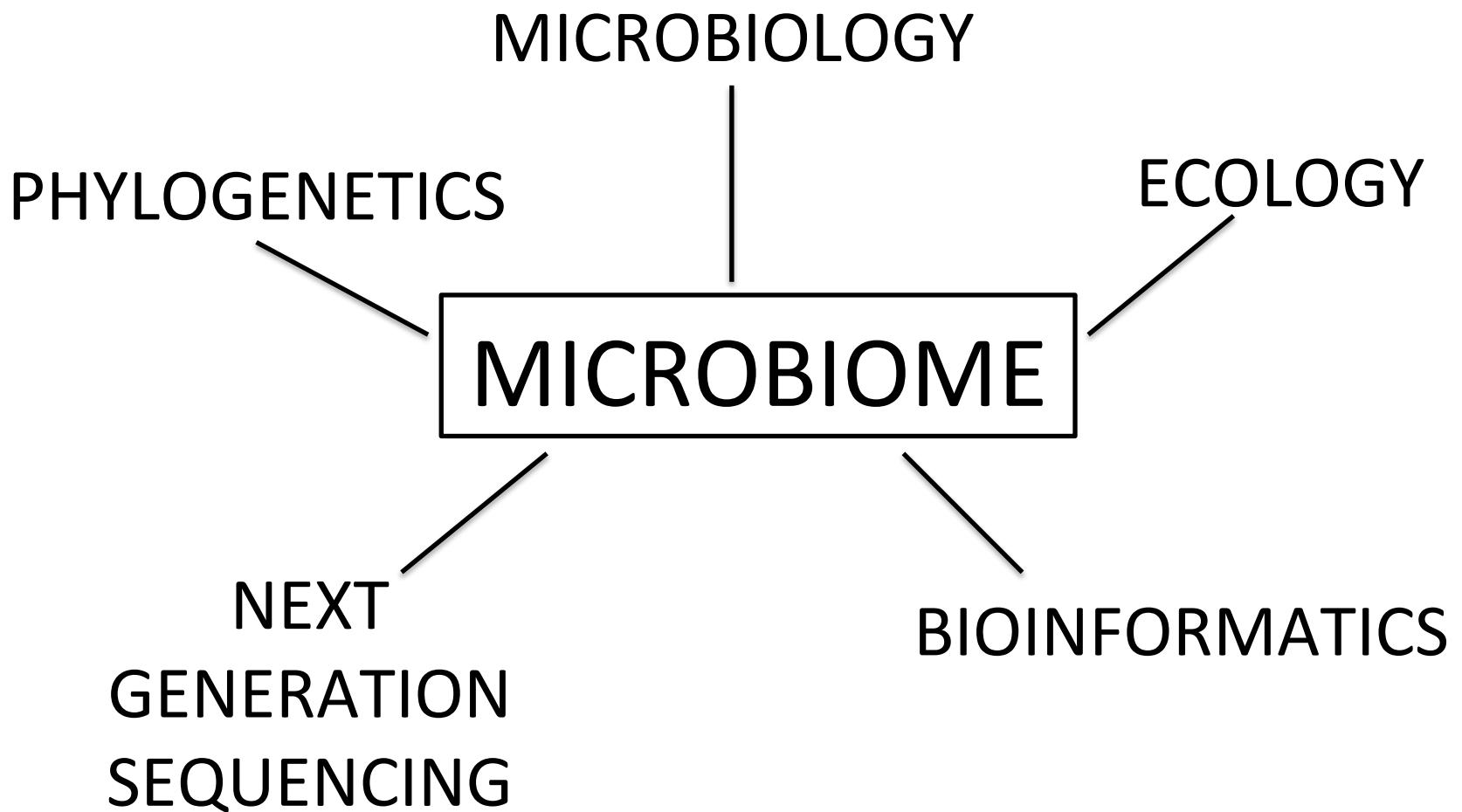


Microbiome 101

Tracy Meiring
Medical Virology
University of Cape Town
Institute of Infectious Diseases and Molecular Medicine





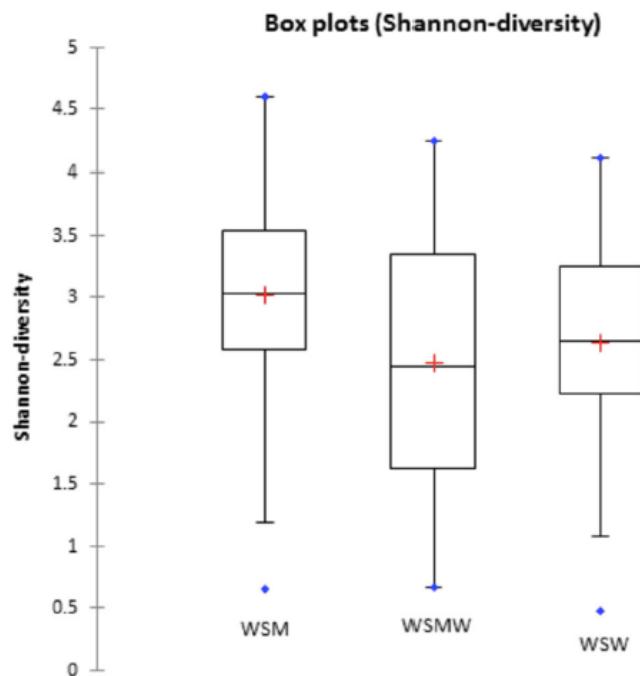
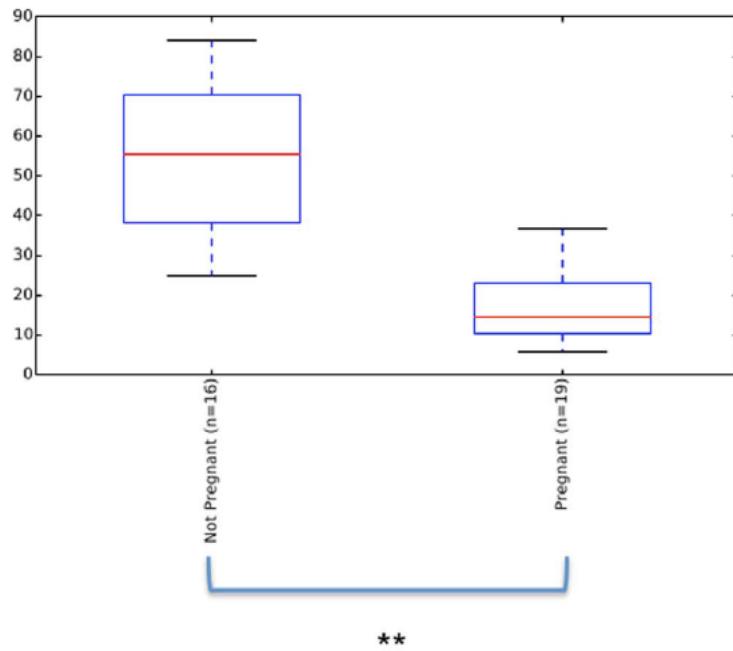


Figure 5. Box plot of Shannon's diversity for each sexual risk behavior group. The red plus signs depict the average value of Shannon's diversity for each sexual risk behavior group; blue dots represent outliers. The WSM group had the highest diversity compared to the other sexual risk behavior groups.

