



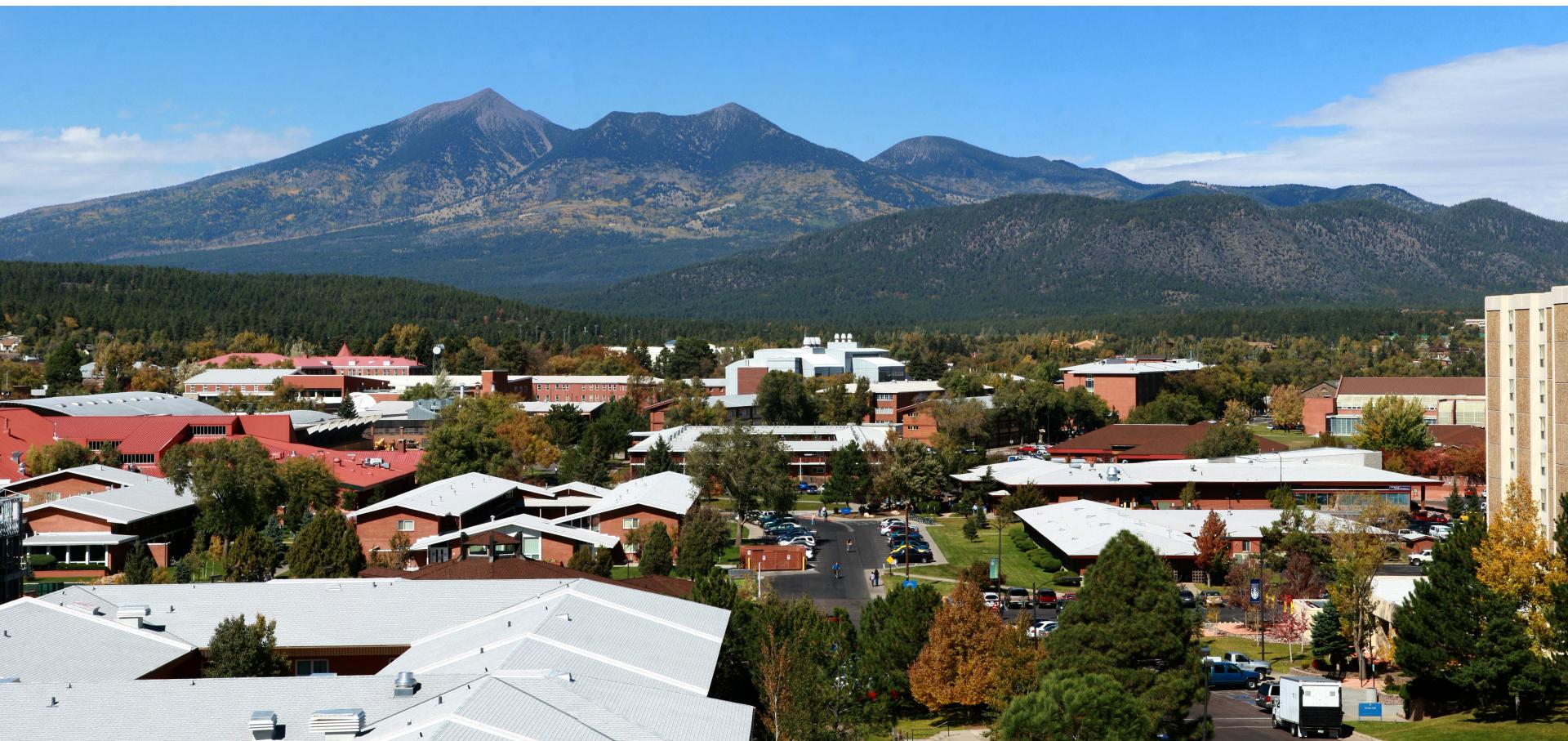
Quantitative Insights Into Microbial Ecology

Mahidol University QIIME Workshop

December 14th-15th, 2015

William (Tony) Walters
Jai Ram Rideout

Northern Arizona University



Agenda

<http://bit.ly/1QfnzBb>

Outline for today's session

December 14th

- 08:00 - 08:30 Intro to microbial community analysis
- 08:30 - 09:00 Intro to QIIME: installation and usage
- 09:00 - 09:30 Preprocessing: metadata and demultiplexing
- 09:30 - 09:45 Break
- 09:45 - 11:15 Preprocessing: OTU picking methods
- 11:15 - 12:00 Diversity and statistical analysis Part I
- 12:00 - 13:30 Lunch break
- 13:30 - 15:30 Diversity and statistical analysis Part II
- 15:30 - 17:00 Advanced topics

First, some background about the microbiome and microbiome studies

A microbe dominated world



The *small subunit ribosomal RNA gene* is frequently used to “fingerprint” different microbial organisms.

Why this gene?

- It's ubiquitous.
- Contains regions conserved across organisms, and regions that are variable across organisms.

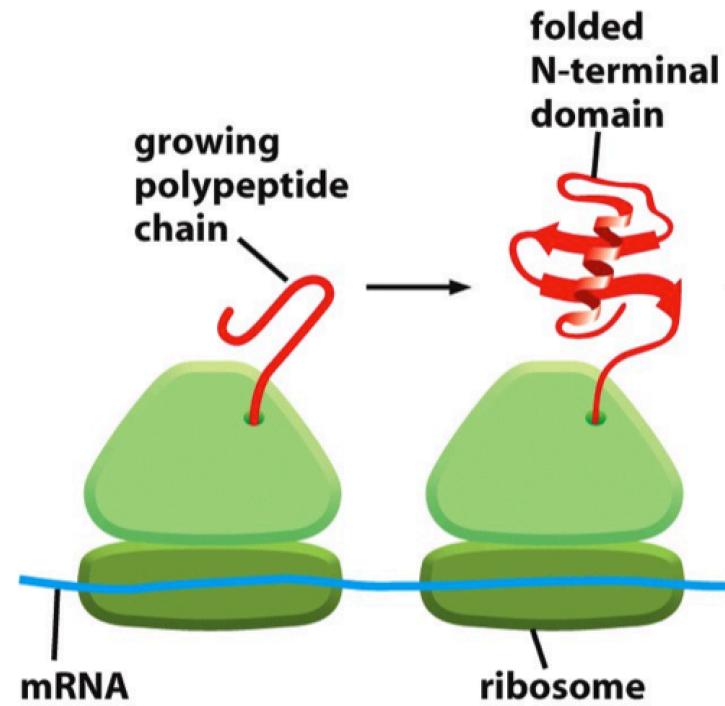
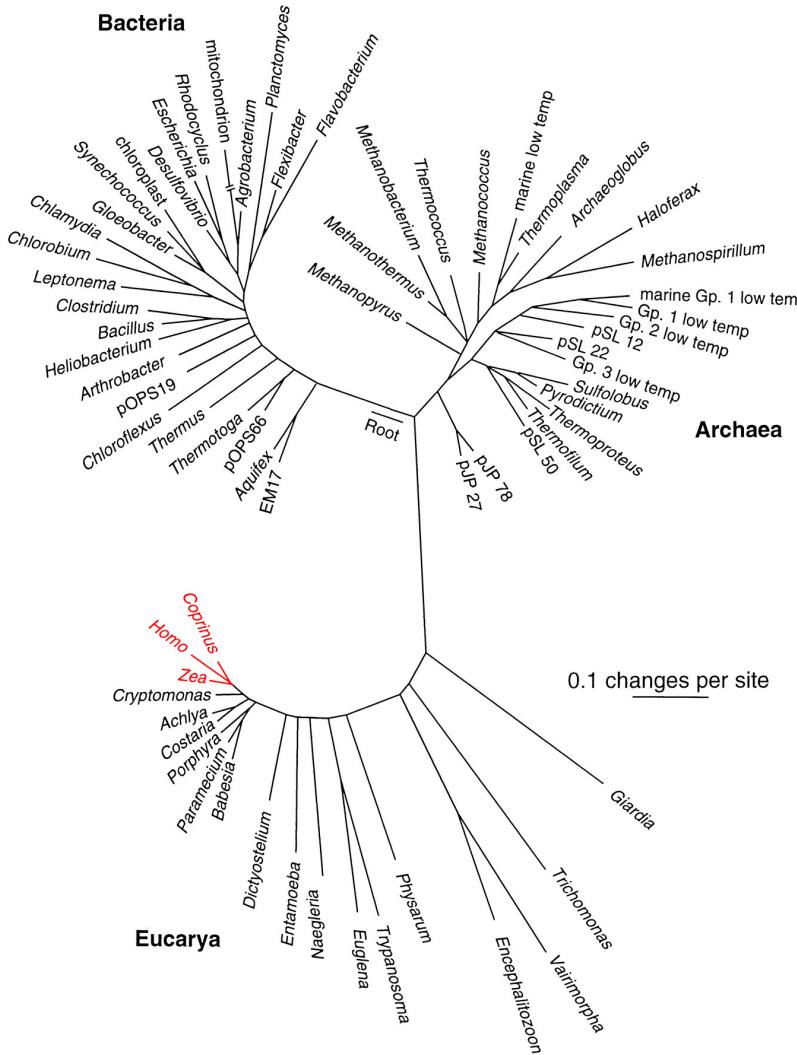


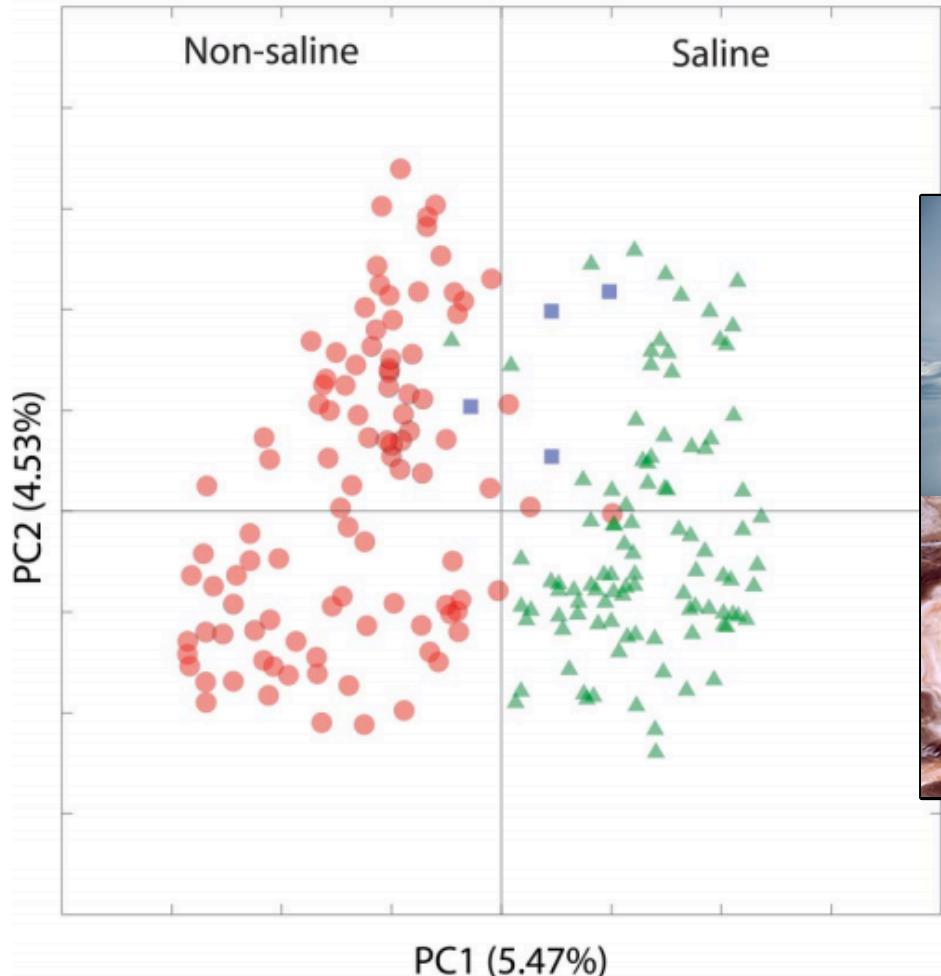
Figure 6-84 *Molecular Biology of the Cell* (© Garland Science 2008)

A microbe dominated world



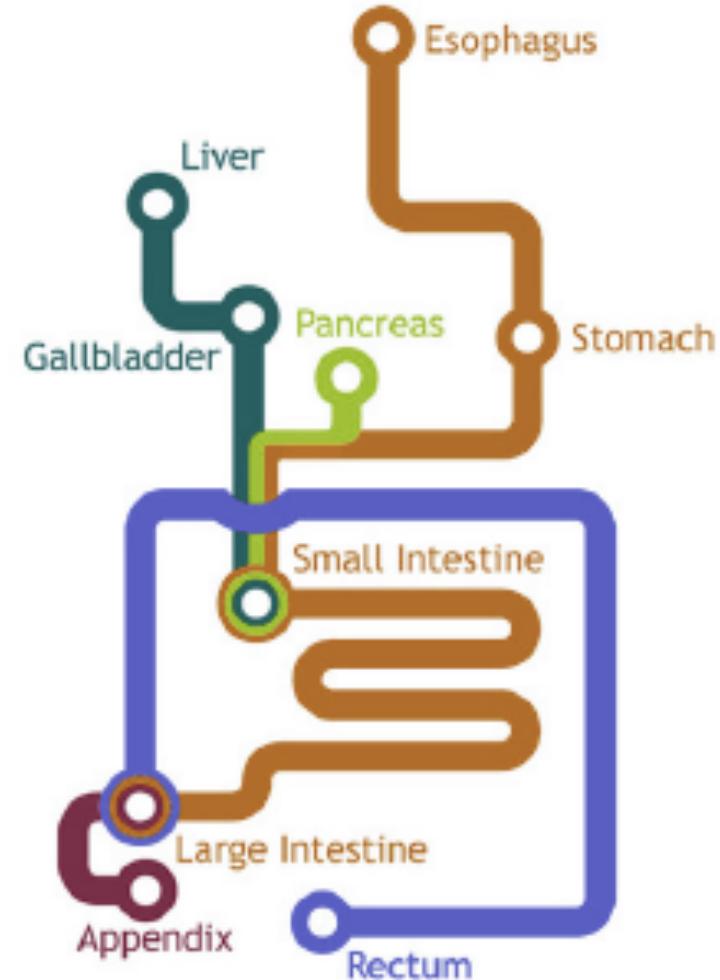
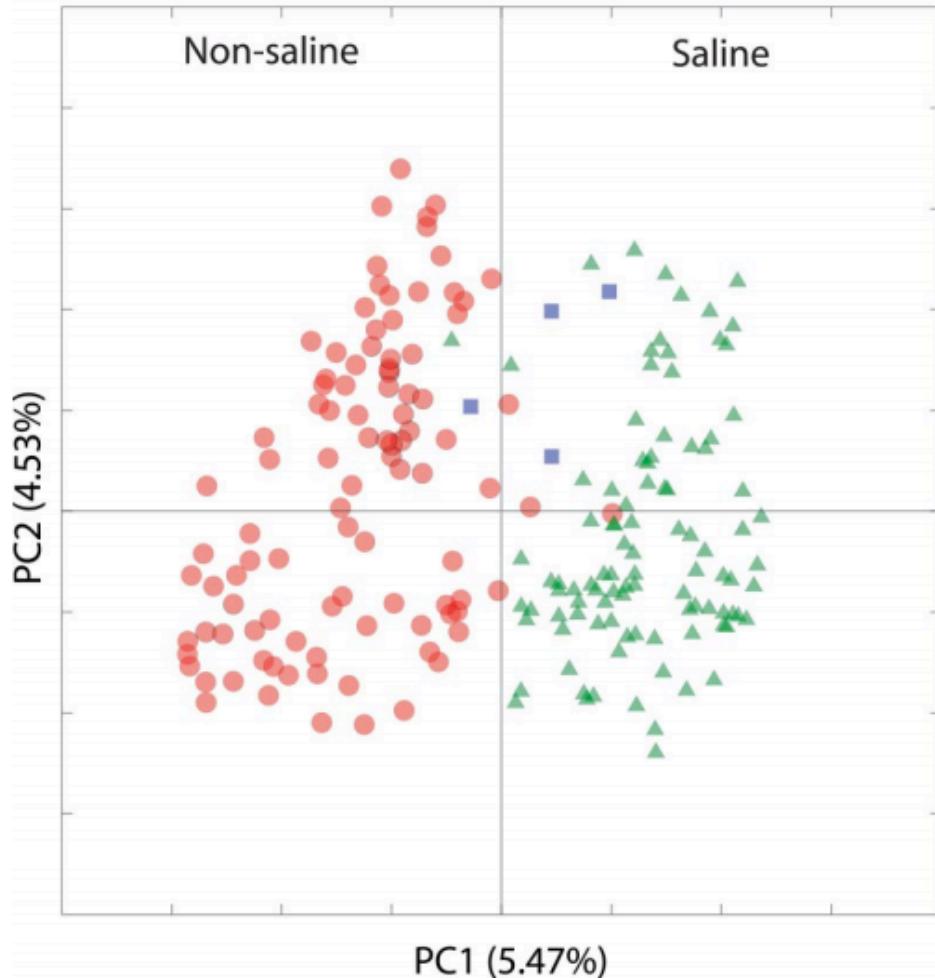
A microbe dominated world

PC1 vs PC2



A microbe dominated world

PC1 vs PC2

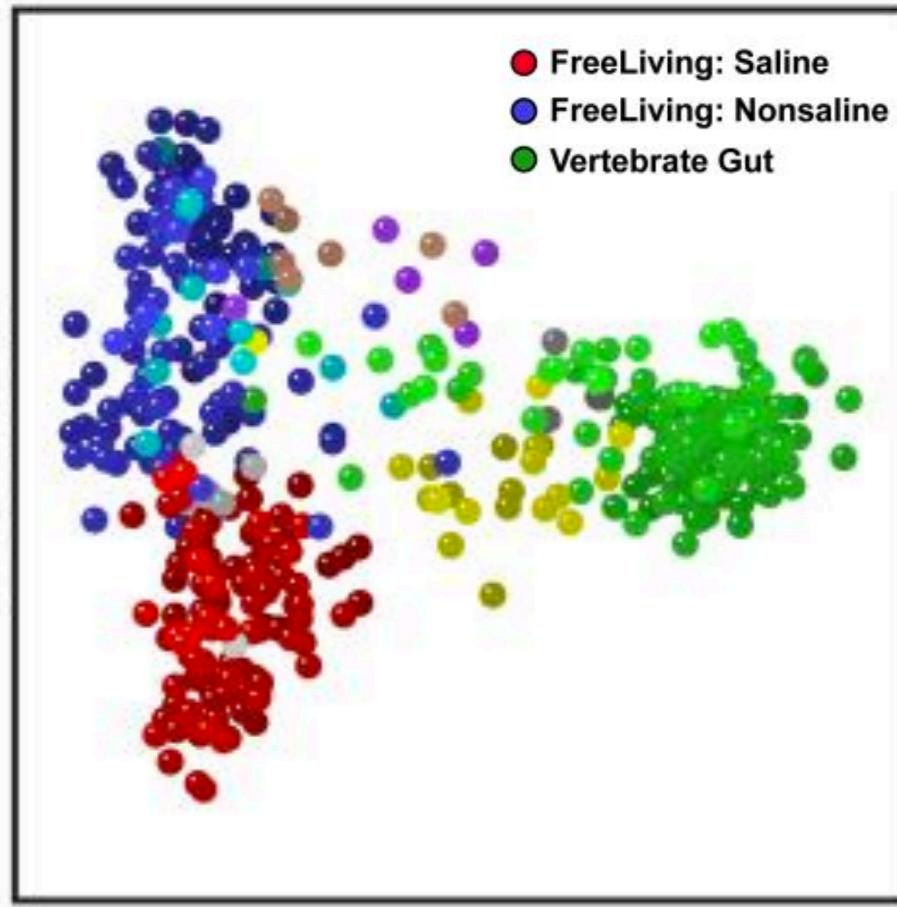


A microbe dominated world

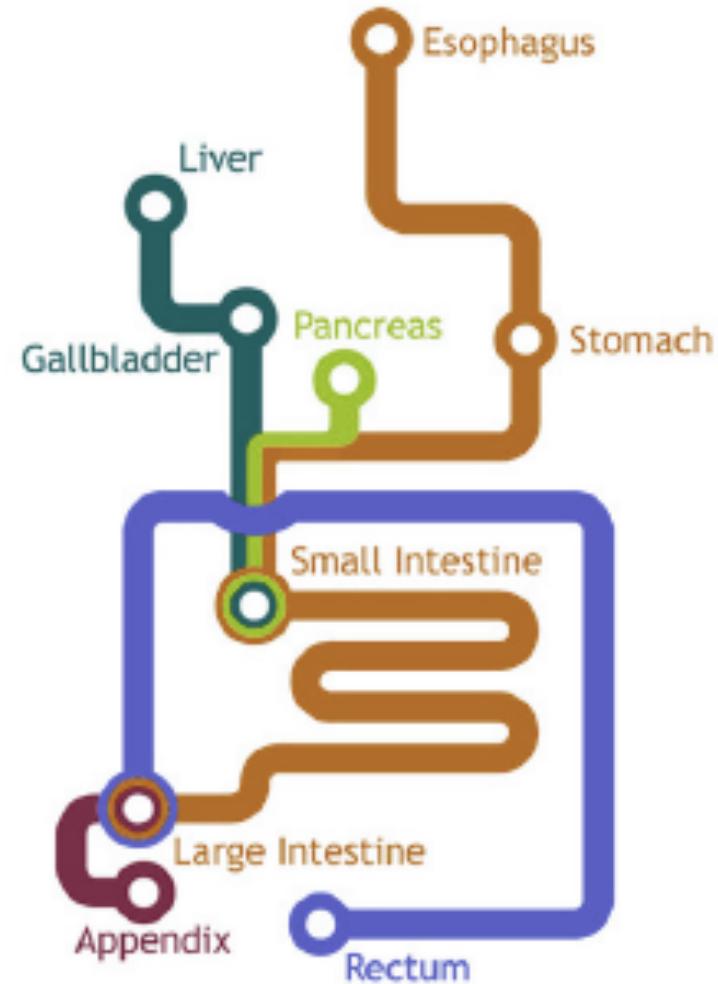
B

All Bacteria

PC3: 2.56%

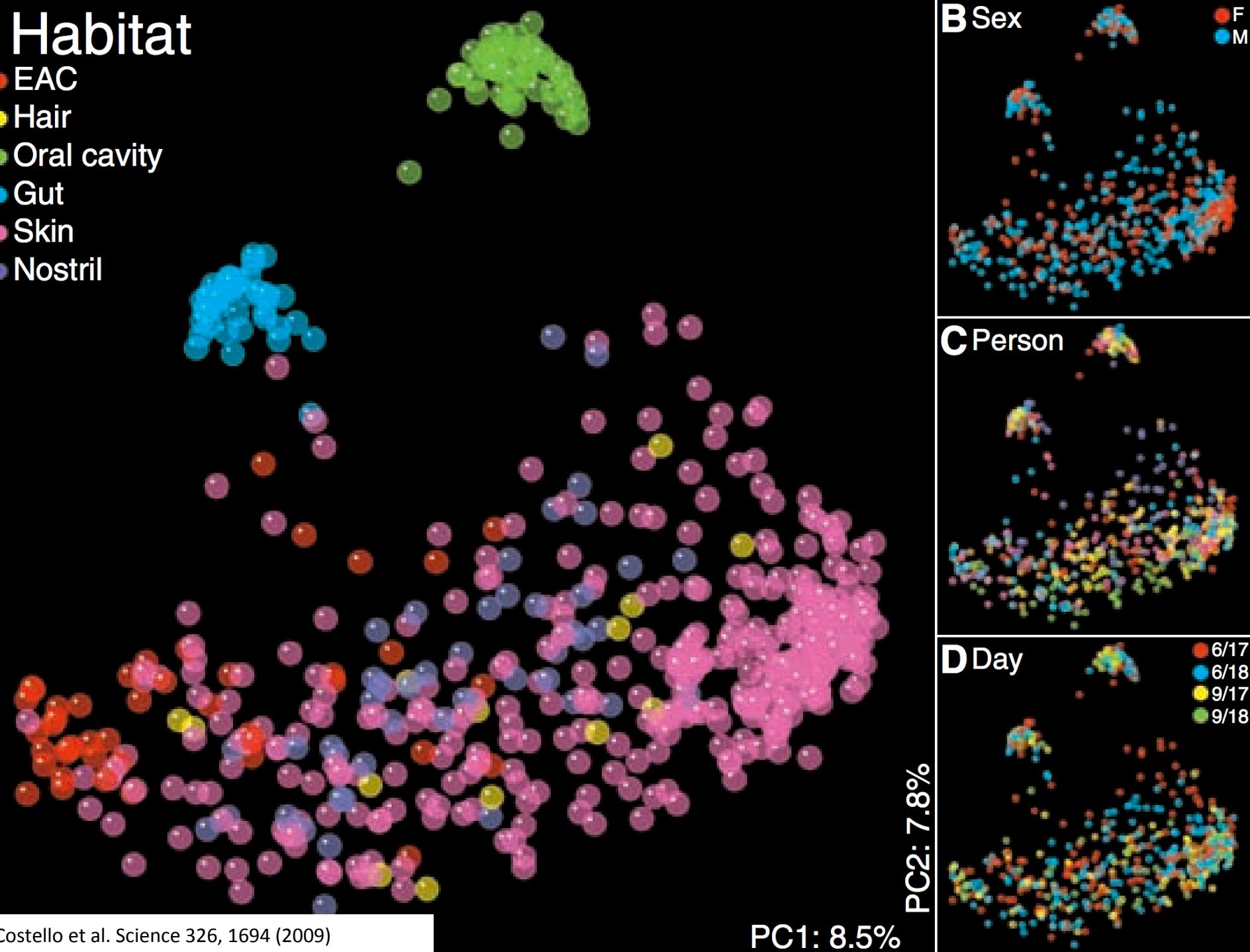


PC1: 7.32%



Habitat

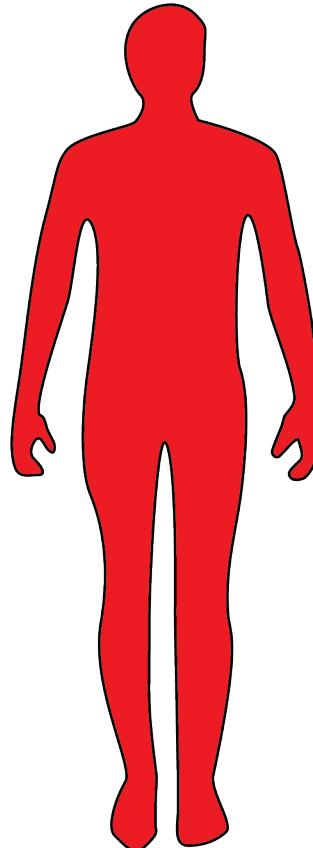
- EAC
- Hair
- Oral cavity
- Gut
- Skin
- Nostril



How “human” are we, really?

Human

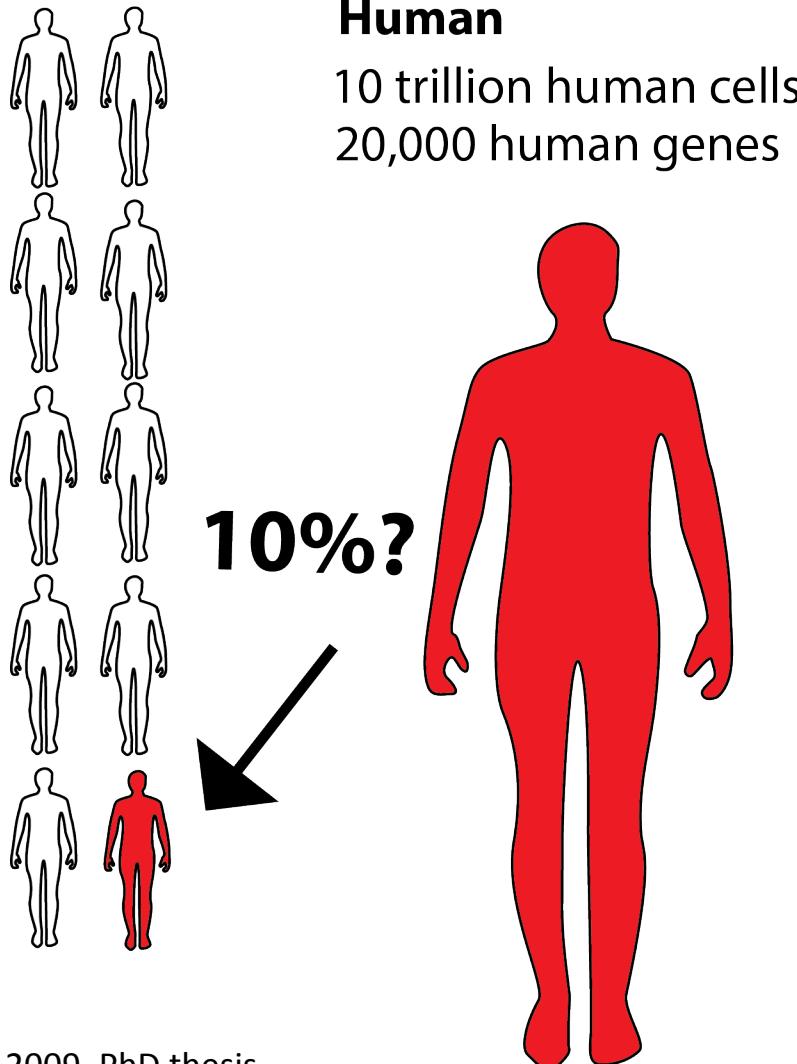
10 trillion human cells
20,000 human genes



How “human” are we, really?

Microbiota

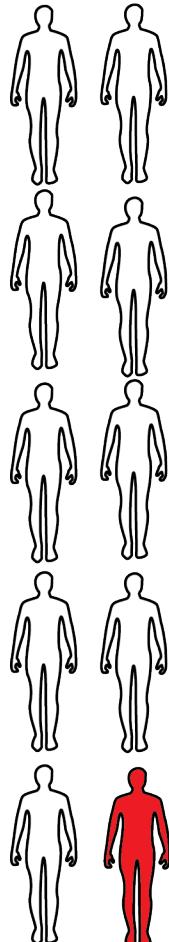
100 trillion microbial cells



How “human” are we, really?

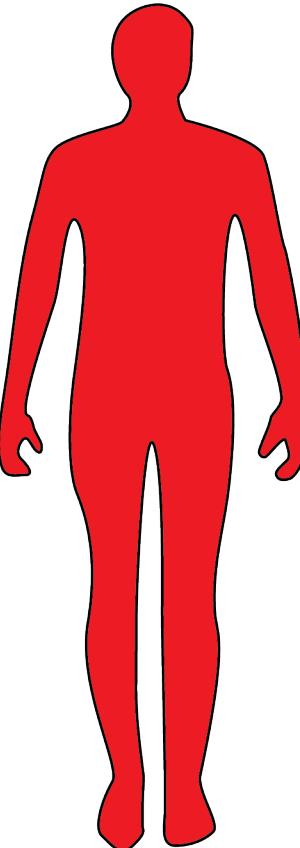
Microbiota

100 trillion microbial cells



Human

10 trillion human cells
20,000 human genes



10%?

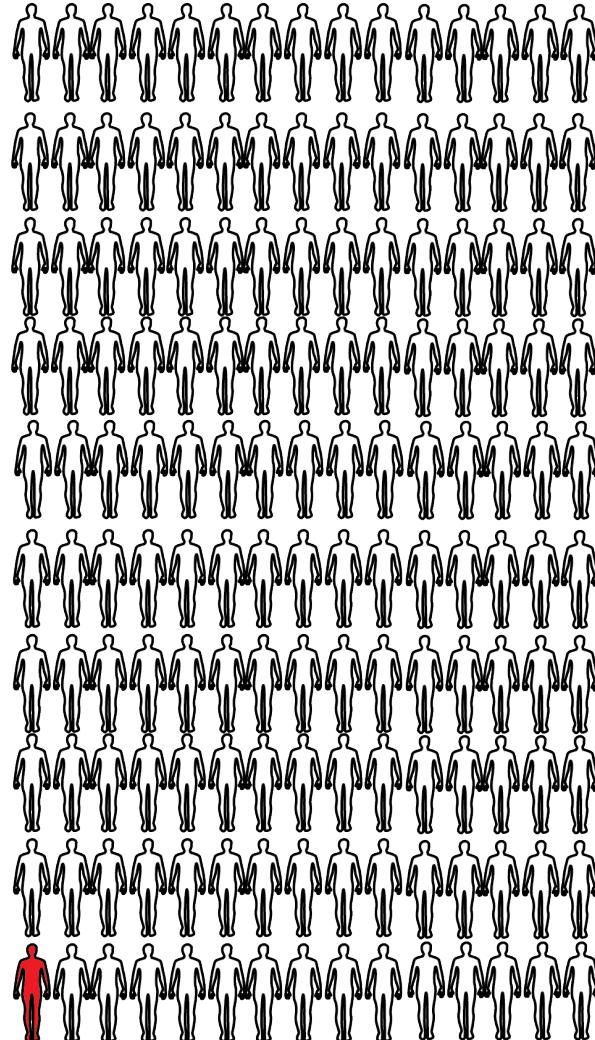


<1%?



Microbiome

3,000,000 microbial genes





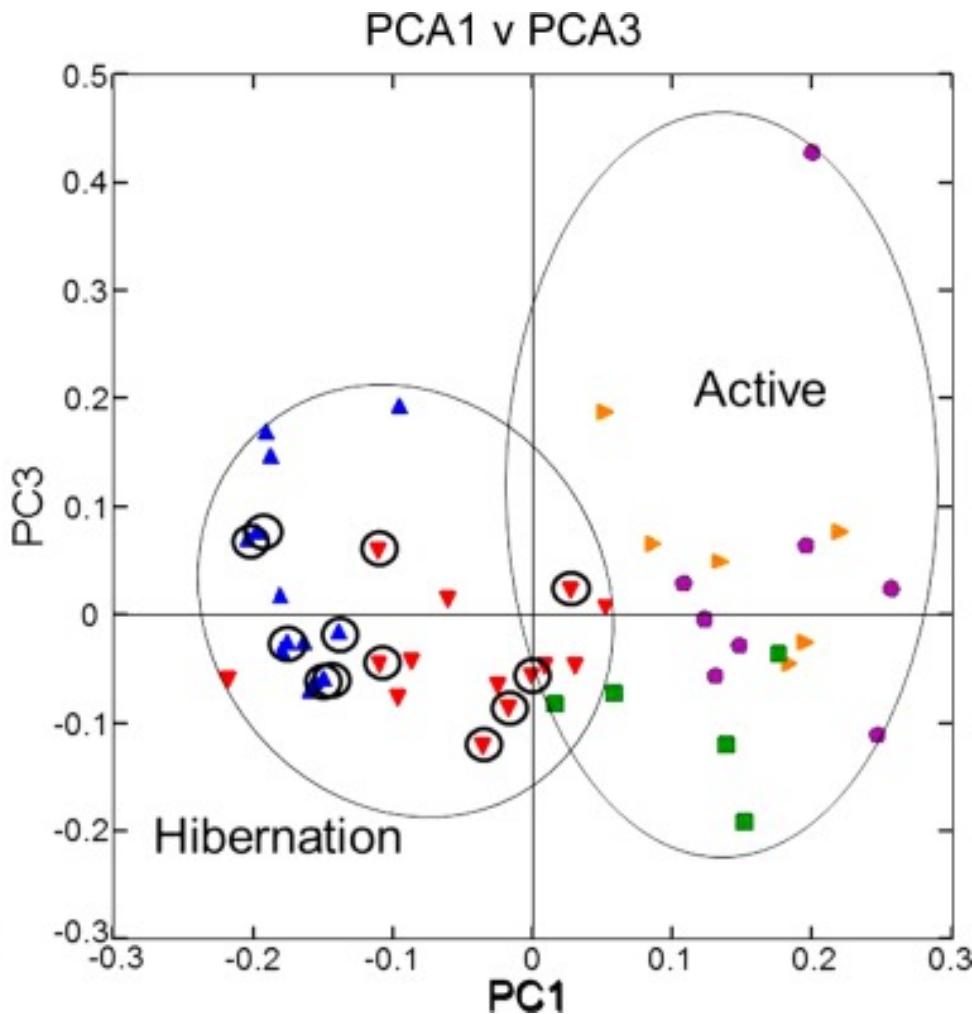
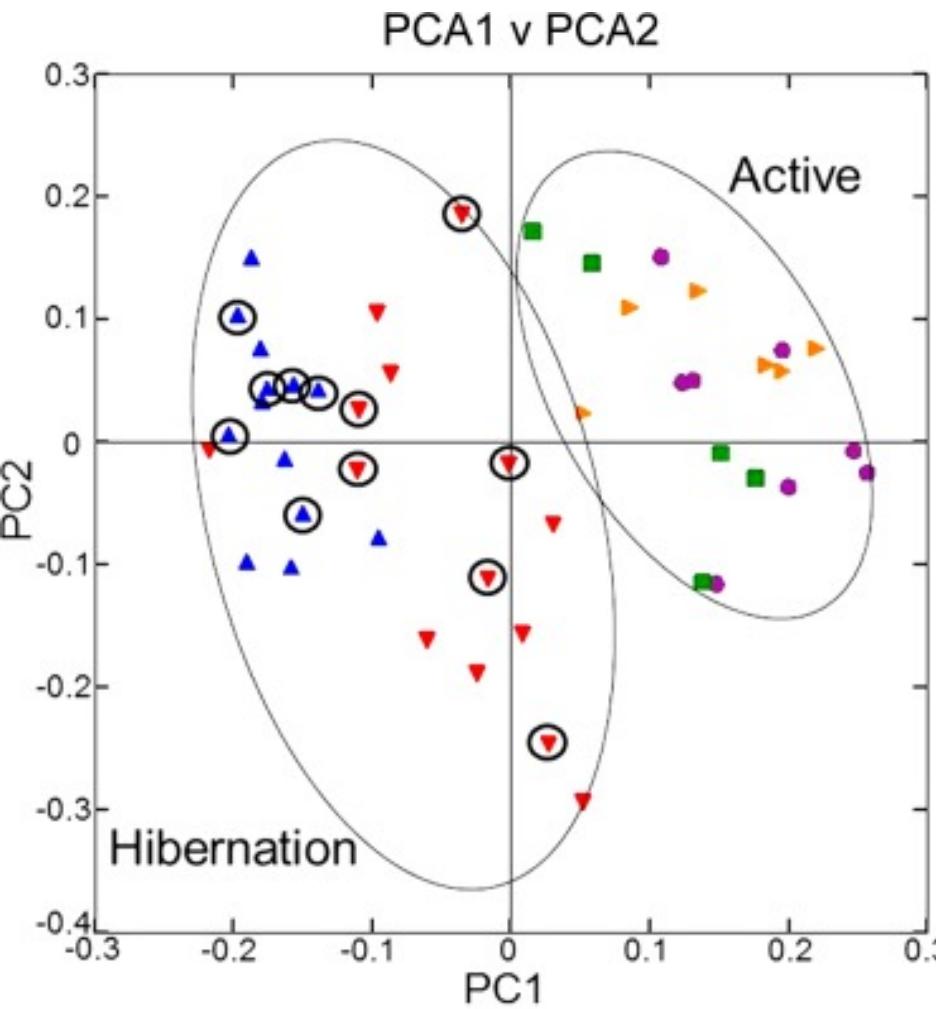
"Your gut is infallible, it is like the pope of your torso."

