

# Lecture 21. Microbiology & Metagenomics

Michael Schatz & Kelly Moffat

April 25, 2017

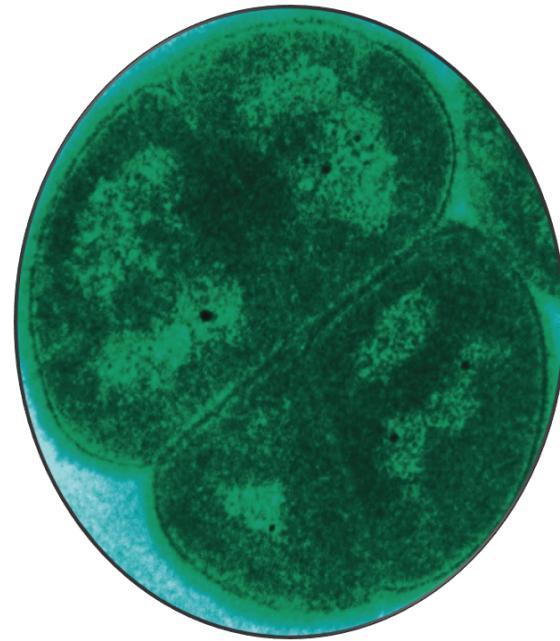
JHU 600.649: Applied Comparative Genomics



# Kelly's background



*A. thaliana*

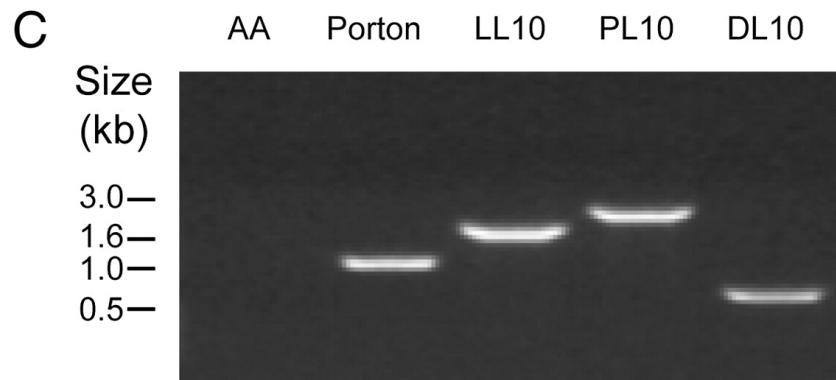
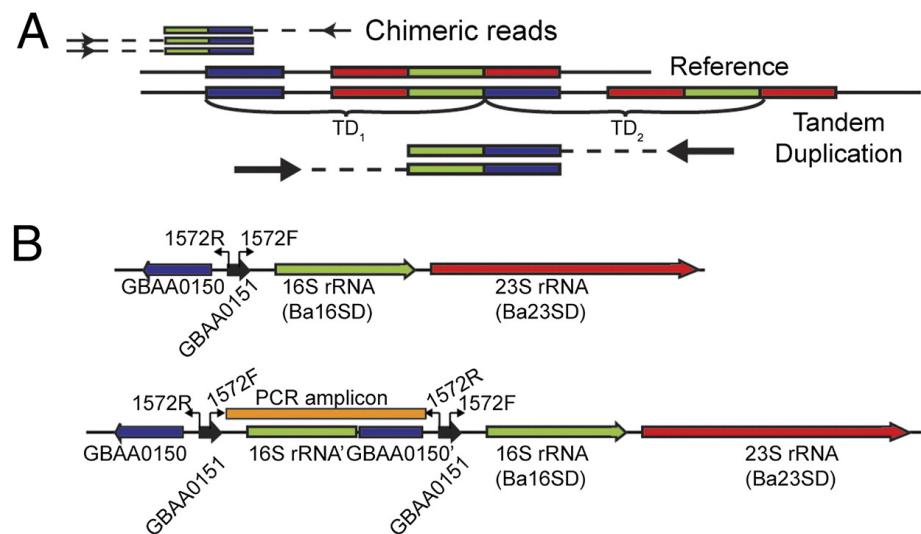
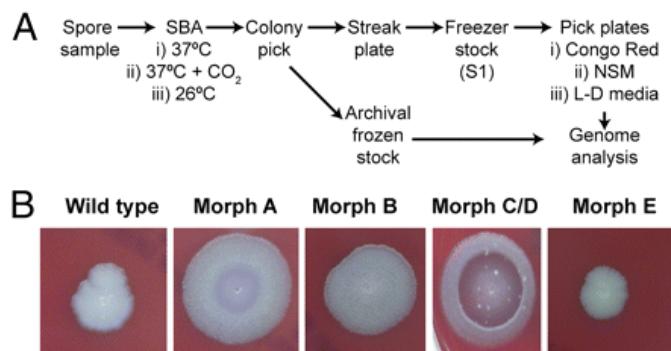
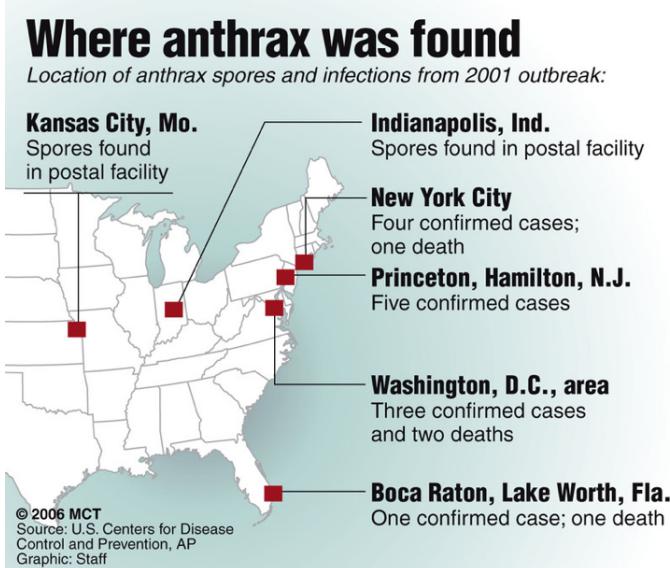