doi: 10.1371/journal.pone.0080254.g005

African-American



Caucasian

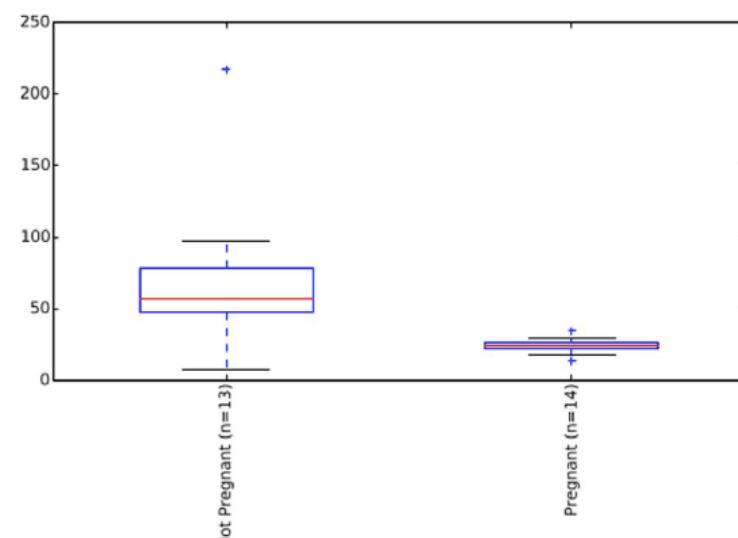


Figure 4. Pregnancy effect in African-American and Caucasian subjects as measured by Chao1 diversity Index. Diversity is significantly reduced during pregnancy in both ethnicities (** $p<0.01$, Monte Carlo analyses, 999 permutations).
doi:10.1371/journal.pone.0098514.g004

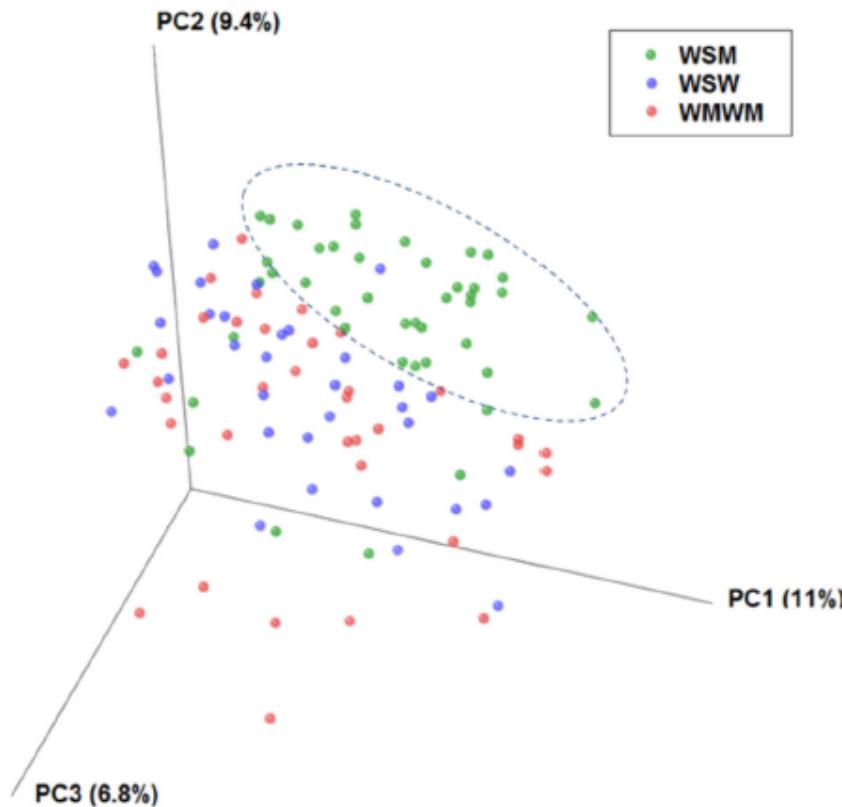
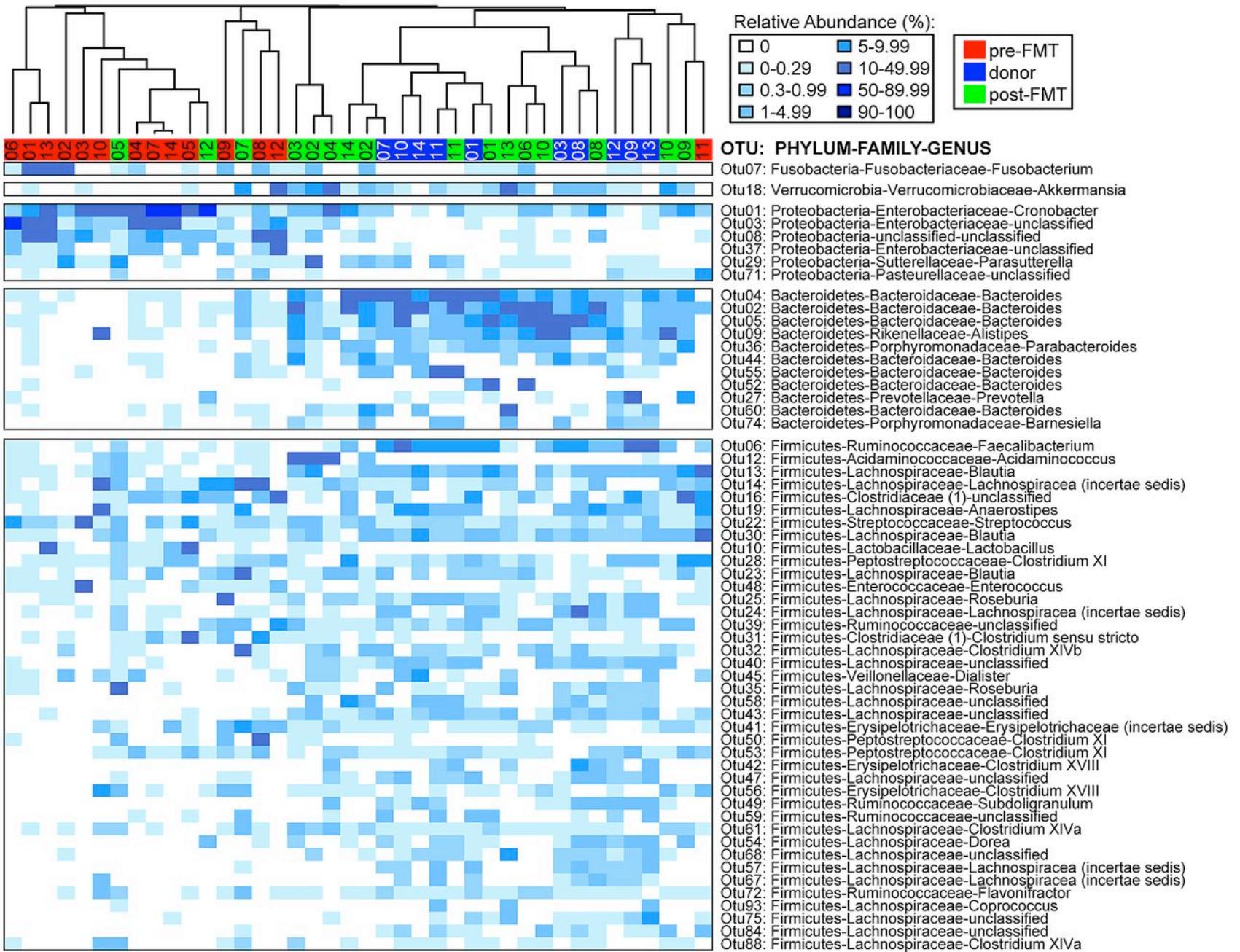


Figure 6. Unweighted beta-diversity Unifrac analysis. This figure depicts a distinct cluster within the WSM group (indicated in green) that is not present within the WSW or WSWM groups (indicated in blue and red, respectively).

doi: 10.1371/journal.pone.0080254.g006

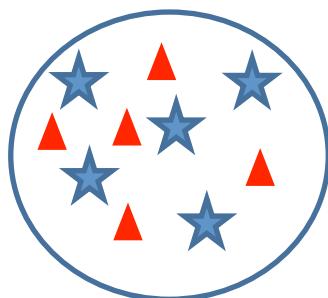


- How do we describe microbiomes?
- How do we compare microbiomes?