—Stephen Colbert

Hibernating Ground Squirrels



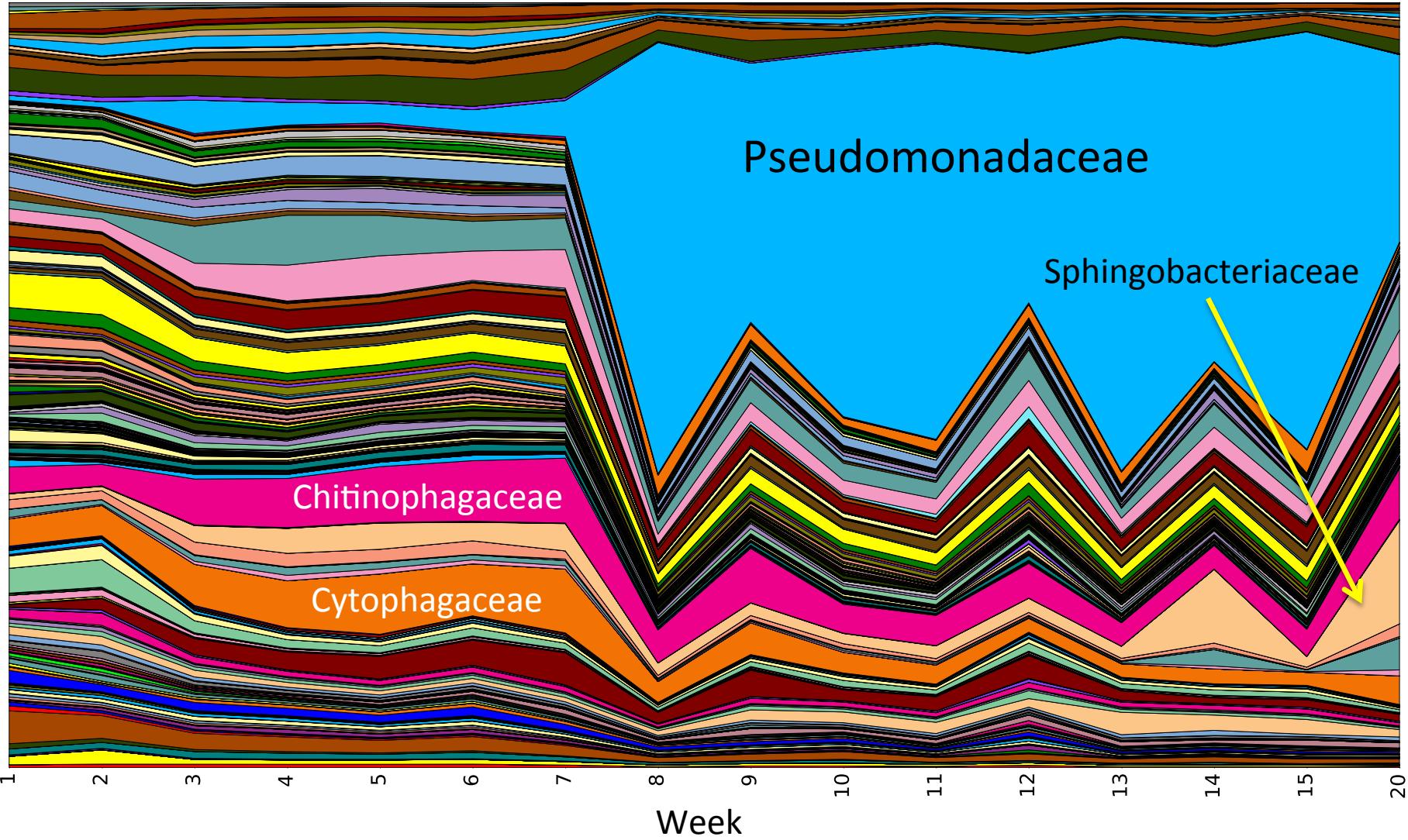
Results of ground squirrel cecal microbiota clustering from Hannah Carey et al, 2012.



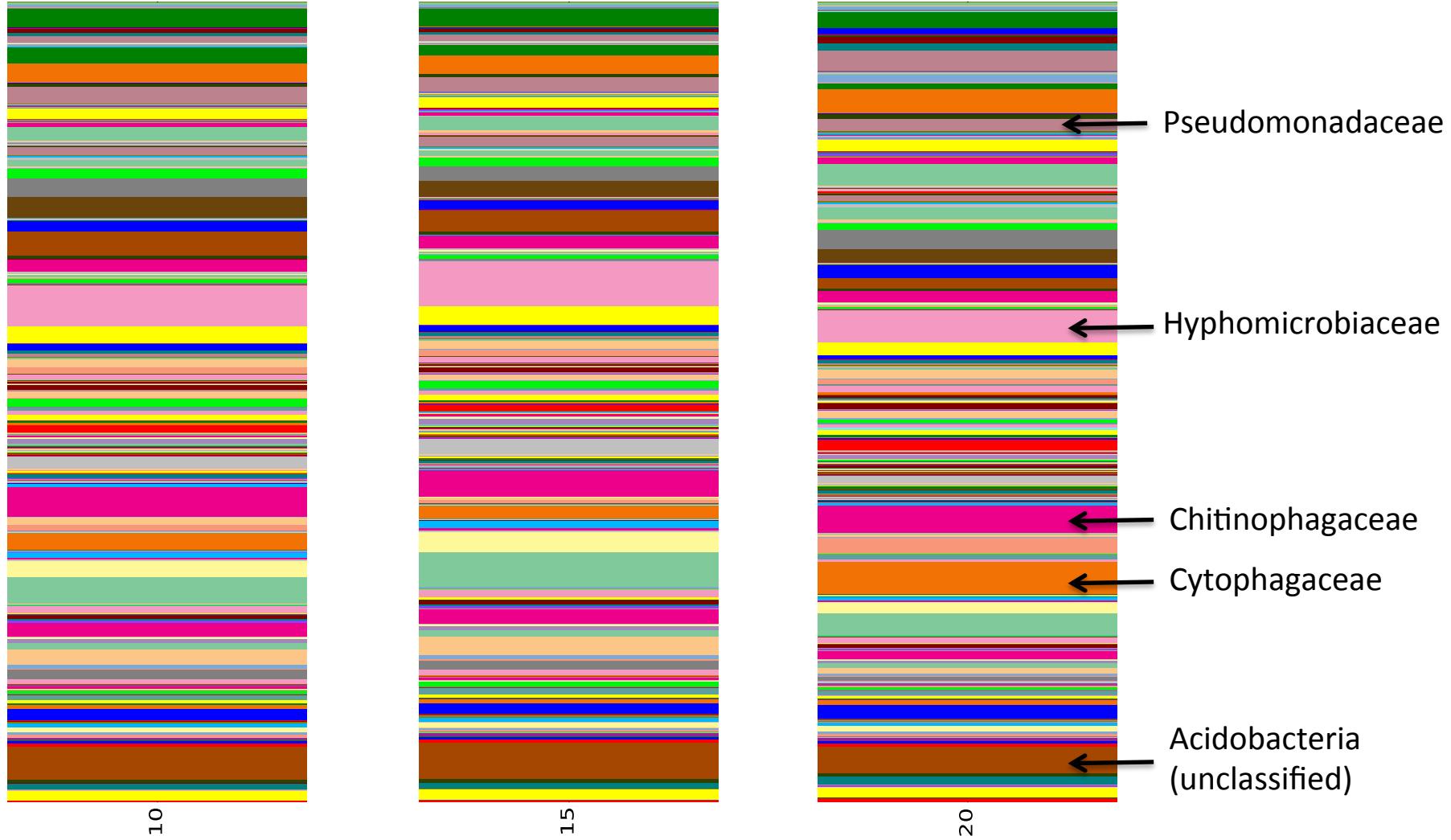
Mothers (yellow triangles) and their pups in Summer (purple circles), Early Winter (red triangles), Late Winter (blue triangles), and Spring (green squares)



Maize rhizo(endo?)sphere taxa shifts over time

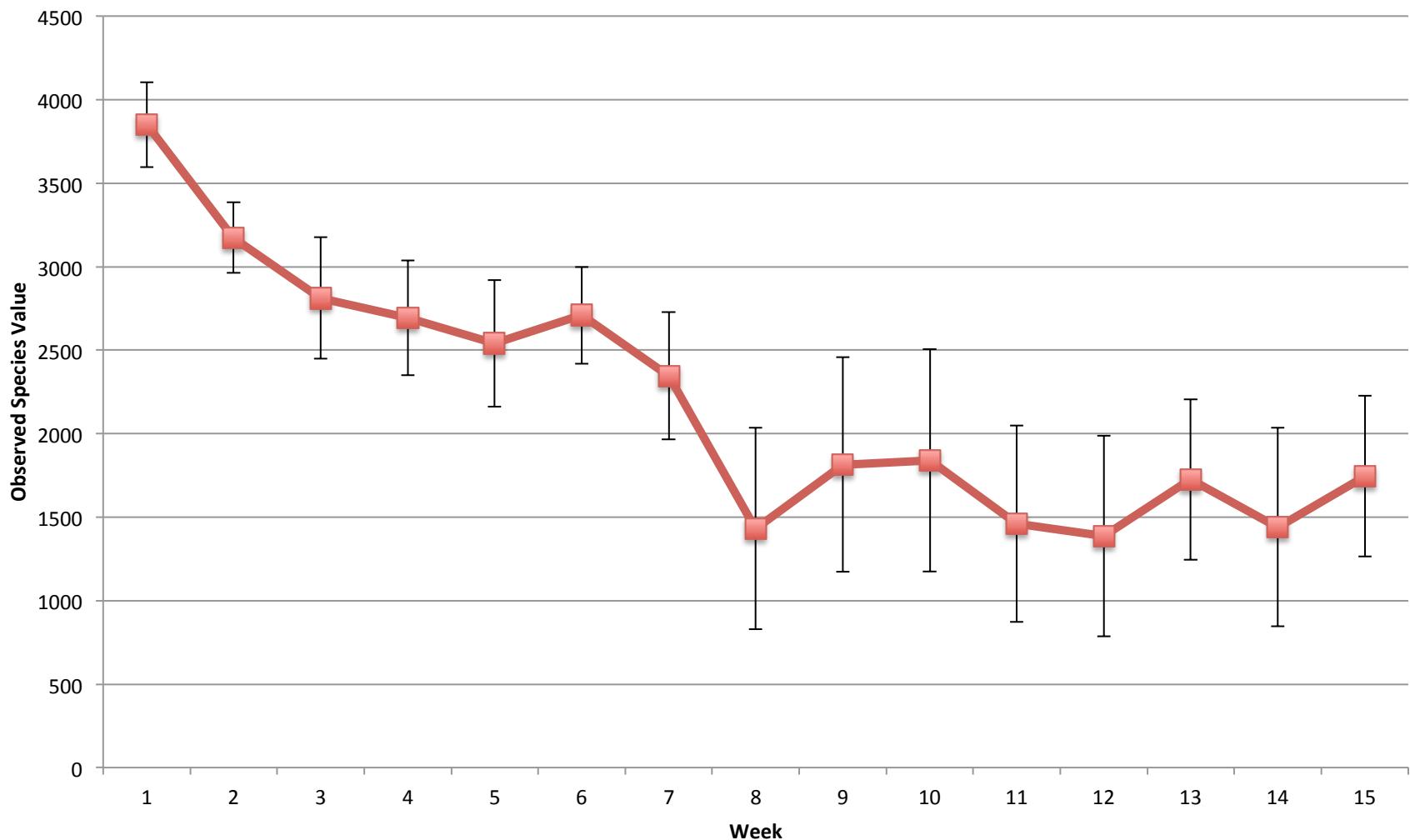


Taxonomy – Bulk Soils by Week

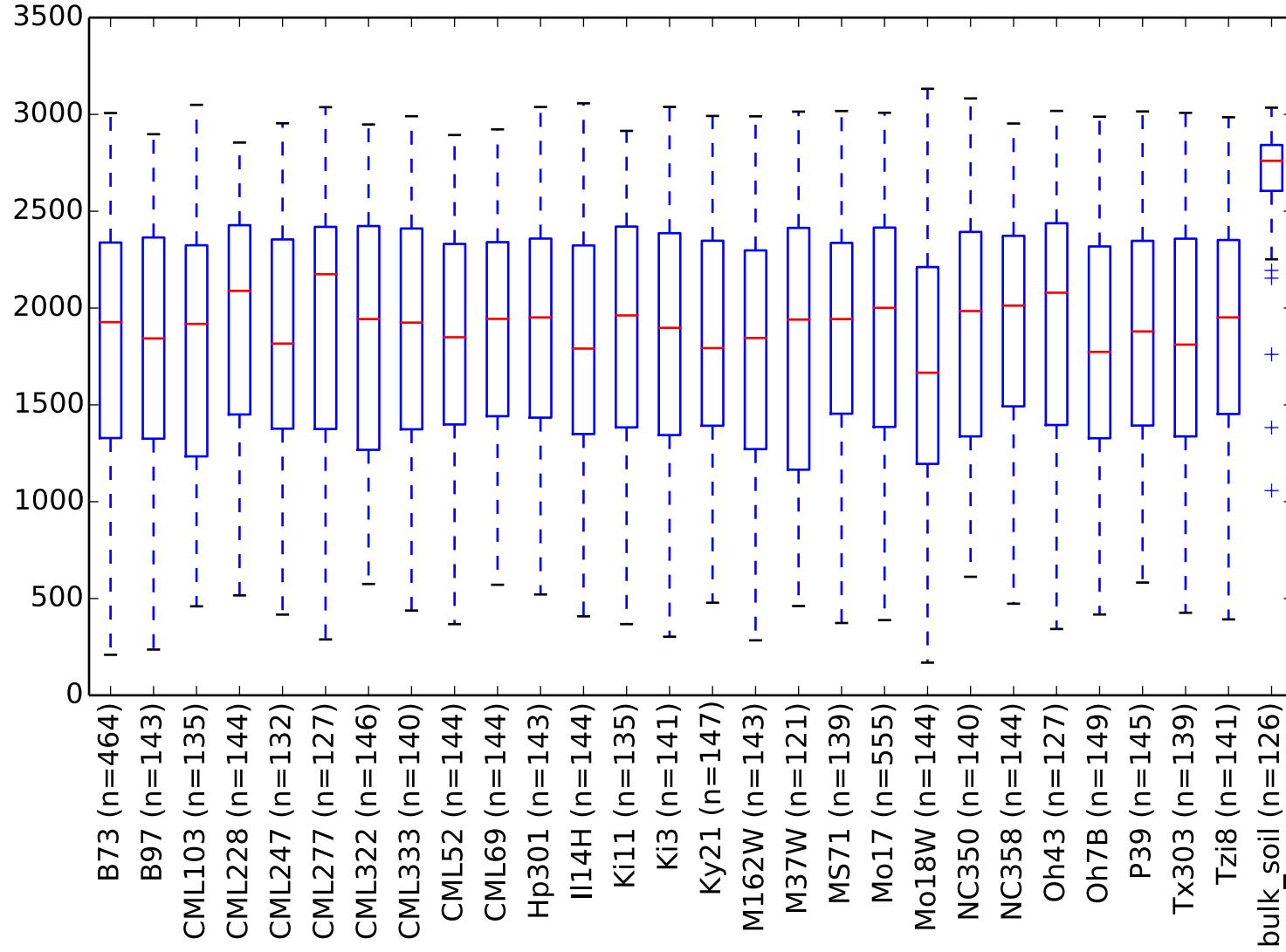


Alpha Diversity, all Maize

Observed Species with stdev



No significant difference between maize strains



Emperor Plots

What can we study with QIIME?

Metagenomics versus marker-gene surveys.

What can we study with QIIME?

Metagenomics versus marker-gene surveys.

What can we investigate?

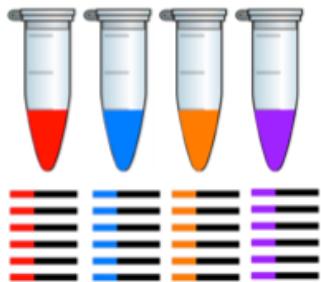
What can we study with QIIME?

Metagenomics versus marker-gene surveys.

What can we investigate?

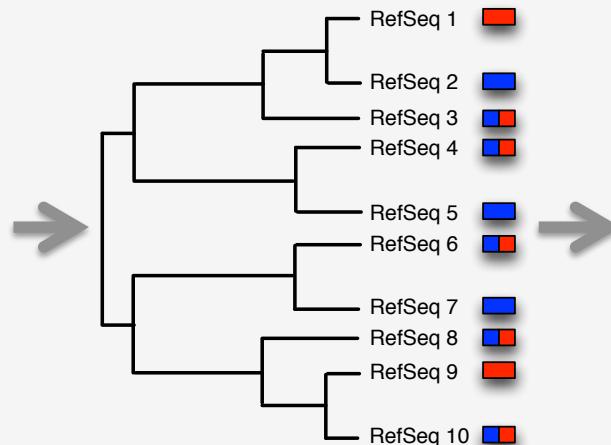
What are the limitations?

Intro to QIIME: installation and usage

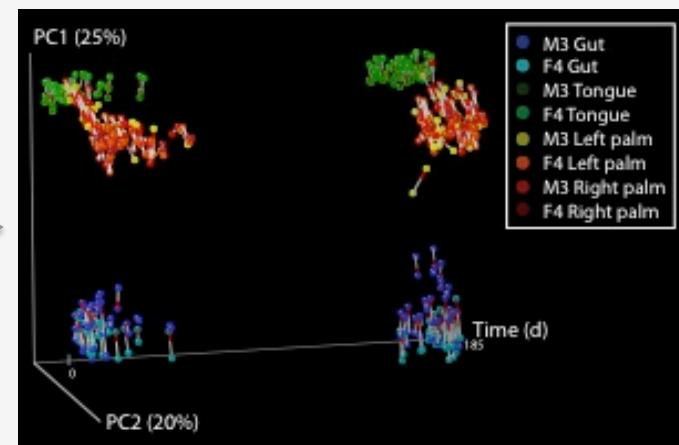


```
>GCACCTGAGGACAGGCATGAGGAA...
>GCACCTGAGGACAGGGGAGGGAGGA...
>TCACATAACCTAGGCAGGACGAA...
>CTACCGGAGGACAGGCATGAGGAT...
>TCACATAACCTAGGCAGGAGGAA...
>GCACCTGAGGACACGCAGGACGAC...
>CTACCGGAGGACAGGCAGGAGGAA...
>CTACCGGAGGACACACAGGAGGAA...
>GAACCTTCACATAGGCAGGAGGAT...
>TCACATAACCTAGGGGCAAGGAA...
>GCACCTGAGGACAGGCAGGAGGAA...
```

Assign reads to samples



Assign millions of
sequences from thousands
of samples to reference



Compare samples
statistically and visually

www.qiime.org



Native installation

OS X or Linux (laptops through 153,408-core compute cluster*)

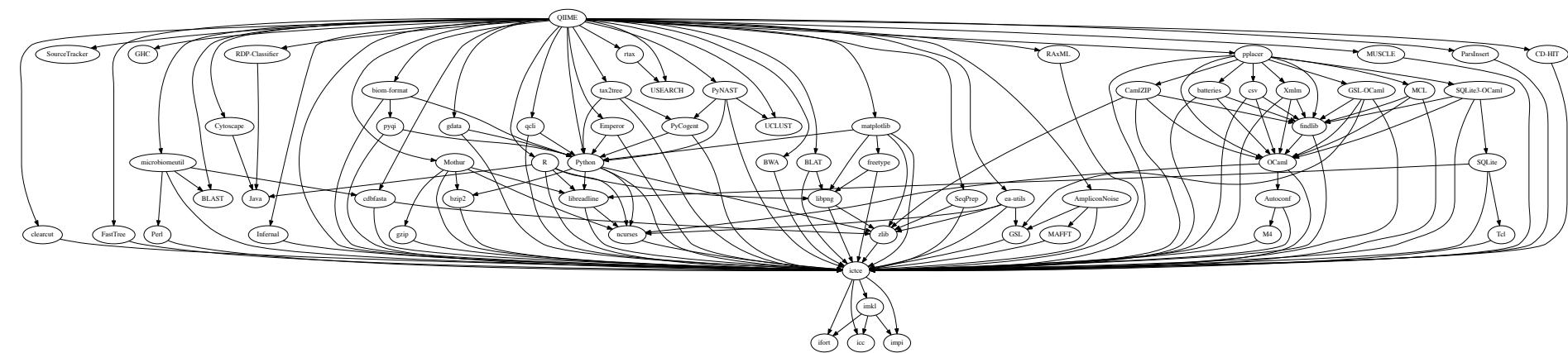
Virtual machines

VirtualBox (local) or
Cloud environments
(e.g., Amazon Web Services, iPlant)

**NEW: Now available as a
BaseSpace app!**

*Hopper (<http://i.top500.org/system/176952>)

QIIME software dependencies



QIIME software dependencies

ONE DOES NOT SIMPLY

INSTALL QIIME

imgflip.com

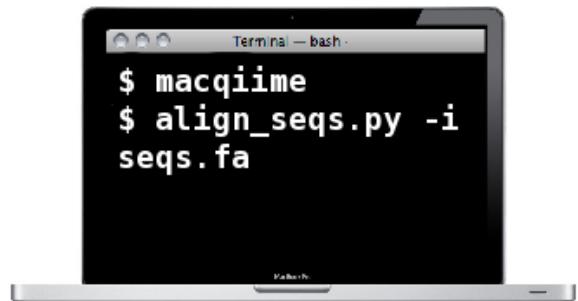
<https://imgflip.com/memegenerator/One-Does-Not-Simply>
<http://users.ugent.be/~kehoste/QIIME.pdf>

Native installation

- Installed directly on operating system / computer
- Pros
 - Best performance
 - Better use of resources (e.g., memory)
- Cons
 - Can be difficult to install

Native installation: MacQIIME

- Mac users: easiest way is MacQIIME
 - Excellent tool maintained by Jeff Werner's lab
 - Easiest way to get a (mostly) complete QIIME installation on a Mac
 - Can be difficult to customize/change the installation (e.g., updating specific pieces of MacQIIME)



Native installation: pip

- Mac/Linux users:

```
pip install qiime
```

- Basic QIIME installation
 - Supports basic upstream analyses and most downstream analyses
 - Install other dependencies as needed

Native installation: qiime-deploy

- Linux users (especially Ubuntu): easiest way to get a **full** QIIME installation
- Commonly used by system administrators to install QIIME in a cluster environment

Virtual machines

- A “guest” operating system running within a “host” operating system
- A software implementation of a computer, that operates like a physical computer.
- A developer can create a virtual machine *image* which contains their tools pre-installed. Users can then *instantiate* that image to work with those tools.

Benefits that virtual machines offer bioinformatics

- Reproducibility: can publish protocols with a virtual machine instance id.
- Updates are burden of developer, not user.
- Coupled with cloud computing, it's the perfect model for users with sporadic compute needs.

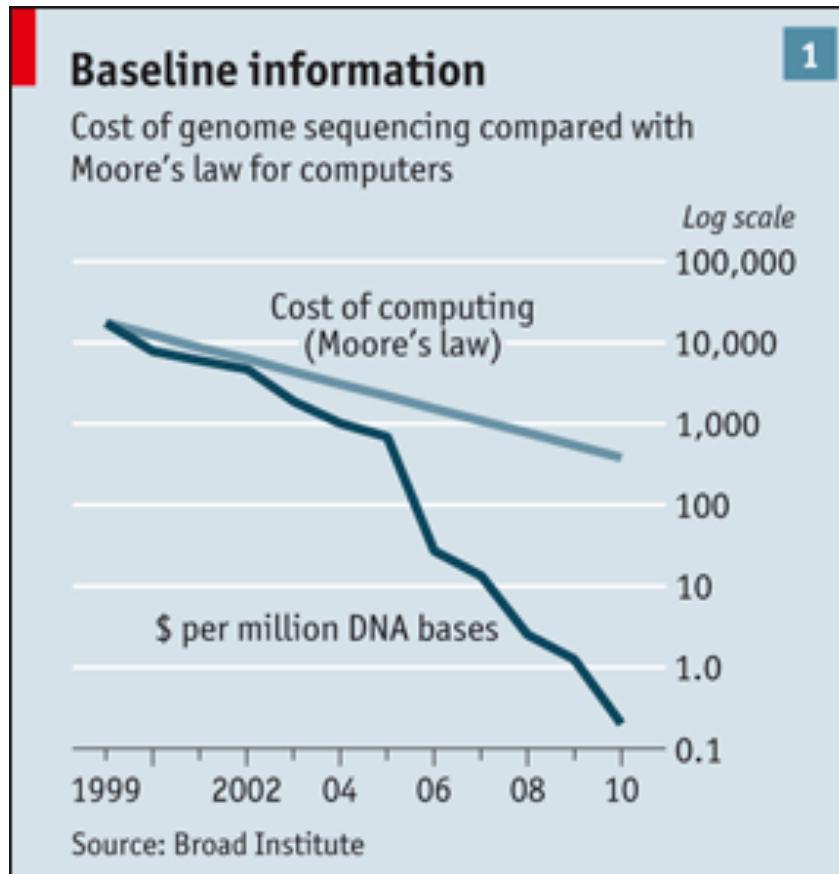
QIIME virtual machines

- VirtualBox (Ubuntu)
 - Easiest way to get a **full** QIIME installation on pretty much **any** platform
 - Especially for Windows users
 - Great for trying out QIIME
 - Slower than a native installation
- Cloud-based installations
 - Amazon Web Services (EC2)
 - iPlant Atmosphere
 - ANL Magellan

Isn't my laptop powerful enough?



Why is parallel computing important in bioinformatics?



Why is parallel computing important in bioinformatics?

Platform	Sanger	454 (Titanium)	Illumina Genome Analyzer II	Illumina HiSeq2000	Illumina MiSeq
Read Length (bases)	~1000	~400	150 (single end)	100 (single end)	150 (single end); 250 soon
Number of reads	96 or 384	~1,000,000	~100,000,000	~1,600,000,000	~10,000,000
Maximum number of samples per run	n/a	1000	12,000 (barcode-limited)	24,000 (barcode-limited)	2500 (barcode-limited)
Sequences per \$1 (sequencing costs only)	0.44	100	5000	200,000	12,500

Cluster computing

- Many computers connected to one another to serve as a larger compute resource.
- Compute-intensive jobs can be split over many systems and run in parallel.
- Similar to desktop compute hardware, but different casing, no (or only few) displays/ keyboards directly connected.
- Owned and maintained “in-house”.

Maintaining hardware is expensive

- Temperature (redundant cooling systems)
- Redundant network connections
- Hardware maintenance (e.g., replacing hard drives)
- Non-water fire suppression
- Backup power
- System administrator (\$\$)

Cloud computing (IaaS model)

- Implemented on a cluster (or grid), but compute power is rented as a service to support arbitrary applications.

Pay-as-you-go compute power

- Public clouds (e.g., Amazon) rent compute resources
- Log in, boot virtual machine image, run analyses, and terminate instance.
- Cheaper for many tasks than buying, maintaining, and supporting a compute cluster.

Interacting with the Amazon Cloud

- Boot virtual machine image via web interface (or a third-party tool like StarCluster).
- Log in and work via terminal (or via web interface with IPython Notebook)
- Move data back and forth via sftp/scp or a graphical sftp client (e.g., Cyberduck [free/cross-platform])

For information on costs, see <http://www.ec2instances.info>

Learning QIIME

- Start with the tutorials
 - <http://qiime.org/tutorials>
- Call any script with –h to get help or see the script usage pages
 - <http://scripts.qiime.org>
- Ask questions on the QIIME Forum
 - <http://help.qiime.org>
 - <http://forum.qiime.org>
- Report bugs on the issue tracker
 - <http://github.com/biocore/qiime/issues>

An Introduction To Applied Bioinformatics

Interactive lessons in bioinformatics.

[View the Project on GitHub](#)

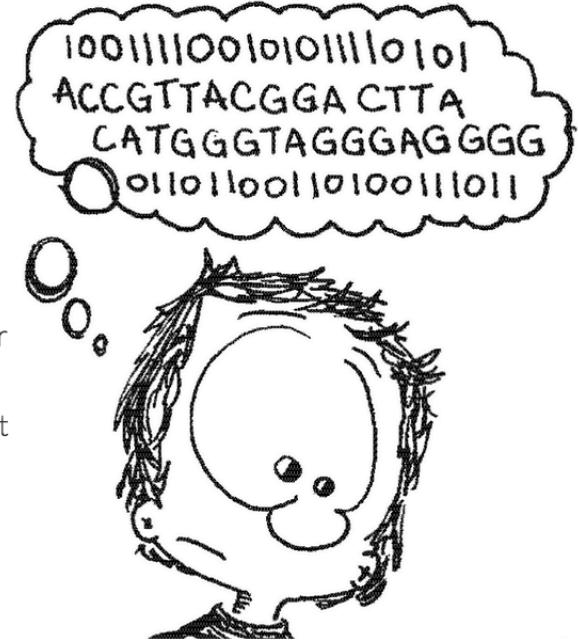
gregcaporaso/An-Introduction-To-Applied-Bioinformatics

Download
ZIP File

Download
TAR Ball

View On
GitHub

Bioinformatics, as I see it, is the application of the tools of computer science (things like programming languages, algorithms, and databases) to address biological problems (for example, inferring the evolutionary relationship between a group of organisms based on fragments of their genomes, or understanding if or how the community of microorganisms that live in my gut changes if I modify my diet). Bioinformatics is a rapidly growing field, largely in response to the vast increase in the quantity of data that biologists now grapple with. Students from varied disciplines (e.g., biology, computer science, statistics, and biochemistry) and stages of their educational careers (undergraduate, graduate, or postdoctoral) are becoming interested in bioinformatics.



I teach bioinformatics at the undergraduate and graduate levels at Northern Arizona University. This repository contains some of the materials that I've developed in these courses, and represents an initial attempt to organize these materials in a standalone way. In some cases, I'm just linking out to other materials for now.

<http://readIAB.org>



ALFRED P. SLOAN
FOUNDATION

Gut Check: Exploring Your Microbiome

by University of Colorado Boulder & University of Colorado System

Course Info

UNIVERSITY OF COLORADO BOULDER & UNIVERSITY OF COLORADO SYSTEM

Gut Check: Exploring Your Microbiome

About this Course

Imagine if there were an organ in your body that weighed as much as your brain, that affected your health, your weight, and even your behavior. Wouldn't you want to know more about it? There is such an organ — the collection of microbes in and on your body, your human microbiome.

Subtitles available in English

Log in to enroll in this course

Log in

Instructors



Professor Rob Knight

Professor

Howard Hughes Medical Institute, and Department of Chemistry & Biochemistry and Computer Science, and Biofrontiers Institute



Dr. Jessica Metcalf

Senior Research Associate

BioFrontiers Institute



Dr. Katherine Amato

Postdoctoral Research Associate

Department of Anthropology, BioFrontiers Institute

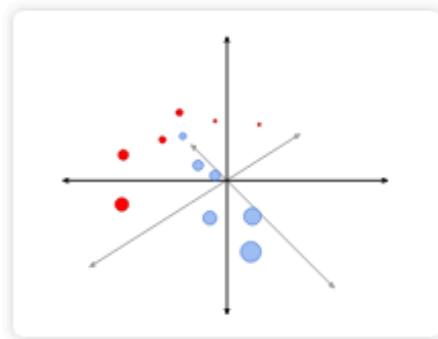
<https://www.coursera.org/learn/microbiome>

Welcome to the GUide to STasitical Analysis in Microbial Ecology (GUSTA ME)!

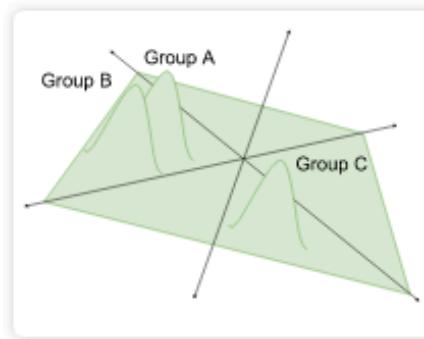
Where would you like to start?

You may start exploring the guide by browsing topics in the sidebar, using the search box at the top right of this page to find a particular method, or by clicking on one of the entry points below...

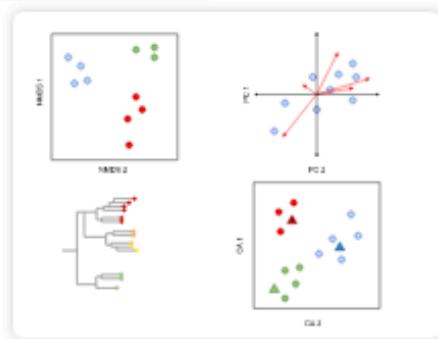
Explore data...



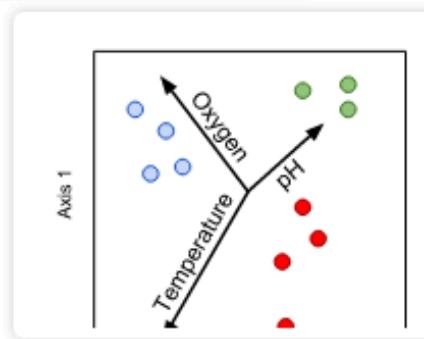
Test a hypothesis...



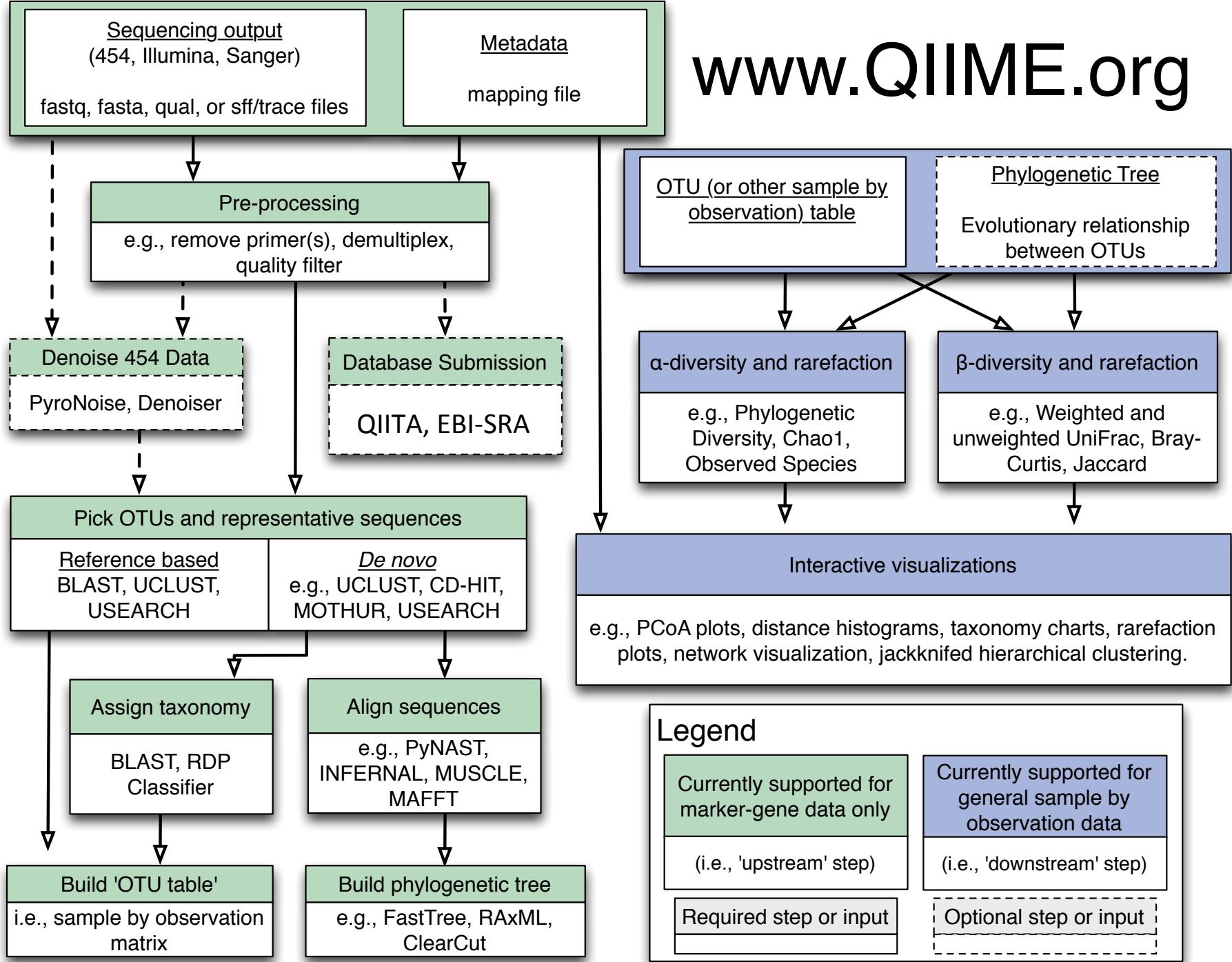
Browse visualisations...

