME 2 version: 2017.6.0
q2cli version: 2017.6.0

Installed plugins
alignment 2017.6.0
composition 2017.6.0
dada2 2017.6.0
```

Untitled - idle

```
[1] import pandas as pd
from qiime2 import Artifact

[2] t = Artifact.load('table.qza')
t.view(pd.DataFrame)

4b5eeb300368260019c1fb7a3c718fc

L1S105 2222.0
L1S140 0.0
L1S208 0.0
L1S257 0.0
L1S281 0.0
```

Python 3 | idle Not saved yet

Galaxy
Cyverse

Galaxy
EBI Metagenomics
Portal
QIITA
NIH Nephele
Cyverse

mothur
QIIME 1

phyloseq

Related Software

Computational Sophistication

Data analyst
(clinician, policy maker,
research subject)

Cancer researchers
and other domain
scientists

Power users

Data Scientists

Plugins

Developing tools to democratize database acquisition, formatting, and curation.

RESEARCH ARTICLE

RESCRIPT: Reproducible sequence taxonomy reference database management

Michael S. Robeson, II¹, Devon R. O'Rourke¹, Benjamin D. Kaehler¹,
Michał Ziemska², Matthew R. Dillon², Jeffrey T. Foster², Nicholas A. Bokulich^{1*}

1 University of Arkansas for Medical Sciences, Department of Biomedical Informatics, Little Rock, Arkansas, United States of America, **2** Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, Arizona, United States of America, **3** School of Science, University of New South Wales, Canberra, Australia, **4** Laboratory of Food Systems Biotechnology, Institute of Food, Nutrition, and Health, ETH Zürich, Switzerland

DOI 10.5281/zenodo.3891931



lint-build-test

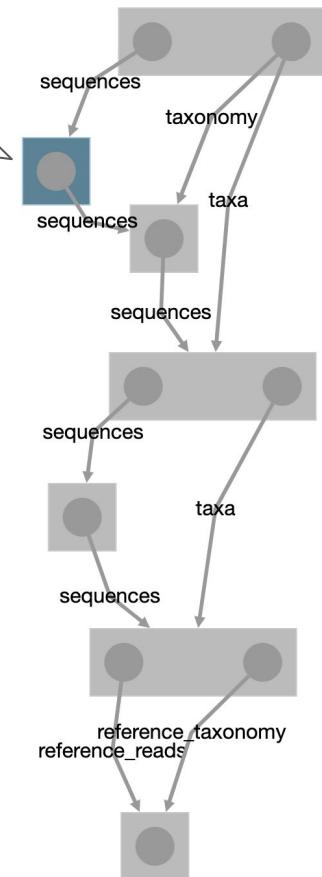
passing

DOI 10.1371/journal.pcbi.1009581

qiime2

▼ parameters:

- ▼ 0:
num_degenerates: 5
- ▼ 1:
homopolymer_length: 8
output-name: "clean_sequences"



RESCRIPT: Reproducible sequence taxonomy reference database management

74

2021

MS Robeson, DR O'Rourke, BD Kaehler, M Ziemska, MR Dillon, JT Foster, ... Paperpile
PLoS computational biology 17 (11), e1009581

- Bolyen *et al.* 2019. "Reproducible, Interactive, Scalable and Extensible Microbiome Data Science Using QIIME 2." *Nature Biotechnology* 37 (8): 852–57.