Important ecological concepts

- How is biodiversity defined and measured?
- Component of biodiversity:
 - RICHNESS
 - EVENNESS
- Species richness: number of different species in a habitat/sample
- Species relative abundance: number of each species relative to total number of all species in a sample (number of reads per OTU in a sample relative to total number of reads in that sample)
- Species evenness: how close in numbers each species in an environment are; distribution

Simple example:

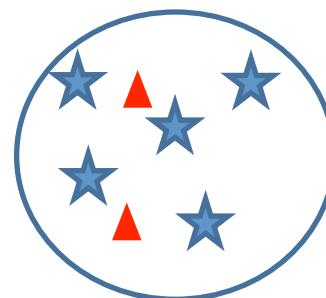


Richness: 2 species

Relative abundance:
 $5/10 = 0.5$ or 50%

★ $5/10 = 0.5$ or 50%

High Evenness



Richness: 2 species

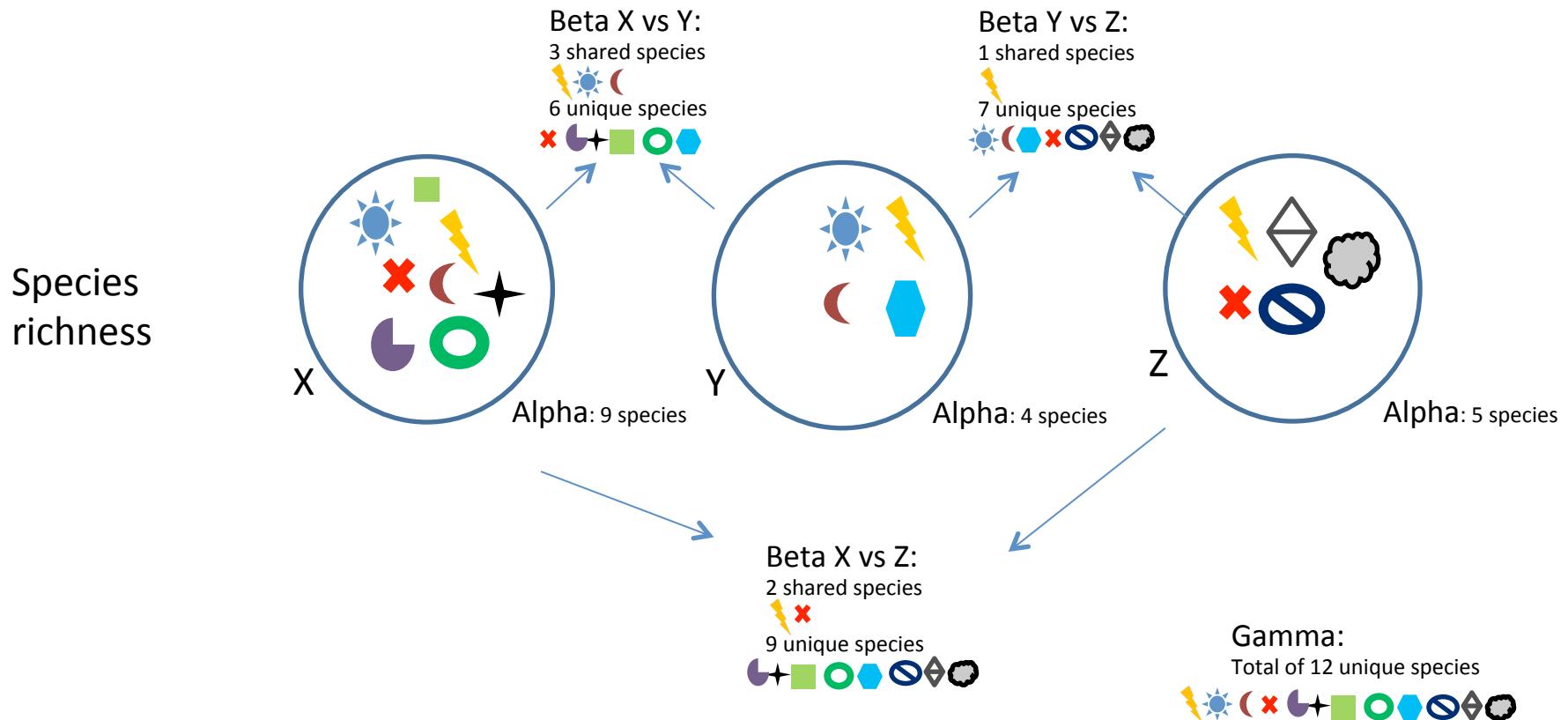
Relative abundance:
 $2/7 = 0.29$ or 29%

★ $5/7 = 0.71$ or 71%

Low Evenness

- Alpha diversity: Diversity within a single sample (Alone)
- Beta diversity: Diversity Between samples
- Gamma diversity: total diversity in a landscape

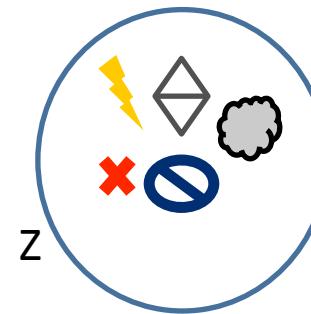
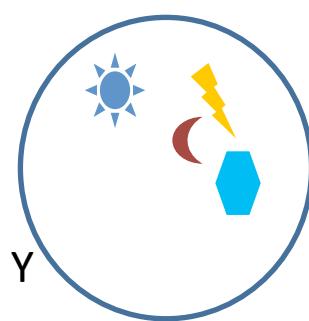
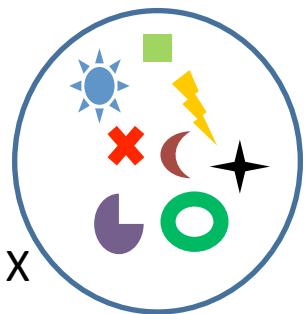
SIMPLE EXAMPLE



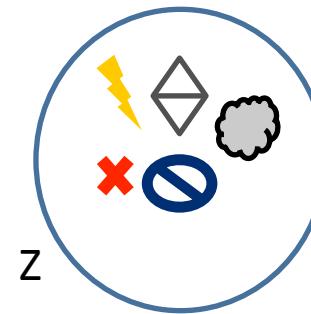
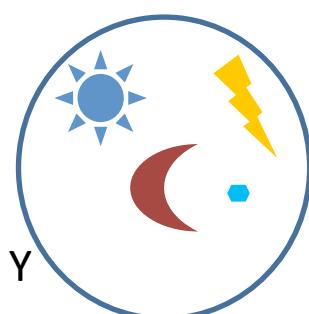
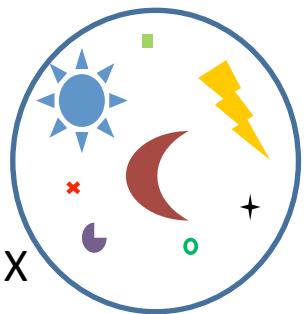
(modified from http://www.webpages.uidaho.edu/veg_measure/Modules/Lessons/Module%209%28Composition&Diversity%29/9_2_Biodiversity.htm)

PROBLEM: Doesn't take abundance of each species OR relatedness of each species into account

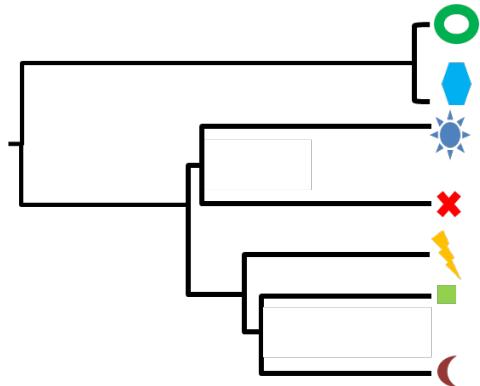
Species richness



Size adjusted
according to
abundance



Phylogenetic
relationship



Metrics used to describe diversity measure different aspects of the community

Rarefaction

50 individuals



2 species

250 individuals



4 species

500 individuals



8 species

Proc. R. Soc. Lond. B (2002) **269**, 2401–2405
DOI 10.1098/rspb.2002.2116

Rarefaction

- Collector's curves
- NGS: individuals = reads
- Evaluate sample size: is sequencing depth(reads per sample) enough?
- Comparing the richness and diversity observed in different samples
- Note rarefaction is not the same as rarefying