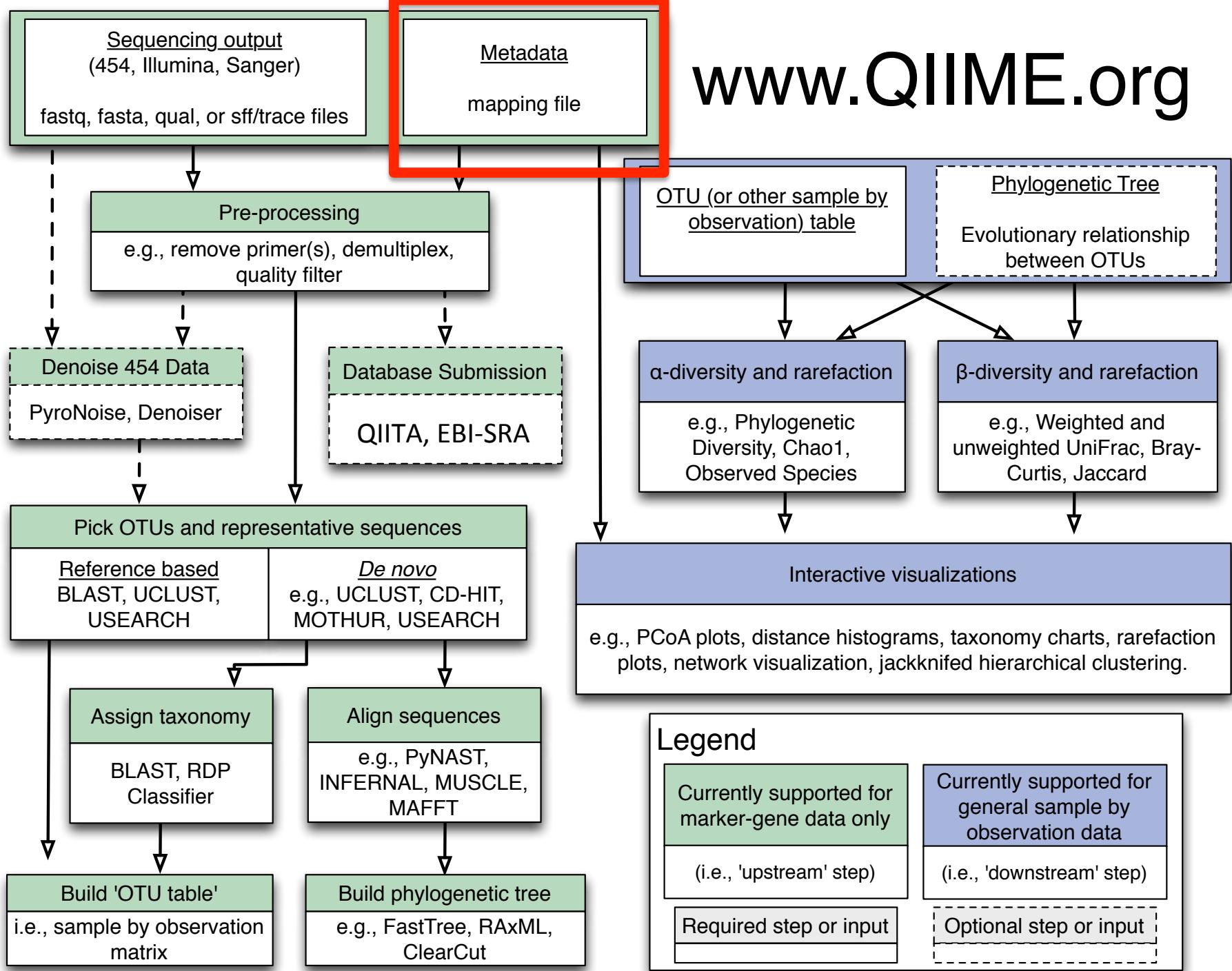


Explore environmental influence...

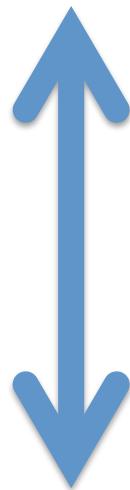
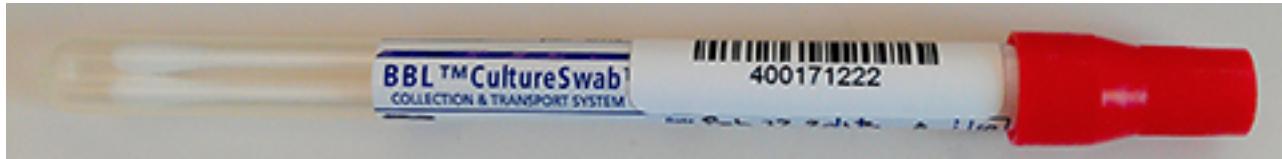


Preprocessing: metadata and demultiplexing





Metadata relates samples to variables



SampleID	Sex	SampleType	Plotting your doom?
400171222	Male	feces	Naturally

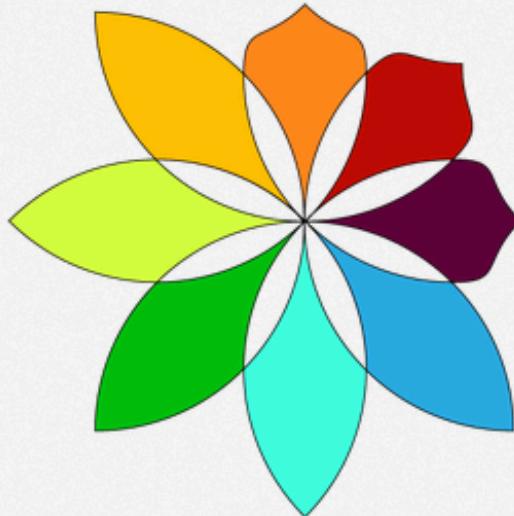
Mapping file

```
Fasting_Map.txt
1 #SampleID> BarcodeSequence>LinkerPrimerSequence> Treatment> DOB>Description>
2 PC.354> AGCAGCGAGCCTA> YATGCTGCCTCCCGTAGGAGT> Control> 20061218> Control_mouse_I.D._354>
3 PC.355> AACTCGTCGATG> YATGCTGCCTCCCGTAGGAGT> Control> 20061218> Control_mouse_I.D._355>
4 PC.356> ACAGACCACTCA> YATGCTGCCTCCCGTAGGAGT> Control> 20061126> Control_mouse_I.D._356>
5 PC.481> ACCAGCGACTAG> YATGCTGCCTCCCGTAGGAGT> Control> 20070314> Control_mouse_I.D._481>
6 PC.593> AGCAGCACTTGT> YATGCTGCCTCCCGTAGGAGT> Control> 20071210> Control_mouse_I.D._593>
7 PC.607> AACTGTGCGTAC> YATGCTGCCTCCCGTAGGAGT> Fast> 20071112> Fasting_mouse_I.D._607>
8 PC.634> ACAGAGTCGGCT> YATGCTGCCTCCCGTAGGAGT> Fast> 20080116> Fasting_mouse_I.D._634>
9 PC.635> ACCGCAGAGTCA> YATGCTGCCTCCCGTAGGAGT> Fast> 20080116> Fasting_mouse_I.D._635>
10 PC.636> ACGGTGAGTGTC> YATGCTGCCTCCCGTAGGAGT> Fast> 20080116> Fasting_mouse_I.D._636>
11
```

Mapping file: always run validate_mapping_file.py

	#SampleID	BarcodeSequence	LinkerPrimerSequence	Treatment	DOB	Description
1	PC.354	AUCACCGAGCTTA	YATGCTGCCTCCCGTAGGAGT	Control	20061218	Control_mouse_I.D._354
2	PC.355	AACTCGTCGATG	YATGCTGCCTCCCGTAGGAGT	Control	20061218	Control_mouse_I.D._355
3	PC.356	ACAGACCACTCA	YATGCTGCCTCCCGTAGGAGT	Control	20061126	Control_mouse_I.D._356
4	PC.481	ACCAGCGACTAG	YATGCTGCCTCCCGTAGGAGT	Control	20070314	Control_mouse_I.D._481
5	PC.593	AGCAGCACTTGT	YATGCTGCCTCCCGTAGGAGT	Control	20071210	Control_mouse_I.D._593
6	PC.607	AACTGTGCGTAC	YATGCTGCCTCCCGTAGGAGT	Fast	20071112	Fasting_mouse_I.D._607
7	PC.634	ACAGAGTCGGCT	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._634
8	PC.635	ACCGCAGAGTCA	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._635
9	PC.636	ACGGTGAGTGTC	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._636
10						
11						

= required field



Keemei: Validate bioinformatics metadata in Google Sheets

Keemei (canonically pronounced *key may*) is an open source [Google Sheets](#) add-on for validating bioinformatics metadata, including [QIIME](#) mapping files.

keemei-demo

File Edit View Insert Format Data Tools Add-ons Help All changes saved in Drive

Validate QIIME mapping file

Clear validation status

About

Help

Sample_Type Sample_Plate Print

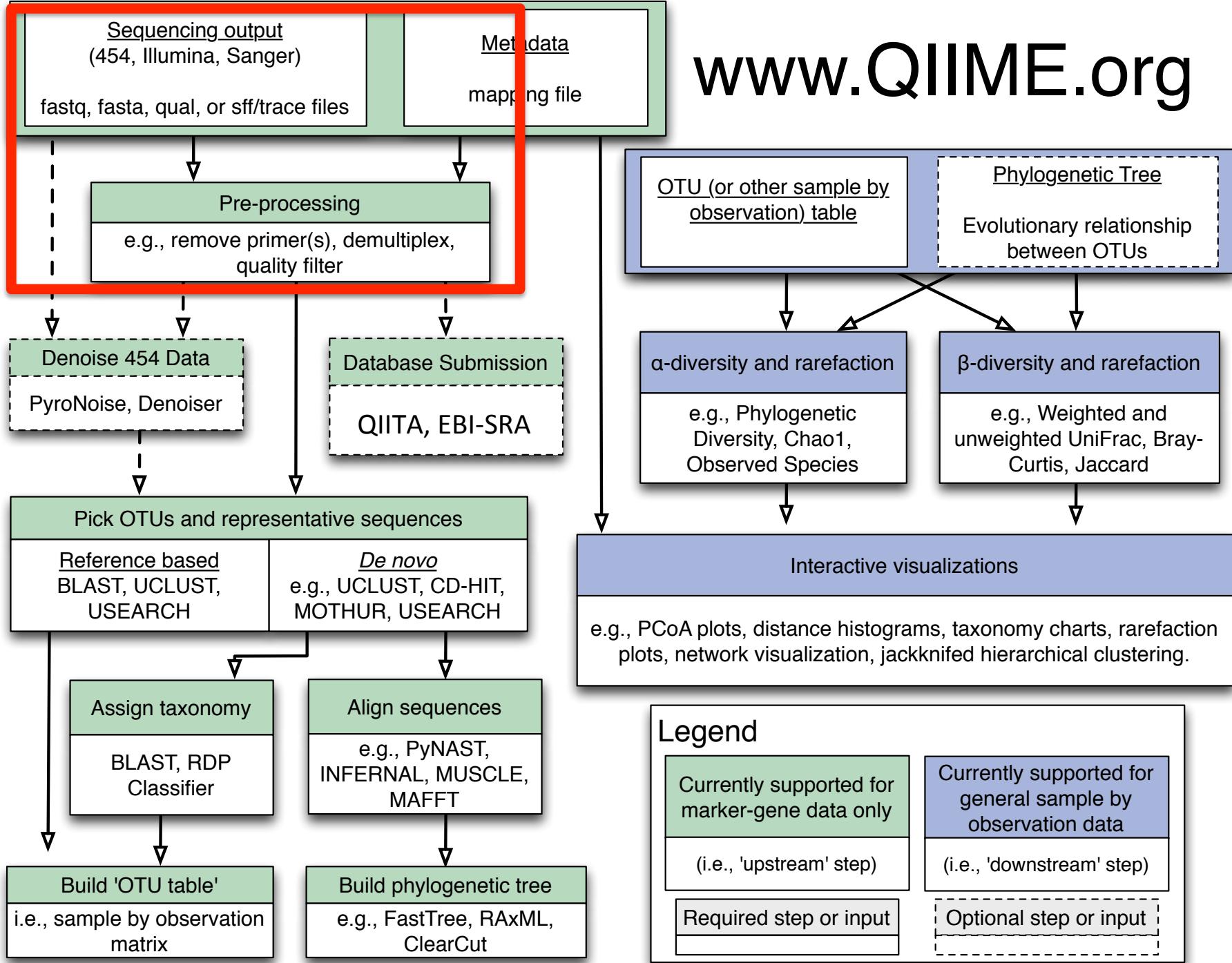
	A	B			H	I	
1	#SampleID	BarcodeSequence	Linker		Stool	ABTX1	
2	AB.Fece.10.21.2008	CATGGCTACACA	CATGCTGCCTCCC	10	21	2008	14173 Skin_Palm ABTX6
3	AB.LPalm.10.21.2008	CTGTTCGTAGAG	CATGCTGCCTCCC	10	21	2008	14173 Tongue ABTX17
4	AB.Tong.10.21.2008	GCTGTAGTATGC	CATGCTGCCTCCC	10	22	2008	14174 Stool ABTX1
5	AB.Fece.10.22.2008	CGAAGACTGCTG	CATGCTGCCTCCC	10	22	2008	14174 Skin_Palm ABTX6
6	AB.LPalm.10.22.2008	GACAGGAGATAG	CATGCTGCCTCCC	10	22	2008	14174 Tongue ABTX17
7	AB.Tong.10.22.2008	GGTCACTGACAG	CATGCTGCCTCCC	10	22	2008	14175 Stool ABTX1
8	AB.Fece.10.23.2008	GGTCACTGACAG	CATGCTGCCTCCC	10	23	2008	14175 Skin_Palm ABTX6
9	AB.LPalm.10.23.2008	GACTCACTCAAT	CATGCTGCCTCCC	10	23	2008	14175 Tongue ABTX17
10	AB.Tong.10.23.2008	GTAGTGTCTAGC	CATGCTGCCTCCC	10	23	2008	14176 Stool ABTX1
11	AB.Fece.10.24.2008	CGGAGTGTCTAT	CATGCTGCCTCCC	10	24	2008	14176 Skin_Palm ABTX6
12	AB.LPalm.10.24.2008	GAGCATTCTCTA	CATGCTGCCTCCC	10	24	2008	14176 Tongue ABTX17
13	AB.Tong.10.24.2008	GTCGCTGTCTTC	CATGCTGCCTCCC	10	24	2008	14177 Stool ABTX1
14	AB.Fece.10.25.2008	CGTGACAATGTC	CATGCTGCCTCCC	10	25	2008	14177 Skin_Palm ABTX6
15	AB.LPalm.10.25.2008	GATCAGAAAGATG	CATGCTGCCTCCC	10	25	2008	14177 Tongue ABTX17
16	AB.Tong.10.25.2008	GTGAGGGTCGCTA	CATGCTGCCTCCC	10	25	2008	14178 Stool ABTX1
17	AB.Fece.10.26.2008	GTACTACAGCTG	CATGCTGCCTCCC	10	26	2008	14178 Skin_Palm ABTX6

keemei-demo

File Edit View Insert Format Data Tools Add-ons Help All changes saved in Drive

fx #SampleID

	A	B	C	D	E	F	G	H	I
1	#SampleID	BarcodeSequence	LinkerPrimerSequer	Month	Day	Year	days_since_epoch	Sample_Type	Sample_Plate
2	AB.Fece.10.21.2008	CATGGCTACACA	CATGCTGCCTCCC	10	21	2008	14173	Stool	ABTX1
3	AB.LPalm.10.21.2008	CTGTTCTAGAG	CATGCTGCCTCCC	10	21	2008	14173	Skin_Palm	ABTX6
4	AB.Tong.10.21.2008	GCTGTAGTATGC	CATGCTGCCTCCC	10	21	2008	14173	Tongue	ABTX17
5	AB.Fece.10.22.2008	CGAAGACTGCTG	CATGCTGCCTCCC	10	22	2008	14174	Stool	ABTX1
6	AB.LPalm.10.22.2008	GACAGGAGATAG	CATGCTGCCTCCC	10	22	2008	14174	Skin_Palm	ABTX6
7	AB.Tong.10.22.2008	GGTCACTGACAG	Error: Duplicate barcode sequence. Duplicates in B7, B8				14174	Tongue	ABTX17
8	AB.Fece.10.23.2008	GGTCACTGACAG		23	2008	14175	Stool	ABTX1	
9	AB.LPalm.10.23.2008	GACTCACTCAAT		23	2008	14175	Skin_Palm	ABTX6	
10	AB.Tong.10.23.2008	GTAGTGTCTAGC		23	2008	14175	Tongue	ABTX17	
11	AB.Fece.10.24.2008	CGGAGTGTCTAT		10	24	2008	14176	Stool	ABTX1
12	AB.LPalm.10.24.2008	GAGCATTCTCTA	CATGCTGCCTCCC	10	24	2008	14176	Skin_Palm	ABTX6
13	AB.Tong.10.24.2008	GTCGCTGTCTTC	CATGCTGCCTCCC	10	24	2008	14176	Tongue	ABTX17
14	AB.Fece.10.25.2008	CGTGACAATGTC	CATGCTGCCTCCC	10	25	2008	14177	Stool	ABTX1
15	AB.LPalm.10.25.2008	GATCAGAAGATG	CATGCTGCCTCCC	10	25	2008	14177	Skin_Palm	ABTX6



Preprocessing of reads

The data that we will be processing tomorrow was generated in the format that comes from Caporaso 2011, PNAS, where one has separate reads and barcodes fastq files.

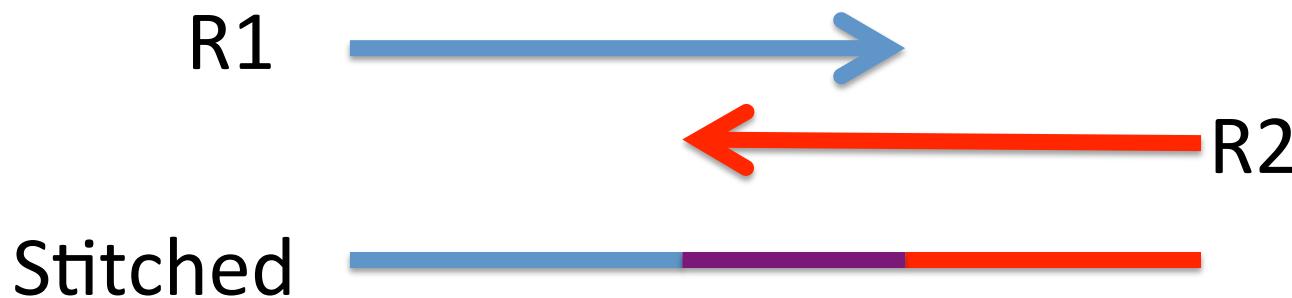
Preprocessing of reads

The data that we will be processing tomorrow was generated in the format that comes from Caporaso 2011, PNAS, where one has separate reads and barcodes fastq files.

For most cases, you will receive paired-end reads (with R1 and R2 in their filenames) for the read 1 and read 2 of Illumina data.

Preprocessing of reads

If one's reads overlap, then they can be stitched together. In QIIME, the `join_paired_ends.py` script can be used to stitch the reads.



Preprocessing of reads

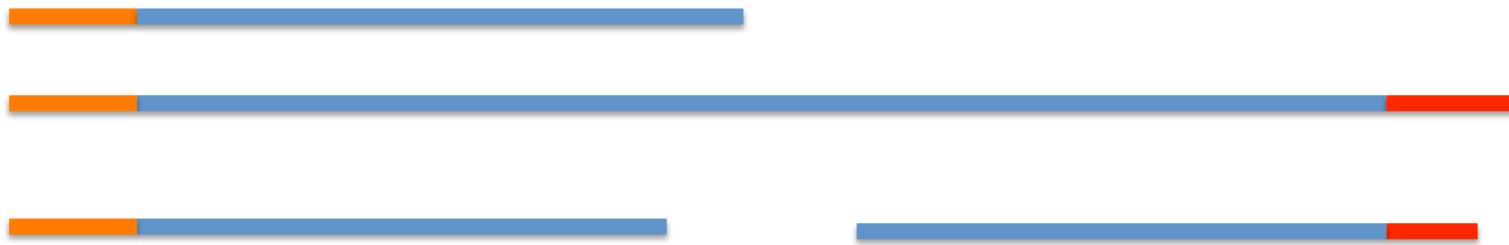
The output of `join_paired_ends.py` will be a set of “joined” reads, as well as unjoined R1 and R2 files. You want to use the joined reads for further processing and analysis.

If you have a separate barcode read, you want to pass this as a parameter when calling `join_paired_ends.py`, so the barcodes will be filtered to match the joined reads.

One consideration: if there is low yields from the stitching process, it may be better to use a single end read, such as R1, than using a low number of stitched reads.

Alternative Barcoding

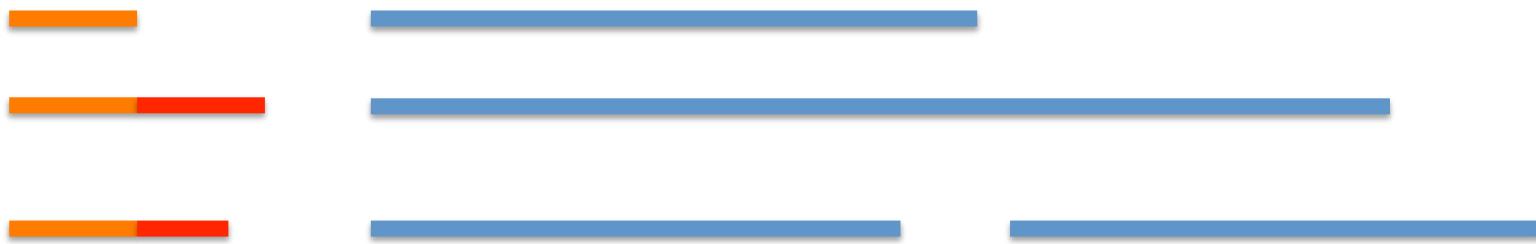
Sometimes (not in the tutorial data tomorrow),
the barcodes will be at the end(s) of your reads.



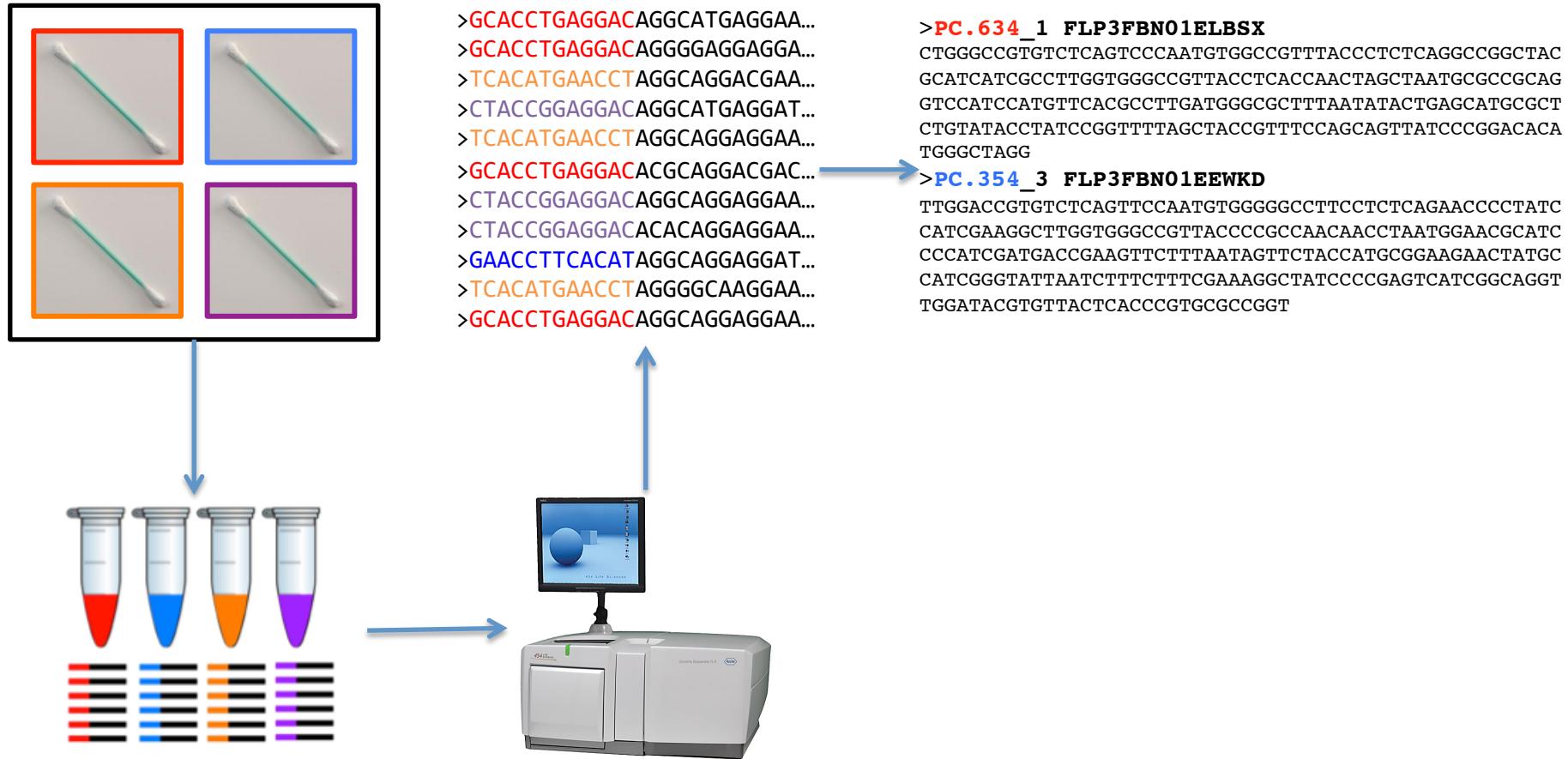
Orange/Red indicates barcodes

Alternative Barcoding

You can use the `extract_barcodes.py` script to generate separate barcodes and reads (with barcodes stripped) fastq files.



Error-correcting codes allow multiplex sequencing



Demultiplexing

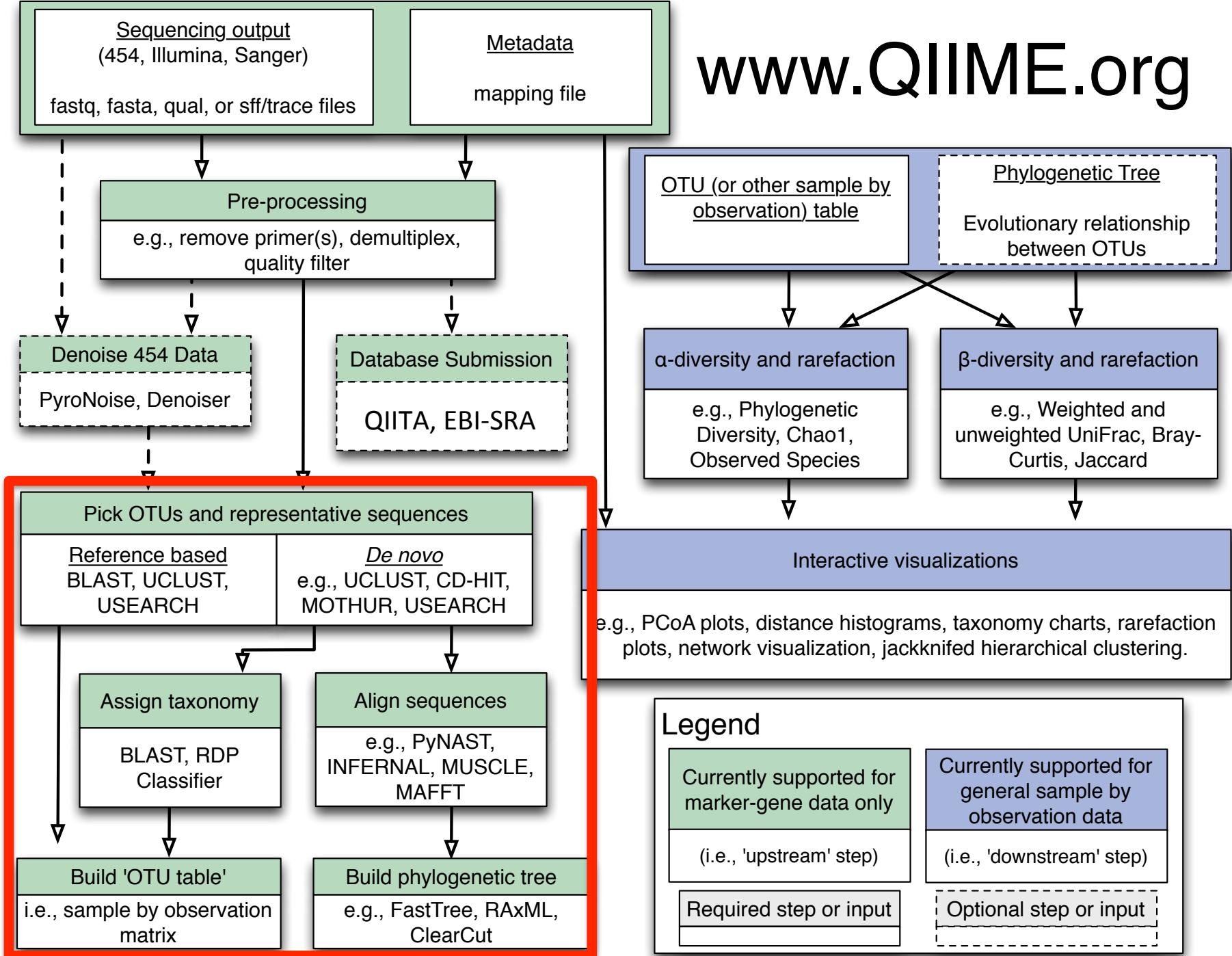
For Illumina data, the script is `split_libraries_fastq.py`.

For 454 and IonTorrent data, use `split_libraries.py`.

A common situation (not in this tutorial): If you get already-split up reads (separated into fastq file(s) by sample), use the script `multiple_split_libraries_fastq.py`

`multiple_join_paired_ends.py` and
`multiple_extract_barcodes.py` also exist.

OTU picking methods



Which microbial organisms are represented by the rRNA gene sequences in each sample?

rRNA reference database
(sequences are available for each ‘tip’ in the tree)

>PC.634_1 FLP3FBN01ELBSX

CTGGGCCGTGTCAGTCCAAATGTGCCGTTACCCCTCAGGCCGG
CTACGCATCATGCCCTGGTGGGCCGTTACCTCACCAACTAGCTAATG
CGCCGCAGGTCCATCCATGTTACGCCCTGATGGCGCTTAATATAC
TGAGCATGCGCTCTGTATAACCTATCCGGTTAGCTACCAGTTCCAGC
AGTTATCCCAGACACATGGGCTAGG

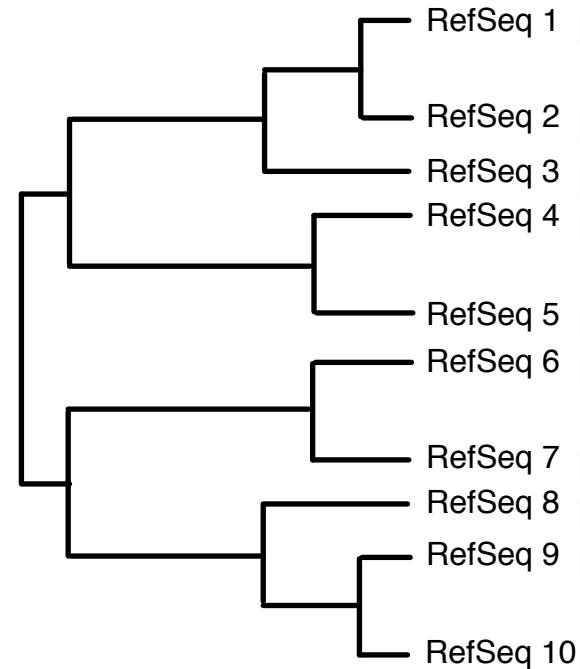
>PC.634_2 FLP3FBN01EG8AX

TTGGACCGTGTCTCAGTCCAATGTGGGGCCTCCTCTCAGAACCCC
TATCCATCGAAGGCTTGGTGGCCGTTACCCCGCCAACAACCTAATGG
AACGCATCCCCATCGATGACCGAAGTTCTTAATAGTTCTACCATGCG
GAAGAACTATGCCATCGGGTATTAATCTTCTTCGAAAGGCTATCCC
CGAGTCATCGGCAGGTTGGATACGTGTTACTCACCCGTGCGCCGGT

>PC.354_3 FLP3FBN01EEWKD

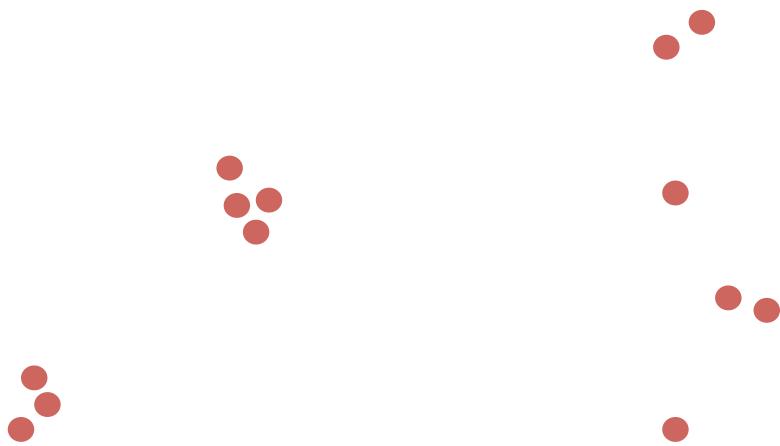
TTGGGCCGTGTCAGTCCAAATGTGCCGATCAGTCTCTTAACTCGG
CTATGCATCATGCCCTGGTAAGCCGTTACCTTACCAACTAGCTAATG
CACCGCAGGTCCATCCAAGAGTGATAGCAGAACCATTTCAAACCTCT
AGACATGCGTCTAGTGTATTCCGGTATTAGCATCTGTTCCAGGT
GTTATCCCAGTCTCTGGG