*D. radiodurans*



Rice

# Amerithrax Analysis



**Bacillus anthracis comparative genome analysis in support of the Amerithrax investigation**  
 Rasko et al (2011) PNAS. doi: 10.1073/pnas.1016657108

**Poorani Sub...**

HMP\_SRS023583\_rm2.fasta.gz

Select database Database statistics File metadata View settings

BETA Bacteria Q1 2016 Number of organisms: 90 File size: 1.57 GiB Filtered results. Display non filtered results Uploaded: Monday, June 6, 2016 3:59 PM

Table Sunburst Bubble Bubble Packing Collapsible Tree Radial Tree

Name	Frequency	Unique Matches %	Total Matches %	Relative Abundance
Bacteroides vulgatus dnLKV7	107645	7.06	9.51	15.07
Bacteroides dorei CL03T12C01	427646	56.29	93.33	12.71
Bacteroides fragilis str 3725 D9 ii	7375	60.43	82.21	9.06
Alistipes putredinis DSM 17216	6879195	63.43	63.43	8.33
Bacteroides 1166 Branch	277	87.5	84.45	8.09
Bacteroides uniformis str 3978 T3 ii	684	10.99	66.38	7.2
Bacteroides ovatus str 3725 D1 iv	410	12.71	78.38	6.97
Bacteroides plebeius DSM 17135	9270784	23.34	23.34	6.95
Bacteroides finegoldii DSM 17565	2529854	90.03	90.74	4.22
Bacteroides caccae ATCC 43185	532352	32.15	81.25	2.41
Bacteroides stercoris CC31F	665518	39.47	60.61	1.93
Alistipes shahii WAL 8301	1178469	78.49	77.31	1.54
Parabacteroides merdae CL09T00C40	158275	29.63	64.46	1.36
Parabacteroides distasonis str 3999B TB 6	2285	35.98	59.11	1.34

Showing 1 to 50 of 90 entries

Save to CSV

Drop files to upload. Click to open file browser.