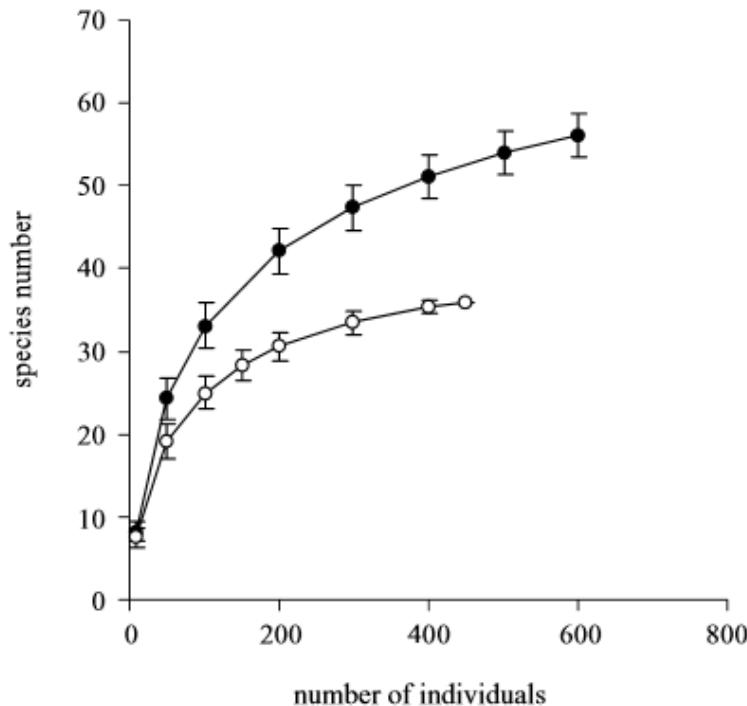
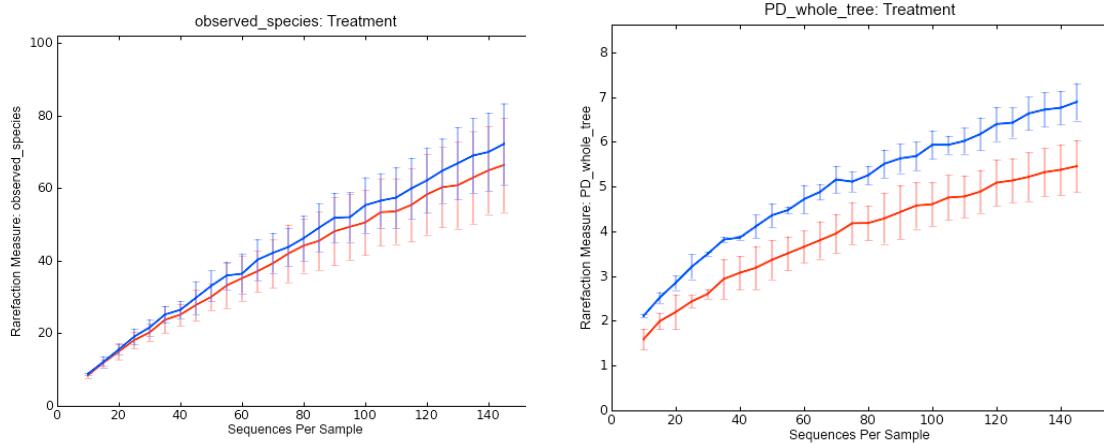
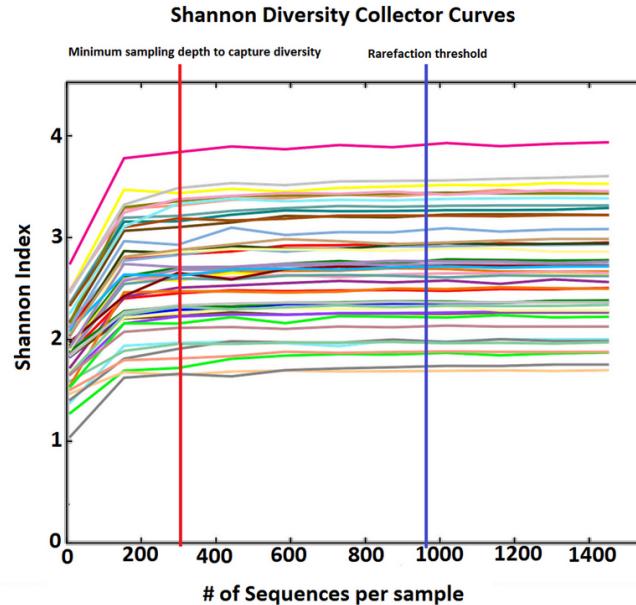


Figure 2. Rarefaction curves for numbers of species of ground dwellers in native savannah (filled circles) and in adjacent agriculture (open circles). For the same sample of individuals (450) there were 30% fewer species in agriculture. The curve for agriculture is approaching an asymptote while the savannah curve climbed to 73 species. Vertical lines indicate one s.d.

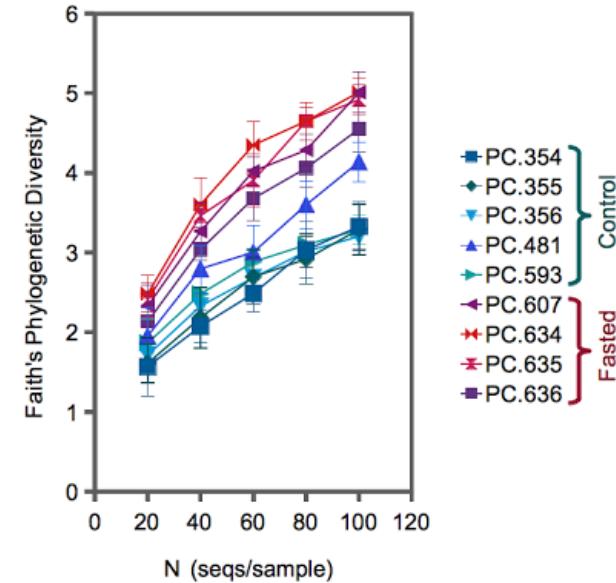
Alpha diversity rarefaction QIIME tutorial and some examples



<http://qiime.org/tutorials/tutorial.html>



Cardenas PA, Cooper PJ, Cox MJ, Chico M, Arias C, et al. (2012) Upper Airways Microbiota in Antibiotic-Naive Wheezing and Healthy Infants from the Tropics of Rural Ecuador. PLoS ONE 7(10): e46803. doi:10.1371/journal.p



<http://www.wernerlab.org/teaching/qiime/overview/e>

Alpha metrics

- Richness: observed species, chao1
- Diversity (Richness and Abundance / Evenness): Shannon, Simpson
- Phylogenetic: PD, Faith's PD or PD_whole tree

Beta diversity

- Diversity between samples
- Single metric to describe difference or similarity between samples
- Non-phylogenetic metrics
- Phylogenetic metrics

Phylogenetic beta diversity: UniFrac distance

APPLIED AND ENVIRONMENTAL MICROBIOLOGY, Dec. 2005, p. 8228–8235
0099-2240/05/\$08.00+0 doi:10.1128/AEM.71.12.8228–8235.2005
Copyright © 2005, American Society for Microbiology. All Rights Reserved.