Search against
reference
sequences

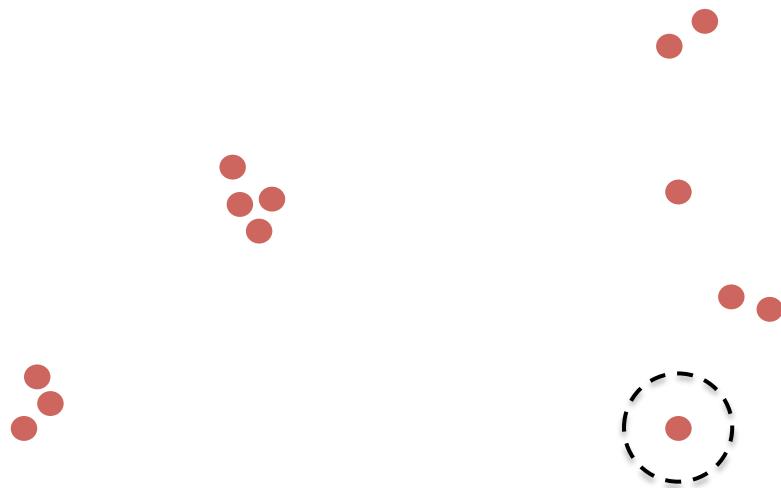


OTU table: counts of OTUs in each sample

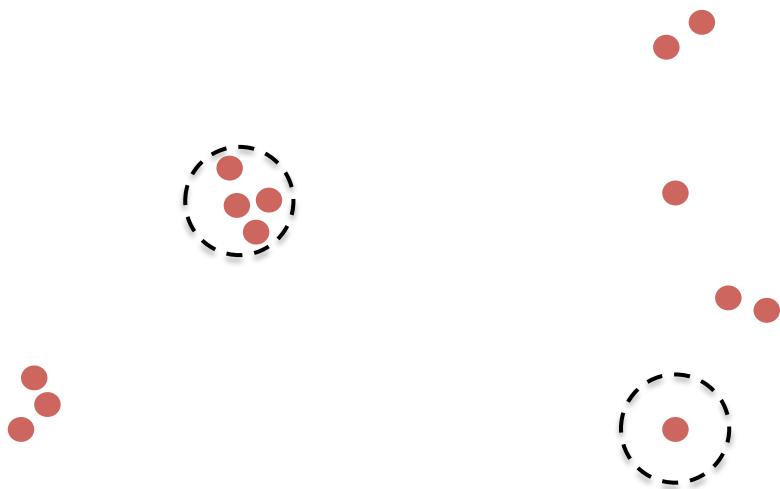
Understanding OTU picking



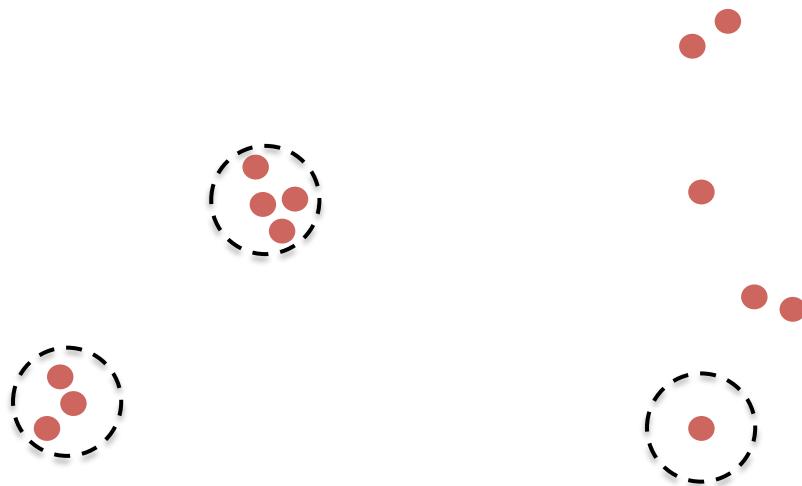
Understanding OTU picking



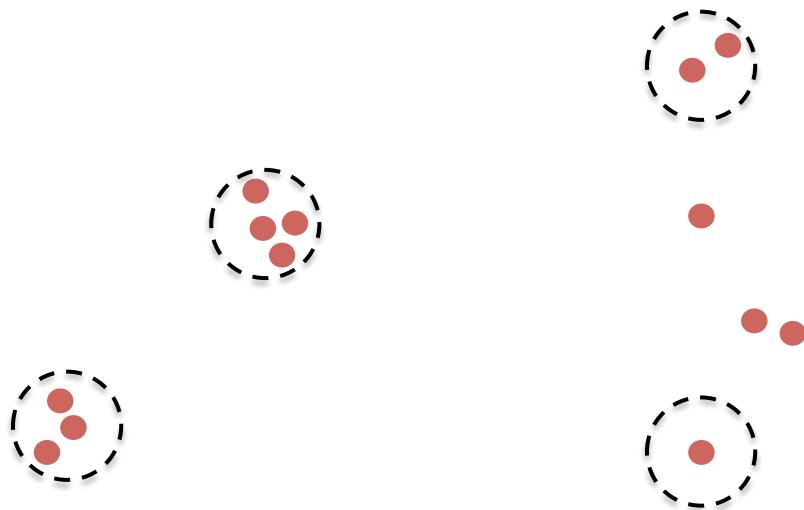
Understanding OTU picking



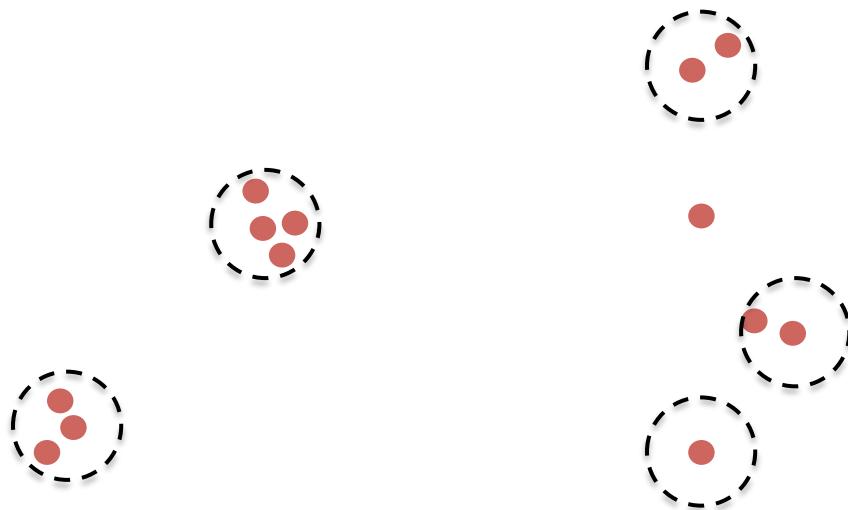
Understanding OTU picking



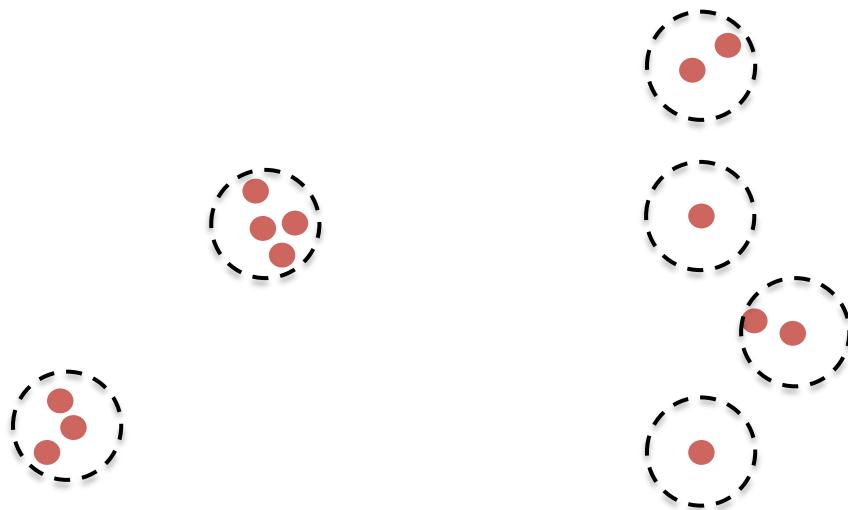
Understanding OTU picking



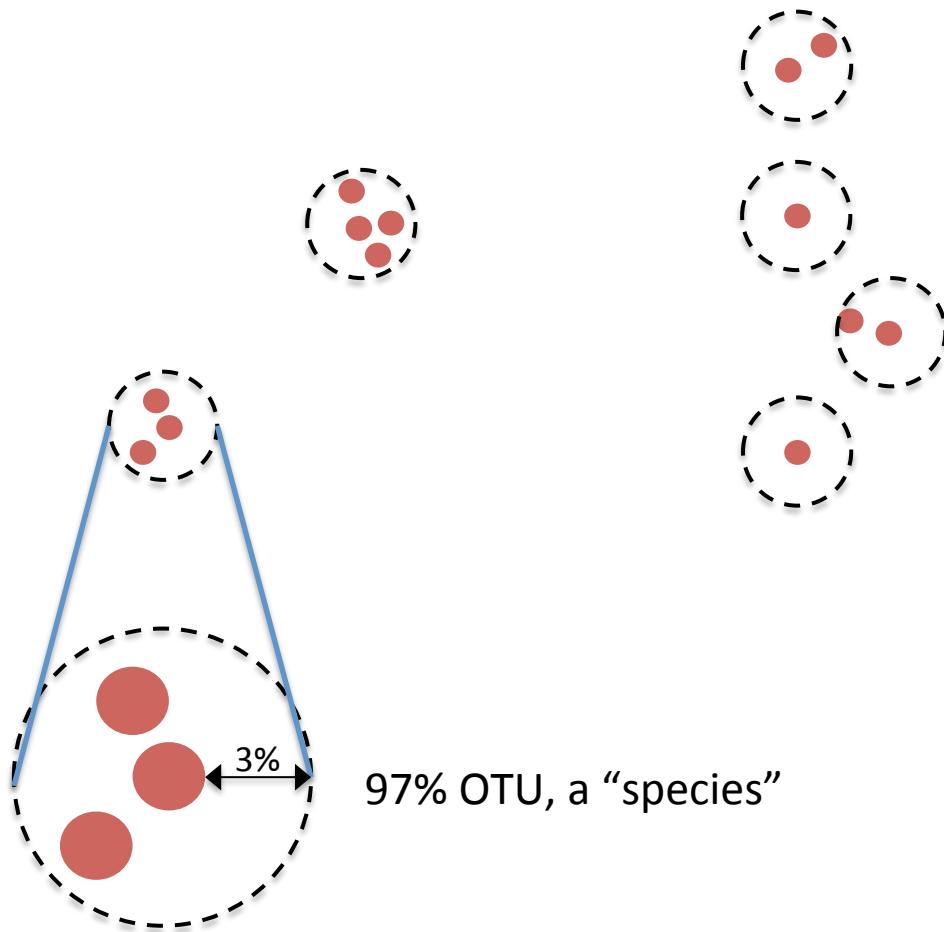
Understanding OTU picking



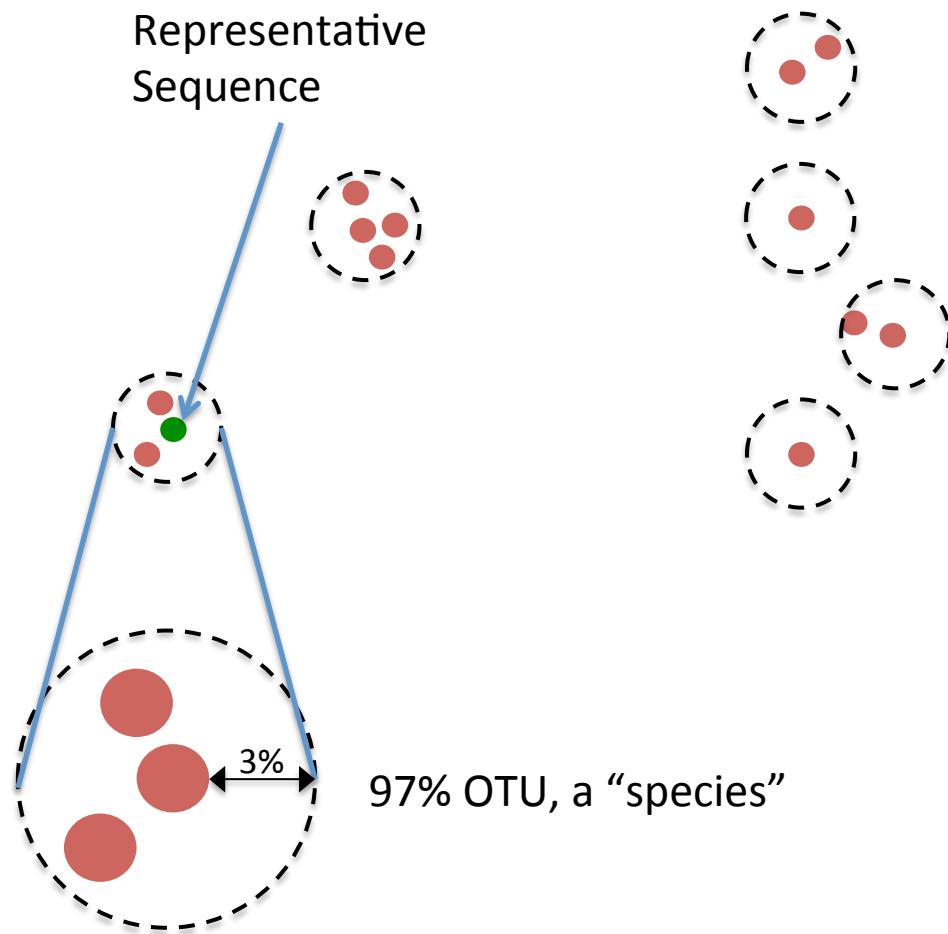
Understanding OTU picking



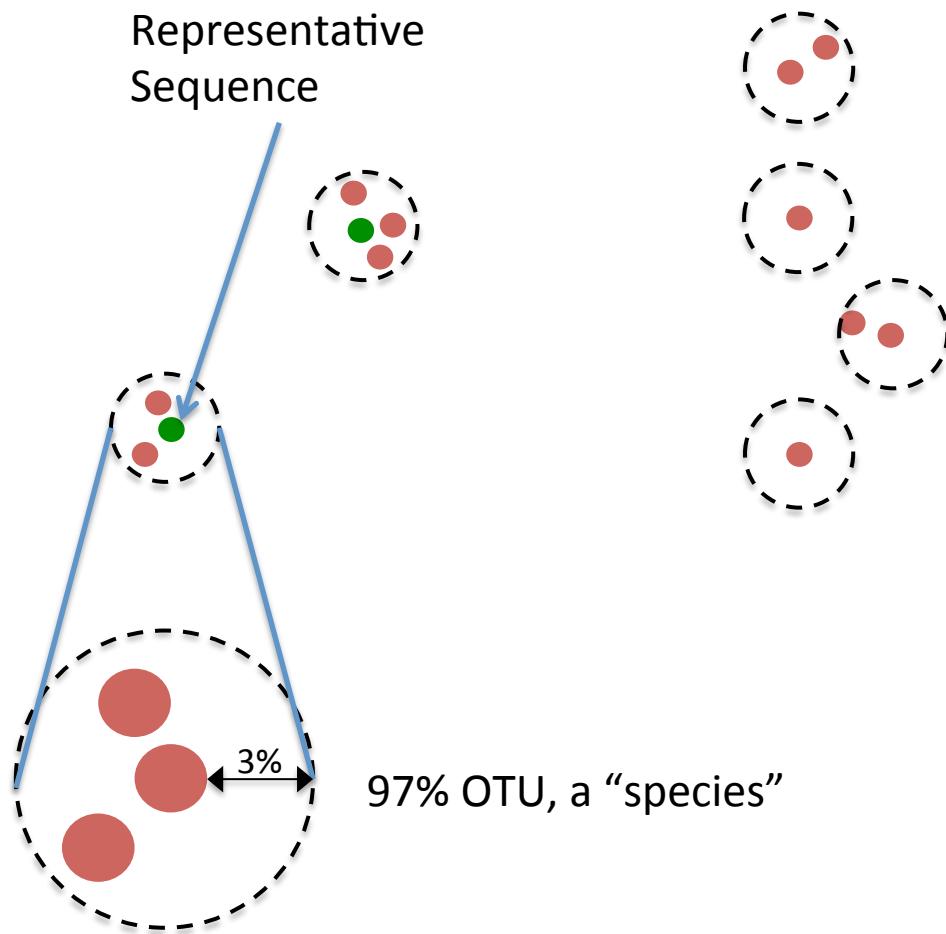
Understanding OTU picking



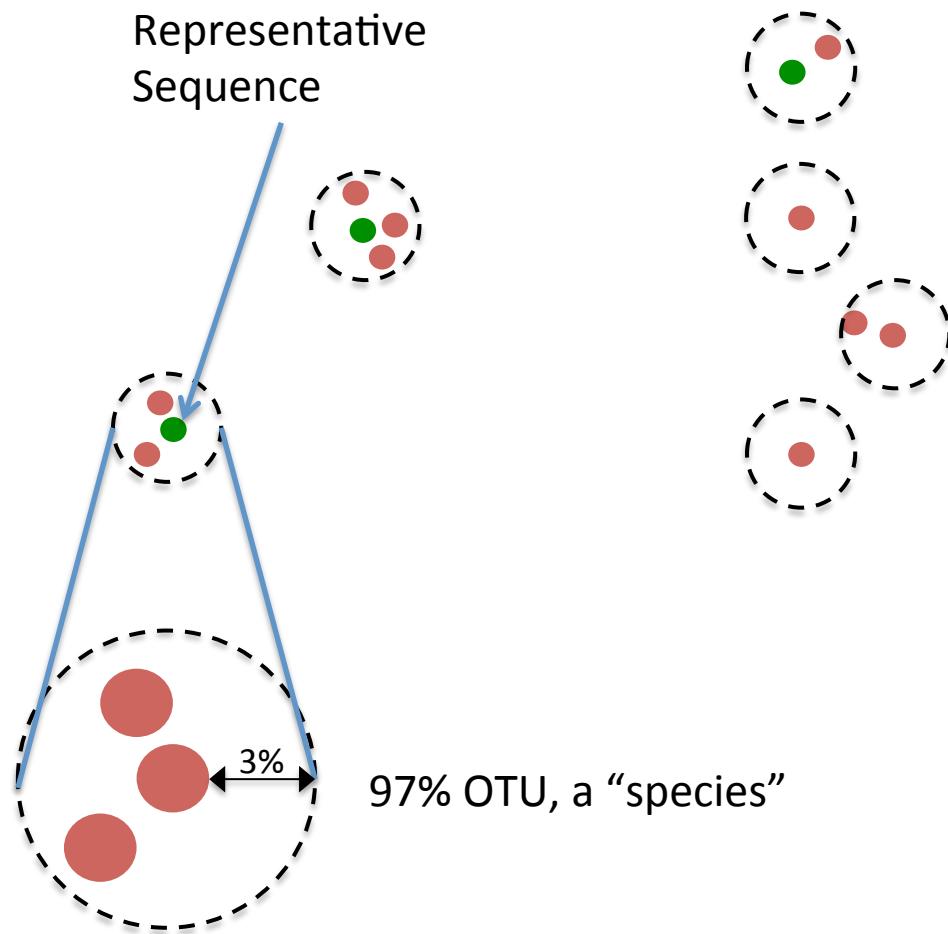
Understanding OTU picking



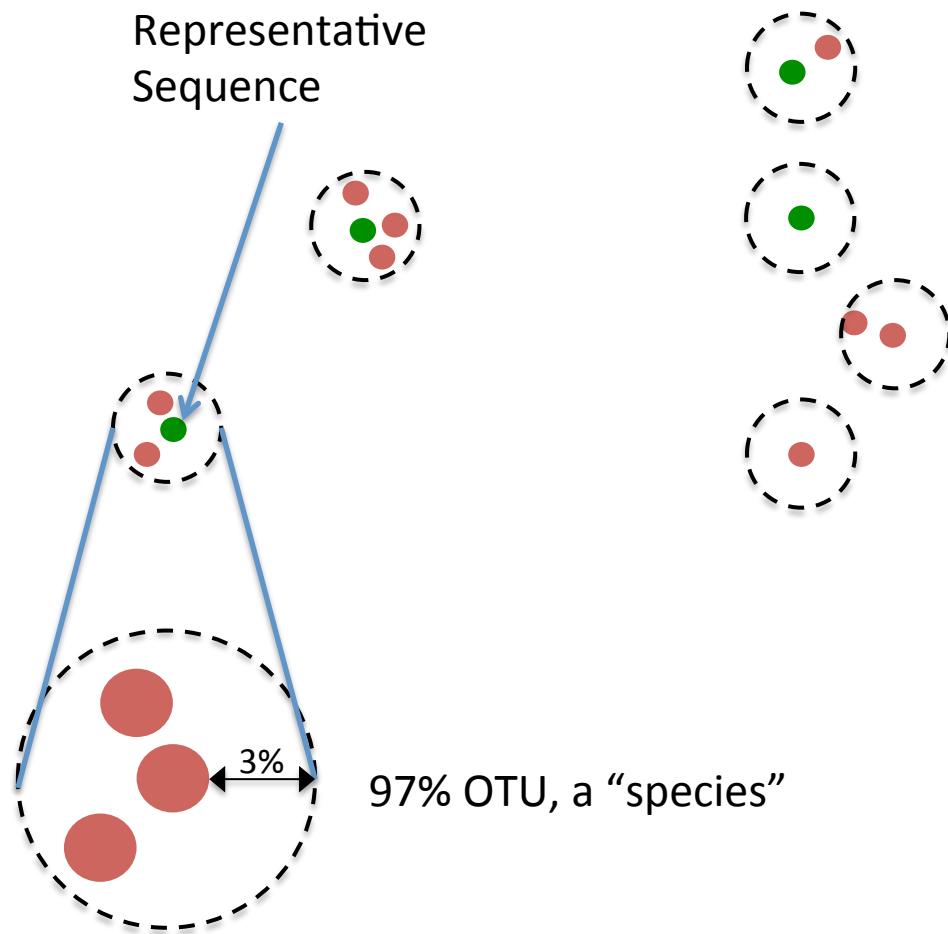
Understanding OTU picking



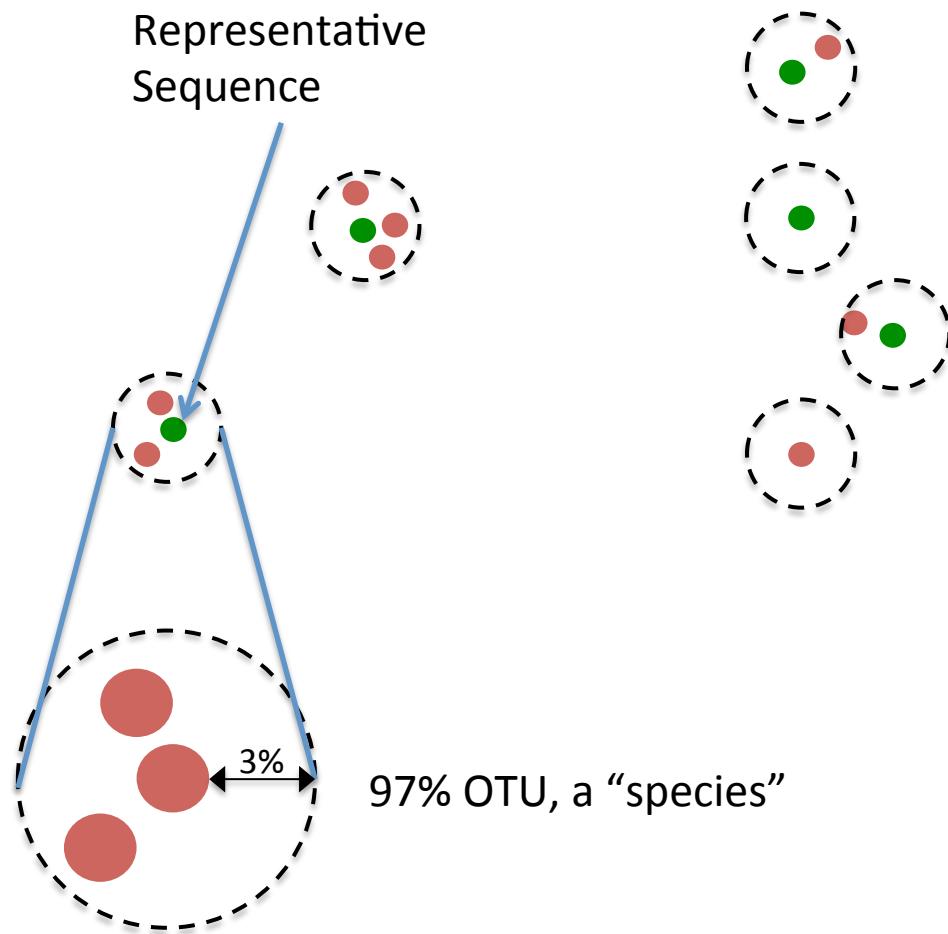
Understanding OTU picking



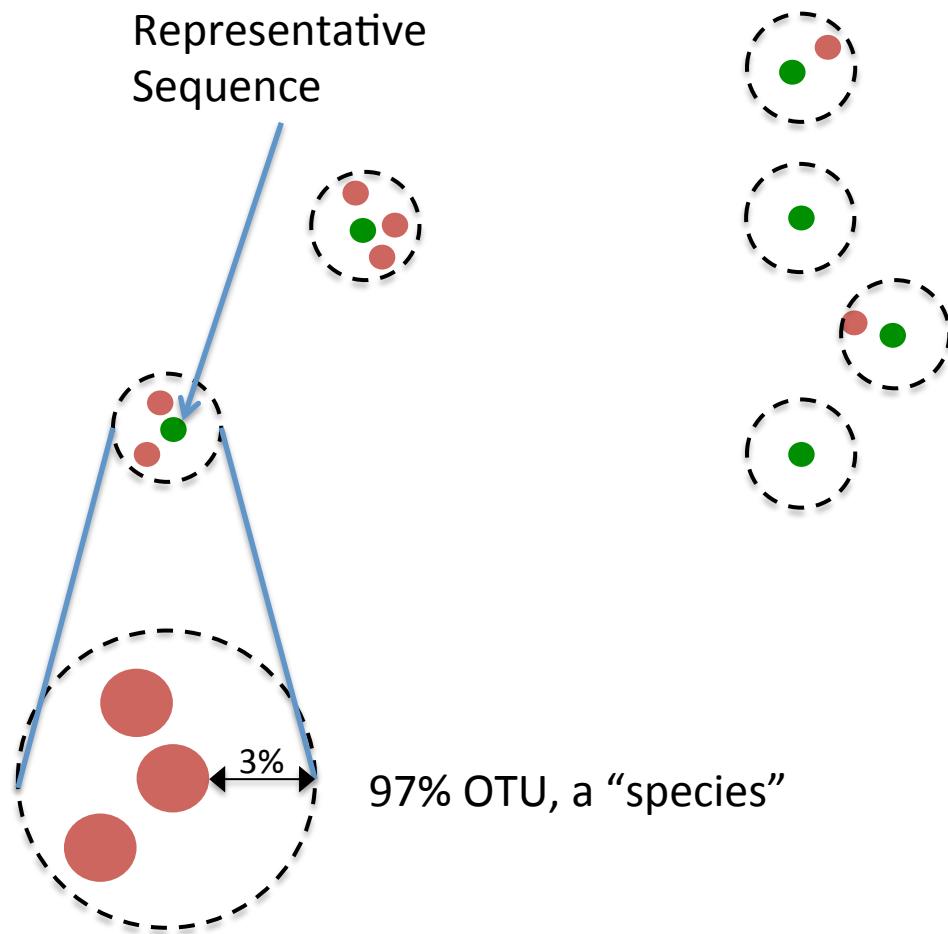
Understanding OTU picking



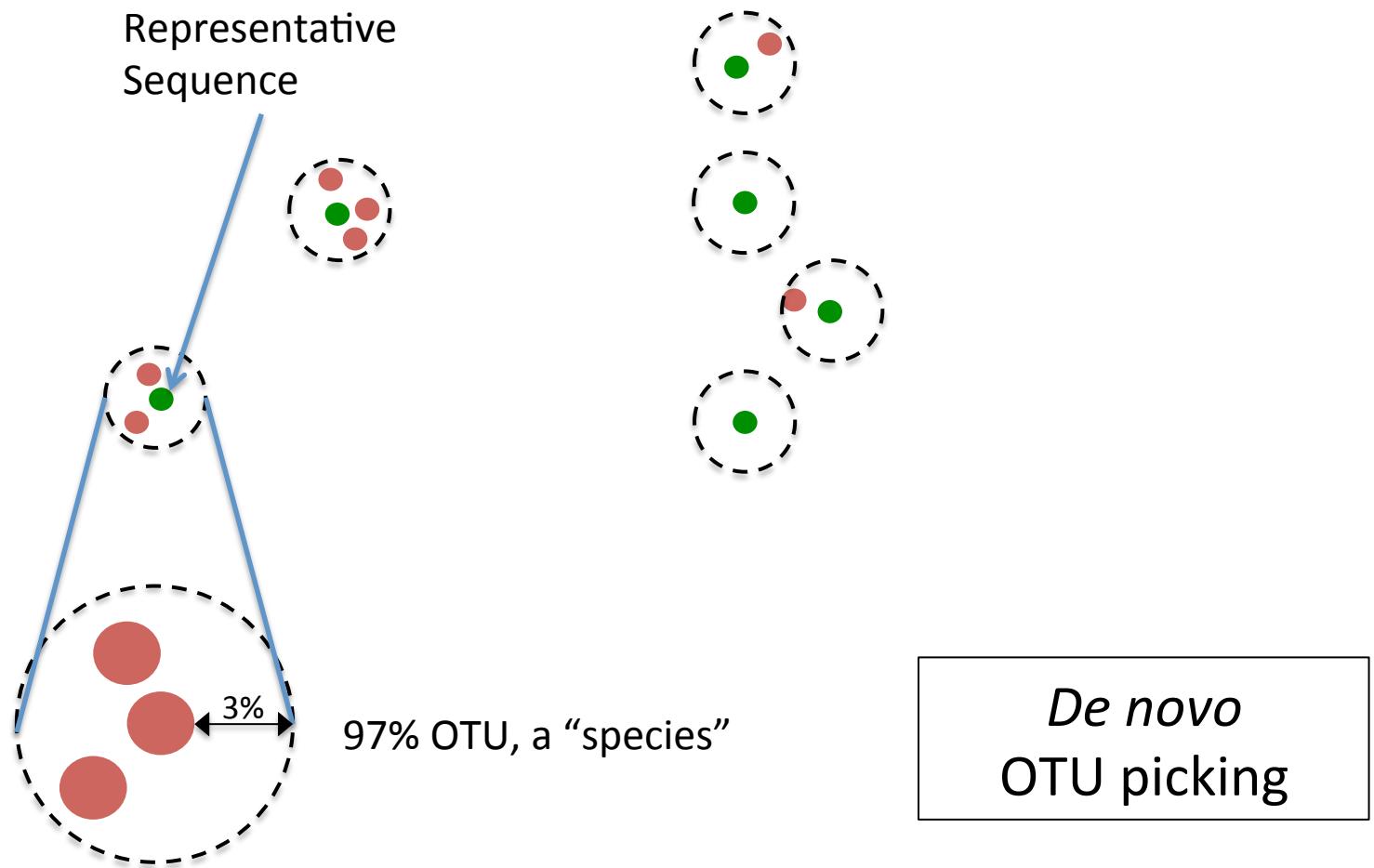
Understanding OTU picking



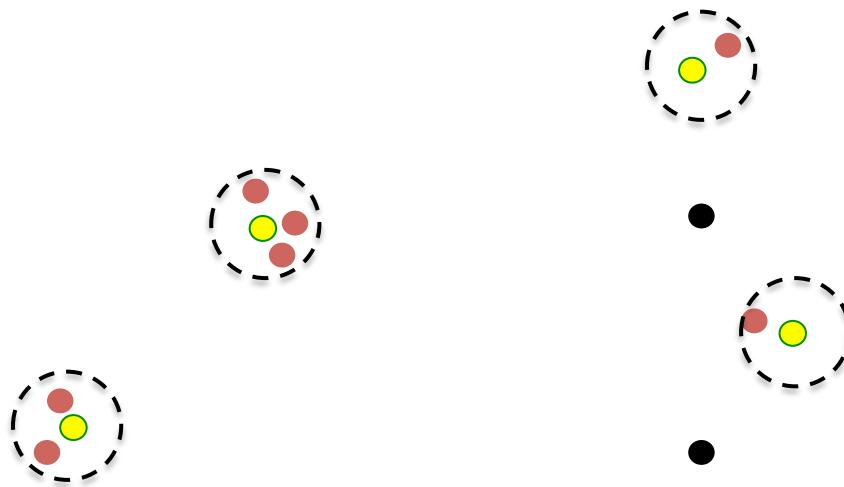
Understanding OTU picking



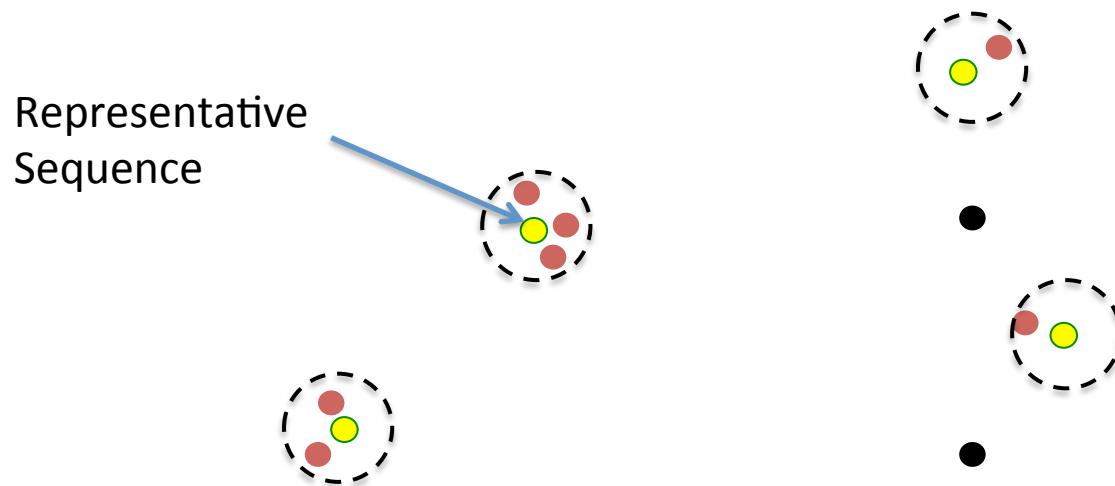
Understanding OTU picking



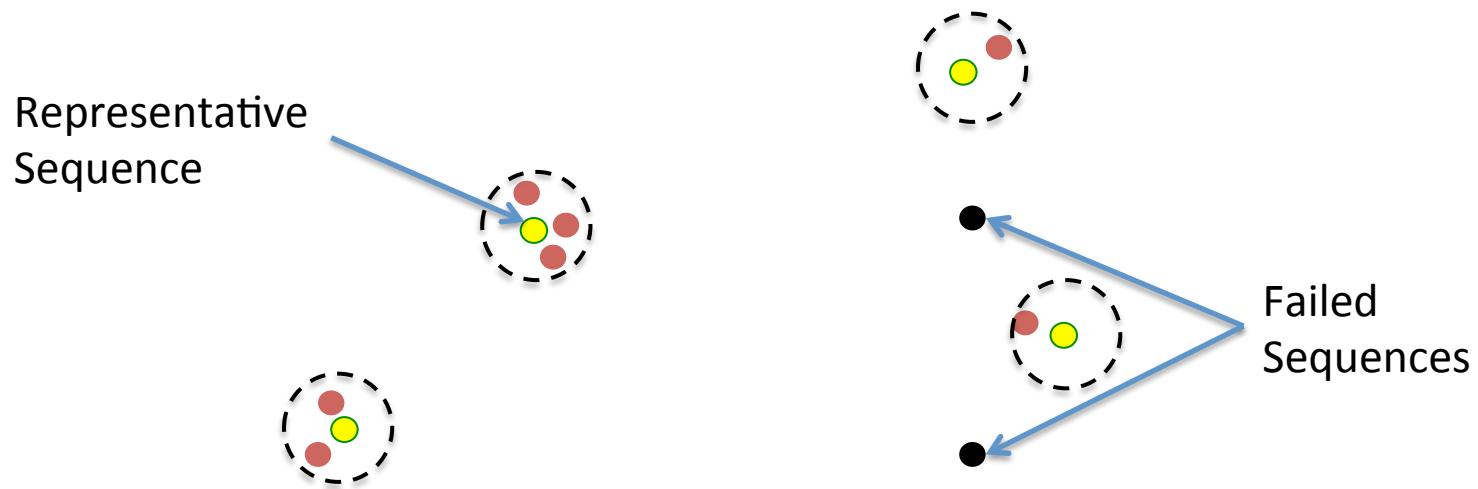
Understanding OTU picking



Understanding OTU picking



Understanding OTU picking

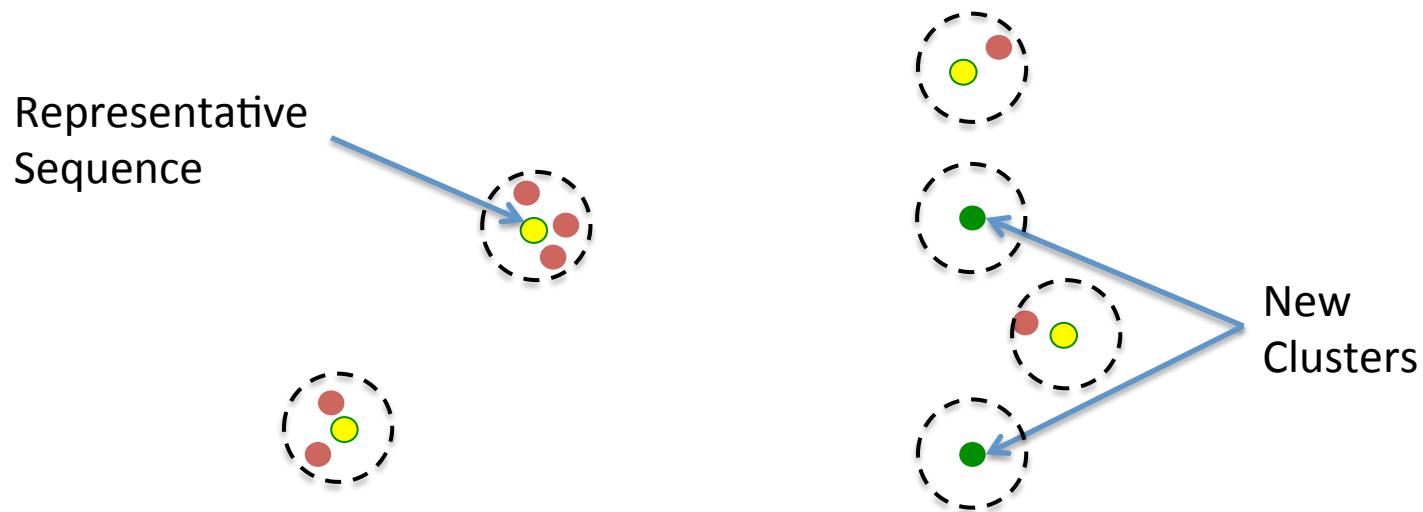


Understanding OTU picking

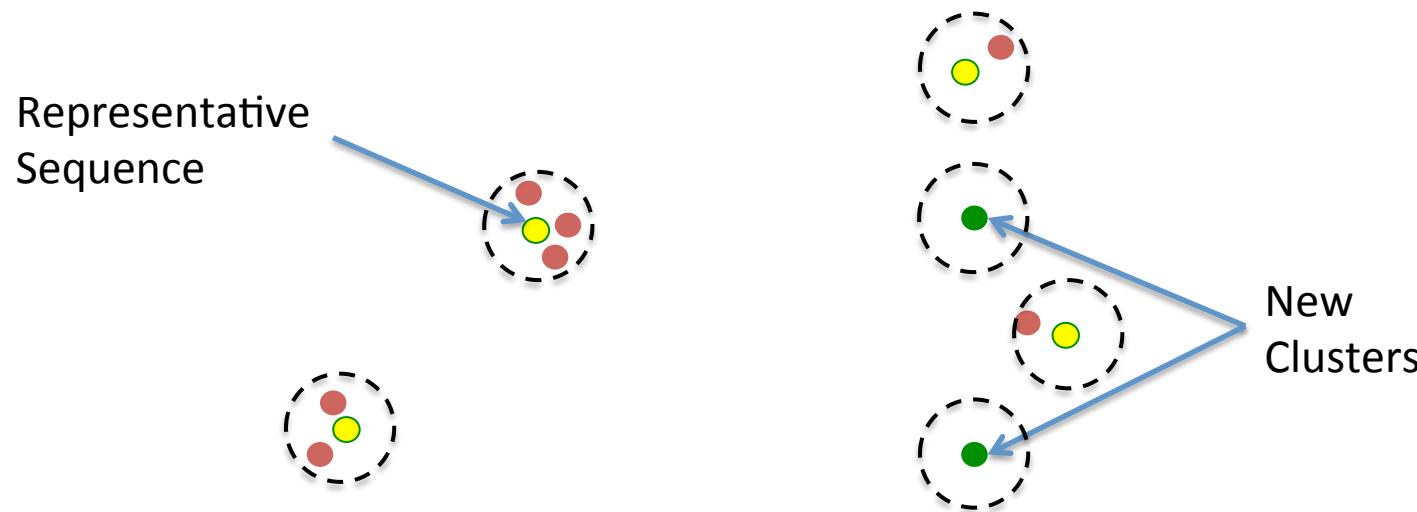


Closed reference
OTU picking

Understanding OTU picking



Understanding OTU picking



Open reference
OTU picking

De novo OTU picking

- Pros
 - All reads are clustered
- Cons
 - Not parallelizable
 - OTUs may be defined by erroneous reads

De novo OTU picking

- You **must** use if:
 - You do not have a reference sequence collection to cluster against, for example because you're working with an infrequently used marker gene.
- You **cannot** use if:
 - You are comparing non-overlapping amplicons, such as the V2 and the V4 regions of the 16S rRNA.
 - You are working with very large data sets, like a full HiSeq 2000 run. (Technically you can, but it will be *really* slow.)