- Robeson *et al.* 2021. "RESCRIPT: Reproducible Sequence Taxonomy Reference Database Management." *PLoS Computational Biology* 17 (11): e1009581. <http://dx.doi.org/10.1371/journal.pcbi.1009581>

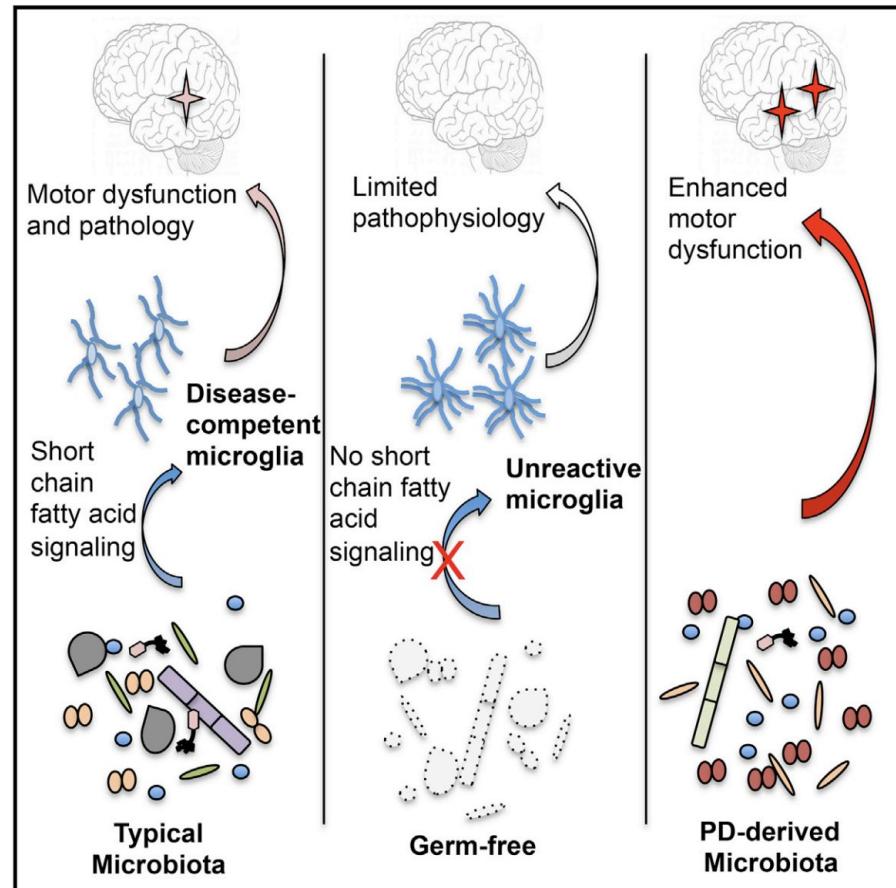
Parkinson's Mouse QIIME 2 Tutorial

<https://docs.qiime2.org/2020.11/tutorials/pd-mice/>

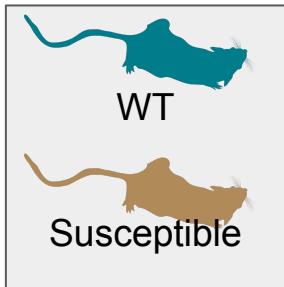
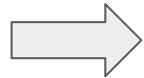
Gut Microbiota Regulate Motor Deficits and Neuroinflammation in a Model of Parkinson's Disease

Timothy R. Sampson,^{1,*} Justine W. Debelius,² Taren Thron,¹ Stefan Janssen,² Gauri G. Shastri,¹ Zehra Esra Ilhan,³ Collin Challis,¹ Catherine E. Schretter,¹ Sandra Rocha,⁴ Viviana Gradinaru,¹ Marie-Francoise Chesselet,⁵ Ali Keshavarzian,⁶ Kathleen M. Shannon,^{7,9} Rosa Krajmalnik-Brown,³ Pernilla Wittung-Stafshede,⁴ Rob Knight,^{2,8} and Sarkis K. Mazmanian^{1,10,*}

<http://dx.doi.org/10.1016/j.cell.2016.11.018>

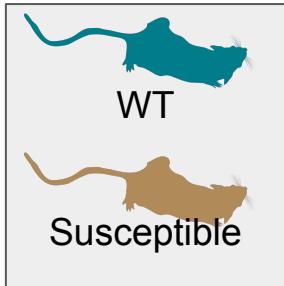


Humanized Mice Experiment



x 3

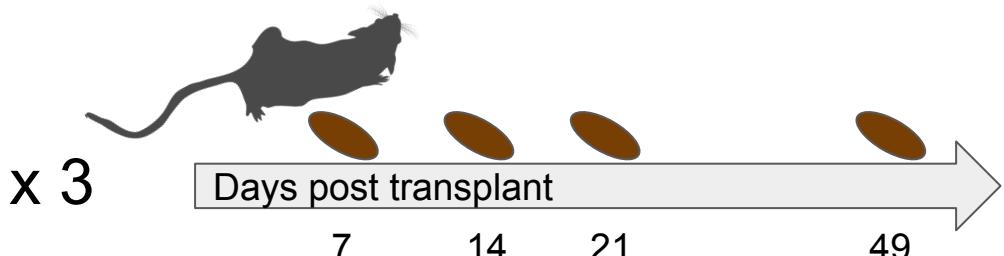
Healthy Control



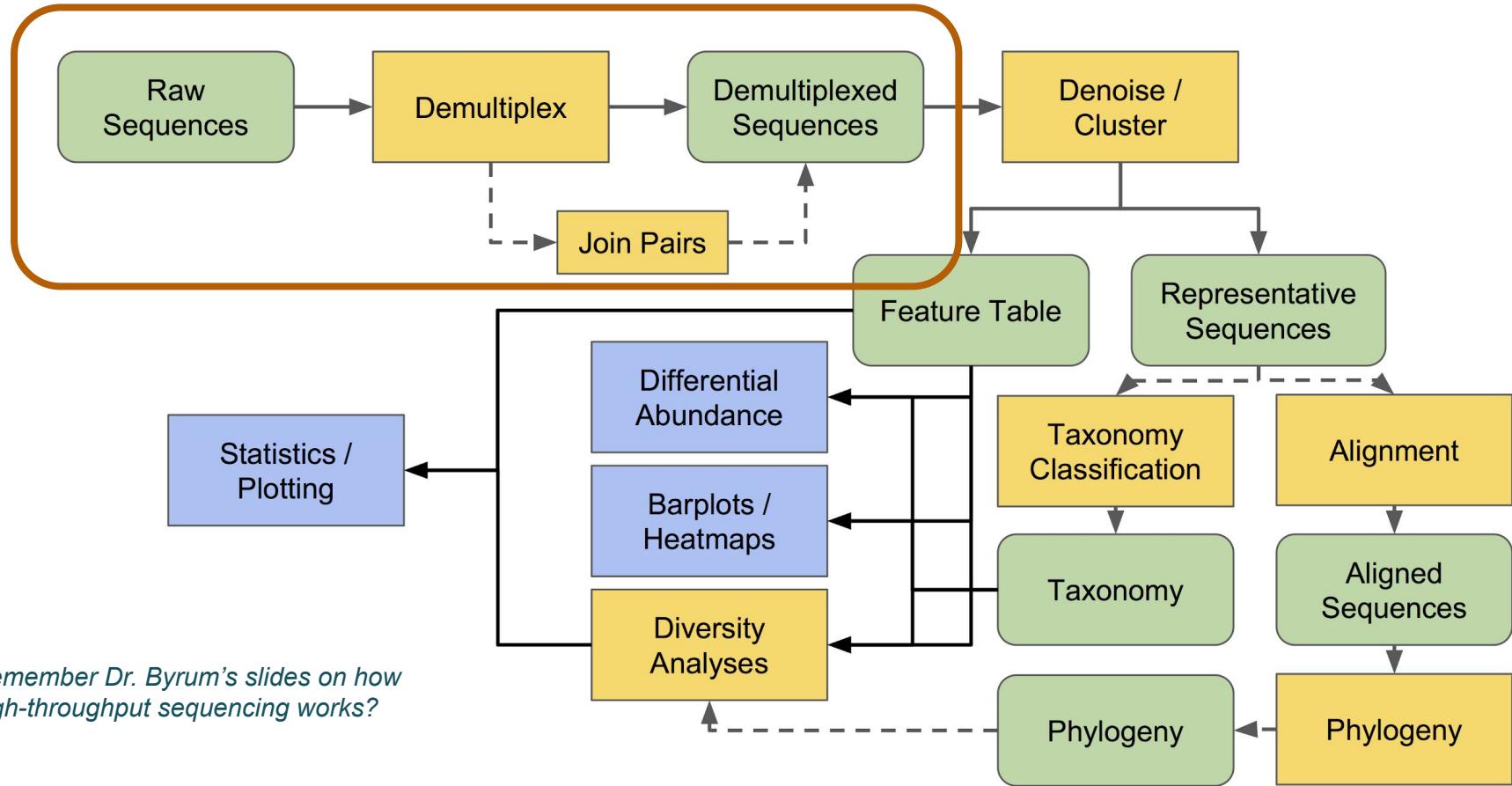
x 3

PD Patient

- 12 donors:
 - 6 healthy
 - with parkinson's disease
- 2 genotypes:
 - Susceptible
 - PD
- 3 cages per donor with mixed mice
- 4 time points per mouse



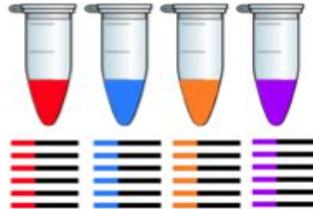
Import and demultiplex data



Remember Dr. Byrum's slides on how high-throughput sequencing works?

Extract DNA, isolate and amplify the rRNA from all samples using barcoded PCR, and sequence.

Barcoded per-sample
rRNA

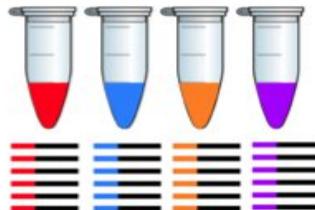


Track per-sample
barcodes (e.g., in
spreadsheet)

| sample-metadata.tsv | |
|---------------------|-----------------|
| SampleID | BarcodeSequence |
| 4ac2 | AACGCAC |
| e375 | AAGAGAT |
| 4gd8 | ACAGCAG |
| 9872 | ACAGCTA |

Extract DNA, isolate and amplify the rRNA from all samples using barcoded PCR, and sequence.

