

Acknowledgements

Susan Holmes	Former postdoc advisor, mentor, co-author
Benjamin Callahan	DADA2 first author, slides, discussions, feedback, etc.
Holmes Group	Helpful advice and feedback
Wolfgang Huber	Helpful advice and feedback, creator of DESeq and DESeq2
BioC and CRAN	Support, Feedback, Distribution of phyloseq and biom
Rob Knight	QIIME, UniFrac, etc.
Huttenhower grp	Biobakery suite, slides, etc.
Hadley Wickham	ggplot2, reshape2, plyr R packages, Rstudio

Schedule for today

Sec	Day	Start	End	Торіс	Lead Instr.
1	Mon	09:00	10:00	Introduction to Metagenomics. Culture independent techniques, 16S rRNA, etc. (60 -75 min)	Joey
2	Mon	10:00	11:00	Introduction to microbiome analysis concepts Exploratory data analysis, Distances, PCoA, Ordination, taxa & sample-level inferences (75 min)	Joey
3	Mon	11:00	11:59	Introduction to microbiome analysis practices: QIIME, phyloseq, reproducible research (30 min)	Joey
	Mon	12:00	14:00	Lunch (120min)	
4	Mon	14:00	17:00	QIIME Lab (180min)	Daniel
	Mon	17:00	19:00	Dinner (120min)	
5	Mon	19:00	22:00	phyloseq Lab (180min)	Joey

An Introduction to Metagenomics

Outline for Today:

- What is metagenomics?
 - What methods, theoretical basis?
 - Why is it useful?
 - Where is it headed?
- How can I use it?
 - wet lab procedures (dry workshop)
 - computational protocols, practices





Why Study Microbiomes?

Environmental Science

- Critical elemental cycles (carbon, nitrogen, sulfur, iron, ...)
- Pollution control, cleanup
- Ecology / Evolution (chloroplasts, mitochondria, genetic evolution, ...)

Industrial Applications

- Wastewater treatment (V. cholera, algal blooms, etc.)
- Bioprospecting (novel enzymes, compounds)
- Novel biosynthesis
- Fermentations: Consortia (yogurt) / wild (kombucha, Belgian ales)

Human Health

- Protection from pathogens (e.g. Clostridium difficile)
- Absorption/Production of nutrients in the gut
- Possible Role in chronic diseases (obesity, Crohn's/IBD, other autoimmune, UTIs, periodontitis, ...)

What is a microbiome?

The totality of microbes in a defined environment, especially their genomes and interactions with each other and surrounding environment.

- A population of a single species/strain is a culture, extremely rare outside of lab, some infections
- A microbiome is a mixed population of different microbial species (microbial ecosystem)

A mixed community is the norm!

What is a microbiome?





Cow Rumen



Human Microbiomes

What is a microbiome?

acid mine biofilm

acid mine biofilm

types

Ancestral strains



Tyson, et al. (2004) Nature, 428(6978), 37-43

What is a microbiome?



E Eukaryotes 4% Sulfobacillus spp. 1% Archaea 10% Leptospirilium gp III 10%



Tyson, et al. (2004) Nature, 428(6978), 37-43





1-10 times more microbial cells than human cells... depends on timing of your last bowel movement

Typical human microbiome < 2 kg





Some provocative oversimplifications...

Microbes can...

- I. "Kill you by acute infection"
- 2. "Prevent same infection"
- 3. "Make you fat(ter)"
- 4. "Give you a heart attack"
- 5. "Give you cancer"
- 6. "Rescue you from cancer"

Can you guess the condition / scenario?



Turnbaugh, et al. (2009). A core gut microbiome in obese and lean twins. Nature

Microbes can make you fat(ter)...

- Lean (n = 10) & obese donors (n=9)
- Colonization of germ-free wild-type mice with microbiota from obese donors causes significant increase in total body fat
- Total body fat content was measured before and after a 2-week colonization
- Confirm that the ob/ob microbiome has an increased capacity for dietary energy harvest



Turnbaugh, et al. (2006). An obesity-associated gut microbiome ... Nature

Gut microbes promote cardiovascular disease



Gut flora required for production of TMAO
Supplementing diet with choline or TMAO promotes atherosclerosis (mouse)
Gut flora suppression (Abx) inhibits dietary choline enhanced atherosclerosis
TMAO is also a renal (kidney) toxin. Fogelman, A. M. (2015). *Circulation Research*.