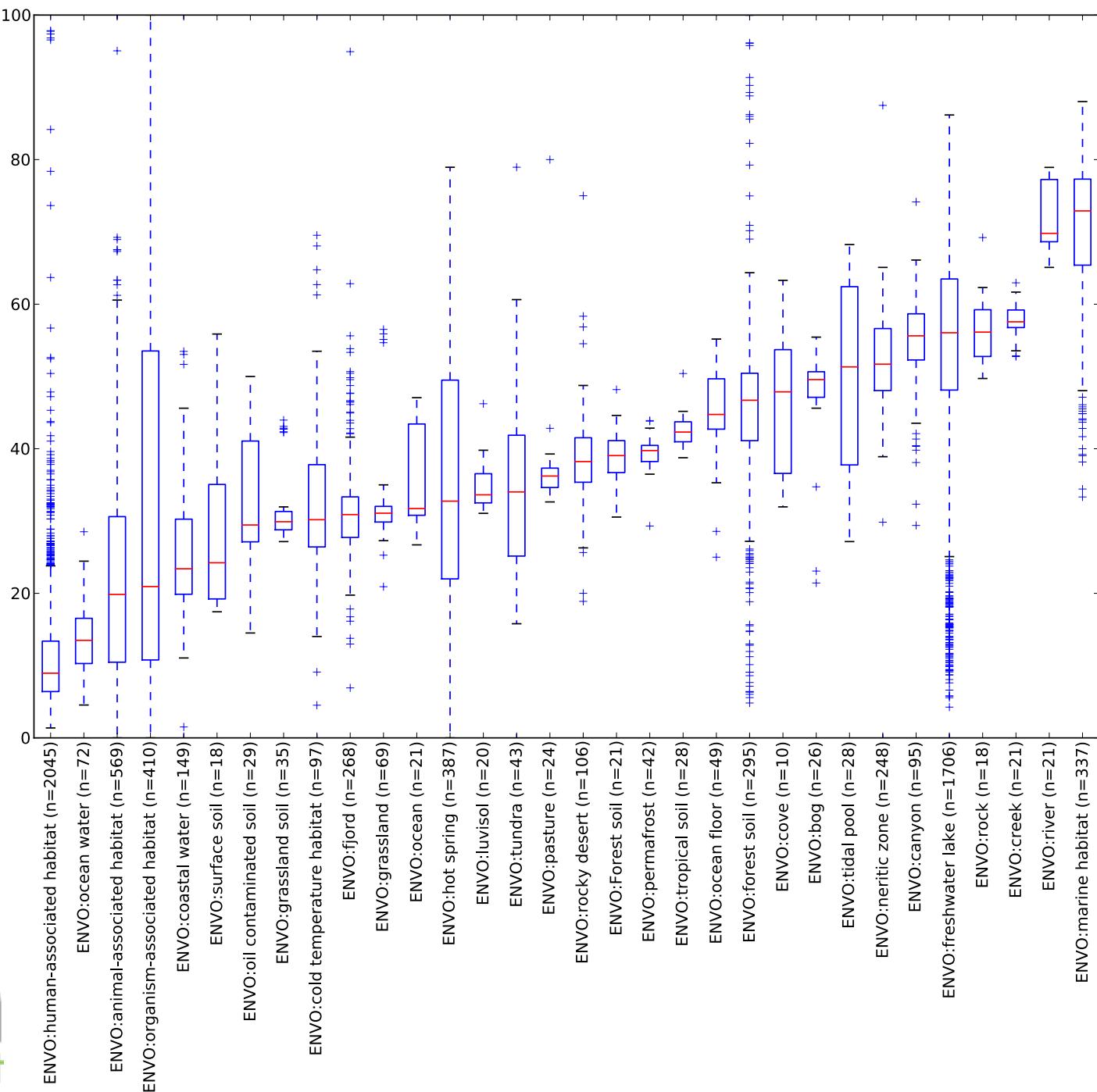
Closed-reference OTU picking

- Pros
 - Built-in quality filter
 - Easily parallelizable
 - OTUs are defined by high-quality, trusted sequences
- Cons
 - Reads that don't hit reference dataset are excluded, so you can never observe new OTUs

Closed-reference OTU picking

- You **must** use if:
 - You are comparing non-overlapping amplicons, such as the V2 and the V4 regions of the 16S rRNA. Your reference sequences must span both of the regions being sequenced.
- You **cannot** use if:
 - You do not have a reference sequence collection to cluster against, for example because you're working with an infrequently used marker gene.

Percentage of reads
that do not hit the
reference collection,
by environment type.



Open-reference OTU picking

- Pros
 - All reads are clustered
 - Partially parallelizable
- Cons
 - Only *partially* parallelizable
 - Mix of high quality sequences defining OTUs (i.e., the database sequences) and possible low quality sequences defining OTUs (i.e., the sequencing reads)

`pick_open_reference_otus.py`

http://qiime.org/tutorials/illumina_overview_tutorial.html

http://qiime.org/tutorials/open_reference_illumina_processing.html

http://qiime.org/tutorials/fungal_its_analysis.html

Open-reference OTU picking

- You **cannot** use if:
 - You are comparing non-overlapping amplicons, such as the V2 and the V4 regions of the 16S rRNA.
 - You do not have a reference sequence collection to cluster against, for example because you're working with an infrequently used marker gene.

`pick_open_reference_otus.py`

http://qiime.org/tutorials/illumina_overview_tutorial.html

http://qiime.org/tutorials/open_reference_illumina_processing.html

http://qiime.org/tutorials/fungal_its_analysis.html

Subsampled open-reference clustering creates consistent, comprehensive OTU definitions and scales to billions of sequences

Jai Ram Rideout^{1,2}, Yan He³, Jose A. Navas-Molina⁴, William A. Walters⁵,
Luke K. Ursell⁶, Sean M. Gibbons^{7,10}, John Chase⁸, Daniel McDonald^{4,9},
Antonio Gonzalez⁹, Adam Robbins-Pianka^{4,9}, Jose C. Clemente², Jack A. Gilbert^{10,11},
Susan M. Huse¹², Hong-Wei Zhou³, Rob Knight^{9,13}, J. Gregory Caporaso  ^{1,8}

Published August 21, 2014

PubMed [25177538](#)

OTU picking methods vs algorithms

- QIIME provides 3 high-level **methods**
- QIIME wraps OTU clustering tools
 - Performs sequence clustering **algorithm**
 - uclust, usearch, SortMeRNA, SumaClust, swarm, ...
 - Be sure to cite these tools too!

OTU table

OTU x sample matrix

Count table

	S1	S2	S3
OTU1	100	0	0
OTU2	100	40	600
OTU3	0	10	0

Relative abundance table

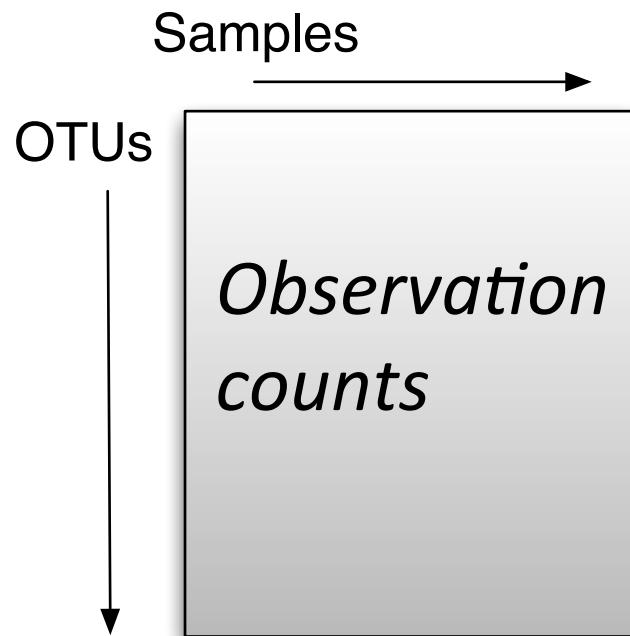
	S1	S2	S3
OTU1	.5	0	0
OTU2	.5	.8	1.0
OTU3	0	.2	0

OTU tables are in Biological
Observation Matrix (*.biom*) format

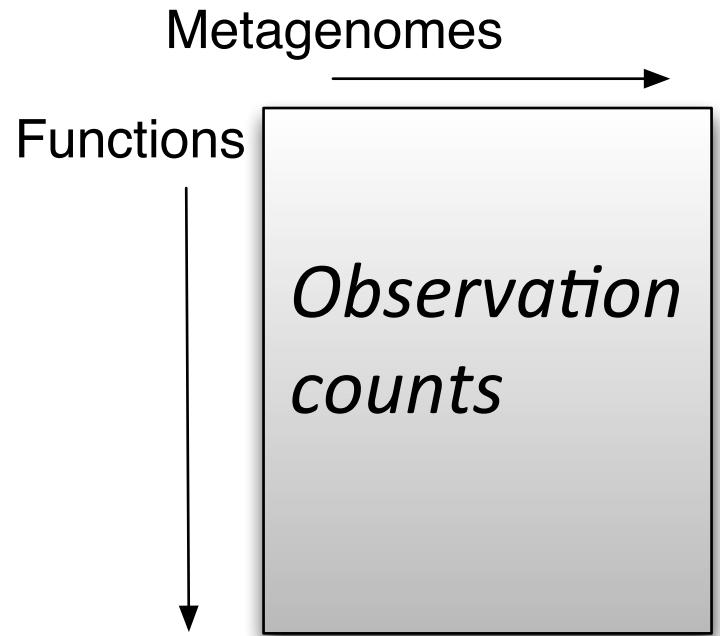
<http://biom-format.org>



sample x observation contingency matrix

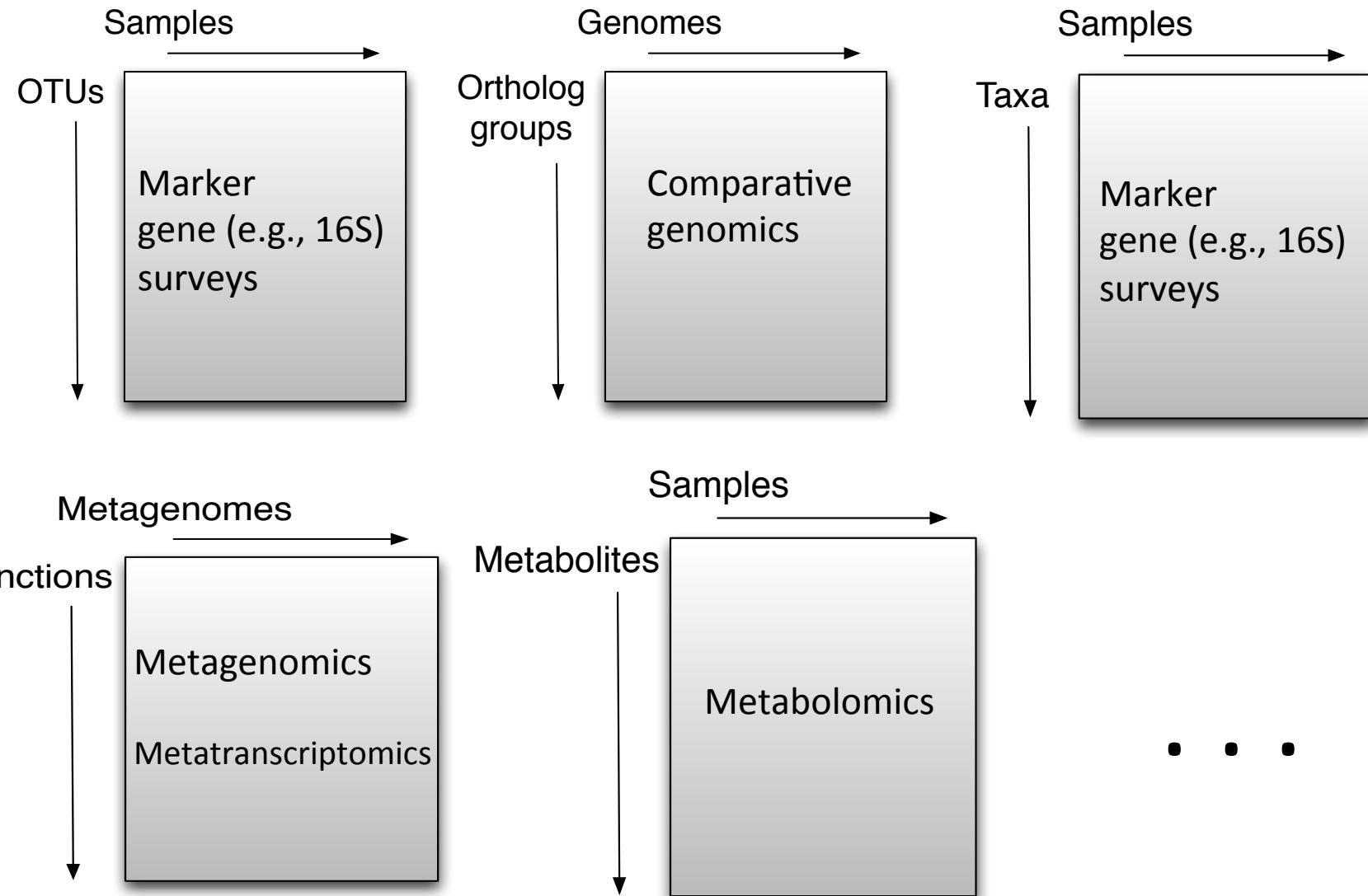


sample x observation contingency matrix



MG-RAST
metagenomics analysis server

sample x observation contingency matrix



The Biological Observation Matrix (BIOM) Format or: How I Learned To Stop Worrying and Love the Ome-ome

Format for representing
arbitrary sample x
observation contingency
tables with optional
metadata



VAMPS
The Visualization and Analysis
of Microbial Population Structures

<http://www.biom-format.org>

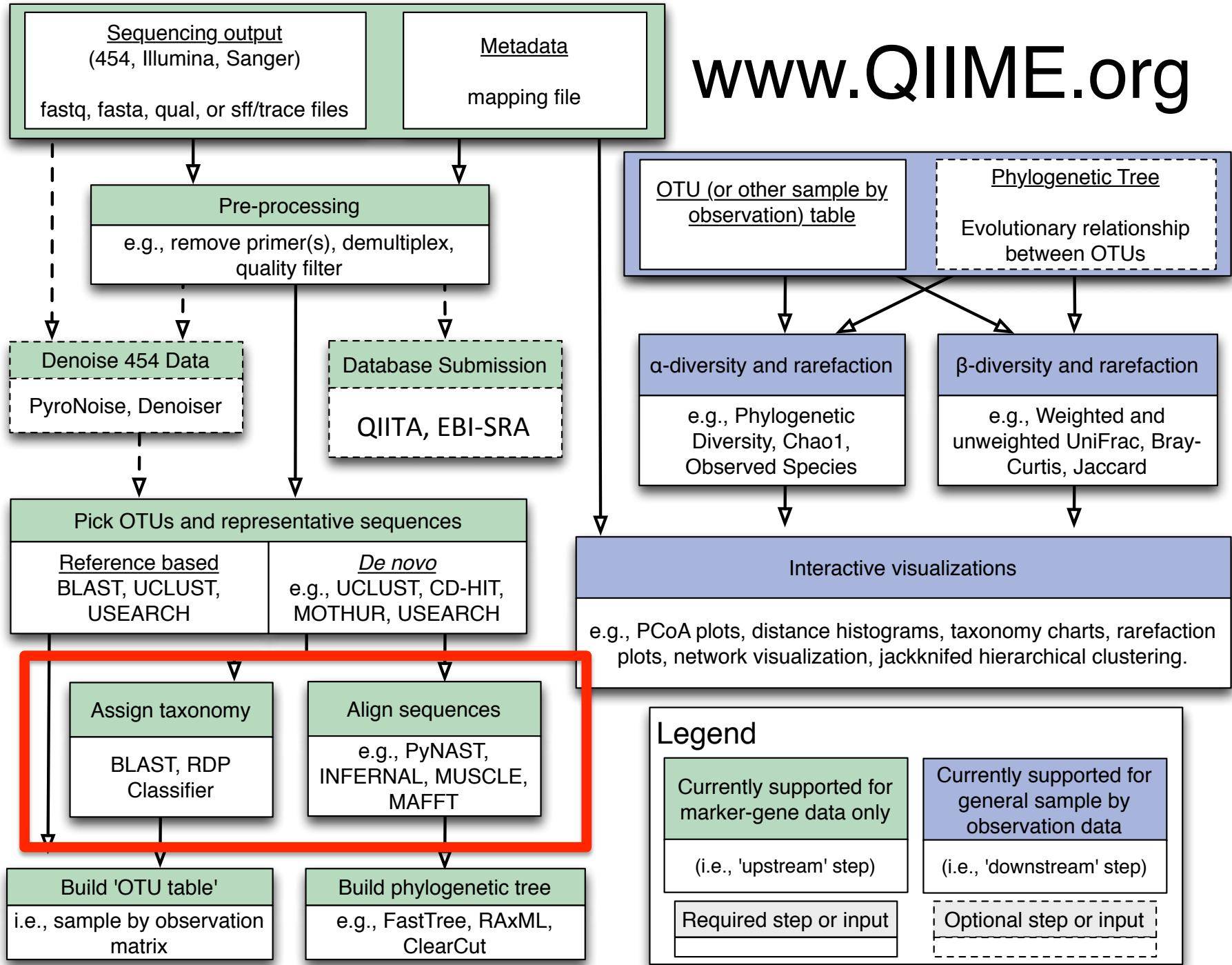
Reference databases

- Greengenes (16S)
 - Default reference database included in QIIME 1.9.0 and later
- SILVA (16S/18S)
- UNITE (ITS)
- IMG (protein sequences)



Reference databases

- Latest QIIME-compatible databases:
 - http://qiime.org/home_static/dataFiles.html
- Tutorials
 - 16S
 - http://qiime.org/tutorials/illumina_overview_tutorial.html
 - <http://qiime.org/tutorials/tutorial.html>
 - ITS
 - http://qiime.org/tutorials/fungal_its_analysis.html
 - 18S
 - http://qiime.org/tutorials/processing_18S_data.html
 - Shotgun data (metagenome data)
 - http://qiime.org/tutorials/shotgun_analysis.html



Representative Sequences

If de novo or open-reference OTU picking is used, a representative sequence will be selected for each OTU. By default, this will be the most abundant read for that OTU. These reads are used for taxonomic assignment and tree building.

Closed-reference OTU picking does not generate a tree or taxonomic assignment. The tree and taxonomies in the reference files are used instead.

Taxonomic Assignment

With workflow scripts, you do not have to manually assign taxonomy, but the underlying script being used is `assign_taxonomy.py`.

Each representative sequence is assigned-some methods can have ambiguous results, while others return the best match.

Phylogenetic Trees

The representative sequences are aligned, then filtered for highly gapped/variable bases, and finally a tree is built.

The sequences that fail to align can be filtered out of the final OTU table.

Scripts used:

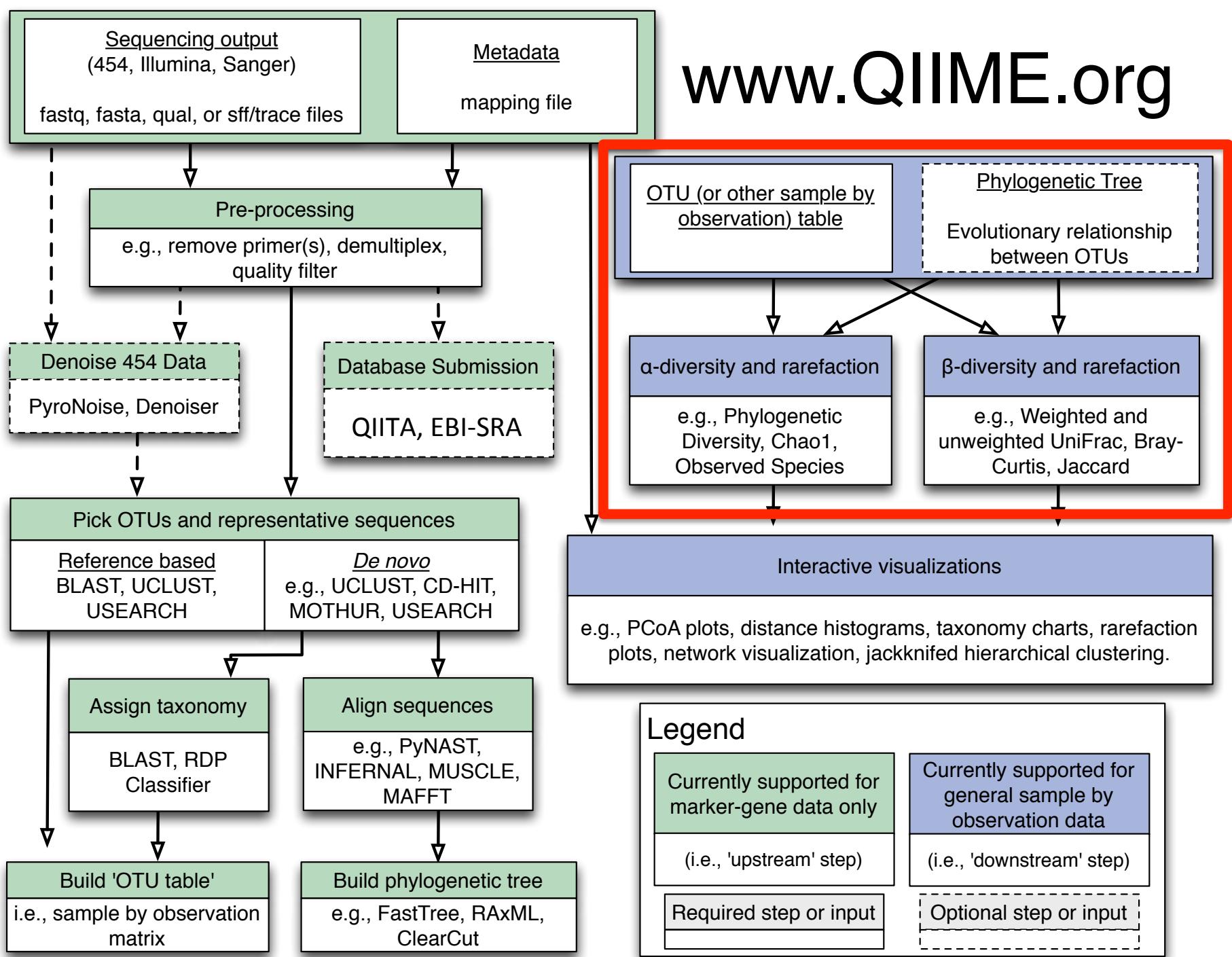
`align_seqs.py`

`filter_alignment.py`

`make_phylogeny.py`

Diversity and statistical analysis

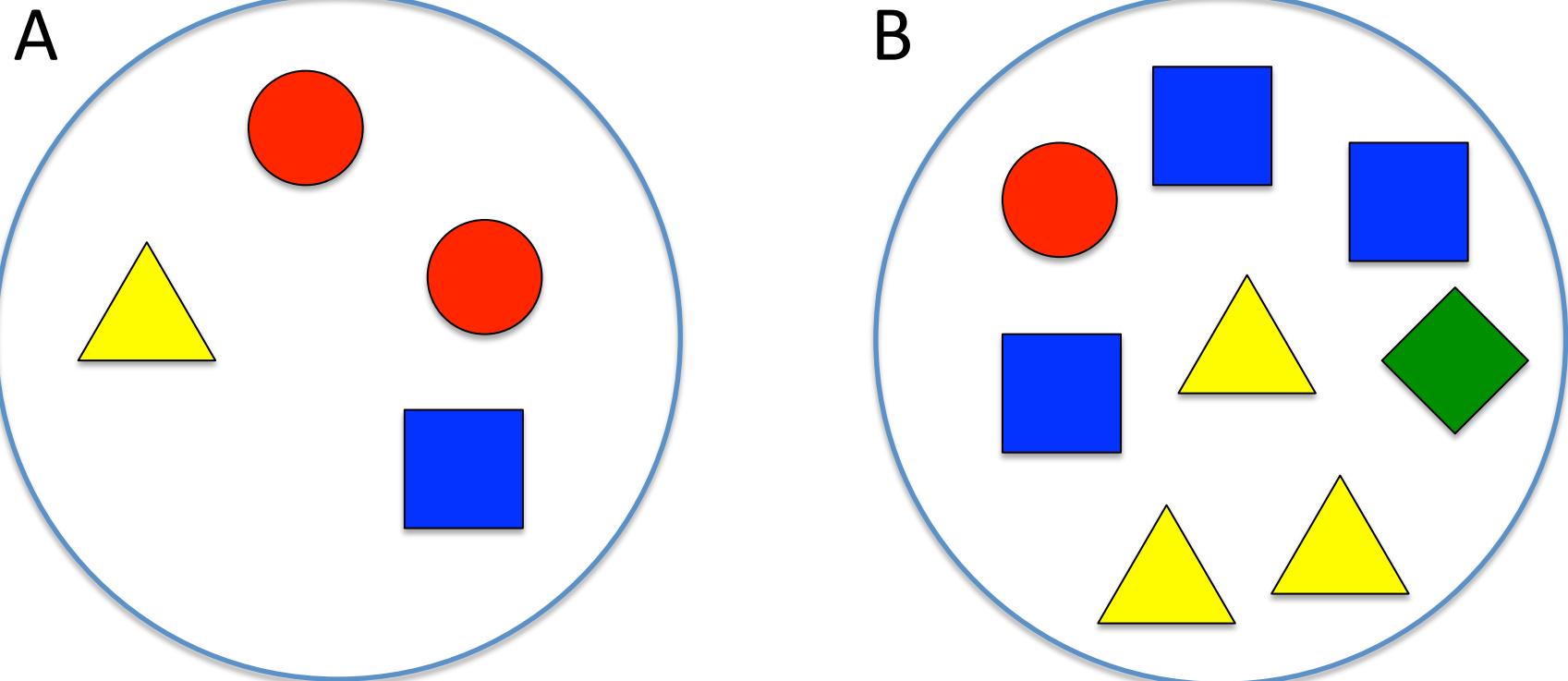
Part I



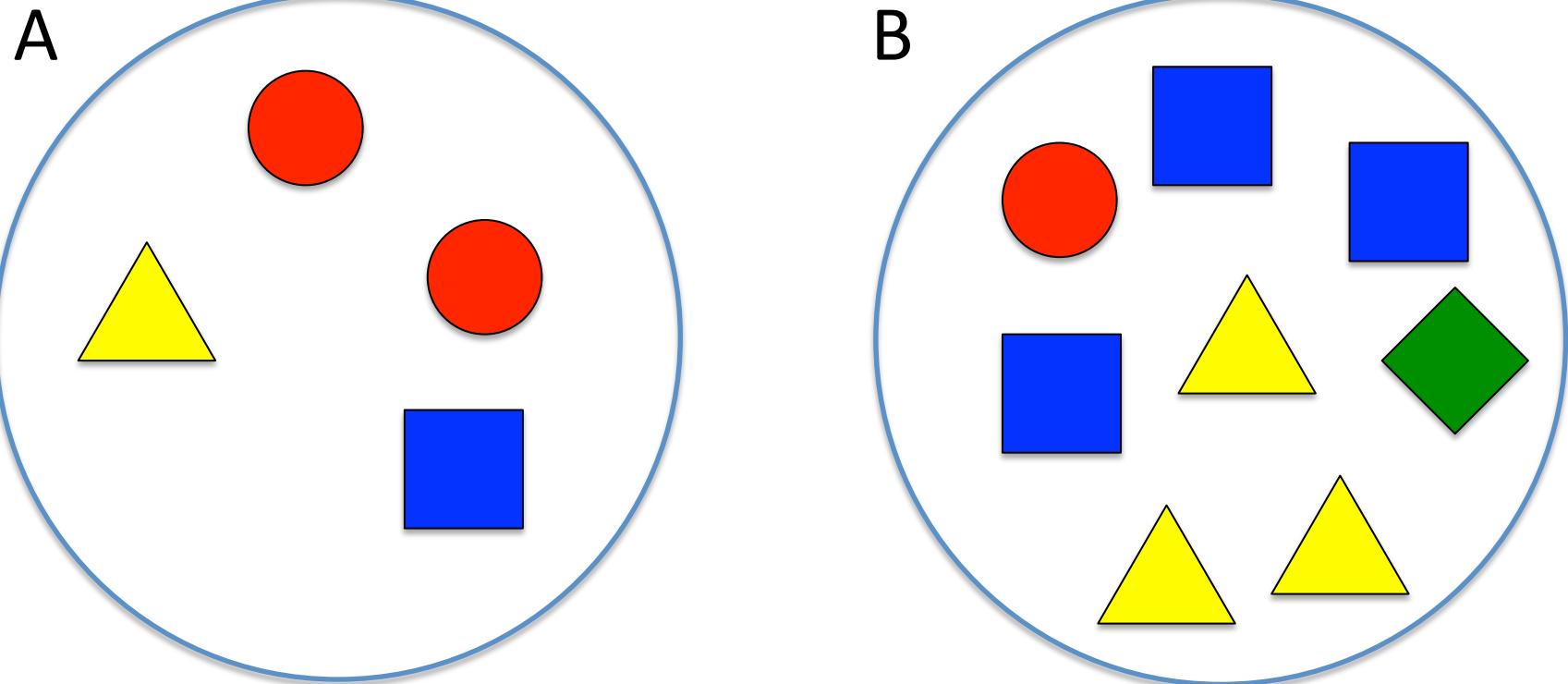
Alpha and beta diversity

- Alpha is within a sample
 - E.g., how many species are in a sample
- Beta is between samples
 - E.g., how similar are two samples
- Lots of ways to calculate these