Barcoded per-sample rRNA



Track per-sample
barcodes (e.g., in
spreadsheet)

Pool and sequence samples



| sample-metadata.tsv | |
|---------------------|-----------------|
| SampleID | BarcodeSequence |
| 4ac2 | AACGCAC |
| e375 | AAGAGAT |
| 4gd8 | ACAGCAG |
| 9872 | ACAGCTA |

sequences.fastq(.gz

@HWI-6X_9267:1:1:25:1051
GACGAAGGTGACGACCGTTGCTCGAATCACTGGGCATAAACGCGCGTAGGTG
GCTTGGTAAGTCATGGTAAATCCCTCGGCTAACCGAGGAAC TG
+

barcodes.fastq(.gz)

aa^__[^_ ^__^_ ^_ ^[^__[^_ _zz [^ @HWI-6X_9267:1:1:25:1051
XUWWURZUYY]XXRZRNVRTNTWUUU^ AACGCAC
@HWT-6X 9267:1:1:25:609 +

TACGTAGGGGCAAGCGTTATCCGGATT +
GATGGACAAAGTCTGATGTGAAAGGCTGG bbbbbb
+ @HWI-6X_9267:1:1:25:267
AAGAGAT

aaab`aaa`aaaaaaaaaaaaaaaaaa^aa
+
[] [I^azz^WW^_ ^`ZZ_T]XY^^\^Z
@HWI-6X 9267:1:1:25:519 | @HWI-6X 9267:1:1:25:609

| | |
|---|------------------------------|
| <pre>GACGGAGGATGCAAGTGTATCCGGAAT GTTTACTAAGTCAACTGTTAAATCTTGA +</pre> | <pre>AACGCAC + bbbbbbb</pre> |
|---|------------------------------|

abaaaaaaaa`aaaaaaaa\aaaaaaaa``^`aa @HWI-6X_9267:1:1:25:519
WY]]_Z_XX\\[]]]^`[\XTVX]^T_V
ACAGCAG
@HWI-6X_9267:1:1:25:1109 +

TACGGAGGGTGCAGCGTTAACGGAAT
GTTAGGTAAAGTCAGATGTGAAAGCCCCG
+
ACAGCTA

```
barcodes.fastq(.gz)
```



Demultiplexed sequence data

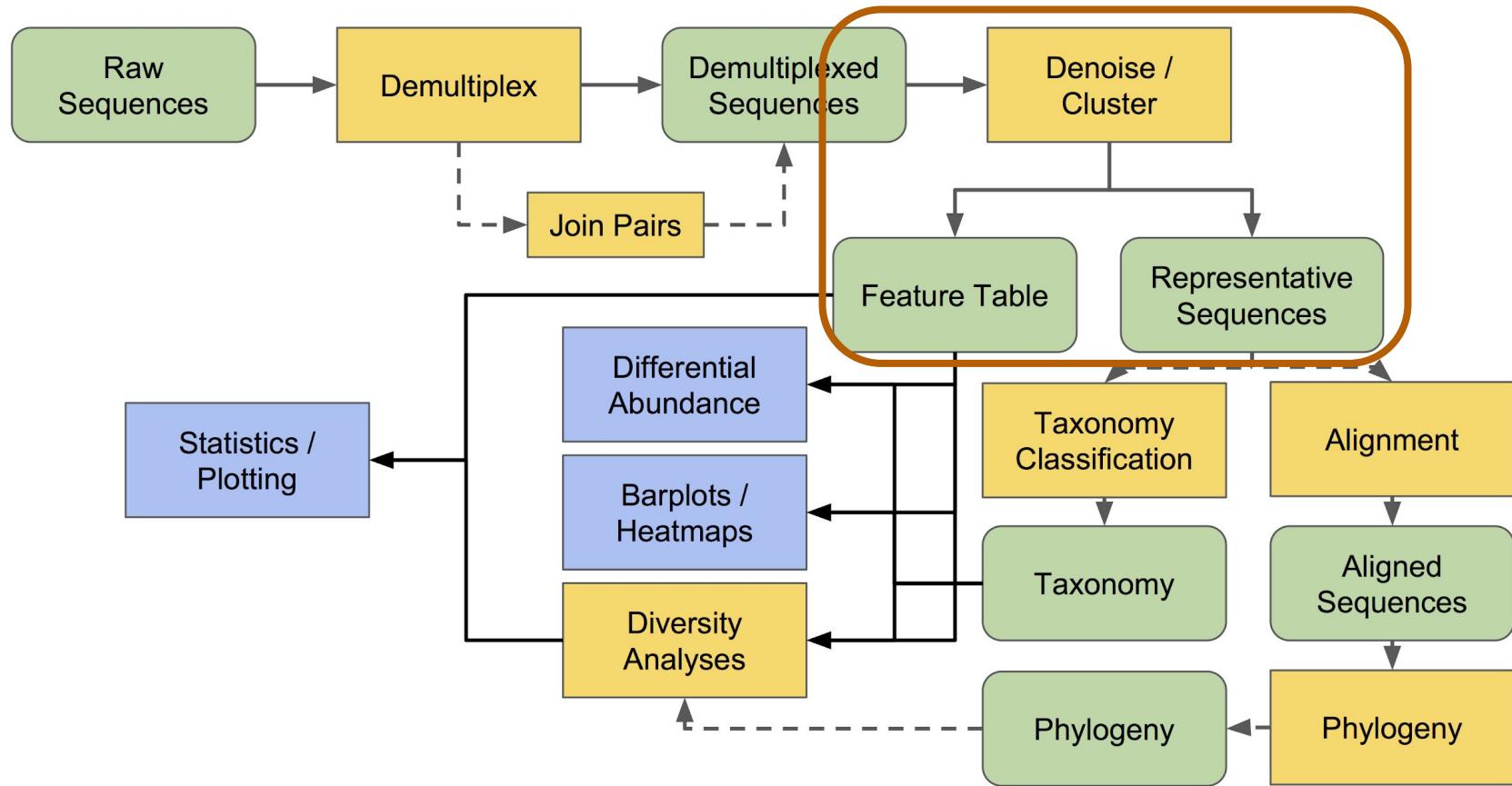
| |
|--|
| 4ac2.fastq(.gz) |
| e375.fastq(.gz) |
| 4gd8.fastq(.gz) |
| 9872.fastq(.gz) |
| <pre> @HWI-6X 9267:1:1:25:1109 TACGGAGGGTGCAGCGTTAACGGAATTACTGGCGTAA AGCGTACGTAGCGGTTAGGTAAAGTCAGATGTGAAAGCCC CGGGCTCCACCTGGGAATGG + aaaba^`a^N `\\ ``a a]Zaa^^\Z`[M]a`[VY a^ X ^ Z]NZ`^]TY\] ^RVH PHOWZM[PTRPTRYUBBBBBB BBBBBBBBBBBBBBBBBBBBBB </pre> |

or

Multiplexed sequence data

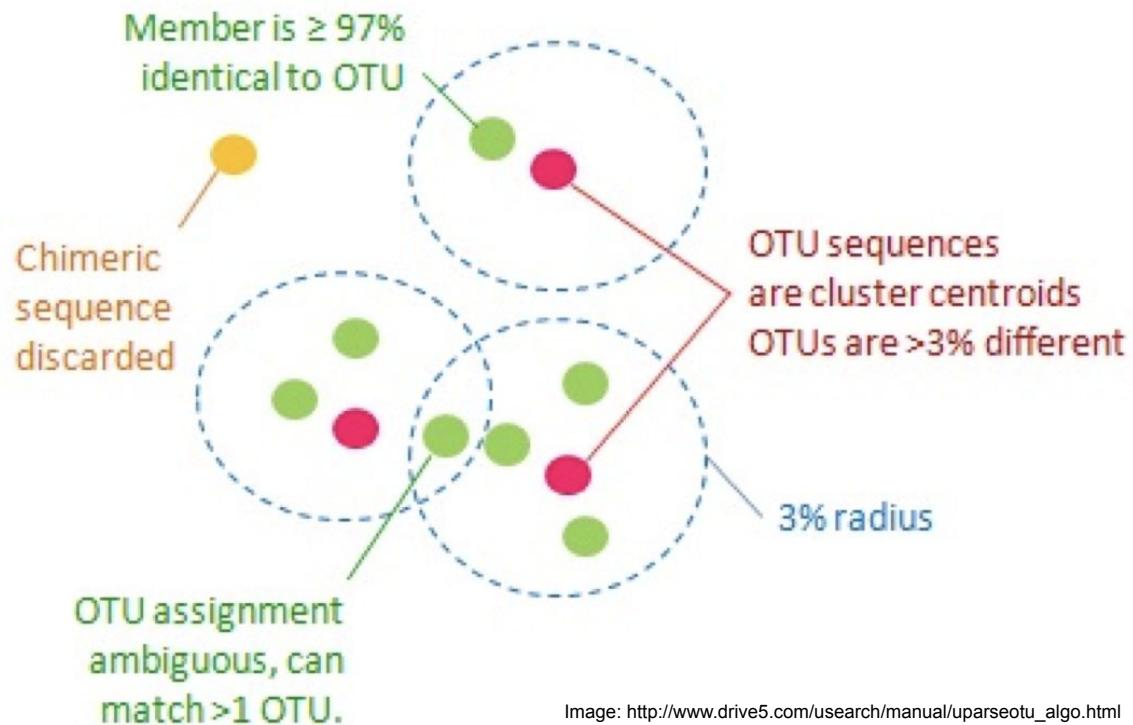
| | | | |
|--|--|----------------------|--|
| sample-metadata.tsv | | | |
| SampleID | BarcodeSequence | | |
| 4ac2 | AACGCAC | | |