ZN Wang, ..., Stanley Hazen. *Nature* **472**, 57-63 (2011) Fogelman, A. M. (2015). TMAO Is Both a Biomarker and a Renal Toxin. *Circulation Research*.

Colorectal Cancer (CRC)

- Microbes affect colonic bile pool exposure, drug metabolism, and mortality-correlated compounds
- Microbe-produced secondary bile acids are among these.
- Gut microbial metabolism may play role in beneficial or detrimental effects of certain foods

Sears, C. L., & Garrett, W. S. (2014). Microbes, Microbiota, and Colon Cancer. Cell Host & Microbe, 15(3), 317–328.



Groundwater: Chlorinated Solvents





McCarty, P. L. (1997). Breathing with chlorinated solvents. Science

Marine picoplankton most abundant organism on Earth?

- Prochlorococcus appears to be the most abundant organism on the planet
- Huge light harvesting proteins
- its density can reach up to 100 million cells per liter
- it can be found down to a **depth of 150 m** in all of the intertropical belt
- picoplankton synchronize cell division at the same time every day -> biological clock

OLIPAC cruise Pacific Ocean 1994 Oligotrophic 16°S





Vertical distribution of the photosynthetic picoplankton populations determined by flow cytometry in the tropical Pacific (OLIPAC cruise, 1994).

Yellowstone National Park



Octopus Spring

- 90° to 93°C
- extremely low in nutrients
- contains abundant biomass
- home to "oldest" known bacteria



Obsidian Pool

- 75° 95°C
- high iron (II) hydrogen sulfide
 - extensive diversity (previously unknown)

Ward, D. M., Weller, R., & Bateson, M. M. (1990). Nature, 345(6270), 63-65. Barns, S. M., Fundyga, R. E., Jeffries, M.W., & Pace, N. R. (1994). PNAS 91(5), 1609–1613.

Symbiosis: sea-floor vent tube worm



Symbiosis: sea-floor vent tube worm





Cavanaugh, C. M. (1983). Nature, 302(5903), 58-61. Cavanaugh, C. M., et al. (1981). Science. 213(4505), 340-342 End: Biological Motivation

Questions before moving on?

Metagenomics Experimental Methods

Exercise: How many species are present?





Confer amongst yourselves. We'll take a poll.

The great "plate count" anomaly

- Cultivation-based cell counts are orders of magnitude lower than direct microscopic observation.
- This is because microbiologists are able to cultivate only a small minority of naturally occurring microbes
- Our nucleic-acid derived understanding of microbial diversity has rapidly outpaced our ability to culture new microbes



Staley, J.T., & Konopka, A. (1985). Measurement of in situ activities of nonphotosynthetic microorganisms in aquatic and terrestrial habitats. *Annual Review of Microbiology*, 39, 321–346.

Why is microbiome research new?

Considering that...

- We have a bacterial endosymbiont in all our cells!
- Humans have always coexisted with bacteria
- We've known about bacteria for a few hundred years





- Historically prokaryotic biology has been focused on microbes that can be grown to large quantities/densities in the lab, especially pathogens; or can be distinguished under the microscope.
- An example of "searching where the light is"...

Why is microbiome research new?

Bias for cultivable microbes, especially pathogens

- Culture-based methods fail to detect most microbes
- Microbes are easy to miss (except pathogens)
- Most microbes are NOT pathogens (even the human-associated)

Availability of tools limited to last 3 decades

- Discovery of culture-independent techniques
- PCR, fast & cheap DNA sequencing, microarrays, etc

Discovery of Culture Independent Techniques

- 1977 rRNA as evolutionary marker Woese & Fox PNAS
- 1985 Polymerase Chain Reaction (PCR) K. Mullis Science
- 1985 "Universal" Primers for rRNA sequencing N. Pace PNAS
- 1989 PCR amplification of 16S rRNA gene Böttger FEMS Microbiol.
- 1996 Large, curated rRNA database (RDP) Maidak Nuc. Acids Res
- 2001 term "microbiome" coined by Joshua Lederberg

Discovery of Culture Independent Techniques

ribosome

- rRNA has both catalytic and structural function.
- The small and large subunits have different lengths, 2nd-structure, 3D shape; but must work together.
- All of the catalytic activity of the ribosome is carried out by the RNA; the proteins reside on the surface and seem to stabilize the structure.