**Poorani Sub...**

HMP\_SRS023583\_rm2.fasta.gz

Select database Database statistics File metadata View settings

BETA Bacteria Q1 2016 Number of organisms: 90 File size: 1.57 GiB Filtered results. Display non filtered results Uploaded: Monday, June 6, 2016 3:59 PM

Table Sunburst Bubble Bubble Packing Collapsible Tree Radial Tree

Max Depth: 6 Collapse: OFF

Node name: Clostridiales  
Relative abundance: 2.23

Bacteria  
Firmicutes  
Clostridia  
Clostridiales 50.00%  
Clostridiaceae 0.90%

**Poorani Sub...**

Comparative Analysis: Report

Name: NICED new  
Date: Tuesday, September 13, 2016 1:16 PM  
Database: BETA Bacteria Q1 2016  
Field: Frequency  
Log Scale: No

Comparative Analysis Results

Principal Component Analysis By attribute: Frequency

PC2

PC3

HMP NICED

Live Chat

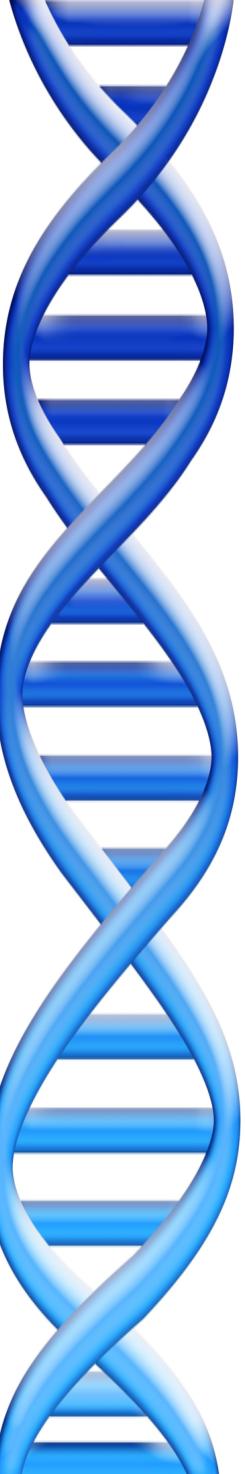
News Settings

Drop files to upload. Click to open file browser.