| e375 | <table border="1"> <tr> <td>sequences.fastq(.gz)</td> </tr> <tr> <td> <pre> @HWI-6X_9267:1:1:25:1051 GACGAAGGTGACGACCGTTGCTCGGAATCACTGGGCATAAAGCGCGTAGGTG GCTTGGTAAGTCCA + abaaaaaa^`a_]`SVYGYVDXOZVT\T @HWI-6X_9267:1:TACGTATGGGGCAA GTGGCTTAAGCGCA + aa^^[___^__^_ _^_XUWWURZUYYY]XXR @HWI-6X_9267:1:TACGTAGGGGCAA GATGGACAAGTCTG + </pre> </td></tr> </table> | sequences.fastq(.gz) | <pre> @HWI-6X_9267:1:1:25:1051 GACGAAGGTGACGACCGTTGCTCGGAATCACTGGGCATAAAGCGCGTAGGTG GCTTGGTAAGTCCA + abaaaaaa^`a_]`SVYGYVDXOZVT\T @HWI-6X_9267:1:TACGTATGGGGCAA GTGGCTTAAGCGCA + aa^^[___^__^_ _^_XUWWURZUYYY]XXR @HWI-6X_9267:1:TACGTAGGGGCAA GATGGACAAGTCTG + </pre> |
| sequences.fastq(.gz) | | | |
| <pre> @HWI-6X_9267:1:1:25:1051 GACGAAGGTGACGACCGTTGCTCGGAATCACTGGGCATAAAGCGCGTAGGTG GCTTGGTAAGTCCA + abaaaaaa^`a_]`SVYGYVDXOZVT\T @HWI-6X_9267:1:TACGTATGGGGCAA GTGGCTTAAGCGCA + aa^^[___^__^_ _^_XUWWURZUYYY]XXR @HWI-6X_9267:1:TACGTAGGGGCAA GATGGACAAGTCTG + </pre> | | | |
| 4gd8 | <table border="1"> <tr> <td>barcodes.fastq(.gz)</td> </tr> <tr> <td> <pre> @HWI-6X_9267:1:1:25:1051 AACGCAC + bbbbbbb @HWI-6X_9267:1:1:25:267 AAGAGAT + bbbbbbb @HWI-6X_9267:1:1:25:609 AACGCAC + </pre> </td></tr> </table> | barcodes.fastq(.gz) | <pre> @HWI-6X_9267:1:1:25:1051 AACGCAC + bbbbbbb @HWI-6X_9267:1:1:25:267 AAGAGAT + bbbbbbb @HWI-6X_9267:1:1:25:609 AACGCAC + </pre> |
| barcodes.fastq(.gz) | | | |
| <pre> @HWI-6X_9267:1:1:25:1051 AACGCAC + bbbbbbb @HWI-6X_9267:1:1:25:267 AAGAGAT + bbbbbbb @HWI-6X_9267:1:1:25:609 AACGCAC + </pre> | | | |
| 9872 | | | |

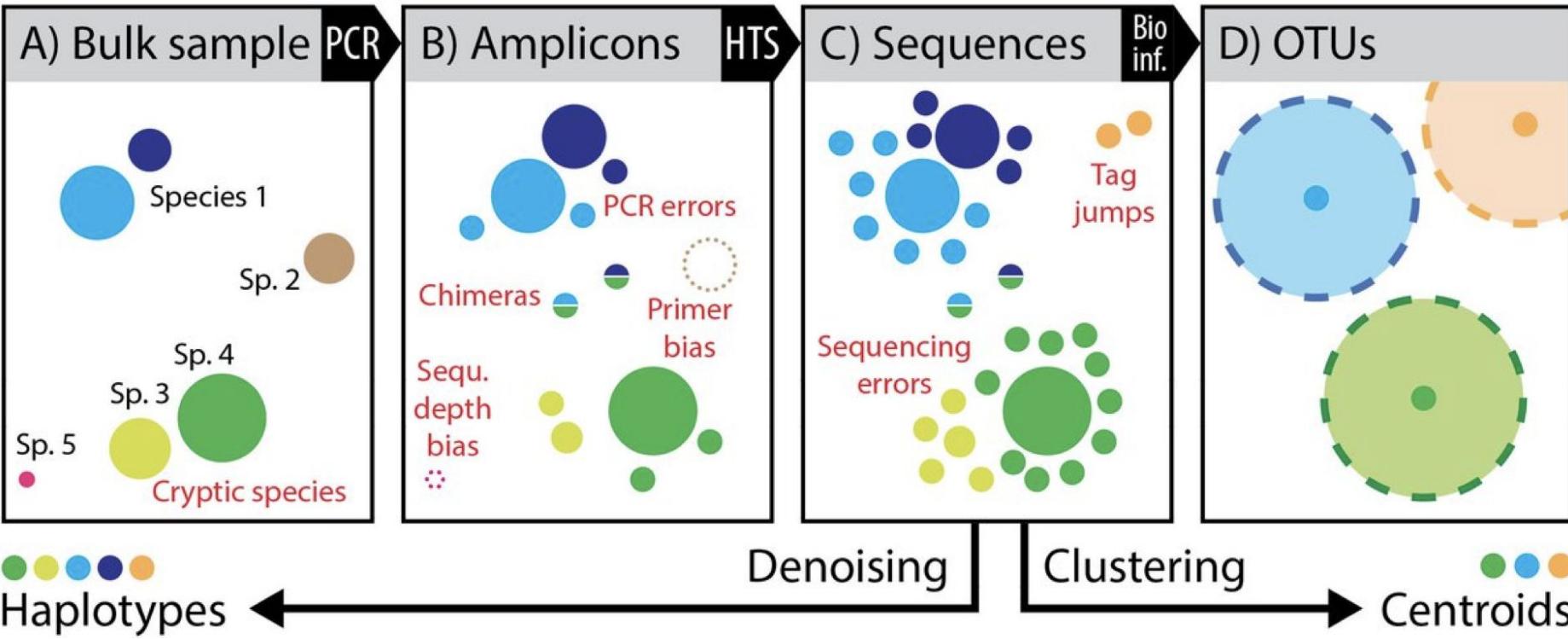
Denoise



- Clustering

- Typical threshold, i.e. 97% similarity, “within species variation”.
- Remove noisy sequences and reduce the amount of sequences to process
- There are different methods (closed or open reference) and algorithms (sortmerna, vclust)

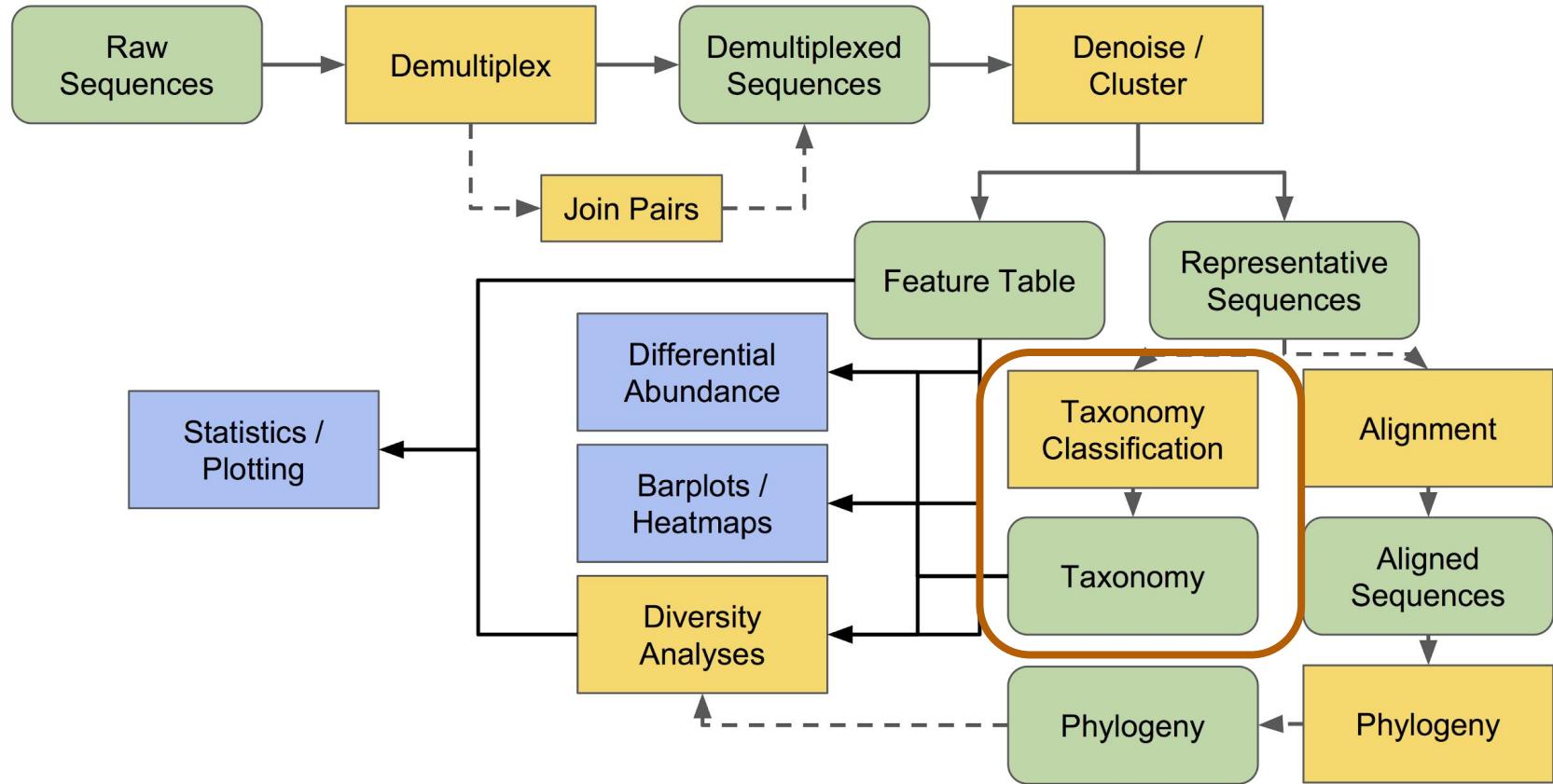




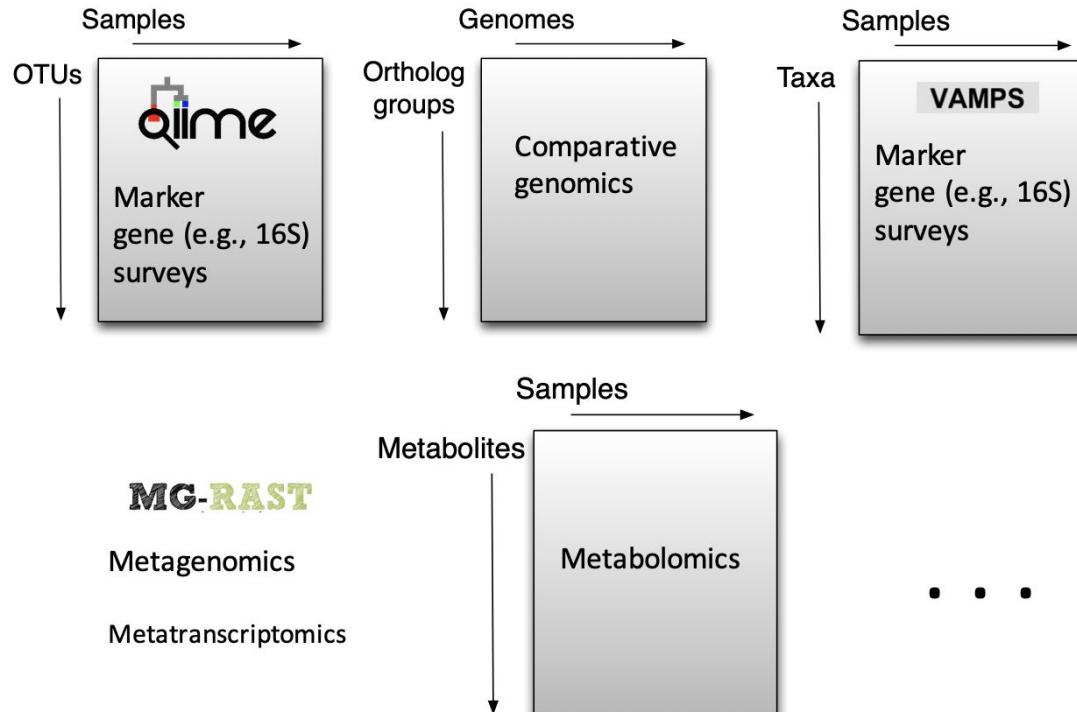
- Remove noise
 - Find the cleanest sequence
 - Correct and/or discard super noisy sequences
 - Examples are: DADA2 and Deblur

Credit: Mehrbod Estaki & QIIME 2 devs.