Vol. 71, No. 12

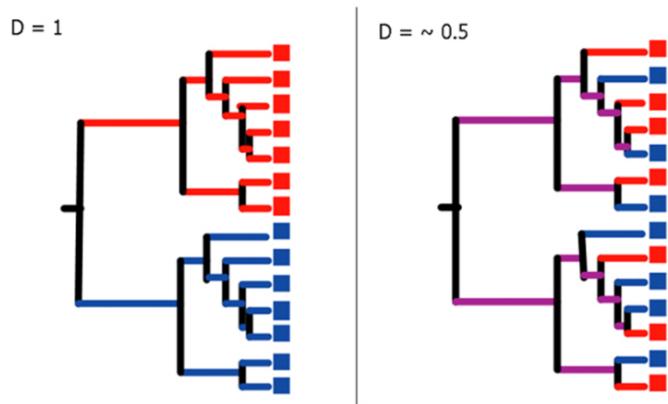
UniFrac: a New Phylogenetic Method for Comparing Microbial Communities

Catherine Lozupone¹ and Rob Knight^{2*}

Department of Molecular, Cellular, and Developmental Biology, University of Colorado, Boulder, Colorado 80309,¹ and Department of Chemistry and Biochemistry, University of Colorado, Boulder, Colorado 80309²

UNIFRAC help: <http://bmf.colorado.edu/unifrac/help.psp>

- Raw unweighted Unifrac: sum of branch length that is unique to one environment or the other

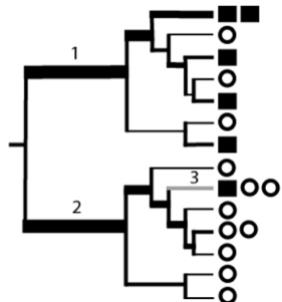


$$u = \frac{\sum_{i=1}^N l_i |A_i - B_i|}{\sum_{i=1}^N l_i \max(A_i, B_i)}$$

l_i is the branch length between node i and its parent, and A_i and B_i are indicators equal to 0 or 1 as descendants of node i are absent or present in communities A and B respectively

A = red, B = blue, branches in common are purple, branches unique to A are red and unique to B are blue. Presence/absence metric.

- Raw weighted Unifrac: Branch lengths are weighted by the relative abundance of sequences

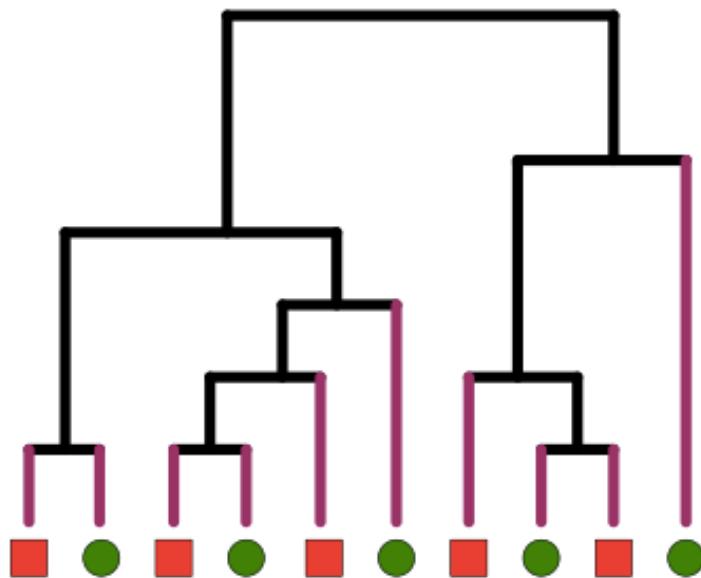


$$u = \sum_i^n b_i \times \left| \frac{A_i}{A_T} - \frac{B_i}{B_T} \right|$$

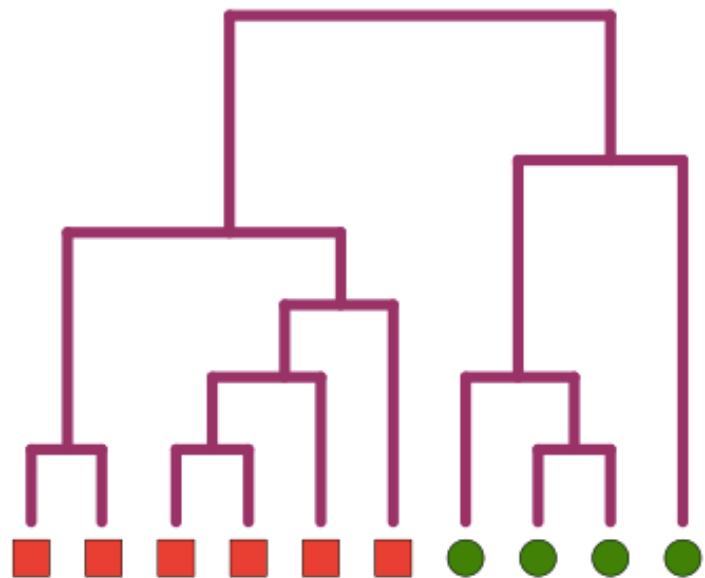
Here, n is the total number of branches in the tree, b_i is the length of branch i, A_i and B_i are the number of descendants of branch i from communities A and B respectively, and A_T and B_T are the total number of sequences from communities A and B respectively. In order to control for unequal sampling effort, A_i and B_i are divided by A_T and B_T .

- Normalised weighted Unifrac: takes abundance and normalises branch length
 - Rapidly evolving lineages (with long branch length can skew unifrac)

Similar Communities



Maximally Different Communities



$$\text{UniFrac Distance Measure} = \frac{\text{Branch length of unique branches}}{\text{Branch length of common branches} + \text{Branch length of unique branches}}$$

Branch length of unique branches

Branch length of common branches

Unweighted Unifrac does not take abundance into account

Weighted Unifrac takes abundance into account – branch lengths weighted by relative abundance

Completely genetically different communities D=1

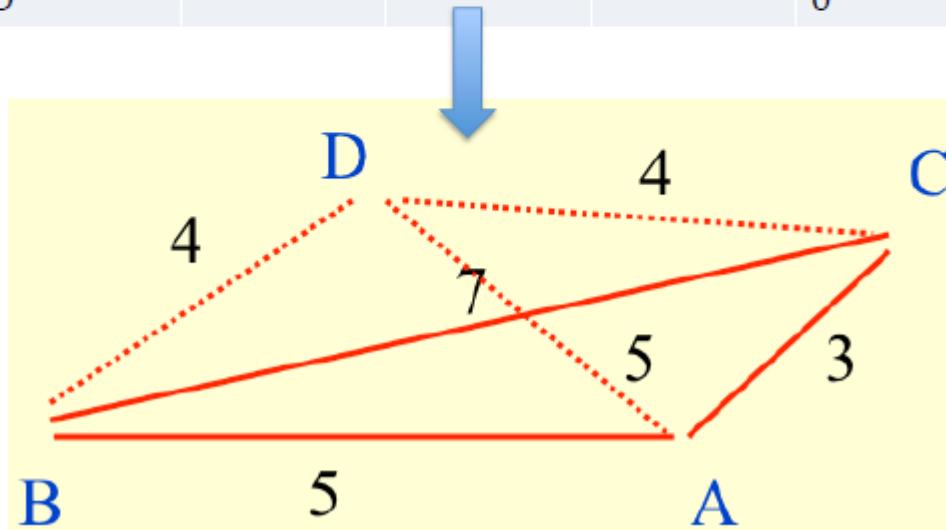
Exactly same communities D=0

Principal Coordinate Analysis (PCoA or PCO)

- Also sometimes called classical MDS (multidimensional scaling)
- Can use any distance matrix (must obey triangle inequality), in this case the Unifrac distance matrix
- Assumes linear relation
- represent distance between samples graphically in multidimensional space (n-1 dimension, n = number samples)
- A new set of reduced variables is derived from the original distances and used to scale samples
- Samples now represented on 2D or 3D plot with these new variables as axes and the relationship between the sample on the plot should reflect their underlying distance
- Ordinates data on plot so that axis 1 (PC1) explains the greatest amount of variance, axis 2 (PC2) explains the next greatest amount of variance, etc.