Live Chat

3D PCA

0.005% 0.05% 0.5% 0.5%



# Part I: Introduction

# Microbial Taxonomy

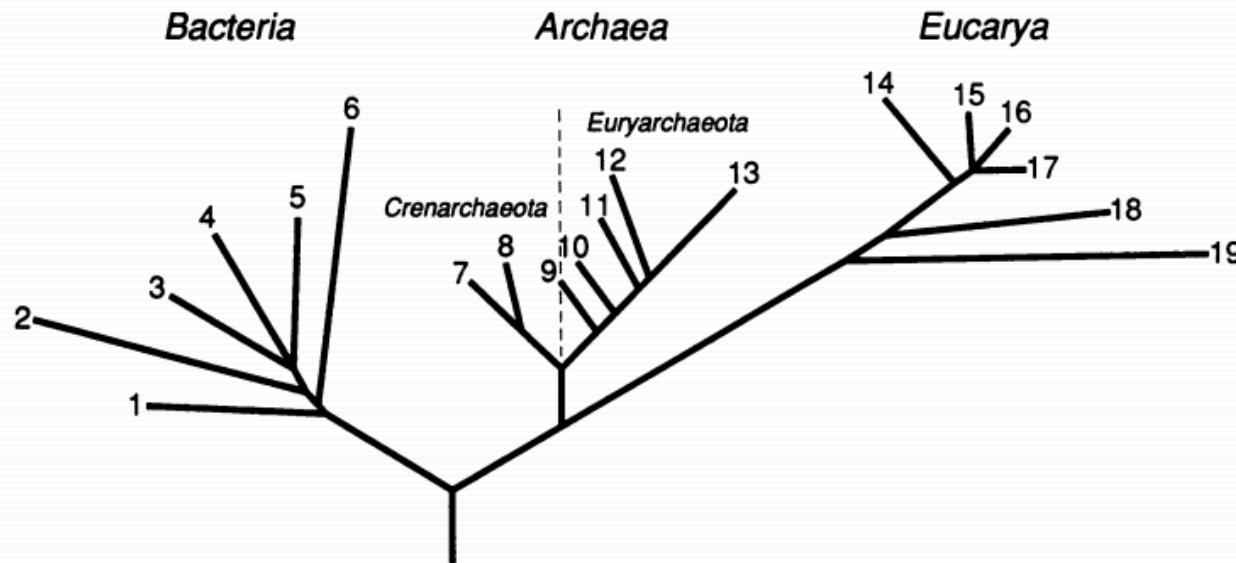
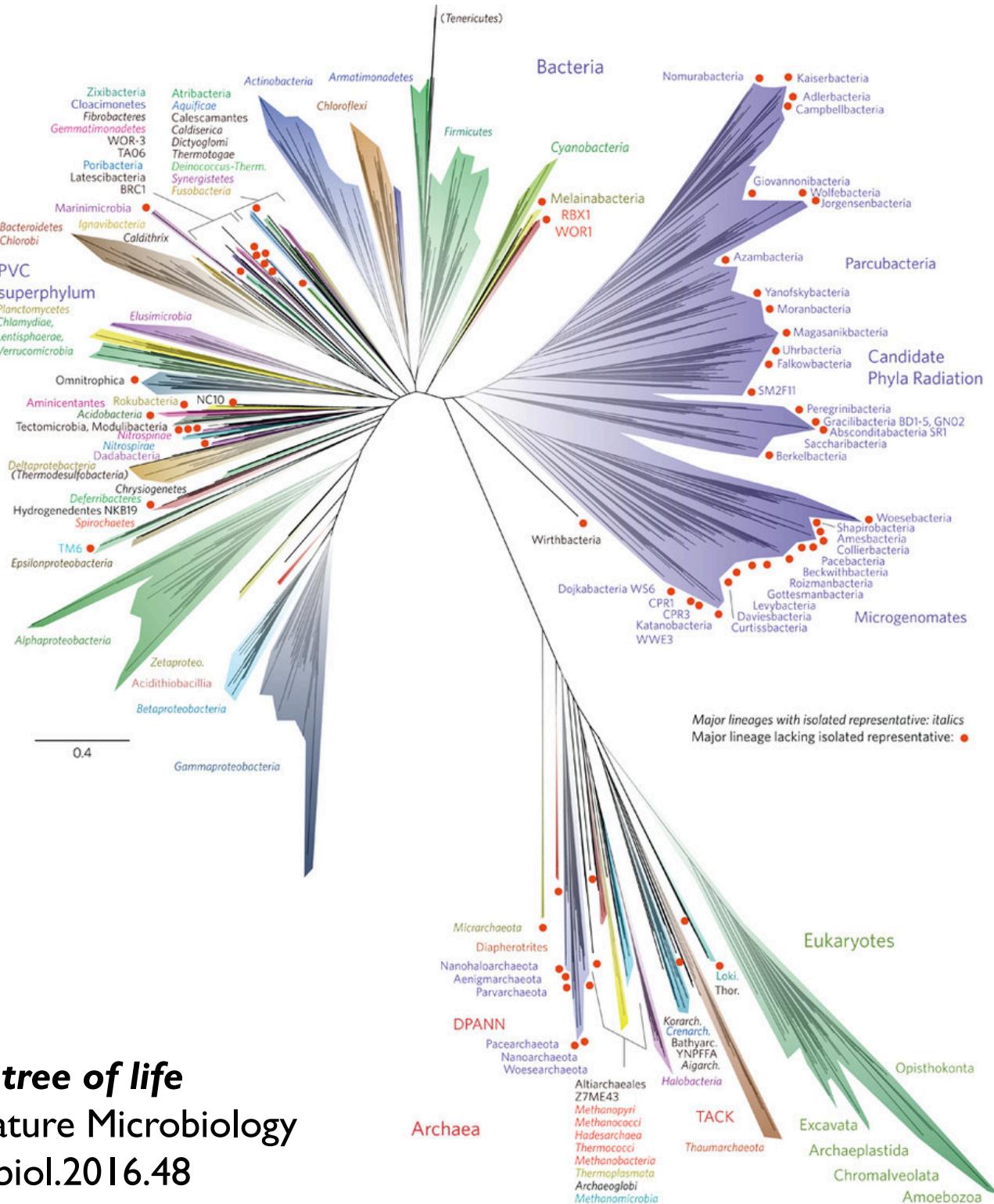