Taxonomy Classification



sample x observation contingency matrix



Credit: Bob Bowers

Taxonomic assignment of observed sequences.

Reference Database
Silva, Greengenes, etc.

FeatureData [Sequence]

```
>feature5
GACGAAGGTGACGCCGTTGCTCGGAATCACTGGGCATAAAGCGCCGTAGGTGGCTTGGTAAGTCATGGTCAA
ATCCCTCGGCTCAACCGAGGAACCTG
>feature4
TACGTAGGGGCAAGCGTTATCCGATTACTGGGTGAAAGGGAGCGTAGACGGATGGACAAGTCTGATGTGAA
AGGCTGGGCTCAACCCGGACGG
>feature2
TACGTATGGGCAAGCGTTATCCGAATTATGGGCGTAAAGAGTCGTAGGTGGCTTAAGCGCAGGGTTA
AGGAATGGCTTAACCTATTGTTCTC
>feature1
GACGGAGGATCCAAGTGTATCCGAATCACTGGGCATAAGCGCTGTAGGTGGTTACTAACGTAACTGTTAA
ATCTTGAGGCTCAACCTCGAAATCG
>feature3
TACGGAGGGTGCAGCGTTAACGAAATTACTGGGCGTAAAGCGTACGTAGGCCTTAGGTAAGTCAGATGTGAA
AGCCCCGGGCTCCACCTGGGATGG
```

FeatureData [Sequence]

```
>reference-sequence-1
TTGAAGGTGGACGCCGTTGCTCGGAATCACTGGGCATAAAGCGCCGTAGGTGGCTTGGTAAGTCACATGGT
GACTCAACCGAGGAACCTGAAAGTGAGGTGGACGCCGTTGCTCGGAATCACTGGGCATAAAGCGCCGTAGGTG
GCTTGTTAAGTCACATGGTACTAACCGAGGAACCTGAA
>reference-sequence-2
AACGTAGGCAAGCGTTATCCGATTACTGGGTGAAAGGGAGCGTAGACGGATGGACAAGTCTGATGTGAAAGG
CTGGGCTCAACCCGGACGGCTTGTGCTCGGAATCACTGGGCATAAAGCGCCGTAGGTG
```

FeatureData [Taxonomy]

| | |
|----------------------|---|
| reference-sequence-1 | Bacteria; Proteobacteria; Gammaproteobact |
| reference-sequence-2 | Bacteria; Bacteroidetes; Flavobacteria; F |
| reference-sequence-3 | Bacteria; Proteobacteria; Deltaproteobact |
| reference-sequence-4 | Archaea; Euryarchaeota; DSEG; 104A5 |



Compare observed sequences to annotated reference sequences to make taxonomic assignments.



FeatureData [Taxonomy]

| | |
|----------|--|
| feature5 | Bacteria; Proteobacteria |
| feature4 | Bacteria; Proteobacteria |
| feature2 | Bacteria; Bacteroidetes; Flavobacteria; Flavobacteriales |
| feature1 | Bacteria; Proteobacteria |
| feature3 | Bacteria; Proteobacteria; Deltaproteobacteria |

???

Classify Taxonomies

RDP Classifier (naive Bayes)

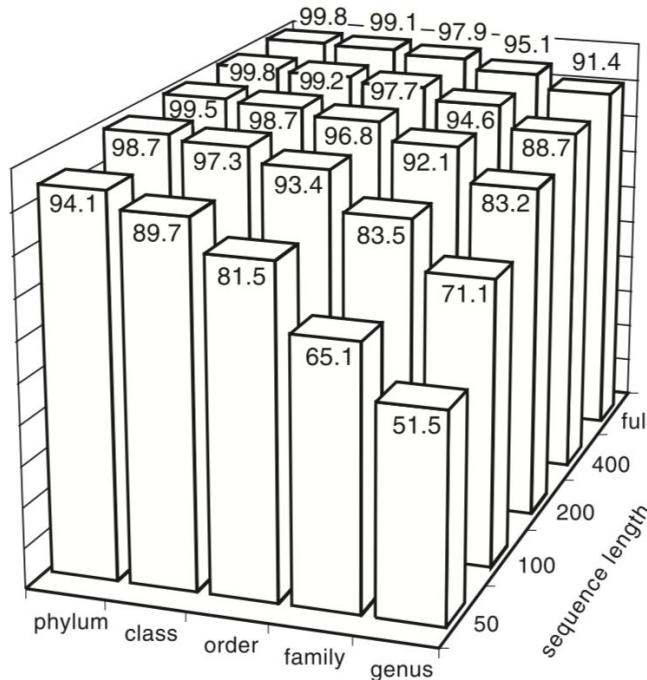
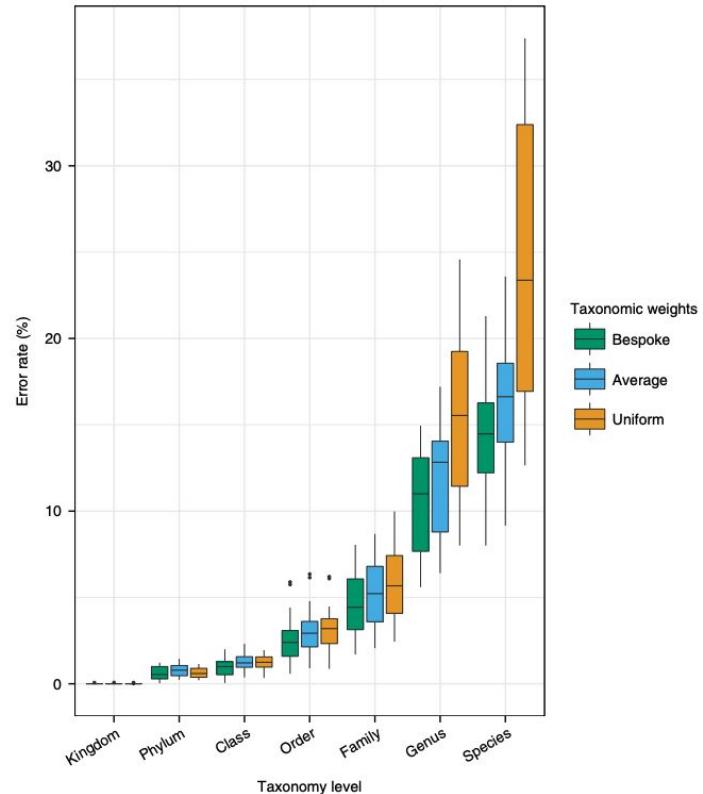


FIG. 1. Overall classification accuracy by query size (exhaustive leave-one-out testing using the Bergey corpus). Numbers are percentages of tests correctly classified.

Naive Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. Wang et al. 2007. Applied and Environmental Microbiology.

Habitat Aware Classification

qiime clawback ...



Kaehler et al. (2019) "Species Abundance Information Improves Sequence Taxonomy Classification Accuracy." Nature Communications 10 (1): 4643.