Principal Coordinates Analysis

	A	B	C	D
A	0	5	3	5
B		0	7	4
C			0	4
D				0

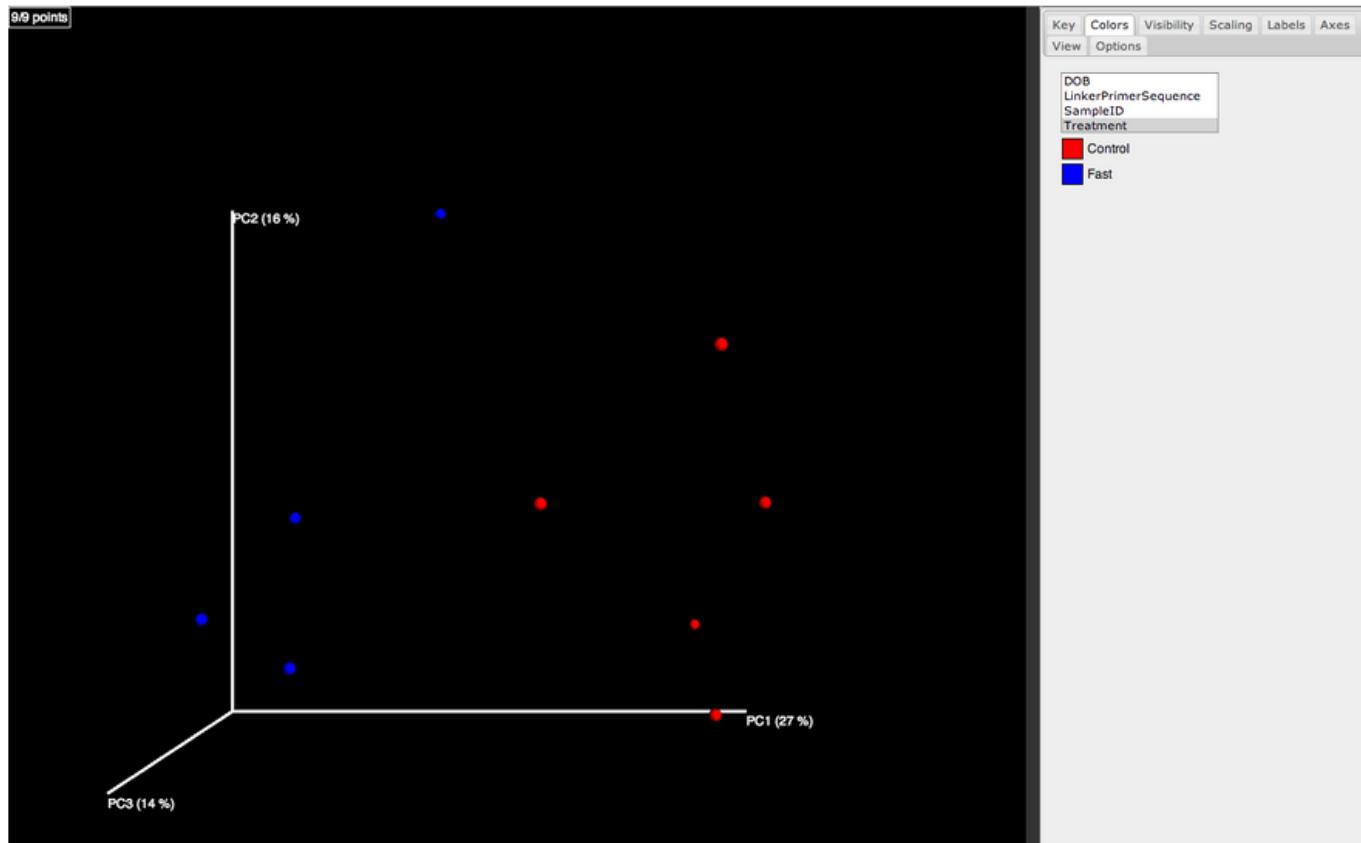


NESCent QIIME Tutorial

Jose Clemente
Daniel McDonald
6.12.12

Adapted from H. Whitehead

[https://www.nescent.org/sites/academy/
File:Nescent_qiimeTutorial_june2012.pdf](https://www.nescent.org/sites/academy/File:Nescent_qiimeTutorial_june2012.pdf)

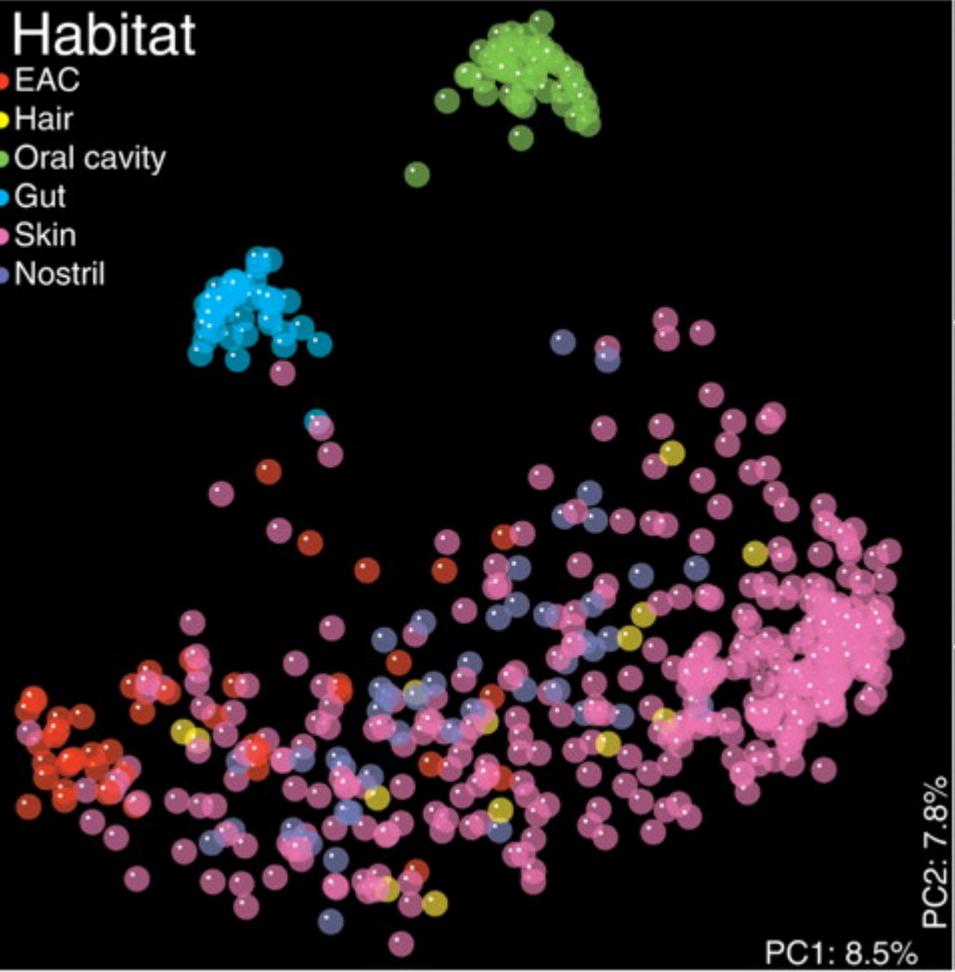


QIIME tutorial data “Five control samples are all red and the four Fast samples are all blue. This lets you easily visualize “clustering” by metadata category. The 3d visualization software allows you to rotate the axes to see the data from different perspectives.” <http://qiime.org/tutorials/tutorial.html>