Discovery of Culture Independent Techniques



Discovery of *Culture Independent* Techniques Small subunit "165" rRNA



- Ubiquitous present in all known life (viruses don't count)
- Functionally constant translation, 2°-structure
- Evolves slowly mutations more rare than for protein-coding genes
- Large information for evolutionary inference
- No exchange Limited examples of rRNA gene-sharing between organisms

Metagenomics: Nucleic acid sequencing as a tool for microbial community analysis

Single microbiome:

- I. Break all cells, extract all DNA (gDNA)
- 2. PCR-amplify a universal gene from gDNA
- 3. DNA sequencing from pool of amplified genes
- 4. Cluster sequences according to species
- 5. Count each species and make a tree

Tringe, S. G., & Rubin, E. M. (2005). Metagenomics: DNA sequencing of environmental samples. *Nature Reviews Genetics*, 6(11), 805–814.







End Metagenomics Lecture I

Questions?

Introduction to Microbiome / Metagenome Analysis Concepts



•Sequence Processing (OTUs)

- •Denoising
- Chimera detection
- •Construction of sequence clusters (OTUs)

•Comparing microbiomes

- •Distances, Diversity
- •Exploratory Data Analysis
- Ordination Methods
- •hierarchical dendrogram
- •extract patterns from a plot
- •clusters gap statistic
- •gradient regression, modeling, etc.
- Identifying important microbes/taxa
- •projected points, coinertia (plots)
- •inferential testing
- modeling

- •Sequence Processing (OTUs)
- •Denoising
- •Chimera detection
- •Construction of sequence clusters (OTUs)
- •Comparing microbiomes
- •Distances, Diversity
- Exploratory Data Analysis
- Ordination Methods
- •hierarchical dendrogram
- extract patterns from a plot
 clusters gap statistic
- •gradient regression, modeling, etc.
- Identifying important microbes/taxa
- •projected points, coinertia (plots)
- •inferential testing
- modeling



Slide graciously provided by Benjamin Callahan, not necessarily with permission O:-)





Slide graciously provided by Benjamin Callahan, not necessarily with permission O:-)







Community Distance Properties

- Range from 0 to 1
- · Distance to self is 0
- \cdot If no shared taxa, distance is 1
- Triangle inequality (metric)
- · Joint absences do not affect distance (biology)
- Independent of absolute counts (metagenomics)

Slide graciously provided by Benjamin Callahan, not necessarily with permission O:-)

Jaccard

 $= ((\mathbf{x}_{A} > 0) \& (\mathbf{x}_{B} > 0))/((\mathbf{x}_{A} > 0) | (\mathbf{x}_{B} > 0))$

AυB

В

А

 $Dist(A, B) = 1 - (A \cap B)/(A \cup B)$

В

A∩B

А

The Distance Spectrum



Slide graciously provided by Benjamin Callahan, not necessarily with permission O:-)







Ordination Methods

Project high-dimensional data onto lower dimensions



Multi-dimensional Scaling



Ordination Methods

Intuition:

Each PC axis is projection that maximizes the area of the shadow Equivalently - max(sum of square of distances between points) Goal: "See" as much variation as possible



MDS Details

Given distances between each observation (sample), MDS finds the closest approximation of that in lower dimensional Euclidean space.

- Algorithm starts from **D** inter-point distances:
 - Center the rows and columns of the distance matrix: $\mathbf{S} = -1/2 \ \mathbf{H} \ \mathbf{D}^{(2)} \ \mathbf{H}$
 - Compute SVD by diagonalizing S: S = U A U^T
 - Extract Euclidean representations: $\mathbf{X} = \mathbf{U} \mathbf{\Lambda}^{1/2}$
- The relative values of diagonal elements of A gives the proportion of variability explained by each of the axes.
- \cdot The valued of Λ should always be looked at in deciding how many dimensions to retain

NMDS is similar, but minimizes a different function (difference in distance ranks) Slide graciously provided by Benjamin Callahan, not necessarily with permission O:-)



Exploratory Data Analysis "Unsupervised Learning" "Ordination Methods"

What we "learn" depends on the data.