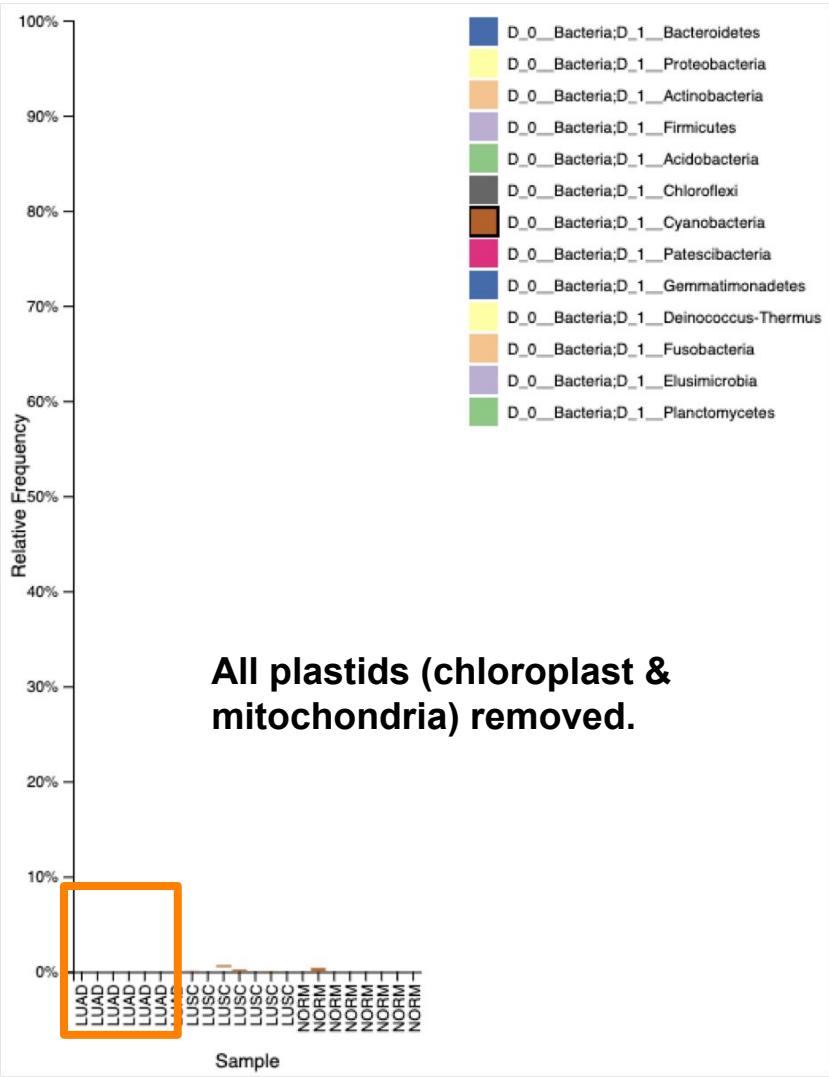
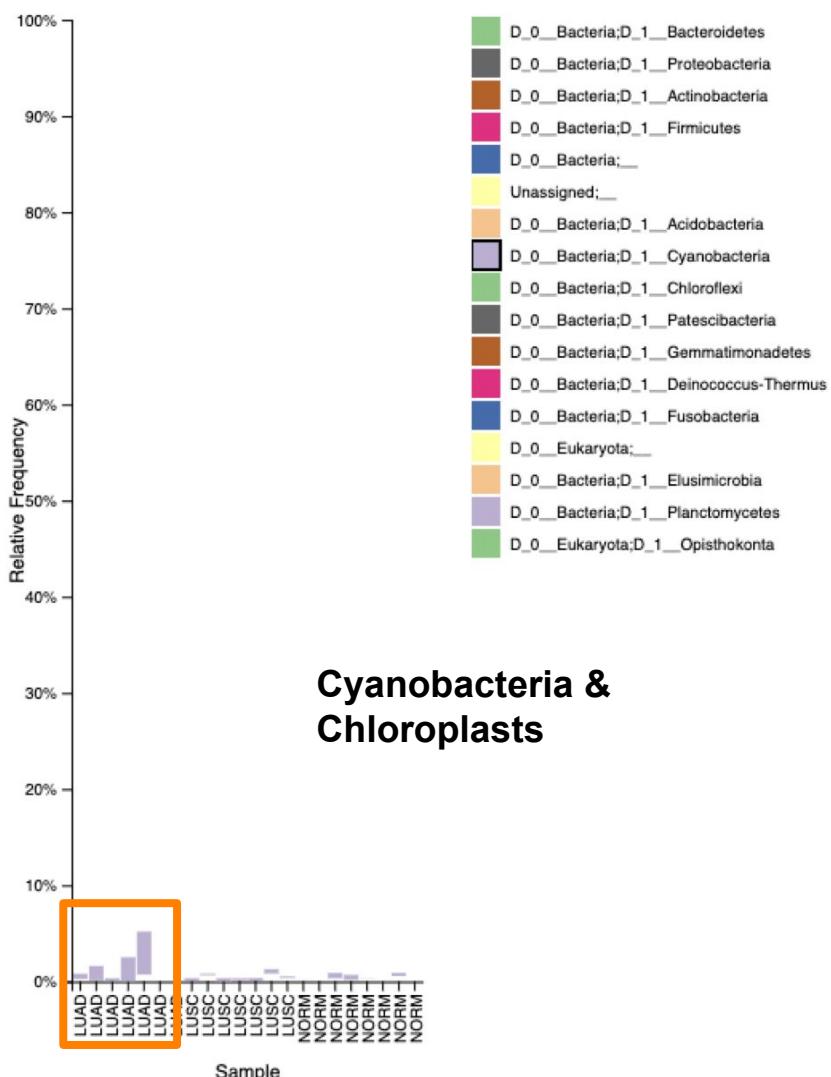
New SILVA taxonomy

d__Bacteria; **p__Proteobacteria**; c__Alphaproteobacteria; o__Rickettsiales; **f__Mitochondria**; **g__Mitochondria**

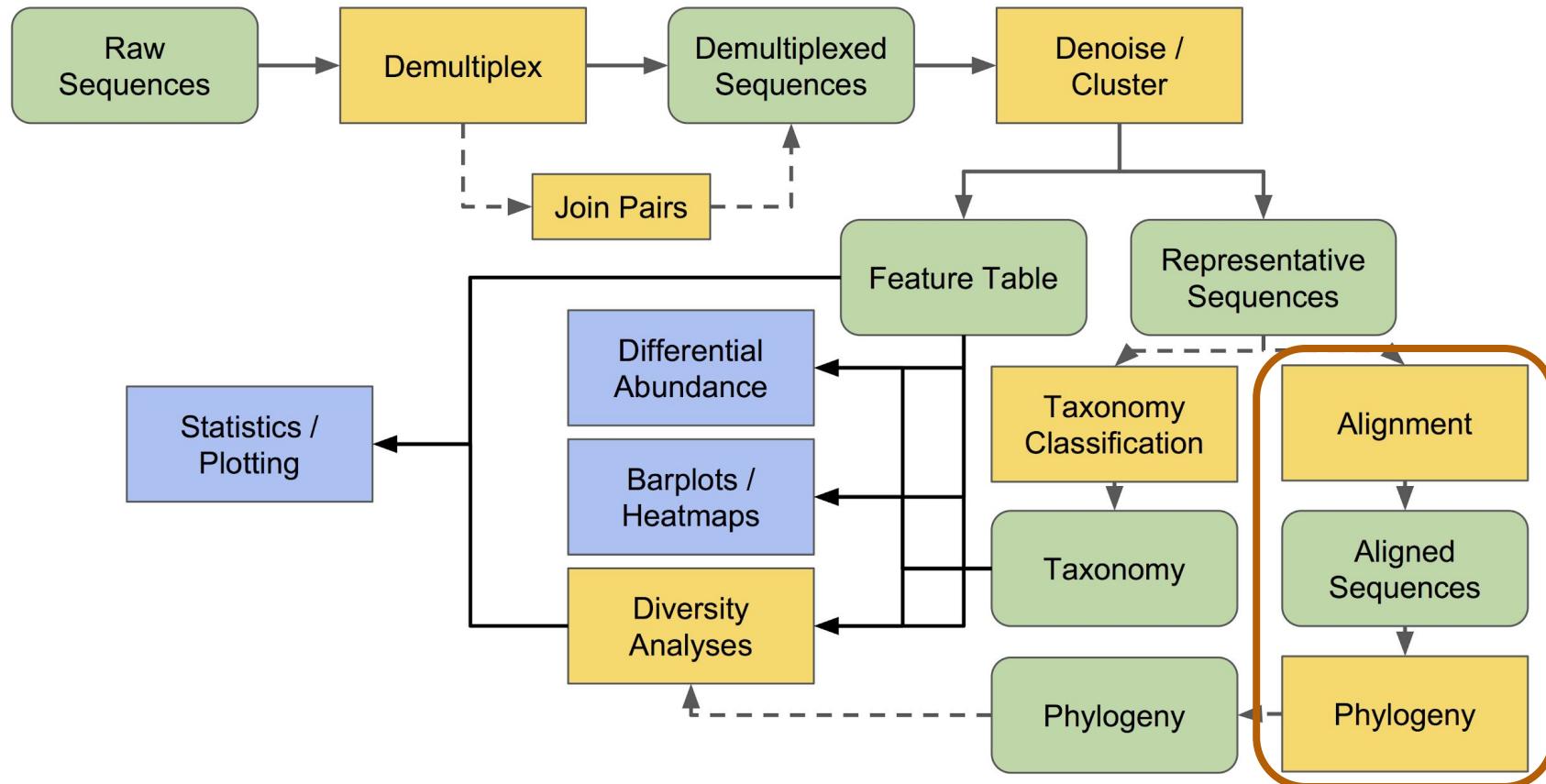
d__Bacteria; **p__Proteobacteria**; c__Alphaproteobacteria; o__Rickettsiales; **f__Mitochondria**; **g__Mitochondria**; s__*Aspergillus_alliaceus*

d__Bacteria; **p__Cyanobacteria**; c__Cyanobacteriia; **o__Chloroplast**; **f__Chloroplast**; **g__Chloroplast**