Alpha diversity

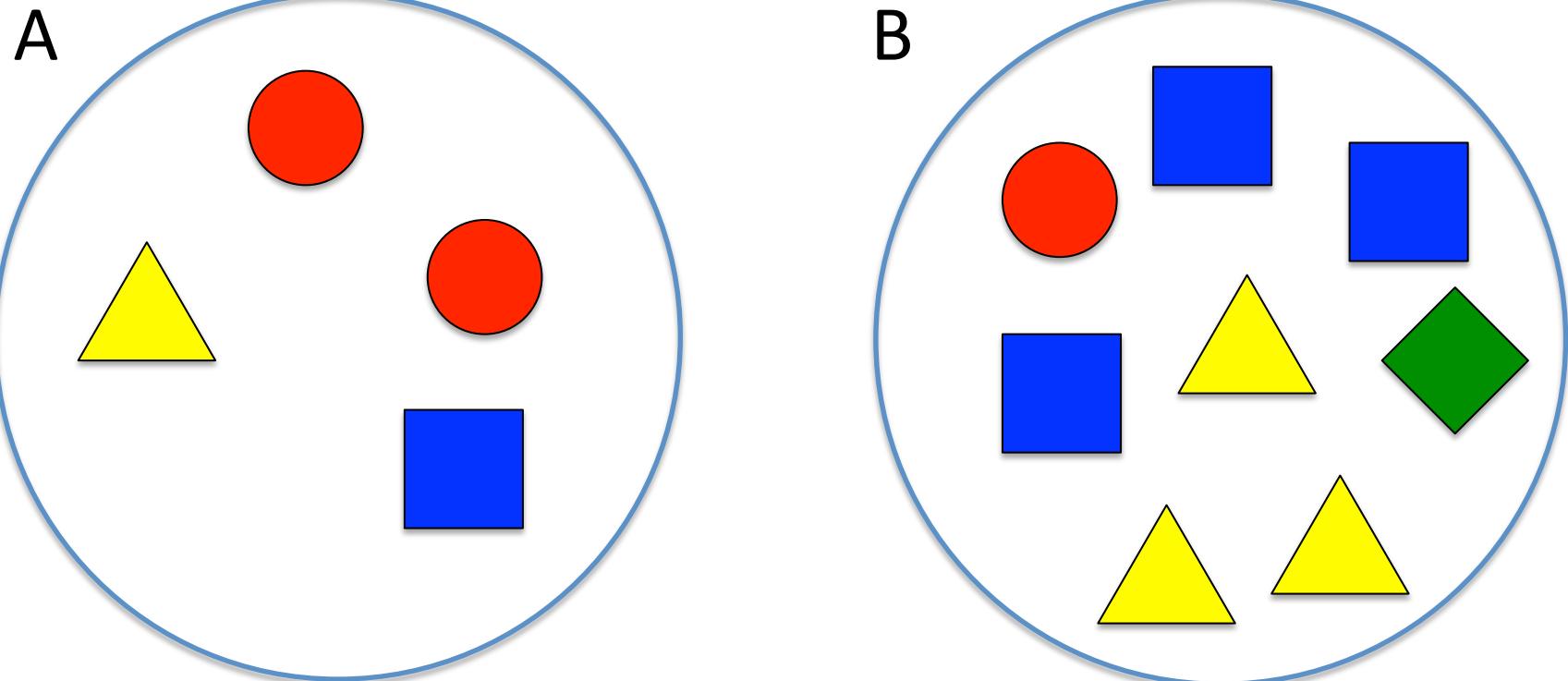


Rarefaction



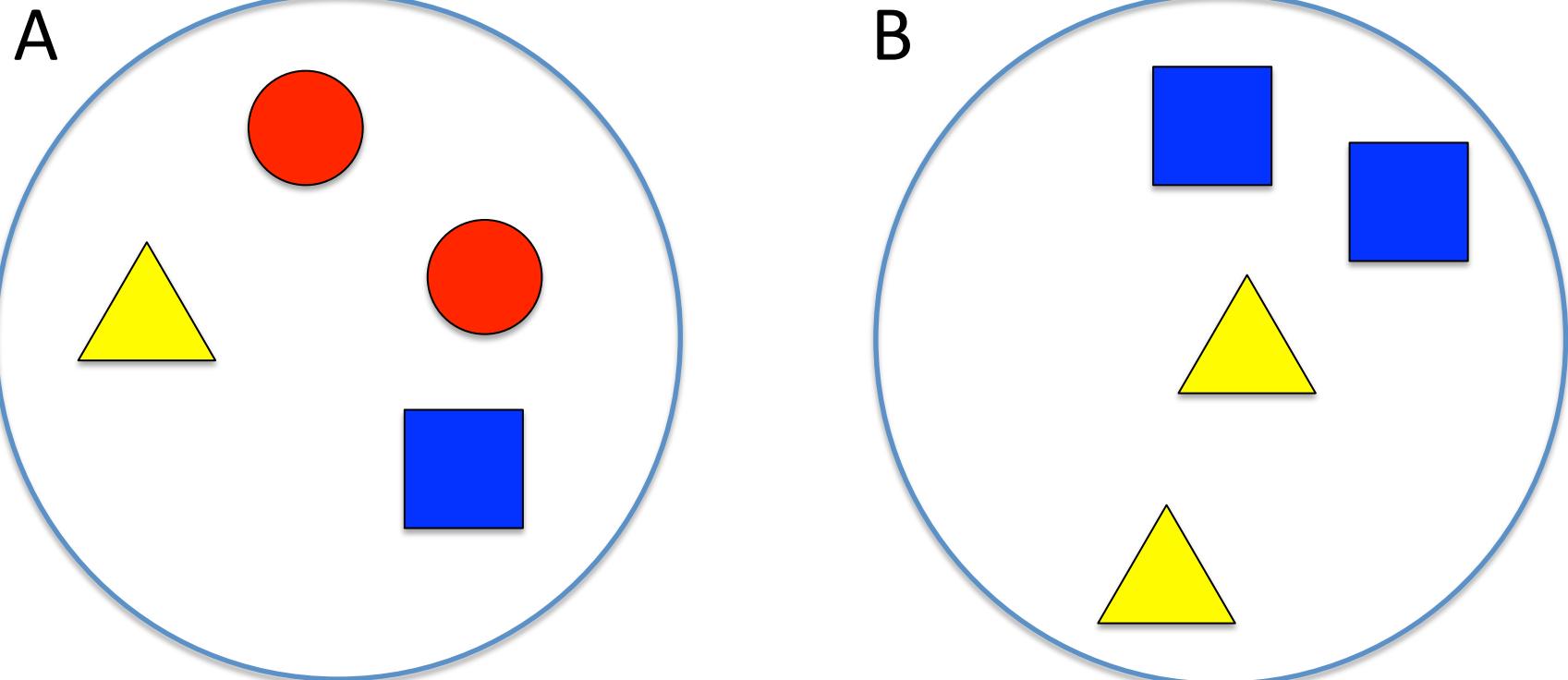
Randomly select 4 sequences from B

Rarefaction

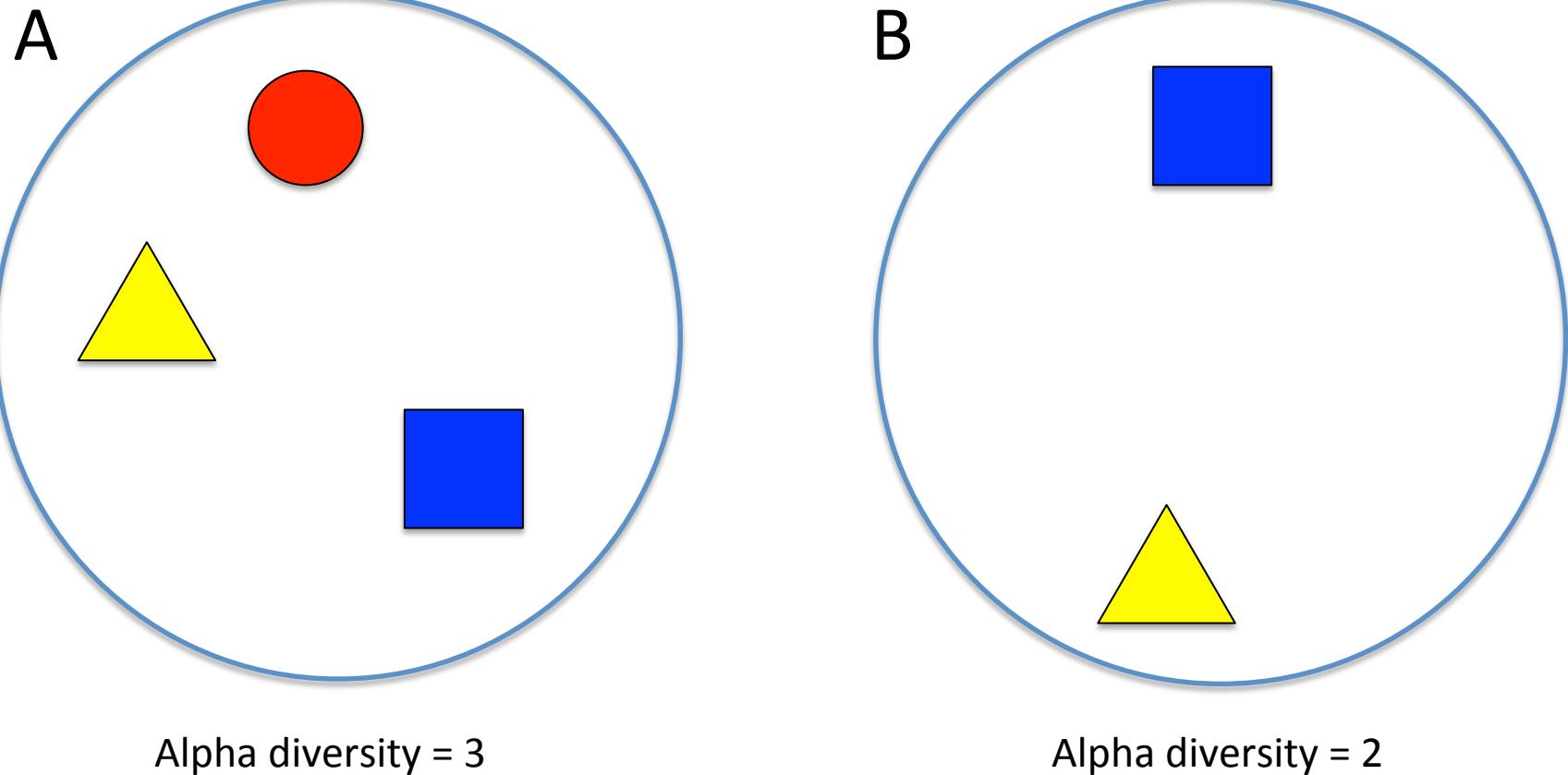


Rarefy to 4 sequences

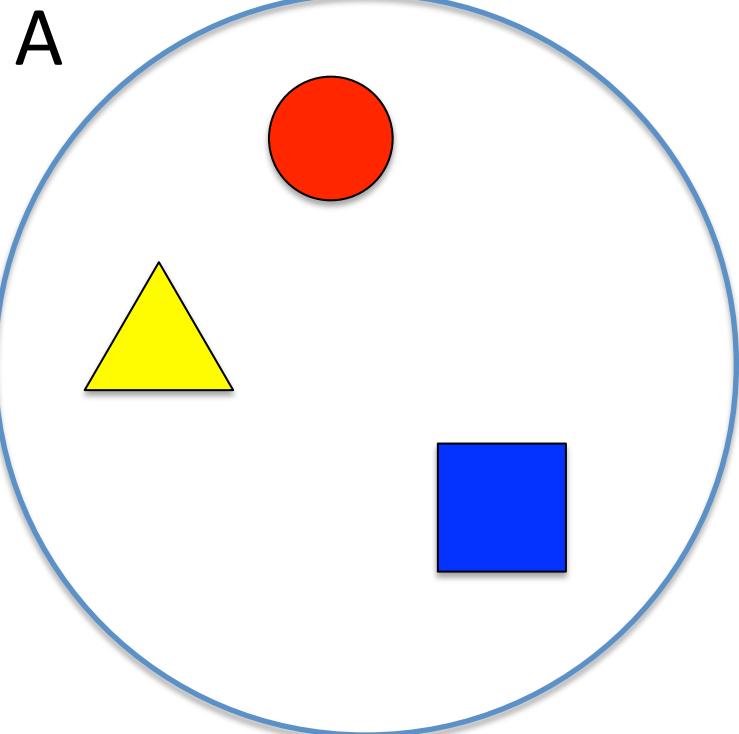
Rarefaction



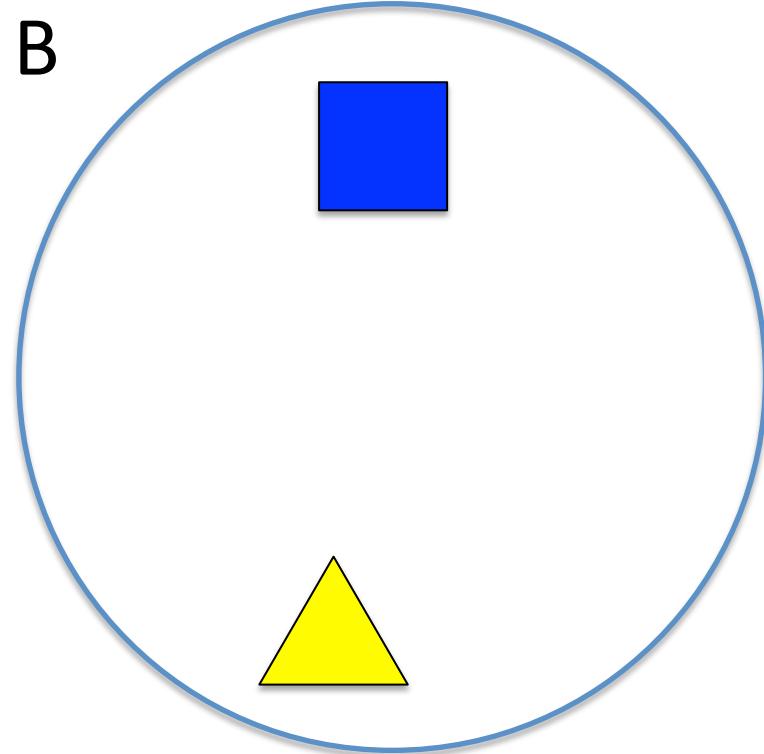
Rarefaction



Rarefaction



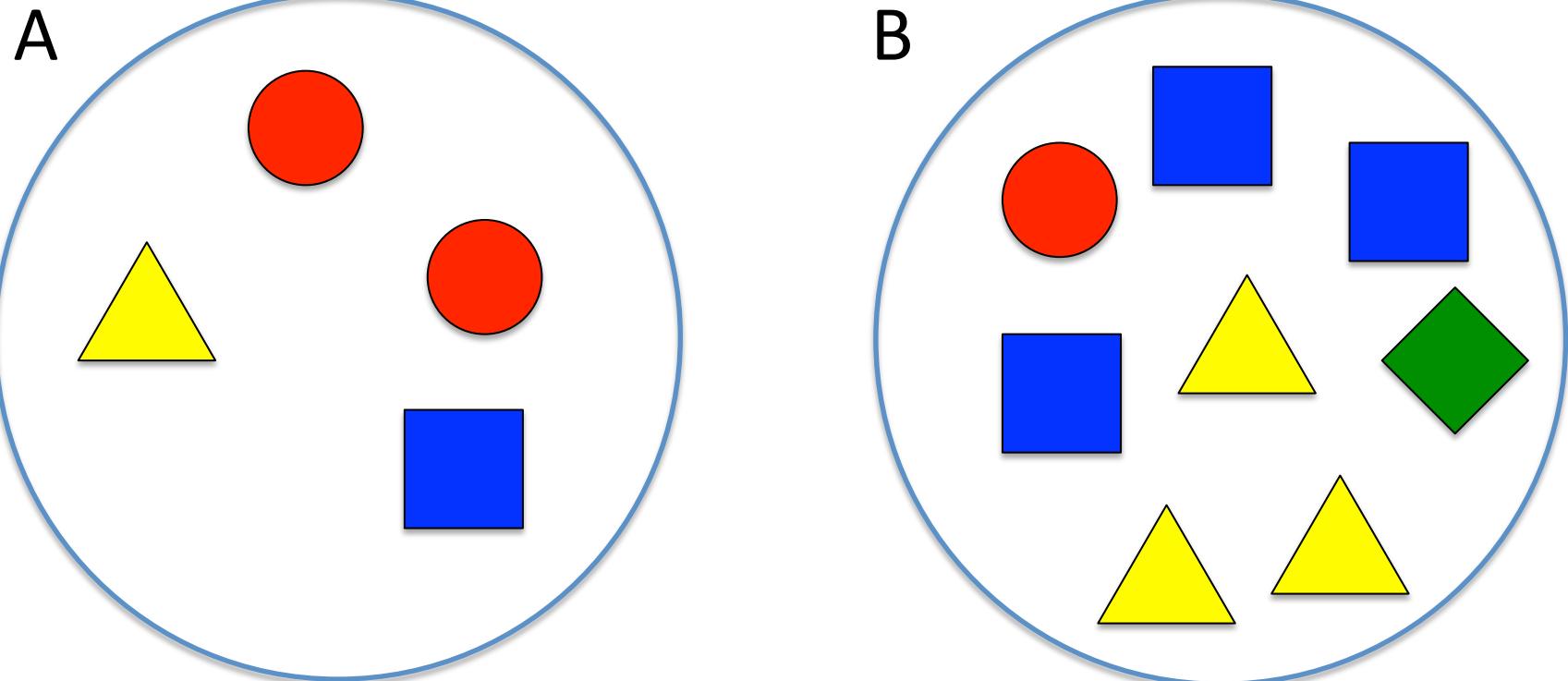
Alpha diversity = 3



Alpha diversity = 2

Sample A is more diverse than sample B

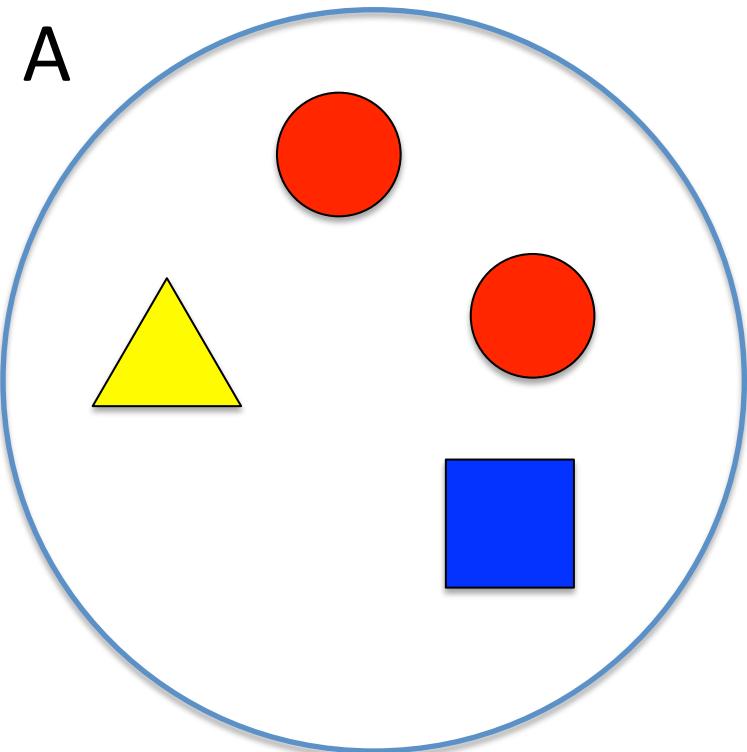
Rarefaction



Rarefy to 4 sequences

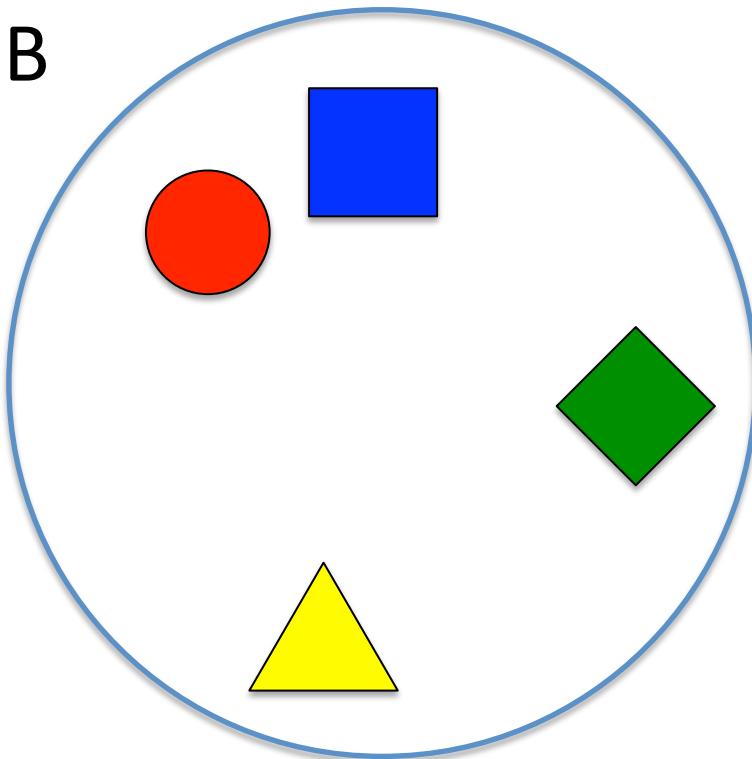
Alpha diversity

A



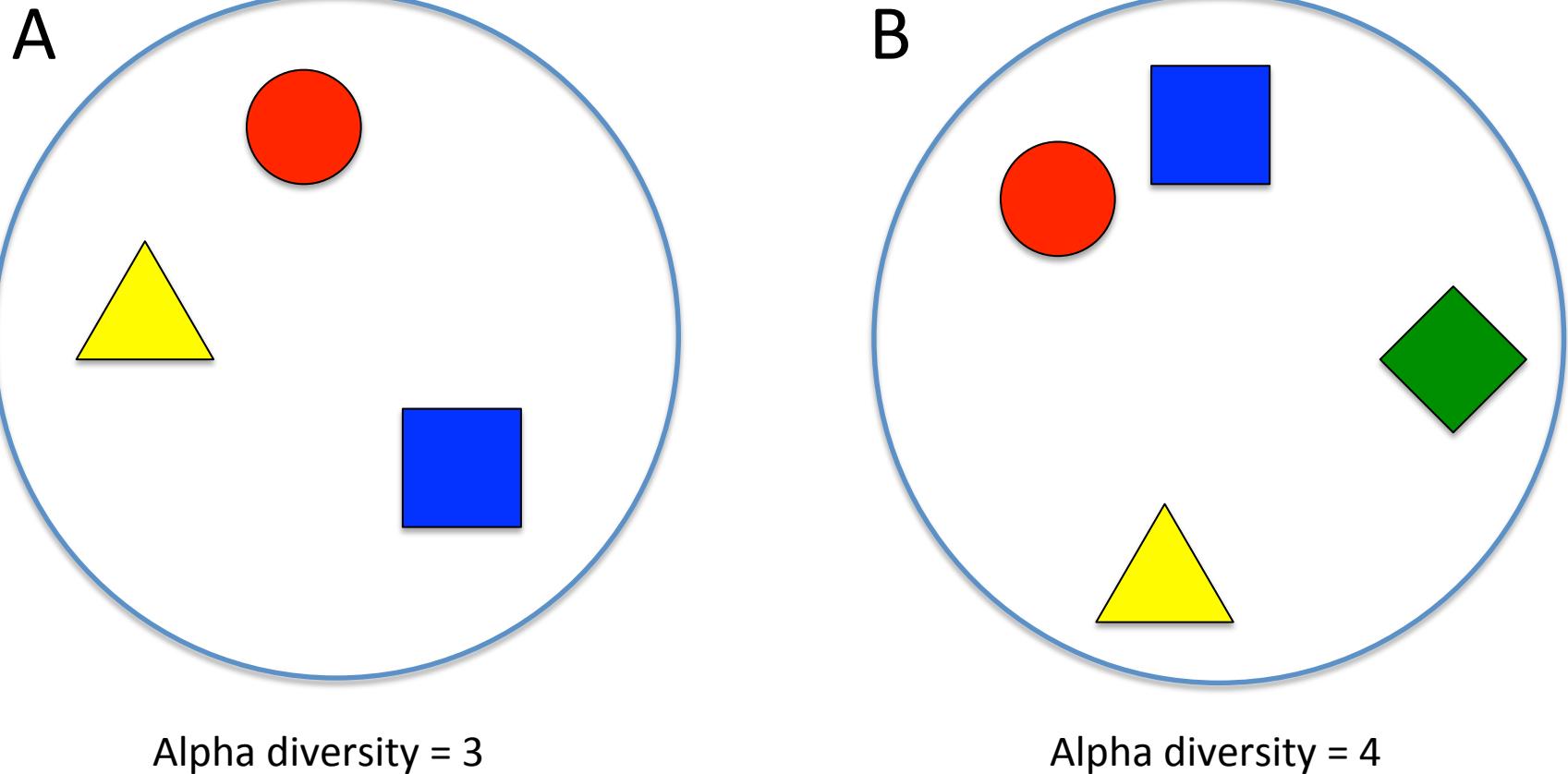
Alpha diversity = 3

B



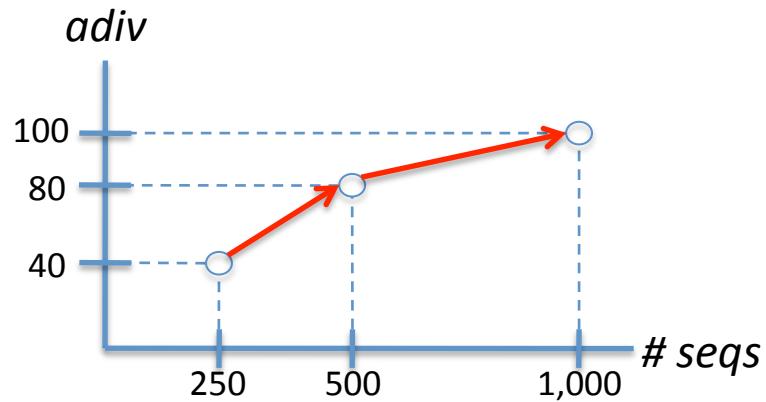
Alpha diversity = 4

Alpha diversity



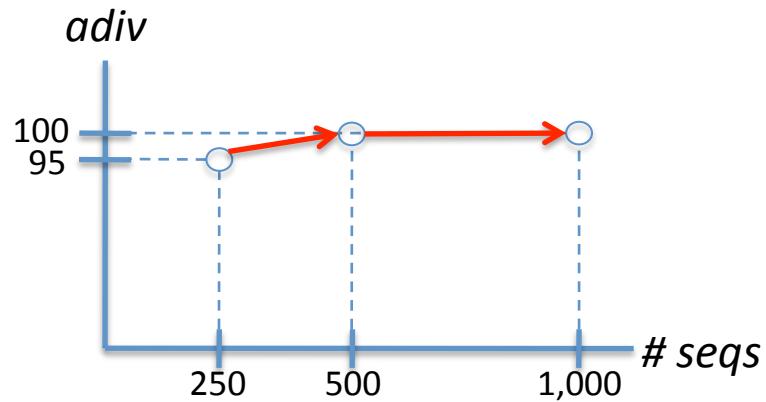
Sample B is more diverse than sample A

Multiple Rarefaction



Higher sequencing effort will probably result in higher observed diversity

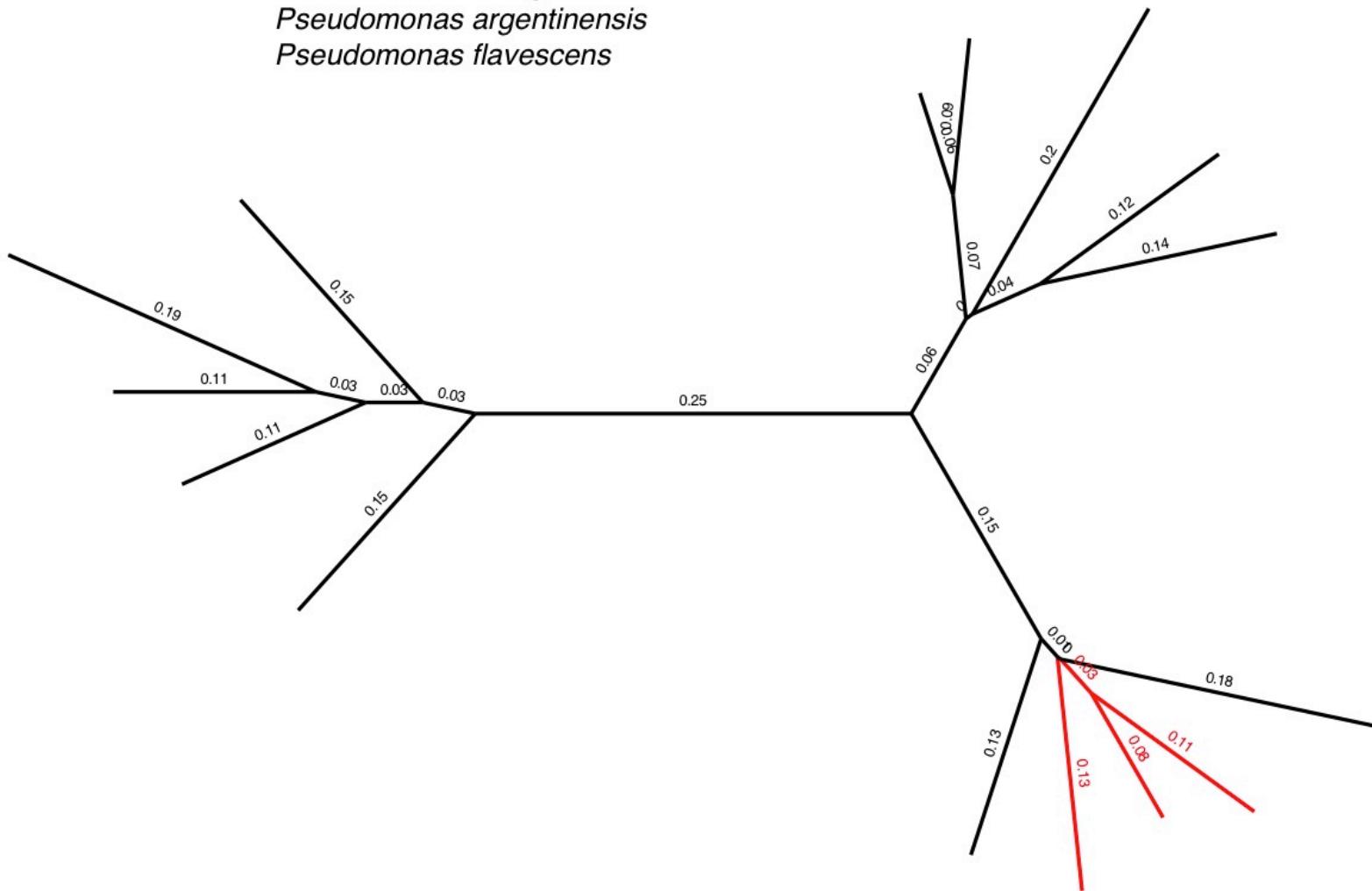
Multiple Rarefaction



Higher sequencing effort will probably
not add to observed diversity

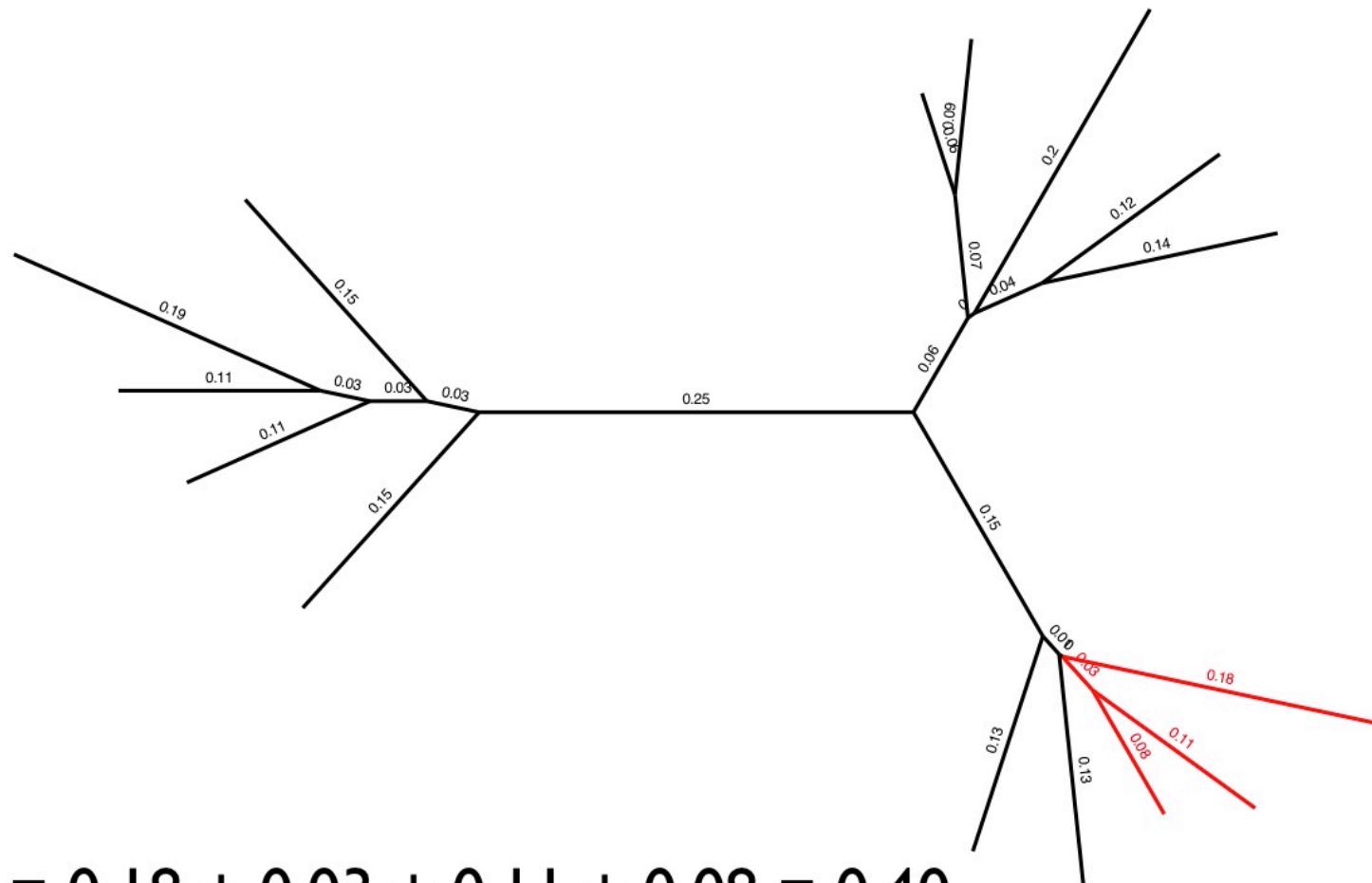
Sample A

Pseudomonas aeruginosa
Pseudomonas argentinensis
Pseudomonas fluorescens



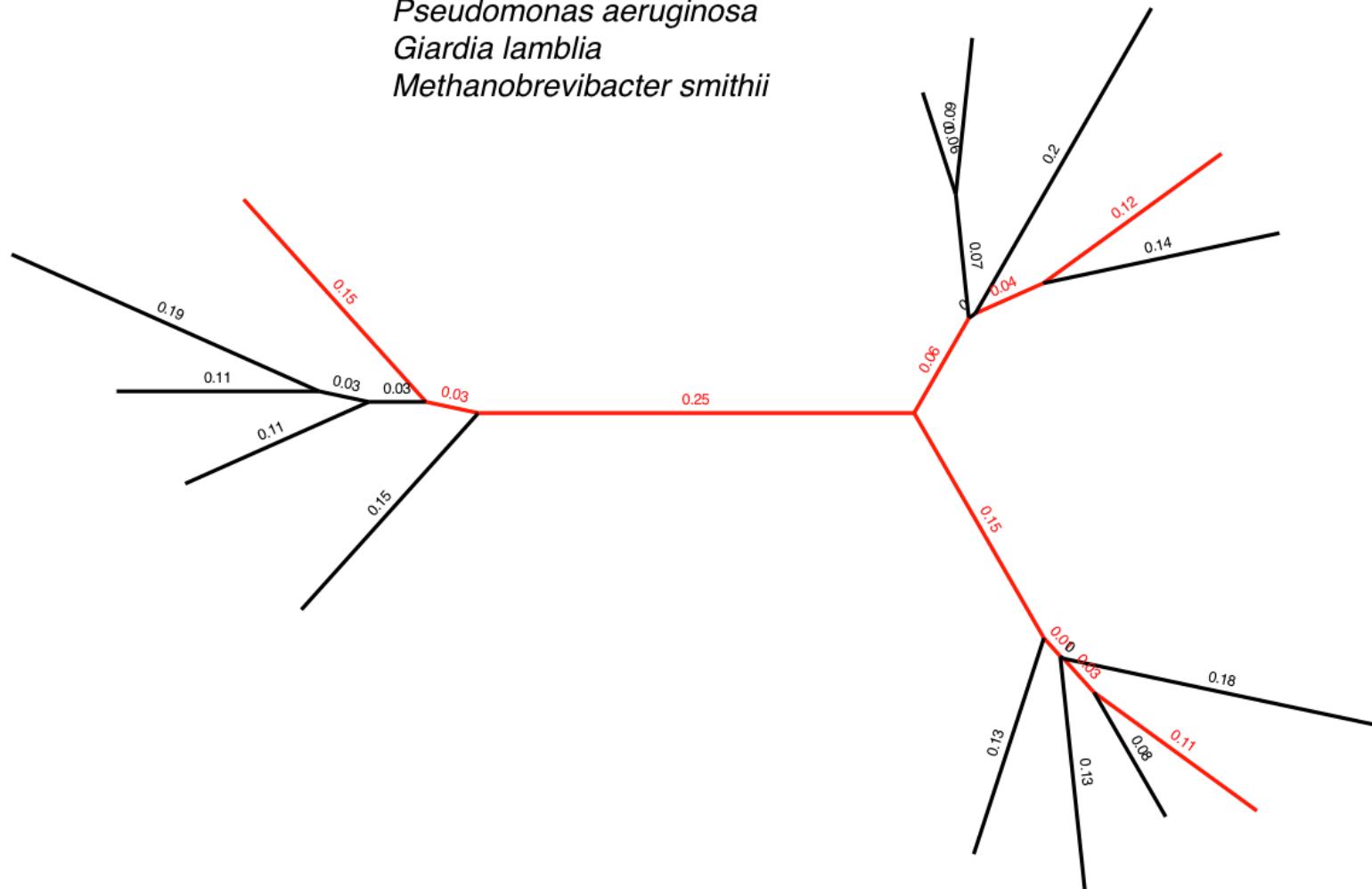
Sample B

Pseudomonas aeruginosa
Pseudomonas argentinensis
Escherichia coli



Sample C

Pseudomonas aeruginosa
Giardia lamblia
Methanobrevibacter smithii



$$PD = 0.15 + 0.03 + 0.25 + 0.06 + 0.04 + 0.12 + 0.15 + 0.01 + 0.03 + 0.11 = 0.95$$

Phylogenetic Diversity (PD)

Sample A

Pseudomonas aeruginosa
Pseudomonas argentinensis
Pseudomonas fluorescens

PD = 0.35

Sample B

Pseudomonas aeruginosa
Pseudomonas argentinensis
Escherichia coli

PD = 0.40

Sample C

Pseudomonas aeruginosa
Giardia lamblia
Methanobrevibacter smithii

PD = 0.95

Phylogenetic Diversity (PD)

Sample A

Pseudomonas aeruginosa
Pseudomonas argentinensis
Pseudomonas fluorescens

Sample B

Pseudomonas aeruginosa
Pseudomonas argentinensis
Escherichia coli

Sample C

Pseudomonas aeruginosa
Giardia lamblia
Methanobrevibacter smithii

$$\text{PD} = 0.35 \quad < \quad \text{PD} = 0.40 \quad < \quad \text{PD} = 0.95$$

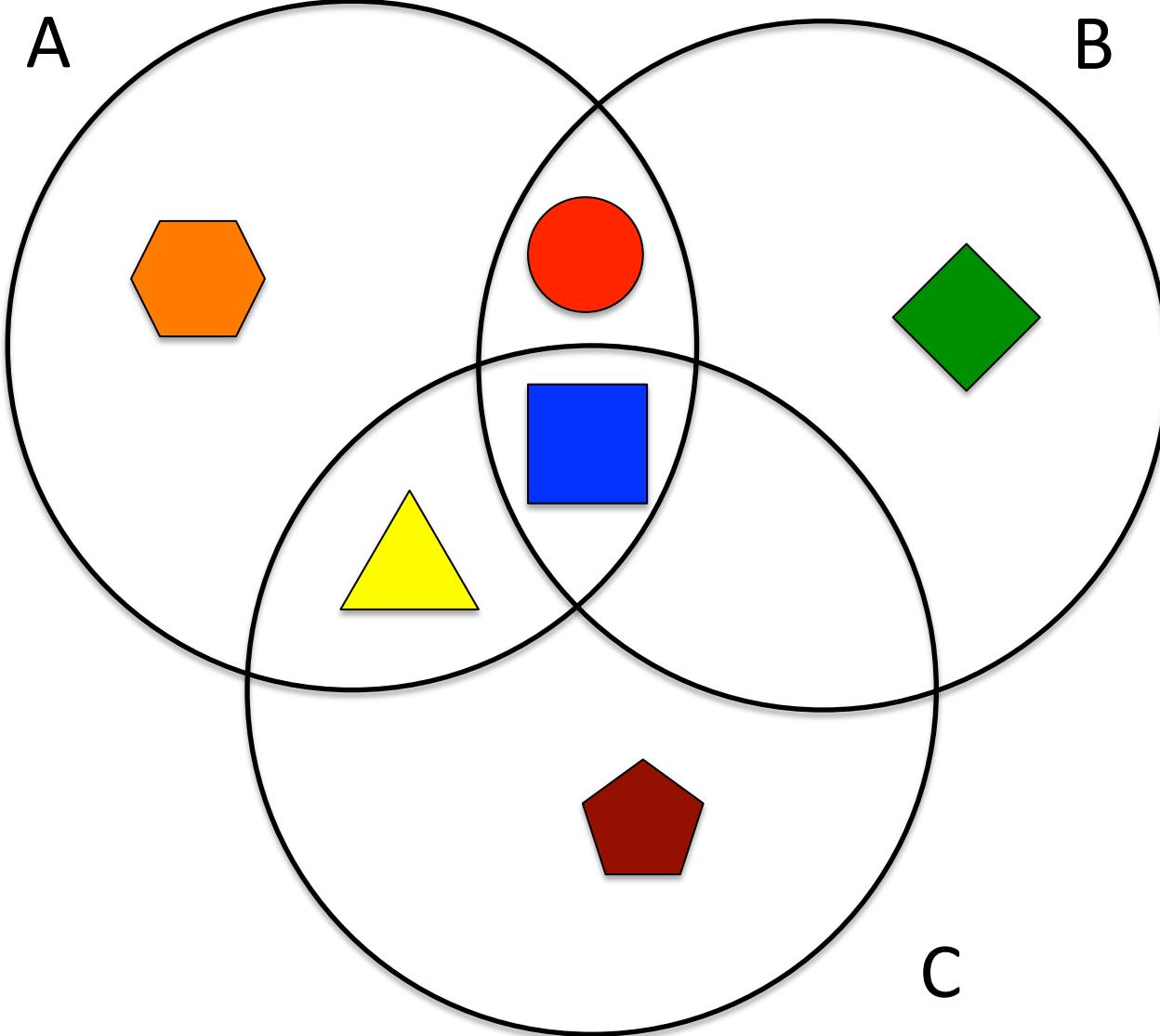
Conclusion:

Sample C is more diverse than sample B,
which is more diverse than sample A.