- How many axes are probably useful?
- Are their clusters? How many?
- Are their gradients?
- Are the patterns consistent with covariates
- (e.g. sample observations)
- How might we test this?



Exploratory Data Analysis "Unsupervised Learning" "Ordination Methods"

• Are their clusters? How many?

Technique: Gap Statistic

Exploratory Data Analysis	Exploratory Data Analysis		
"Unsupervised Learning"	"Unsupervised Learning"		
"Ordination Methods"	"Ordination Methods"		
 Are their gradients? Are they explained by one or	 Are the patterns consistent with covariates? Technique:		
more sample covariates? Technique:	Permutational Multivariate ANOVA		
PC regression (statistics' PCR)	vegan::adonis()		
End: Introduction to Microbiome / Metagenome Analysis Concepts Questions?	<section-header></section-header>		

Introduction to Microbiome / Metagenome Analysis Tools and Practices

- I. Probably-not-comprehensive summary of metagenomic tools
- 2. Short sermon on the virtues of reproducible analysis
- 3. Introduction to phyloseq & send-off this afternoon's lab





I6S rRNA Databases

- · GreenGenes http://greengenes.secondgenome.com
- Silva <u>www.arb-silva.de</u>
- Ribosomal Database Project (RDP) <u>https://rdp.cme.msu.edu</u>
 - •~100Ks millions of unique 16S rRNA genes
 - Curated taxonomy
 - •Classification tools (e.g. RDP classifier, ARB, etc.)

(16S rRNA) Amplicon Sequence Processing Tools:

- · QIIME (Soon 'Qiita') http://qiime.org/
- mothur <u>www.mothur.org/</u>
- usearch <u>www.drive5.com/usearch</u>
- DADA2 <u>https://github.com/benjjneb/dada2</u>

Afternoon will be spent using QIIME Daniel has much more to say about it...



- Looking for variants?
 - Clustering:
 Mapping:
- +specificity, -difficulty
- +sensitivity, -novelty
- What else?

Slide graciously provided by Curtis Huttenhower, not necessarily with permission O:-)



MetaPhIAn: Taxonomic profiling using unique marker genes X is a core gene for clade Y X is a unique marker gene for clade Y Gene X Gene X Gene X Must representative markers used for identification 184±45 markers per species (target 200) ~7,100 species (excludes incomplete annotations, spp., etc.) False positive/False negative rates of ~1 in 10⁶ Profiles all domains of life: bacteria, viruses, euks, archaea Strain level profiling using marker barcodes and SNPs Quasi-markers used to resolve ambiguity in postprocessing

p://huttenhower.sp5libagraciously provided by Curtis Huttenhower, not necessarily with permission O:-)



Reproducible analysis of microbiome / metagenome data

- Why make the effort?
- What if I don't want someone else reproducing my analysis?
- What if I don't know how?
- Isn't it enough to provide a cursory description in the methods section with a light sprinkling of literature citations?

illustrative example favoring reproducible analysis: "Enterotypes of the human genome"



illustrative example favoring reproducible analysis: "Enterotypes of the human genome"

MDS on supported distance metrics: enterotype data















Schedule for today

Sec	Day	Start	End	Торіс	Lead Instr.
1	Mon	09:00	10:00	Introduction to Metagenomics. Culture independent techniques, 16S rRNA, etc. (60 -75 min)	Joey
2	Mon	10:00	11:00	Introduction to microbiome analysis concepts Exploratory data analysis, Distances, PCoA, Ordination, taxa & sample-level inferences (75 min)	Joey
3	Mon	11:00	11:59	Introduction to microbiome analysis practices: QIIME, phyloseq, reproducible research (30 min)	Joey
	Mon	12:00	14:00	Lunch (120min)	
4	Mon	14:00	17:00	QIIME Lab (180min)	Daniel
	Mon	17:00	19:00	Dinner (120min)	
5	Mon	19:00	22:00	phyloseq Lab (180min)	Joey