d__Bacteria; **p__Cyanobacteria**; c__Cyanobacteriia; **o__Chloroplast**; **f__Chloroplast**; **g__Chloroplast**; s__*Phytophthora_lateralis*

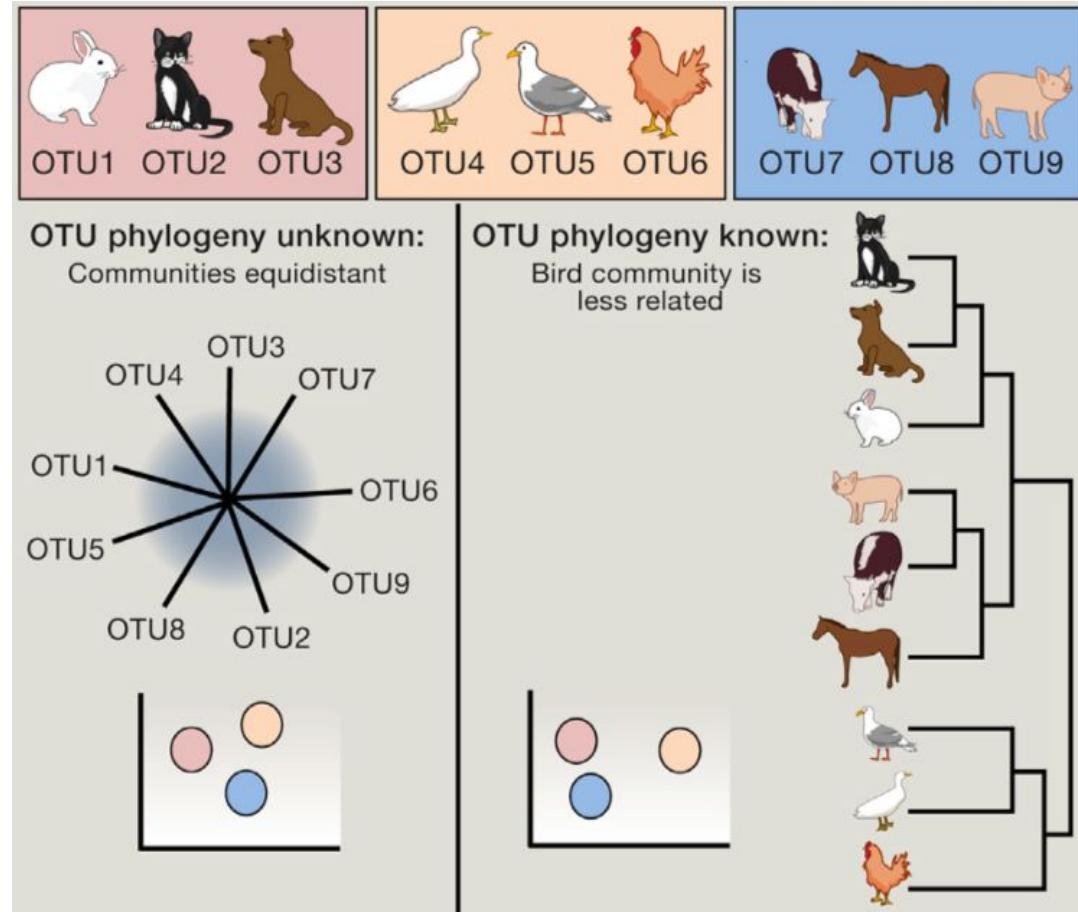


Phylogenetics

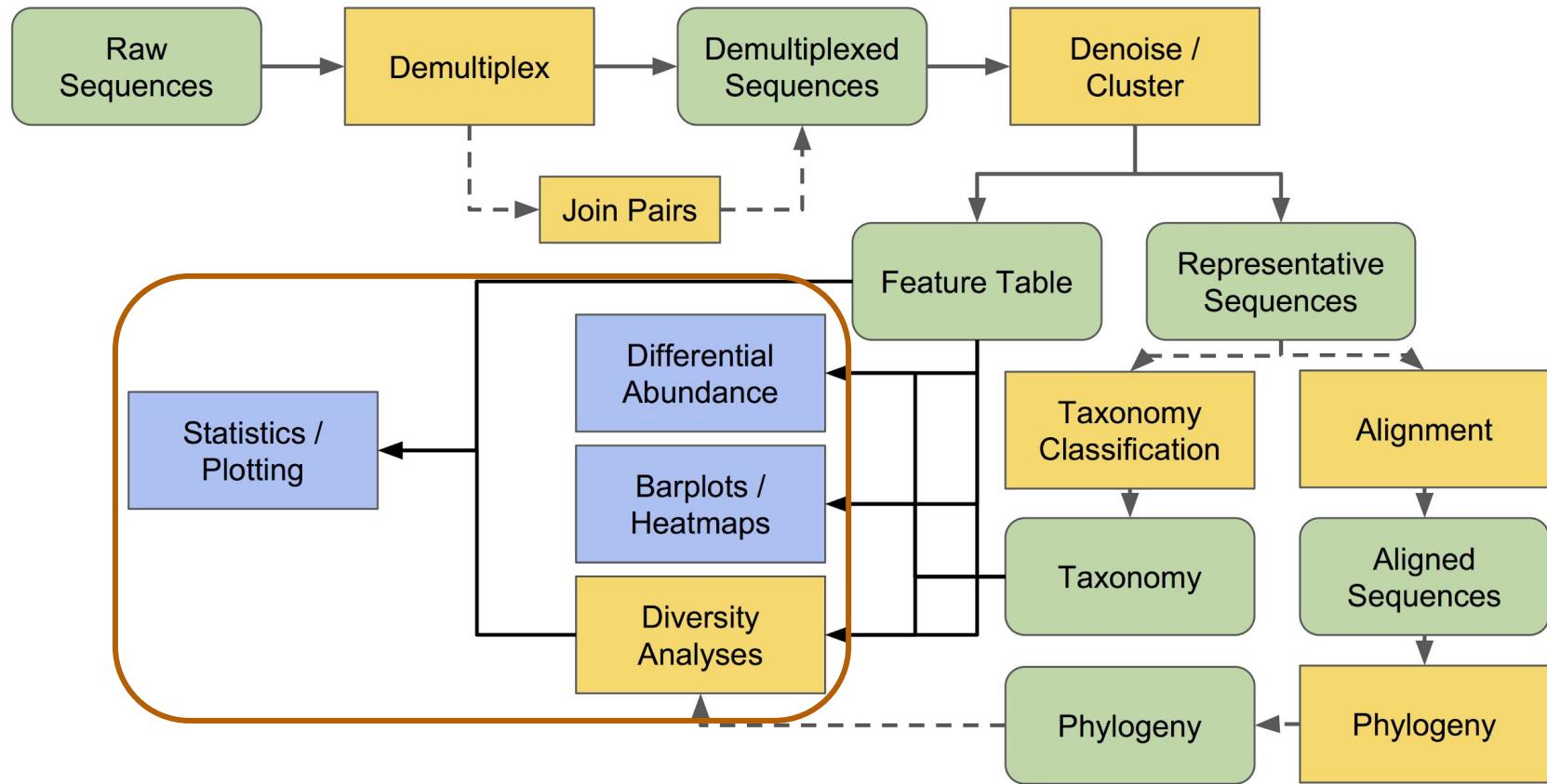


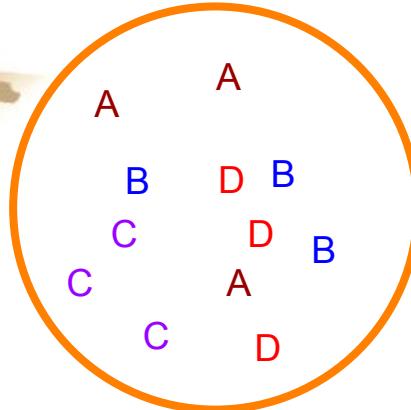
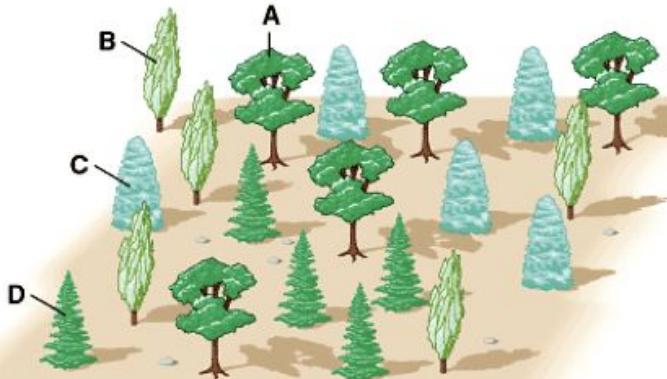
Non-phylogenetic diversity metrics assume that all taxa are equally related, so don't make assumptions about evolutionary relationships.