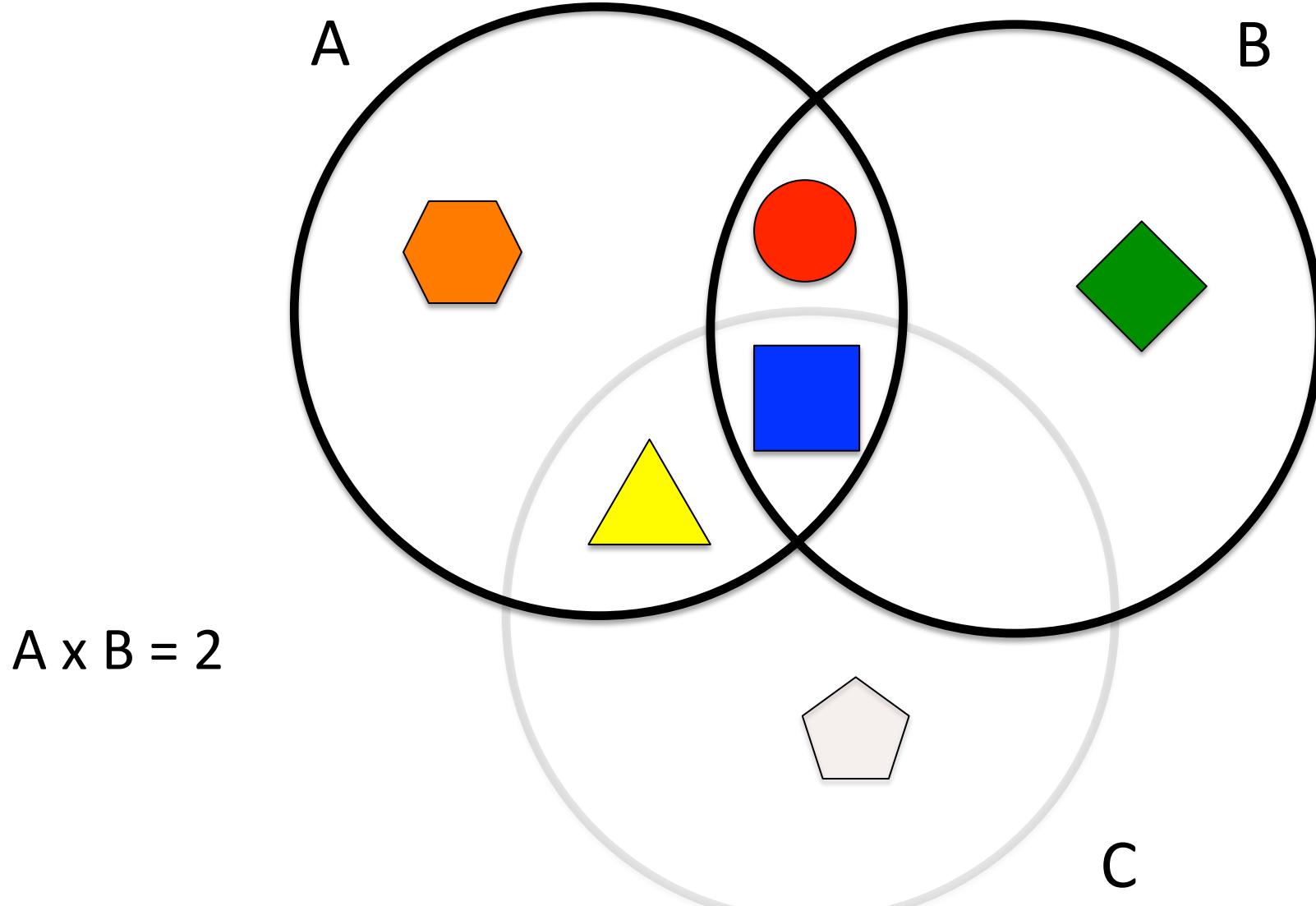
Diversity and statistical analysis

Part II

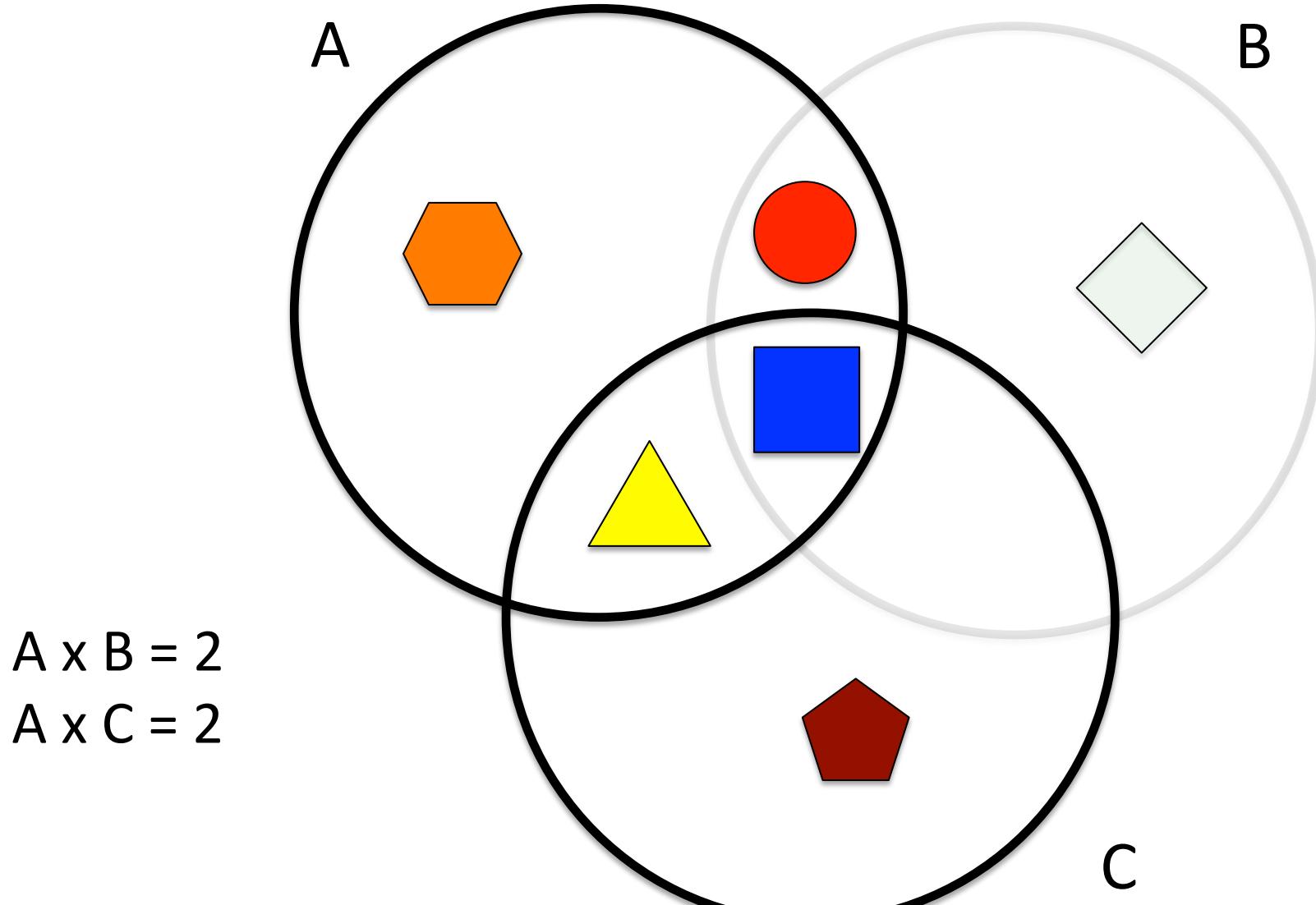
Beta diversity



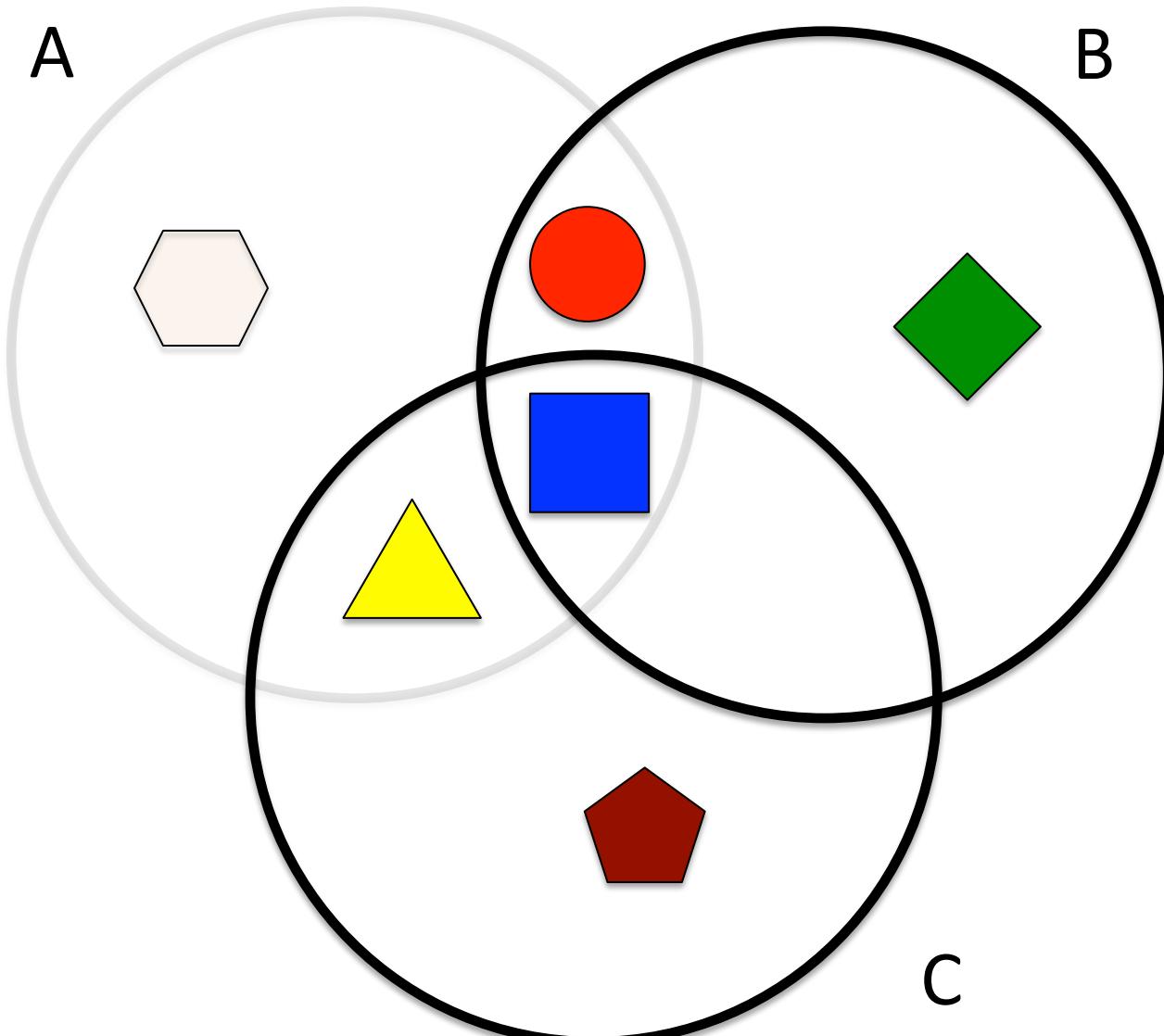
Beta diversity



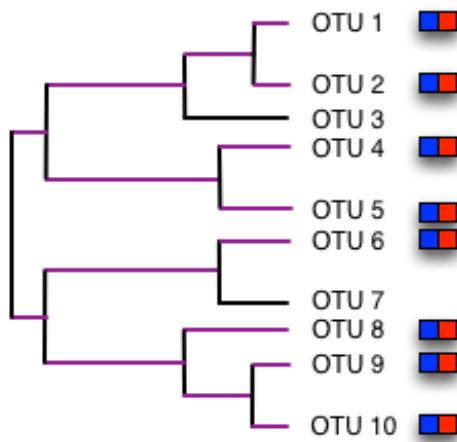
Beta diversity



Beta diversity

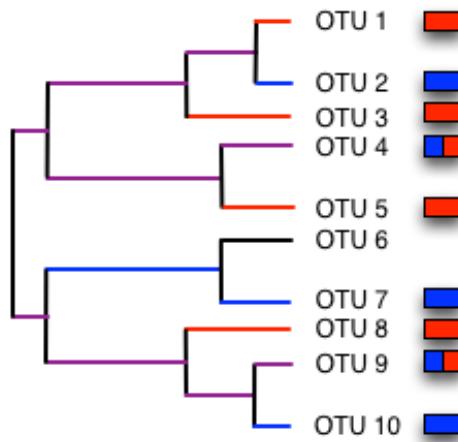


Unweighted UniFrac: a phylogenetic measure of the dissimilarity of microbial communities



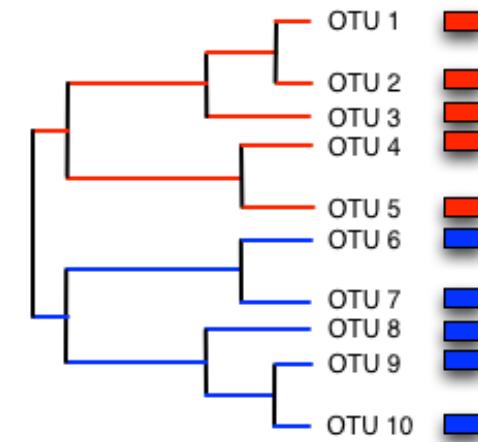
$$U = 0.0$$

Identical communities



$$U \approx 0.5$$

Related communities

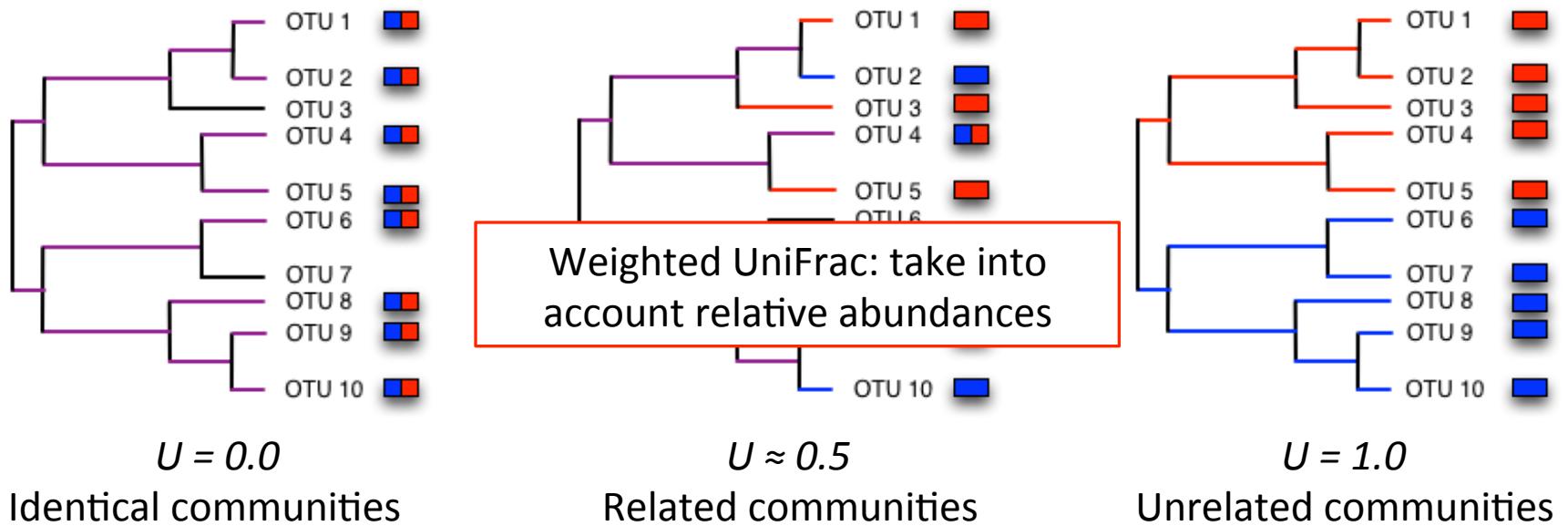


$$U = 1.0$$

Unrelated communities

Percent of observed branch length that is unique to either sample

Unweighted UniFrac: a phylogenetic measure of the dissimilarity of microbial communities



Percent of observed branch length that is unique to either sample

Beta diversity comparison:
visually with ordination plots (e.g., PCoA, NMDS)
statistically (e.g., PERMANOVA, ANOSIM*)

Unweighted UniFrac distance matrix:

	A	B	C	D	E	F
A	0.00	0.35	0.83	0.83	0.90	0.90
B	0.35	0.00	0.86	0.85	0.92	0.91
C	0.83	0.86	0.00	0.25	0.88	0.87
D	0.83	0.85	0.25	0.00	0.88	0.88
E	0.90	0.92	0.88	0.88	0.00	0.50
F	0.90	0.91	0.87	0.88	0.50	0.00

Sample ID	Sample Type
A	Plant (yellow)
B	Plant (yellow)
C	Turtle (red)
D	Turtle (red)
E	Human (green)
F	Dog (blue)

* In QIIME 1.9.1, use `compare_categories.py`

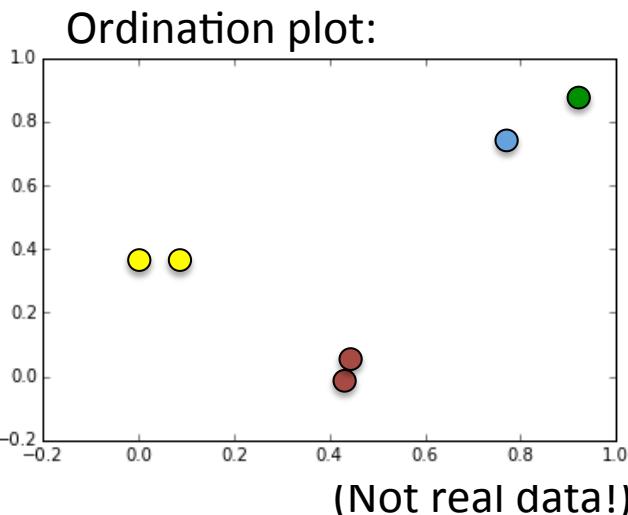
(Not real data!)

Beta diversity comparison:
visually with ordination plots (e.g., PCoA, NMDS)
statistically (e.g., PERMANOVA, ANOSIM*)

Unweighted UniFrac distance matrix:

	A	B	C	D	E	F
A	0.00	0.35	0.83	0.83	0.90	0.90
B	0.35	0.00	0.86	0.85	0.92	0.91
C	0.83	0.86	0.00	0.25	0.88	0.87
D	0.83	0.85	0.25	0.00	0.88	0.88
E	0.90	0.92	0.88	0.88	0.00	0.50
F	0.90	0.91	0.87	0.88	0.50	0.00

Sample ID	Sample Type
A	Plant (yellow)
B	Plant (yellow)
C	Turtle (red)
D	Turtle (red)
E	Human (green)
F	Dog (blue)

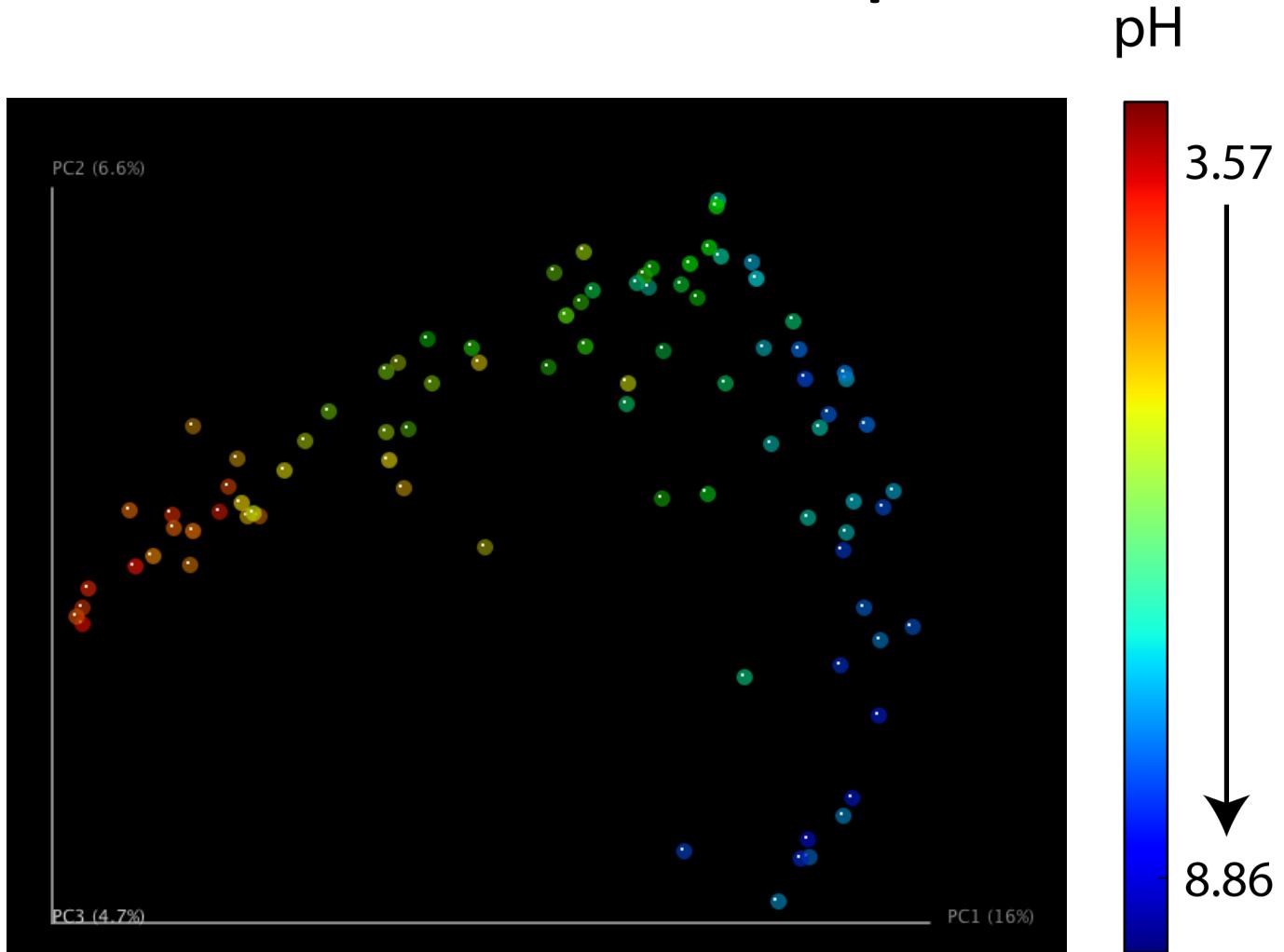


* In QIIME 1.9.1, use `compare_categories.py`

So we've got a distance matrix: now what?

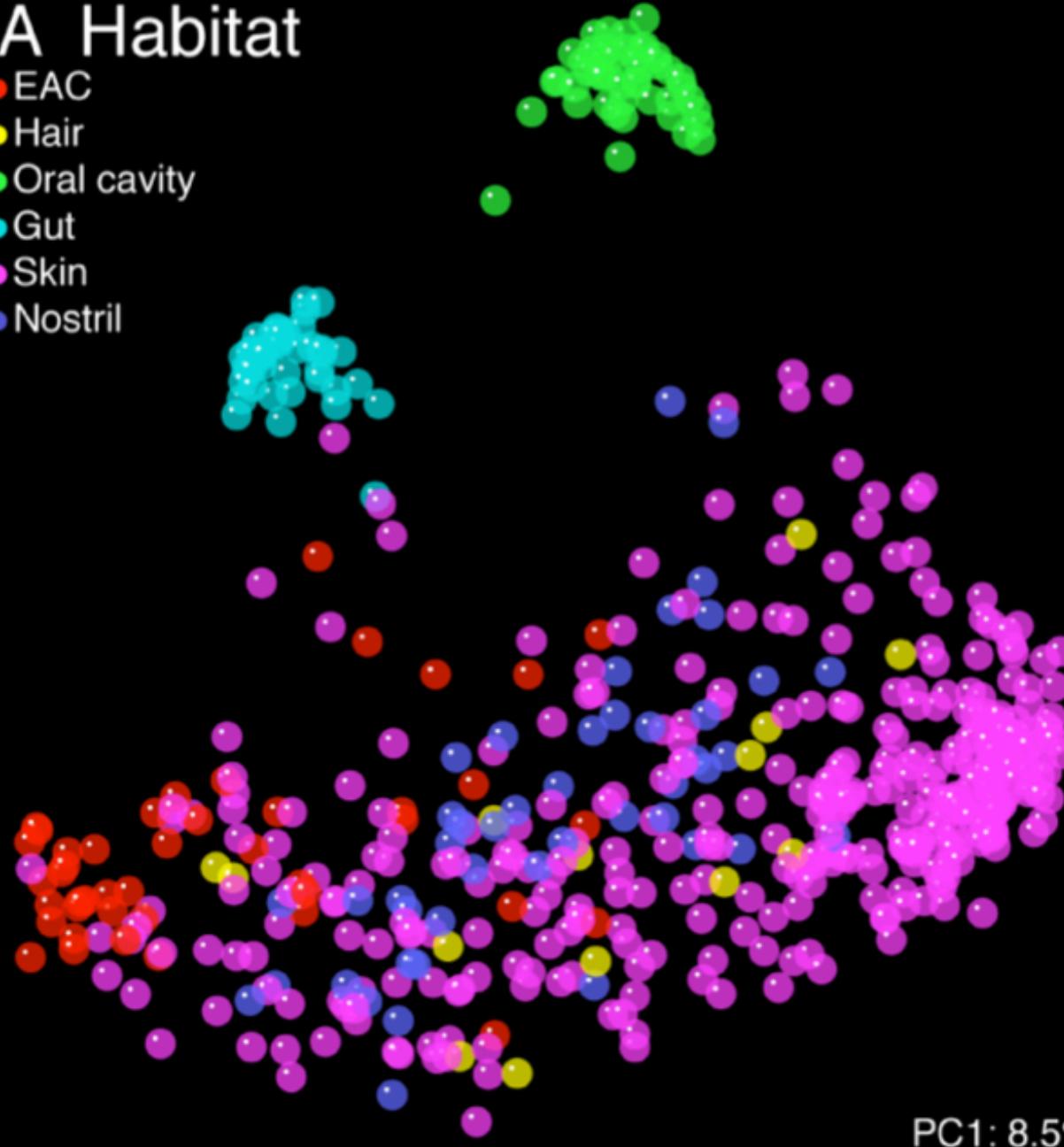
	unweighted_unifrac_dm.txt																		
1	> Stillton4R2	Stillton4R3	Stillton4R1	HCanyon3R3	HCanyon3R2	HCanyon3R1	Halls9R2	HCanyon2R2	HCanyon2R3	HCanyon2R1	Halls9R1	Stillton10R3	HCanyon0R1	HCanyon0R2	HCanyon0R3	HCanyon0R4	HCanyon0R5	HCanyon0R6	
2	Stillton4R2	0.0	0.382273294624	0.391416675288	0.560309484808	0.553938232028	0.566136031815	0.557134987546	0.531719852875	0.53655901824	0.567041909667	0.574831935502	0.474244845718						
3	Stillton4R3	0.382273294624	0.0	0.394399899497	0.586303332083	0.583969320358	0.589363298118	0.588066262733	0.560923487354	0.57767776062	0.601572335069	0.609758664376	0.479910046195						
4	Stillton4R1	0.391416675288	0.394399899497	0.0	0.580628929798	0.583767102601	0.584463513906	0.573828216898	0.550892933103	0.566056393448	0.584611276938	0.607144746402	0.469209424282						
5	HCanyon3R3	0.560309484808	0.586303332083	0.580628929798	0.0	0.354600316745	0.345731612479	0.458450145834	0.412488628393	0.385127601086	0.384525100123	0.46823652766	0.585402372425						
6	HCanyon3R2	0.553938232028	0.583969320358	0.583767102601	0.354600316745	0.0	0.36256949665	0.462707750398	0.414528760004	0.385453766442	0.380087766632	0.46109426257	0.574488221279						
7	HCanyon3R1	0.566136031815	0.589363298118	0.584463513906	0.345731612479	0.36256949665	0.0	0.452351806146	0.415483559719	0.389786424875	0.381981707704	0.485596583752	0.580926122772						
8	Halls9R2	0.557134987546	0.588066262733	0.573828216898	0.458450145834	0.462707750398	0.452351806146	0.0	0.447883445295	0.429943464459	0.409064513124	0.344264504725	0.561110815583						
9	HCanyon2R2	0.531719852875	0.560923487354	0.550892933103	0.412488628393	0.414528760004	0.415483559719	0.447883445295	0.0	0.404179520995	0.388727659604	0.468345488157	0.549328940024						
10	HCanyon2R3	0.53655901824	0.57767776062	0.566056393448	0.385127601086	0.385453766442	0.398796424875	0.429943464459	0.404179520995	0.0	0.361444528983	0.439367245599	0.566330894525						
11	HCanyon2R1	0.567041909667	0.601572335069	0.584611276938	0.384525100123	0.380087766632	0.381981707704	0.409064513124	0.388727659604	0.361444528983	0.0	0.454633280595	0.582635621853						
12	Halls9R1	0.574831935502	0.609758664376	0.607144746402	0.46823652766	0.46109426257	0.485596583752	0.344264504725	0.468345488157	0.43962745599	0.454633280595	0.0	0.578254150137						
13	Stillton10R3	0.474244845718	0.479910046195	0.469209424282	0.5045802372425	0.574488221279	0.580926122772	0.561110815583	0.549328940024	0.566330894525	0.582635621853	0.578254150137							
14	HCanyon0R1	0.776189708931	0.779315452075	0.781074100354	0.69071344042	0.692644612288	0.683086513009	0.727520276229	0.709578414217	0.711736423222	0.707912554104	0.745189206089	0.7						
15	HCanyon0R2	0.797241309582	0.795032626325	0.797809039119	0.71380444578	0.714717508733	0.697490417772	0.75335887058	0.723508997775	0.728837315629	0.722257776622	0.761081079305	0.7						
16	HCanyon0R3	0.784124036272	0.791711071574	0.795527048923	0.693476495058	0.699412891734	0.688360016409	0.732489014053	0.713925200799	0.708354096579	0.708462144382	0.747778043238	0.7						
17	HCanyon7R3	0.539194289149	0.563745794923	0.557216919336	0.457709184042	0.467073058252	0.473888308307	0.42997132225	0.455576236399	0.448029116277	0.441566044862	0.444788573194	0.5						
18	HCanyon7R2	0.665877547192	0.683833494738	0.674471865383	0.553082921975	0.563316168524	0.562611936812	0.546938312206	0.558741865535	0.562533096516	0.534453165883	0.558141564634	0.6						
19	HCanyon7R1	0.554857668962	0.578652540309	0.570402483951	0.475303182122	0.474309745536	0.462587318718	0.43076590681	0.463222699451	0.452021582021	0.42664191335	0.445527585683	0.5						
20	HCanyon1R1	0.746066596617	0.754645795641	0.759778994683	0.654674045733	0.651121732227	0.640930093945	0.703134531793	0.67682938985	0.675957137688	0.676955919301	0.714539249748	0.7						
21	HCanyon1R3	0.751522854919	0.764074726397	0.767226018039	0.641176135477	0.65043468158	0.629924888988	0.702785457767	0.675191090098	0.668258446703	0.66343689429	0.713256418398	0.7						
22	HCanyon1R2	0.748585453295	0.753698061339	0.759152313976	0.666322115311	0.666043304584	0.654075859985	0.711448710816	0.680563093126	0.679964026606	0.678341781539	0.728364004942	0.7						
23	HCanyon10R2	0.56835786215	0.584149104488	0.584183423897	0.466055758676	0.464168854592	0.470680372129	0.458866733827	0.4814190318	0.460896823727	0.442543616102	0.4619876362	0.5						
24	HCanyon10R3	0.570798042641	0.588843077275	0.591942476159	0.50320590177	0.498220269939	0.50309452212318	0.460532212318	0.496962587514	0.490816424808	0.497919239387	0.463277414447	0.5						
25	HCanyon10R1	0.505926071682	0.528347698837	0.529350507225	0.472413250039	0.479553556468	0.478832889376	0.455918418773	0.458027517992	0.428733419477	0.456719821741	0.457712834343	0.5						
26	HCanyon11R3	0.582831143875	0.594413158276	0.60681715155	0.463621293703	0.45842280462	0.486656215192	0.430493752745	0.47439497384	0.445815433526	0.43201382772	0.437061115175	0.5						
27	HCanyon11R2	0.528552555453	0.55774407831	0.555171886694	0.432763917311	0.429877059678	0.440567014981	0.41332040453	0.426696806769	0.398447088329	0.396059530935	0.433778641456	0.5						
28	HCanyon11R1	0.555205897098	0.579991048797	0.575126874355	0.453486254954	0.44588780311	0.456712334902	0.41923496891	0.446869791255	0.430325756526	0.428966840449	0.431167398777	0.5						
29	HCanyon6R2	0.568527068157	0.585132465080	0.583868790939	0.430174126838	0.419369519662	0.432708499988	0.416828113391	0.445815626173	0.426928785146	0.410839210322	0.422532640701	0.5						
30	HCanyon6R3	0.5663346578975	0.587780110737	0.586940661236	0.438634245182	0.425883025257	0.449519180577	0.424862337935	0.447100684715	0.425434475646	0.426218627557	0.424693270252	0.5						
31	HCanyon6R1	0.583972679946	0.586778207344	0.581534515565	0.434912987895	0.424307243199	0.441765941119	0.415236907762	0.428454635151	0.428113957425	0.406954252133	0.422536839493	0.5						
32	HCanyon5R1	0.558910457622	0.57470972671	0.562821643814	0.440681773256	0.451596106791	0.444216261284	0.424754580111	0.4357631361753	0.417737787339	0.416631599644	0.455877139948	0.5						
33	HCanyon5R3	0.559987492698	0.5820196982	0.584250873735	0.450336948189	0.426574204559	0.448524323812	0.439745925494	0.440643510478	0.417613041188	0.413563495656	0.464820696561	0.5						
34	HCanyon5R2	0.577951660858	0.593568671618	0.596389572645	0.422545765979	0.409515246412	0.434403209926	0.429247944328	0.453870629905	0.431590973836	0.4366014112	0.436175614058	0.5						
35	HCanyon4R1	0.516189922282	0.555158817988	0.544402740823	0.416449154884	0.420465768689	0.437380878857	0.460153016099	0.459537309642	0.42519263132	0.4397711880651	0.474121489459	0.5						
36	HCanyon4R2	0.601827528628	0.632734373589	0.628005145571	0.414666304996	0.417716489314	0.443390188238	0.497188855337	0.473766452031	0.447416361784	0.454050692681	0.502442274303	0.6						
37	HCanyon4R3	0.624113459554	0.653406187495	0.649892943176	0.438848993112	0.452700670776	0.456666604112	0.506388612937	0.493725468949	0.479594328128	0.462327913958	0.510888415471	0.6						
38	HCanyon12R1	0.513376789194	0.547194076841	0.519066521533	0.471737648119	0.464638931219	0.461841527684	0.430362902794	0.439277983407	0.4191092078971	0.41992903626	0.466361101825	0.5						
39	HCanyon12R2	0.552892298156	0.560516044793	0.54925756132	0.504271040435	0.504749255647	0.517681811493	0.453014756694	0.480827841616	0.489487879426	0.485996535935	0.465436661211	0.5						
40	HCanyon12R3	0.553185666616	0.585465060164	0.58565605024134	0.439712101711	0.431336959109	0.463612059381	0.409326408365	0.435088785460	0.414920465304	0.407480817122	0.422739541765	0.5						
41	HCanyon8R1	0.5552014316498	0.559504417239	0.562836483618	0.483602818412	0.470838841272	0.485053828711	0.446654297414	0.465192835869	0.464892486774	0.44173317446	0.466837938413	0.5						
42	HCanyon8R2	0.547616166832	0.576958110793	0.578527645361	0.458907582912	0.452778768072	0.456800424019	0.420176245854	0.449481849055	0.430019204227	0.413022684081	0.432050808518	0.5						
43	HCanyon8R3	0.549015567252	0.573444281355	0.580500957664	0.446205572204	0.446996909932	0.453171155433	0.43275915163	0.439870648001	0.425623148959	0.408822278598	0.446705079363	0.5						
44	Halls8R2	0.535370779107	0.569114464296	0.574510728774	0.465288199139	0.452053418068	0.460087502347	0.379587194176	0.458489263366	0.421181481291	0.440297005744	0.366881097176	0.5						
45	Stillton10R1	0.46229629442	0.4684474471913	0.470923059335	0.575228561838	0.580538958921	0.5846812234	0.564450575535	0.557710461145	0.559334675708	0.5881723633	0.584771138785	0.5						
46	Stillton10R2	0.5387982217																	