A

Habitat

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



Communities clustered using PCoA of the unweighted UniFrac distance matrix

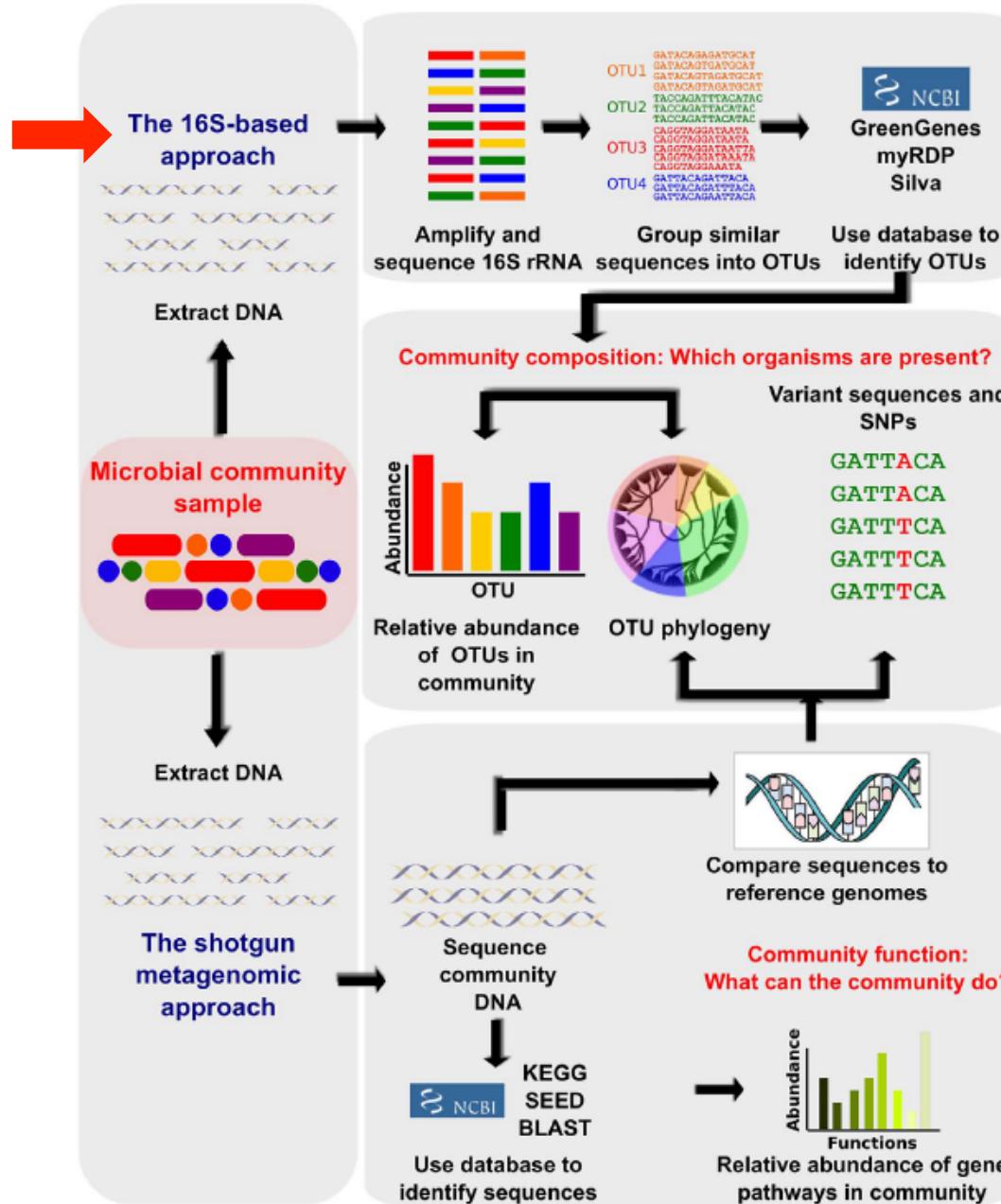


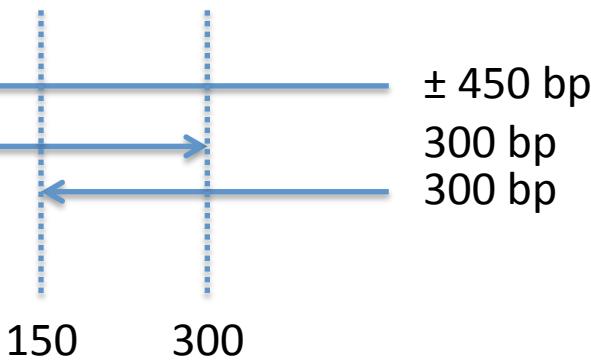
Figure 1. Bioinformatic methods for functional metagenomics. Studies that aim to define the composition and function of uncultured microbial communities are often referred to collectively as “metagenomic,” although this refers more specifically to particular sequencing-based assays. First, community DNA is extracted from a sample, typically uncultured, containing multiple microbial members. The bacterial taxa present in

Which 16S region(/s)?

- Which sequencing platform?



16S V3-V4



MiSeq Series

NextSeq Series

Max Output
15 Gb

Max Read Number
25 M

Max Read Length
2 x 300 bp



HiSeq Series

HiSeq X Series

Max Output
1500 Gb

Max Read Number
5000 M

Max Read Length
2 x 150 bp

ion torrent

Sequencing for all.TM



Max Output

15 Gb

Max Read Number

25 M

Max Read Length

2 x 300 bp

Ion PGM™ System Performance Specifications

	Ion 314™ Chip v2		Ion 316™ Chip v2		Ion 318™ Chip v2	
Output*	200 base	30-50 Mb	300-600 Mb	600 Mb-1 Gb	600 Mb-1 Gb	1.2-2 Gb
	400 base	60-100 Mb				
Reads	400-550 thousand		2-3 million		4-5.5 million	
Run time	200 base	2.3 hr	3.0 hr	4.4 hr	7.3 hr	7.3 hr
	400 base	3.7 hr	4.9 hr			

PCR and sequencing controls

- Controls to detect possible contamination
- Contamination can occur at multiple stages:
 - Sample collection
 - DNA extraction
 - 16S PCR
 - Barcoding and Adaptor PCR
 - Amplicon purification
 - Sequencing

Reagent and laboratory contamination can critically impact sequence-based microbiome analyses

Susannah J Salter^{1*}, Michael J Cox², Elena M Turek², Szymon T Calus³, William O Cookson², Miriam F Moffatt², Paul Turner^{4,5}, Julian Parkhill¹, Nicholas J Loman³ and Alan W Walker^{1,6*}

* Corresponding authors: Susannah J Salter sb18@sanger.ac.uk - Alan W Walker alan.walker@abdn.ac.uk

▼ Author Affiliations

¹ Pathogen Genomics Group, Wellcome Trust Sanger Institute, Hinxton, UK

² Molecular Genetics and Genomics, National Heart and Lung Institute, Imperial College London, London, UK