Phylogenetic diversity metrics incorporate evolutionary relationships between taxa, but assume that we know what those relationships are.

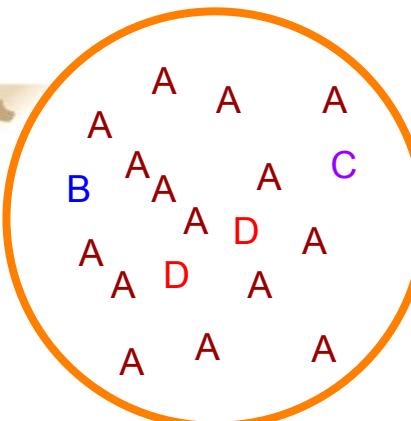


Analysis



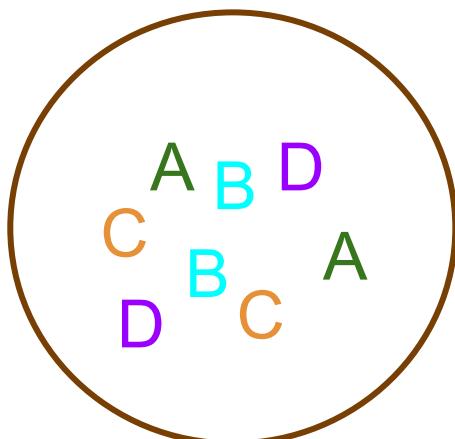


- Evenness.
- Approx. equal representation of taxa within the community.

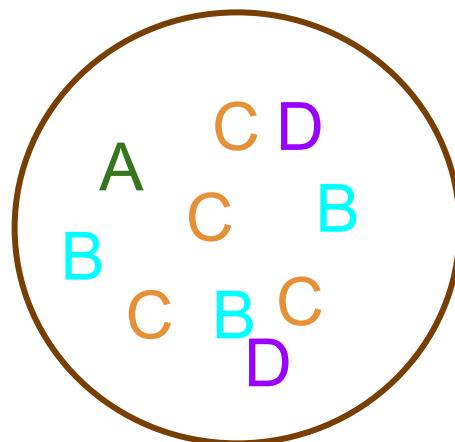


- Not-even / dominance.
- Unequal distribution of taxa.

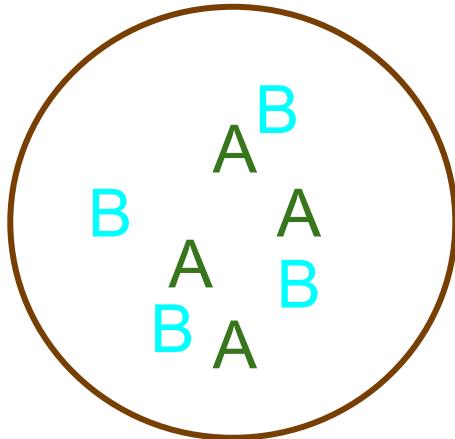
**(A) Rich
and even**



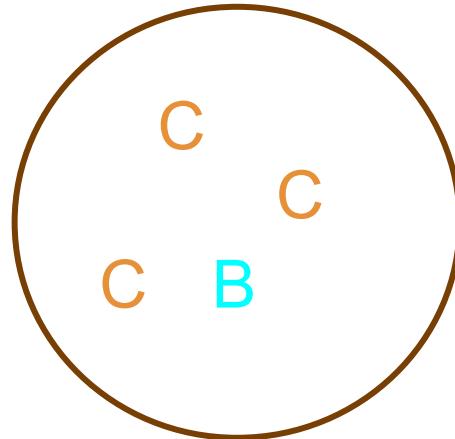
**(B) Rich but
not even**

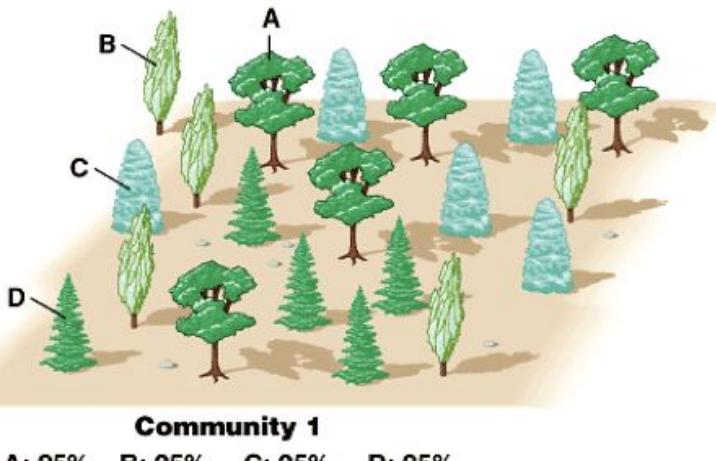


**(C) Not rich
but even**



**(D) Not rich
and not even**





Communities 1 & 2 can be said to have

- A. **same richness**
 - a. both contain same 4 species (presence / absence)
 - b. qualitative metrics

- B. **different abundance / evenness**
 - a. both contain same 4 species, with varying counts
 - b. quantitative metrics

Lozupone et al. 2007. "Quantitative and Qualitative Beta Diversity Measures Lead to Different Insights into Factors That Structure Microbial Communities." *Applied and Environmental Microbiology* 73 (5): 1576–85.

Table 1. Categories of diversity measurements

| | Measurement of diversity within a single community (α -diversity) | Measurement of diversity shared among communities (β -diversity) |
|--|--|--|
| Only presence/absence of taxa considered | <i>Qualitative α-diversity</i> (Richness) <i>Species-based:</i> Chao 1 ACE Rarefaction <i>Divergence-based:</i> Phylogenetic diversity (PD) | <i>Qualitative β-diversity</i> Species-based: Sørensen index Jaccard index <i>Divergence-based:</i> Unweighted UNIFRAC Taxonomic similarity (Δ_S) |
| Additionally accounts for the number of times that each taxon was observed | <i>Quantitative α-diversity</i> (Richness and/or Evenness) <i>Species-based:</i> Shannon's index Simpson's index <i>Divergence-based:</i> θ | <i>Quantitative β diversity</i> <i>Species-based:</i> Sørensen quantitative index Morisita-Horn measure <i>Divergence-based:</i> Weighted UNIFRAC F_{ST} DPCoA |

Also see this wonderful community post on diversity measures:

<https://forum.qiime2.org/t/alpha-and-beta-diversity-explanations-and-commands/2282>

Lozupone & Knight. 2008. "Species Divergence and the Measurement of Microbial Diversity." *FEMS Microbiology Reviews* 32 (4): 557–78.

Variation in sequencing depth can be a problem.



<https://qiime2.org/>

<https://www.youtube.com/c/QIIME2>

TED Talks about the microbiome

Jonathan Eisen, Ph.D.:

https://www.ted.com/talks/jonathan_eisen_meet_your_microbes?language=en



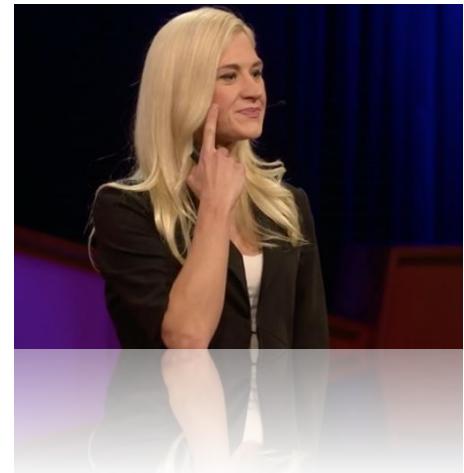
Rob Knight, Ph.D.:

https://www.ted.com/talks/rob_knight_how_our_microbes_make_us_who_we_are?language=en



Anne Madden, Ph.D.:

https://www.ted.com/talks/anne_madden_meet_the_microscopic_life_in_your_home_and_on_your_face?language=en



Some great reading!