Ordination Techniques



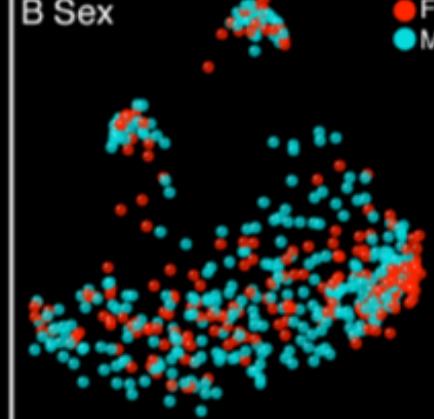
A Habitat

- EAC
- Hair
- Oral cavity
- Gut
- Skin
- Nostril

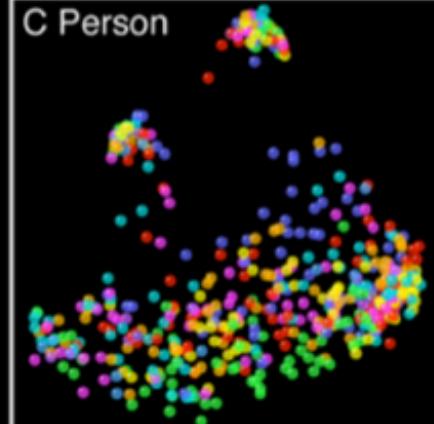


B Sex

- F
- M

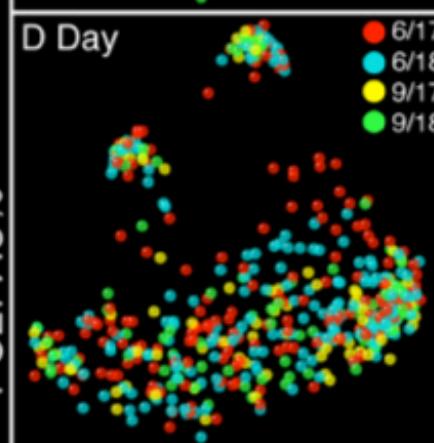


C Person



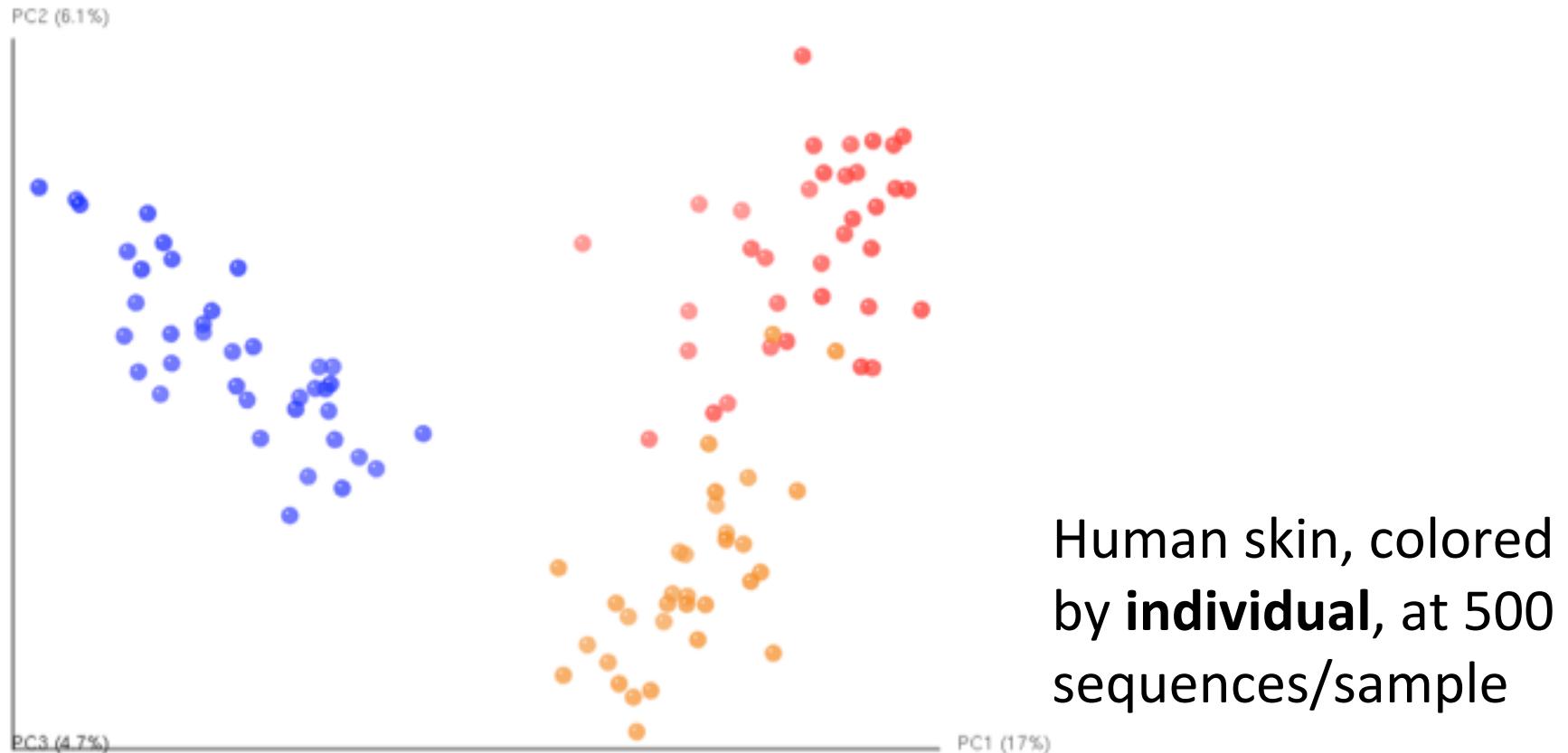
D Day

- 6/17
- 6/18
- 9/17
- 9/18



Variation in sampling depth also needs to be controlled for beta diversity!

Variation in sampling depth is an important consideration



Image/analysis credit: Justin Kuczynski

Data reference:

Forensic identification using skin bacterial communities. Fierer N, Lauber CL, Zhou N, McDonald D, Costello EK, Knight R. Proc Natl Acad Sci U S A. 2010 Apr 6;107(14):6477-81.

Variation in sampling depth is an important consideration



Image/analysis credit: Justin Kuczynski

Data reference:

Forensic identification using skin bacterial communities. Fierer N, Lauber CL, Zhou N, McDonald D, Costello EK, Knight R. Proc Natl Acad Sci U S A. 2010 Apr 6;107(14):6477-81.

Variation in sampling depth is an important consideration



Human skin, colored by **sampling depth**, at either 50 (blue) or 500 (red) sequences/sample

Image/analysis credit: Justin Kuczynski

Data reference:

Forensic identification using skin bacterial communities. Fierer N, Lauber CL, Zhou N, McDonald D, Costello EK, Knight R. Proc Natl Acad Sci U S A. 2010 Apr 6;107(14):6477-81.

How deep is deep enough?

It depends on the question...

- Differences between community types: not many sequences.
- Rare biosphere: more (but be careful about sequencing noise!)

How deep is deep enough?

100 sequences/sample

10 sequences/sample

1 sequence/sample

PC2 (8.4%)



PC2 (11%)



PC1 (13%)

PC3 (8.1%)

PC1 (8.6%)

PC3 (6.2%)

PC2 (17%)



PC1 (2.4%)

PC3 (9.7%)

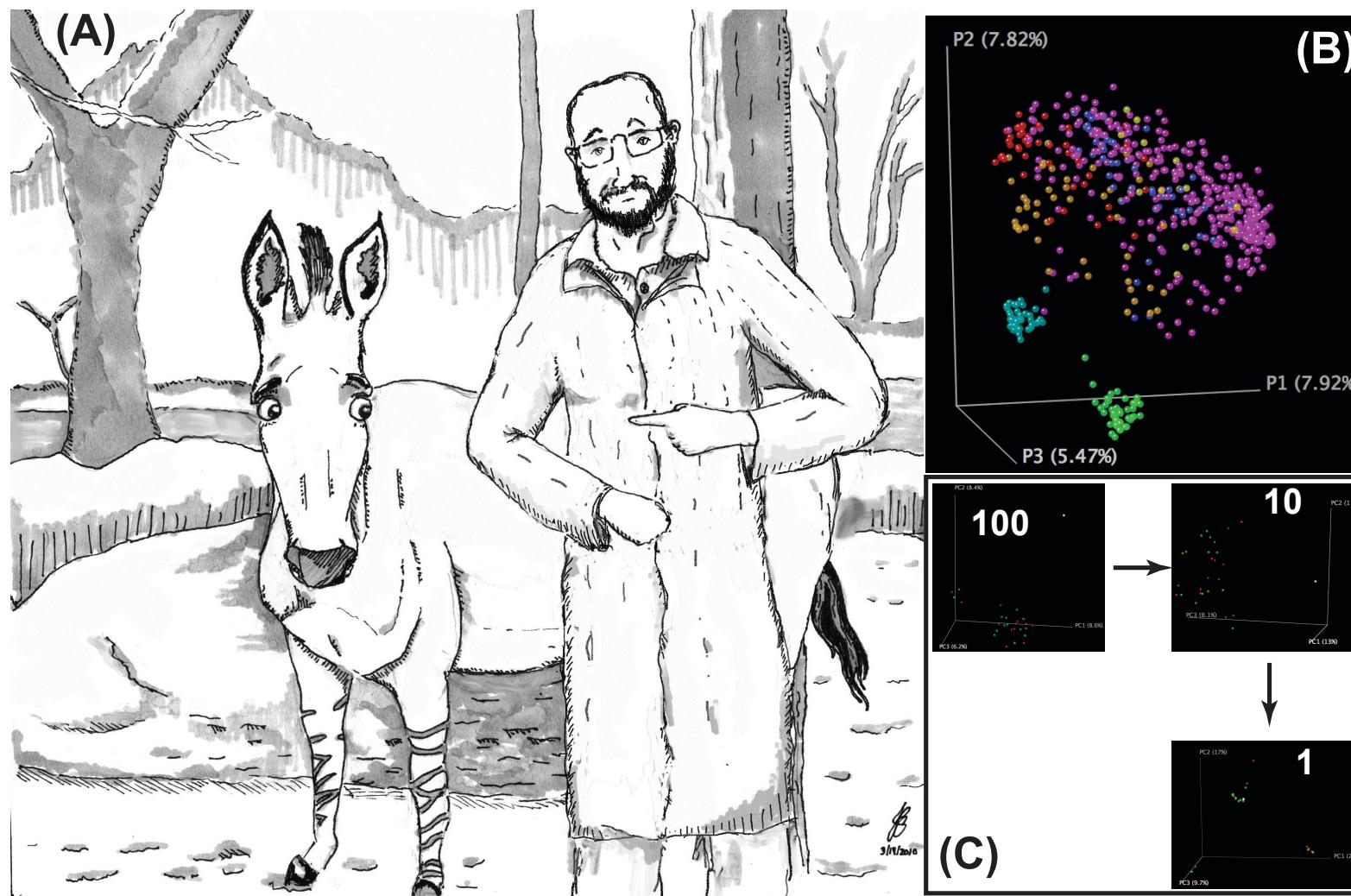
Direct sequencing of the human microbiome readily reveals community differences.

J Kuczynski et al. Genome Biology (2011).

Direct sequencing of the human microbiome readily reveals community differences.

Kuczynski J, Costello EK, Nemergut DR, Zaneveld J, Lauber CL, Knights D, Koren O, Fierer N, Kelley ST, Ley RE, Gordon JI, Knight R.

Figure 1



Differentially abundant OTUs

Which OTUs have statistically different abundances between groups of samples?

Example: which OTUs are differentially abundant across human body sites?

Differentially abundant OTUs

group_significance.py

- Kruskal-Wallis (non-parametric, QIIME default)
 - Groups must contain > 4 samples
 - Commonly used for marker gene surveys
- ANOVA (parametric, previous QIIME default)
 - Equal variance
 - Normality of residuals
 - Independence
 - Usually violated by marker gene surveys
- *Many others*

Differentially abundant OTUs

Currently an active area of research, recent evidence suggests other approaches are better:

Waste not, want not: why rarefying microbiome data is inadmissible. McMurdie and Holmes, 2014.

<http://www.ncbi.nlm.nih.gov/pubmed/24699258>

(Available in QIIME: differential_abundance.py)

Analysis of composition of microbiomes... Mandal et al., 2015.

<http://www.ncbi.nlm.nih.gov/pubmed/26028277>

(Not yet available in QIIME, but coming soon.)

core_diversity_analyses.py



Quantitative Insights Into Microbial Ecology

News and Announcements » • QIIME 1.7.0 is live! • QIIME 1.6.0 is live! • UNITE/QIIME 12_11 ITS reference OTUs now available (alpha release!)

Home »

index

Site index

- Home
- Install
- Documentation
- Tutorials
- Blog
- Developer

Quick search

Go

Enter search terms or a module, class or function name.

core_diversity_analyses.py - A workflow for running a core set of QIIME diversity analyses.

Description:

This script plugs several QIIME diversity analyses together to form a basic workflow beginning with a BIOM table, mapping file, and optional phylogenetic tree.

The included scripts are those run by the workflow scripts `alpha_rarefaction.py`, `beta_diversity_through_plots.py`, `summarize_taxa_through_plots.py`, plus the (non-workflow) scripts `make_distance_boxplots.py`, `compare_alpha_diversity.py`, and `otu_category_significance.py`. To update parameters to the workflow scripts, you should pass the same parameters file that you would pass if calling the workflow script directly.

Usage: `core_diversity_analyses.py [options]`

Input Arguments:

```
[REQUIRED]
-i, --input_biom_fp
    The input biom file [REQUIRED]

-o, --output_dir
    The output directory [REQUIRED]

-m, --mapping_fp
    The mapping filepath [REQUIRED]

-e, --sampling_depth
    Sequencing depth to use for even sub-sampling and maximum rarefaction depth. You should review the output of print_biom_table_summary.py to decide on this value.

[OPTIONAL]
-p, --parameter_fp
    Path to the parameter file, which specifies changes to the default behavior. For more information, see www.qiime.org/documentation/qiime\_parameters\_files.html [if omitted, default values will be used]

-a, --parallel
    Run in parallel where available. Specify number of jobs to start with -O or in the parameters file. [default: False]

--nonphylogenetic_diversity
    Apply non-phylogenetic alpha (chaol and observed_species) and beta (bray_curtis) diversity calculations. This is useful if, for example, you are working with non-amplicon BIOM tables, or if a reliable tree is not available (e.g., if you're working with ITS amplicons) [default: False]

--suppress_taxa_summary
    Suppress generation of taxa summary plots. [default: False]
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
                      -o core_output  
                      -m map.txt  
                      -c SampleType  
                      -t rep_set.tre  
                      -e 20
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
                           -o core_output  
                           -m map.txt  
                           -c SampleType  
                           -t rep_set.tre  
                           -e 20
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
          -o core_output  
          -m map.txt  
          -c SampleType  
          -t rep_set.tre  
          -e 20
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
    -o core_output  
    -m map.txt  
    -c SampleType  
    -t rep_set.tre  
    -e 20
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
                      -o core_output  
                      -m map.txt  
                      -c SampleType  
                      -t rep_set.tre  
                      -e 20
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
                      -o core_output  
                      -m map.txt  
                      -c SampleType  
                      -t rep_set.tre  
                      -e 20
```

core_diversity_analyses.py

```
$ core_diversity_analyses.py -i otu_table.biom  
                      -o core_output  
                      -m map.txt  
                      -c SampleType  
                      -t rep_set.tre  
                      -e 20
```

Hands-on with core_diversity_analyses.py output



Genome **Biology**

Research

Highly accessed

Open Access

Moving pictures of the human microbiome

J Gregory Caporaso¹, Christian L Lauber², Elizabeth K Costello³, Donna Berg-Lyons², Antonio Gonzalez⁴, Jesse Stombaugh¹, Dan Knights⁴, Pawel Gajer⁵, Jacques Ravel⁵, Noah Fierer²⁶, Jeffrey I Gordon⁷ and Rob Knight^{18*}

* Corresponding author: Rob Knight rob.knight@colorado.edu

► Author Affiliations

For all author emails, please [log on](#).

Genome Biology 2011, **12**:R50 doi:10.1186/gb-2011-12-5-r50

The electronic version of this article is the complete one and can be found online at:
<http://genomebiology.com/2011/12/5/R50>

Full tutorial available at:

http://qiime.org/tutorials/illumina_overview_tutorial.html

Moving Pictures of the Human Microbiome

- Two human subjects provided daily microbiome samples, one for 6 months and the other for 18 months.
- Sampled tongue, left and right palms, and gut (via feces).
- Tutorial data is a subset of this full data set.
- *Illumina HiSeq 2000* 16S sequencing with Earth Microbiome Project protocols (earthmicrobiome.org)

core_diversity_analyses.py precomputed results demonstration

<http://bit.ly/1IL5WpX>

Advanced topics

Chimera checking

- Usage in QIIME described here:
[http://qiime.org/tutorials/
chimera checking.html](http://qiime.org/tutorials/chimera_checking.html)
- Should not be used if one is using closed-reference OTU picking.
- Demonstration of chimera checking on the QIIME Illumina tutorial data with usearch61.

Supervised Learning

This approach is useful if one needs to categorize samples based upon microbial data. QIIME implementation is a random forest approach.

The script used is `supervised_learning.py`.

Demonstration of supervised learning with QIIME tutorial data.

SourceTracker

QIIME tutorial page: http://qiime.org/tutorials/source_tracking.html

Describes an older command for converting one's OTU table to tab-delimited format. A command that works with the new version of the biom software is:

```
biom convert --to-tsv -i INPUT -o OUTPUT
```

Plates from two fields of my maize study were mixed up for a particular week (Aurora and Lansing fields, week 9). Can see predictions from SourceTracker which indicate the source of the microbes from the correct field.

PICRUSt

- Predicts metagenomes from 16S data
- Useful for exploring metagenomic data without shotgun sequences
 - Save money and time
- PICRUSt workflow
- Website: <http://picrust.github.io>

Microbial Study Design Considerations

You probably already have sequencing data, ready to analyze, but...

PCR primer choice

Taxa

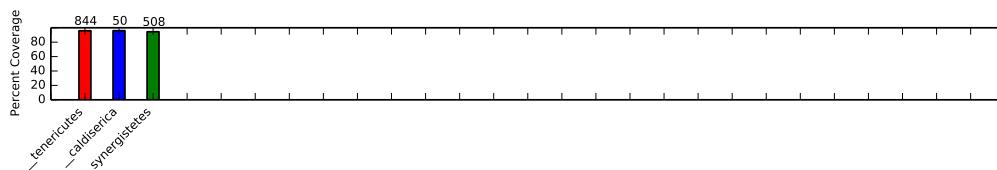
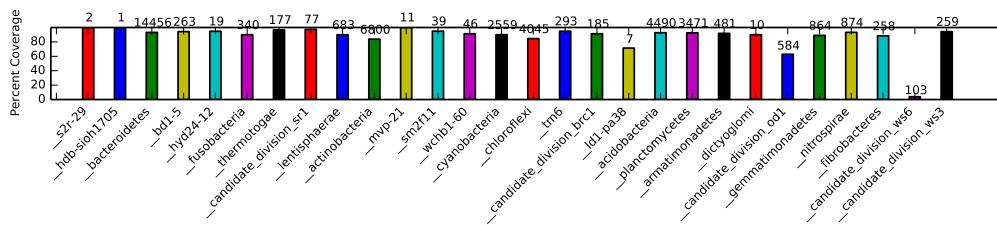
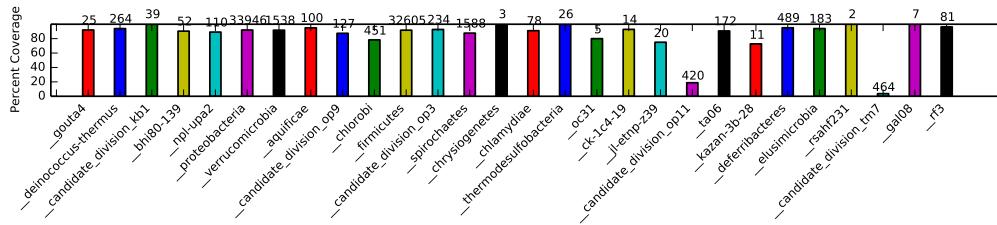
Amplicon Length

Target gene

We can address these with Primer Prospector:
<http://pprospector.sourceforge.net/>

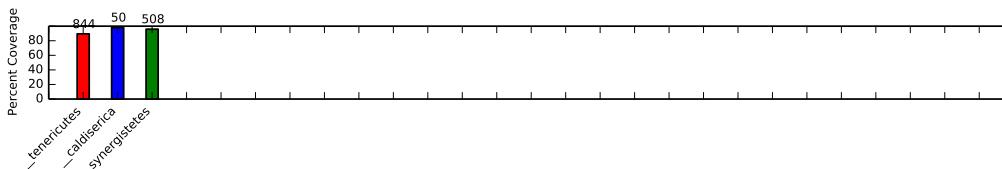
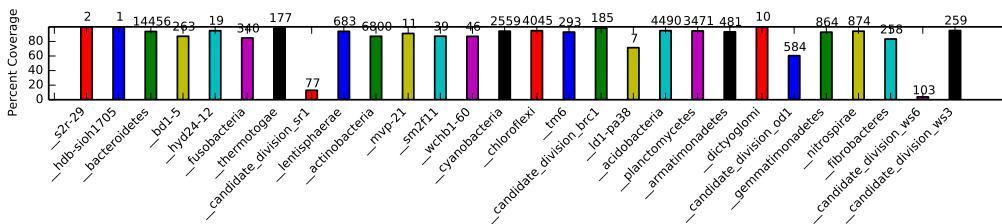
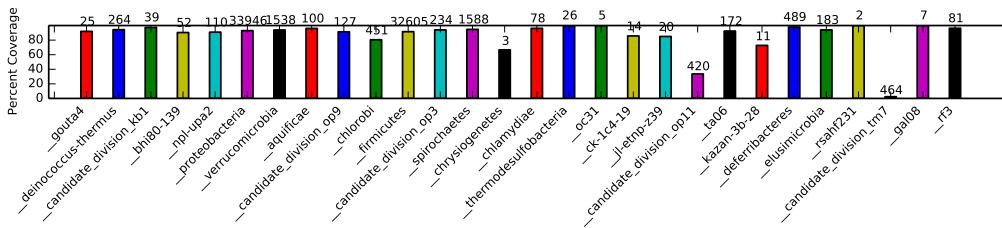
Predicted Taxonomic Coverage
515f_97_Silva_111_rep_set_hits_806r_97_Silva_111_rep_set_hits
Sequences In Category Bacteria
Taxonomy Level 1
 Numeric values above bins represent
 total sequence counts for each set

V4 Bacteria

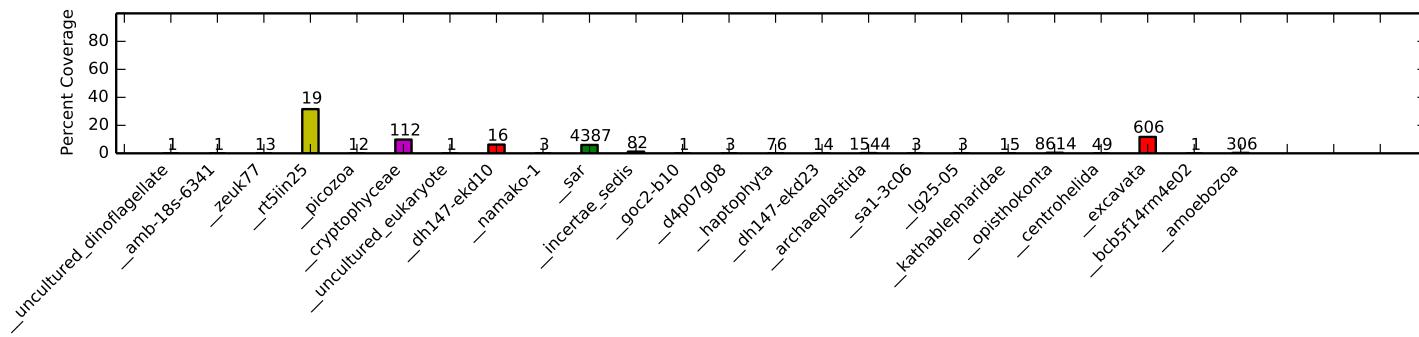


Predicted Taxonomic Coverage
515f_97_Silva_111_rep_set_hits_926r_97_Silva_111_rep_set_hits
Sequences In Category Bacteria
Taxonomy Level 1
 Numeric values above bins represent
 total sequence counts for each set

V4-5 Bacteria

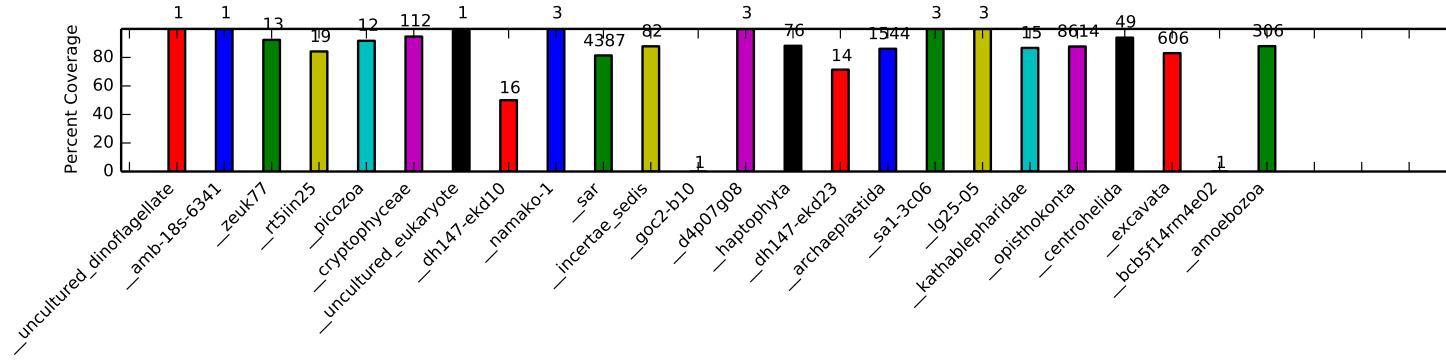


Predicted Taxonomic Coverage
515f_97_Silva_111_rep_set_hits_806r_97_Silva_111_rep_set_hits
Sequences In Category Eukaryota
Taxonomy Level 1
Numeric values above bins represent
total sequence counts for each set



V4 (515f/806r) predicted Eukaryotic amplification

Predicted Taxonomic Coverage
515f_97_Silva_111_rep_set_hits_926r_97_Silva_111_rep_set_hits
Sequences In Category Eukaryota
Taxonomy Level 1
Numeric values above bins represent
total sequence counts for each set



V4-5 (515f/926r) predicted Eukaryotic amplification

Validating metadata

- validate_mapping_file.py demonstration
- Keemei demonstration
 - <http://keemei.qiime.org>

Reproducible bioinformatics

Why is it important for a scientific experiment to be reproducible?

- Differentiate real results from experimental artifacts

Reproducible versus replicable

- Replicable: exact same conditions lead to concordant results
- Reproducible: some experimental variation is allowed, but results are concordant

What does it mean for a bioinformatics experiment to be replicable?

- Our experimental methods are not as ‘noisy’ as most.
- Same commands on the same system should give you the same results (if the algorithm is *deterministic*).

Deterministic versus non-deterministic

- Deterministic algorithm: a given input produces the same series of internal states and results in the same output.
- Non-deterministic algorithm: a given input may produce different internal states and/or result in a different output.
 - Commonly probabilistic algorithms in bioinformatics

Deterministic

- Smith-Waterman alignment: if properly implemented, aligning two sequences will always give the same result

Non-deterministic/Probabilistic

- Sub-sampling a data set (e.g., rarefaction of an OTU table)
- Jackknifed analyses
- Permutation-based p-values (Monte Carlo)

How can we develop software that supports reproducibility?

Open source software

- Consider making your software open source
 - Avoid “black box”
 - Benefits other researchers/community
 - Encourages collaboration
 - Extendibility
- Place it under public hosted revision control
 - GitHub
 - Bitbucket
 - SourceForge



Version control systems

- Git
- Subversion (svn)
- Mercurial (hg)
- CVS (ancient history)
- Allow for viewing history of changes, obtaining previous versions.
- Example: <https://github.com/biocore/qiime>



Virtual Machines

- Publish virtual machine images
 - Gives access to exact software, configuration, and data used in analyses

What's in a good log file?

- Ideally will supplement your lab notebook (for successful runs)
 - Version information
 - Exact commands that were run
 - Details on input files (path, md5)
 - System configuration details
- Publish these as supplementary material

MD5

- A cryptographic hash function: deterministic function which takes some input and returns a fixed-size string – changing the input should change the return value

From [Wikipedia](#):

- it is easy (but not necessarily quick) to compute the hash value for any given message
- it is infeasible to generate a message that has a given hash
- it is infeasible to modify a message without changing the hash
- it is infeasible to find two different messages with the same hash

IPython Notebook

- Interactive, executable documents: code, text, images, etc.
- Easy to share, publish, convert
- Great for keeping track of analysis commands, code, descriptions/comments, etc.
- Make your methods section *executable!*

IP[y]:

<http://ipython.org/>

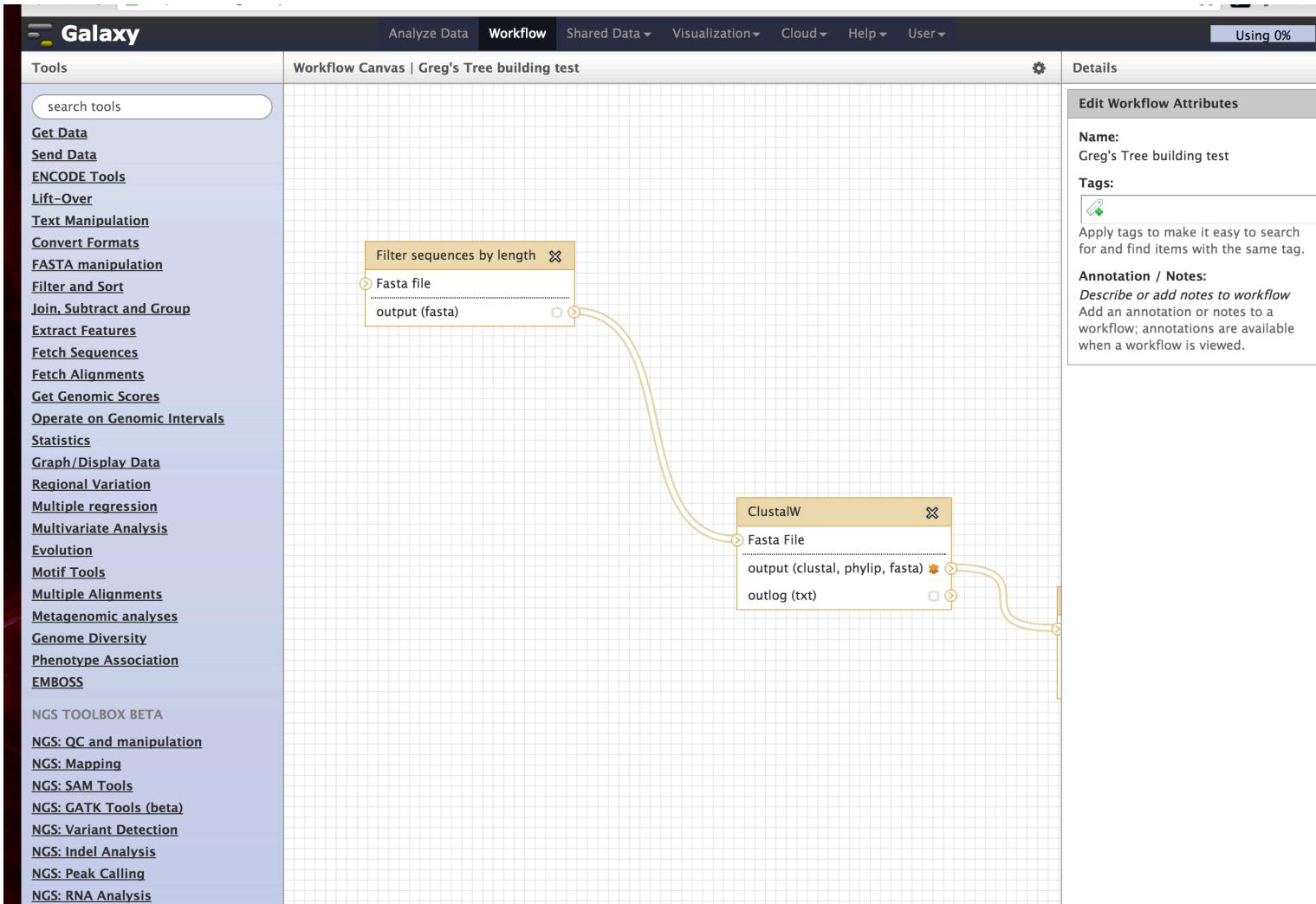
IPython Notebook

- Workshop exercises are written in IPython Notebooks
- An Introduction to Applied Bioinformatics
 - Online interactive bioinformatics book
 - <http://readIAB.org>

IP[y]:

<http://ipython.org/>

Reproducible computing through workflow engines, e.g. Galaxy



Future directions: QiiTA and QIIME 2

QiiTA pre-history

- Previous iteration known as the “QIIME Database” (QIIME-DB)
- Suffered database crash
 - Efforts focused on rewrite instead of recover

QiiTA: QIIME-DB Reboot

- System for:
 - Depositing/archiving microbiome data
 - Performing meta-analysis
 - Combine data from a variety of sources (marker gene, metagenomic, metabolomic, etc.)
- Goals
 - Easy-to-use web interface
 - User-deployable in a variety of environments (e.g., laptops to clusters)
 - Powerful meta-analysis capabilities

QiiTA

- Currently supports:
 - Querying studies
 - Closed-reference OTU picking
 - Submission of data to EBI-SRA (needed for publishing studies)
- Demonstration of querying and downloading study data

Moving toward QIIME 2

- QIIME is currently:
 - Command-line only (very limited Galaxy support)
 - Simplistic execution of workflows
 - Hard to extend and maintain (for both users and devs)
 - Can be difficult to install

Moving toward QIIME 2

- Most requested feature: graphical interface
- Most support efforts: command-line issues

Users spend too much time grappling with the command line and less time performing awesome microbiome research.

Devs spend too much time helping users with installation and command-line issues, and less time answering users' research questions.

QIIME 2 Overview

- Complete redesign/rebuild of QIIME
- Powered by scikit-bio (<http://scikit-bio.org>)
- Graphical web-based interface
- Command-line interface and Python API
- Interactive visualizations
- Deployable on laptops -> clusters
- Extendable by users/devs via plugin system

QIIME 2 Overview

- Currently being developed
 - Initial alpha release scheduled for 2016
- Follow updates at <http://blog.qiime.org> and on Twitter by following @qiime_

Preparation for day 2

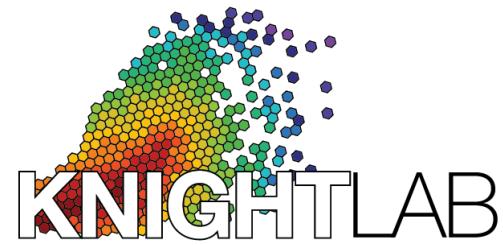
- Download and install Cyberduck
 - <http://cyberduck.io>
- Download and install Google Chrome
 - <https://www.google.com/chrome/>

Note: We'll upload the latest version of today's slides this evening.

Acknowledgements

Evan Bolyen
Nick Bokulich
Katy Califf
Greg Caporaso
Jose Clemente
John Chase
Kevin Cohen
Antonio Gonzalez
Crystal Hepp
Rob Knight
Bruce Hungate
Larry Hunter
Paul Keim
Scott Kelley
Justin Kuczynski

Cathy Lozupone
Daniel McDonald
David Mills
Norm Pace
Fernando Perez
Jai Ram Rideout
Egbert Schwartz
Karen Schwarzberg
Jeffrey Siegel
Jesse Stombaugh
Yoshiki Vazquez
Tony Walters



ALFRED P. SLOAN
FOUNDATION



william.a.walters@gmail.com
jai.rideout@gmail.com



NORTHERN
ARIZONA
UNIVERSITY

Slides compiled by:

Greg Caporaso

John Chase

Jose Clemente

Antonio Gonzalez Peña

Rob Knight

Cathy Lozupone

Daniel McDonald

Jai Ram Rideout

Yoshiki Vázquez Baeza

Tony Walters



This work is licensed under the Creative Commons Attribution 3.0 United States License. To view a copy of this license, visit

<http://creativecommons.org/licenses/by/3.0/us/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

Feel free to use or modify these slides, but please credit us by placing the following attribution information where you feel that it makes sense:

Slides derived from QIIME educational materials www.qiime.org.