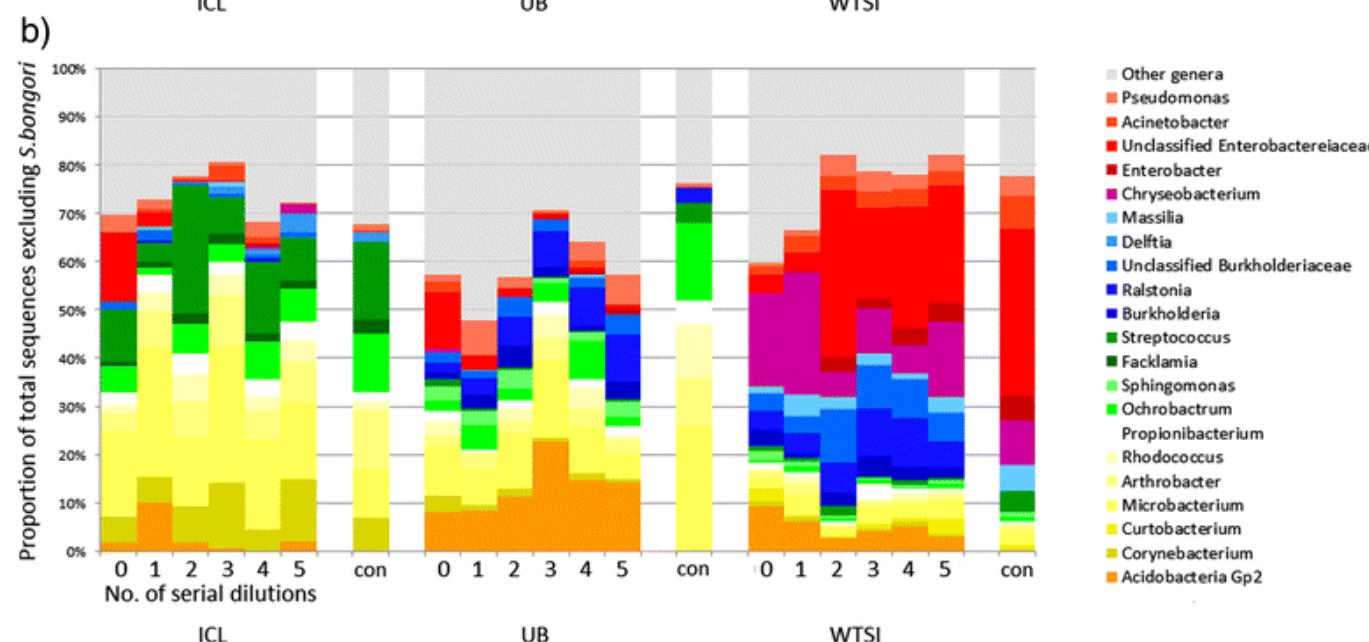
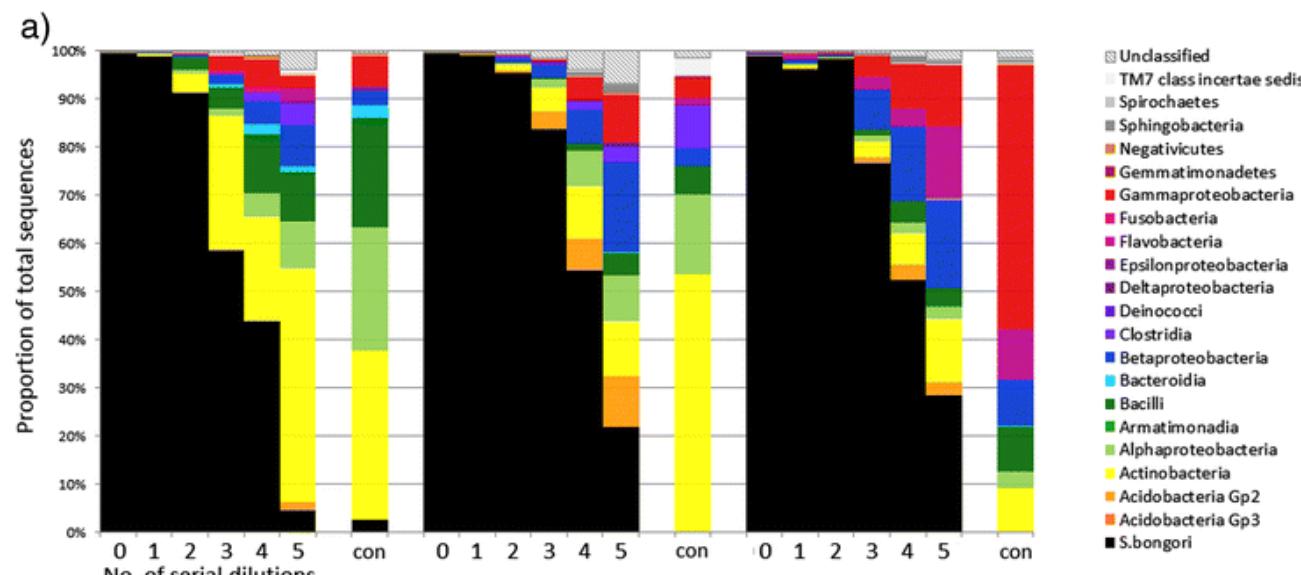
³ Institute of Microbiology and Infection, University of Birmingham, Birmingham, UK

⁴ Shoklo Malaria Research Unit, Mahidol-Oxford Tropical Medicine Research Unit, Faculty of Tropical Medicine, Mahidol University, Mae Sot, Thailand

⁵ Centre for Tropical Medicine, Nuffield Department of Medicine, University of Oxford, Oxford, UK

⁶ Microbiology Group, Rowett Institute of Nutrition and Health, University of Aberdeen, Aberdeen, UK

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



Summary of 16S rRNA gene sequencing taxonomic assignment from ten-fold diluted pure cultures and controls. Undiluted DNA extractions contained approximately 108 cells, and controls (annotated in the Figure with 'con') were template-free PCRs.

Research highlight

Highly accessed

Open Access

Tracking down the sources of experimental contamination in microbiome studies

Sophie Weiss¹, Amnon Amir², Embriette R Hyde², Jessica L Metcalf², Se Jin Song² and Rob Knight^{2,3,4*}

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

▼ Author Affiliations

¹ Department of Chemical and Biological Engineering, University of Colorado at Boulder, Boulder 80309, CO, USA

² BioFrontiers Institute, University of Colorado at Boulder, Boulder 80309, CO, USA

³ Department of Chemistry and Biochemistry, University of Colorado at Boulder, Boulder 80309, CO, USA

⁴ Howard Hughes Medical Institute, Boulder 80309, CO, USA

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

Positive controls

- Mock microbial communities
- BEI Resources (ATCC)

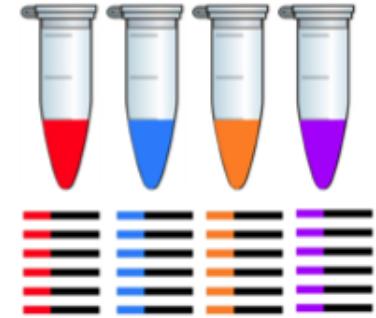


Genomic DNA from Microbial Mock
Community B (Staggered, Low
Concentration), v5.2L, for 16S rRNA Gene
Sequencing

Catalog No. HM-783D

Table 1: Microbial Mock Community B

Organism	NCBI Reference Sequence	Organism	NCBI Reference Sequence
<i>Acinetobacter baumannii</i> , strain 5377	NC_009085	<i>Pseudomonas aeruginosa</i> , strain PAO1-LAC	NC_002516
<i>Actinomyces odontolyticus</i> , strain 1A.21	NZ_AAYI02000000	<i>Rhodobacter sphaeroides</i> , strain ATH 2.4.1	NC_007493, NC_007494
<i>Bacillus cereus</i> , strain NRS 248	NC_003909I	<i>Staphylococcus aureus</i> , strain TCH1516	NC_010079
<i>Bacteroides vulgatus</i> , strain ATCC® 8482™	NC_009614	<i>Staphylococcus epidermidis</i> , FDA strain PCI 1200	NC_004461
<i>Clostridium beijerinckii</i> , strain NCIMB 8052	NC_009617	<i>Streptococcus agalactiae</i> , strain 2603 V/R	NC_004116
<i>Deinococcus radiodurans</i> , strain R1 (smooth)	NC_001263, NC_001264	<i>Streptococcus mutans</i> , strain UA159	NC_004350
<i>Enterococcus faecalis</i> , strain OG1RF	NC_17316	<i>Streptococcus pneumoniae</i> , strain TIGR4	NC_003028
<i>Escherichia coli</i> , strain K12, substrain MG1655	NC_000913		
<i>Helicobacter pylori</i> , strain 26695	NC_000915		
<i>Lactobacillus gasseri</i> , strain 63 AM	NC_008530		
<i>Listeria monocytogenes</i> , strain EGDe	NC_003210		
<i>Neisseria meningitidis</i> , strain MC58	NC_003112		
<i>Propionibacterium acnes</i> , strain KPA171202	NC_006085		



Sequencing

```
>GCACCTGAGGACAGGCATGAGGAA_
>GCACCTGAGGACAGGGGAGGAGGA_
>TCACATGAACCTAGGCAGGACGAA_
>CTACCGGAGGACAGGCATGAGGAT_
>TCACATGAACCTAGGCAGGAGGA_
>GCACCTGAGGACACGCAGGACGAC_
>CTACCGGAGGACAGGCAGGAGGA_
>CTACCGGAGGACACACAGGAGGA_
>GAACCTTCACATAGGCAGGAGGAT_
>TCACATGAACCTAGGGCAAGGAA_
>GCACCTGAGGACAGGCAGGAGGA_
```

Demultiplex

Merge paired end reads

Quality control:
Remove primer sequences
Remove chimeras
Trim reads based on Phred score
Remove singlettons (and doubletions)

High quality sequences

Associate sequence data with samples and sample metadata

Clustering reads

- Sequence clustering into OTUs (Operational Taxonomic Units)
- Pick representative sequence for each OTU

OTU picking (clustering)

- QIIME has three methods for OTU picking
 - *de novo*
 - closed-reference
 - open reference
 - Sequences grouped together based on sequence identity (*de novo*) or alignment to reference sequence (closed = reads discarded if they don't match a reference, open = reads that don't hit reference from *de novo* cluster)
- http://qiime.org/tutorials/otu_picking.html

Assign taxonomy

Tools

- Blast
- RDP classifier
- UCLUST is now preferred method

16S database