FIG. 1. Universal phylogenetic tree in rooted form, showing the three domains. Branching order and branch lengths are based upon rRNA sequence comparisons (and have been taken from figure 4 of ref. 2). The position of the root was determined by comparing (the few known) sequences of pairs of paralogous genes that diverged from each other before the three primary lineages emerged from their common ancestral condition (27). [This rooting strategy (28) in effect uses the one set of (aboriginally duplicated) genes as an outgroup for the other.] The numbers on the branch tips correspond to the following groups of organisms (2). Bacteria: 1, the Thermotogales; 2, the flavobacteria and relatives; 3, the cyanobacteria; 4, the purple bacteria; 5, the Gram-positive bacteria; and 6, the green nonsulfur bacteria. Archae: the kingdom Crenarchaeota: 7, the genus *Pyrodictium*; and 8, the genus *Thermoproteus*; and the kingdom Euryarchaeota: 9, the Thermococcales; 10, the Methanococcales; 11, the Methanobacteriales; 12, the Methanomicrobiales; and 13, the extreme halophiles. Eucarya: 14, the animals; 15, the ciliates; 16, the green plants; 17, the fungi; 18, the flagellates; and 19, the microsporidia.

**Towards a natural system of organisms: proposal for the domains Archaea, Bacteria, and Eucarya.**  
Woese et al (1990) PNAS. doi: 10.1073/pnas.87.12.4576



**A new view of the tree of life**  
Hug et al. (2016) Nature Microbiology  
doi:10.1038/nmicrobiol.2016.48

# Your second genome?



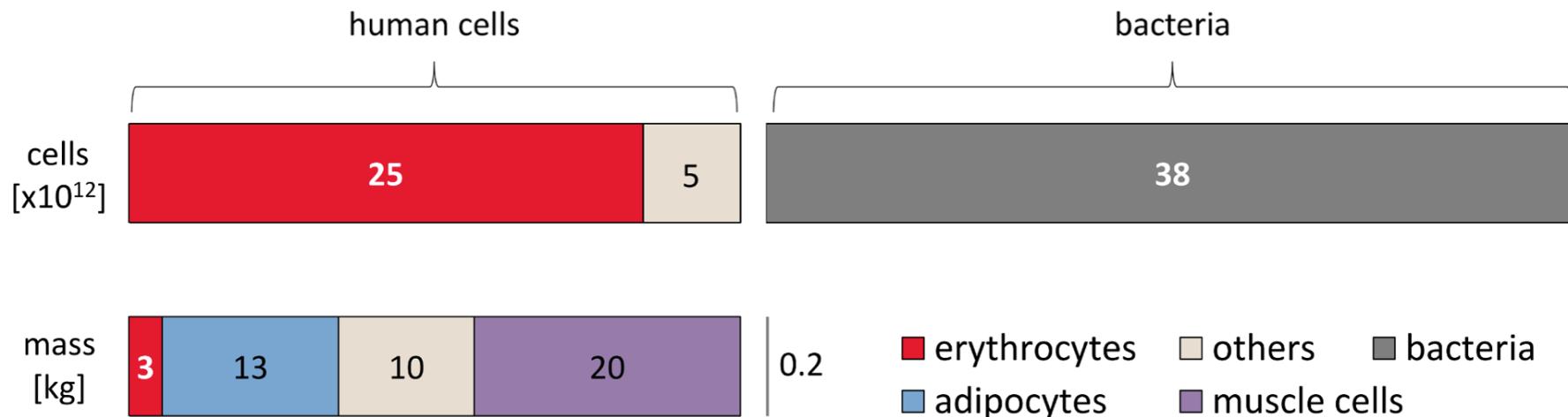
**Human body:**  
**~10 trillion cells**

**Human brain:**  
**~3.3 lbs**

**Microbiome**  
**~100 trillion cells**

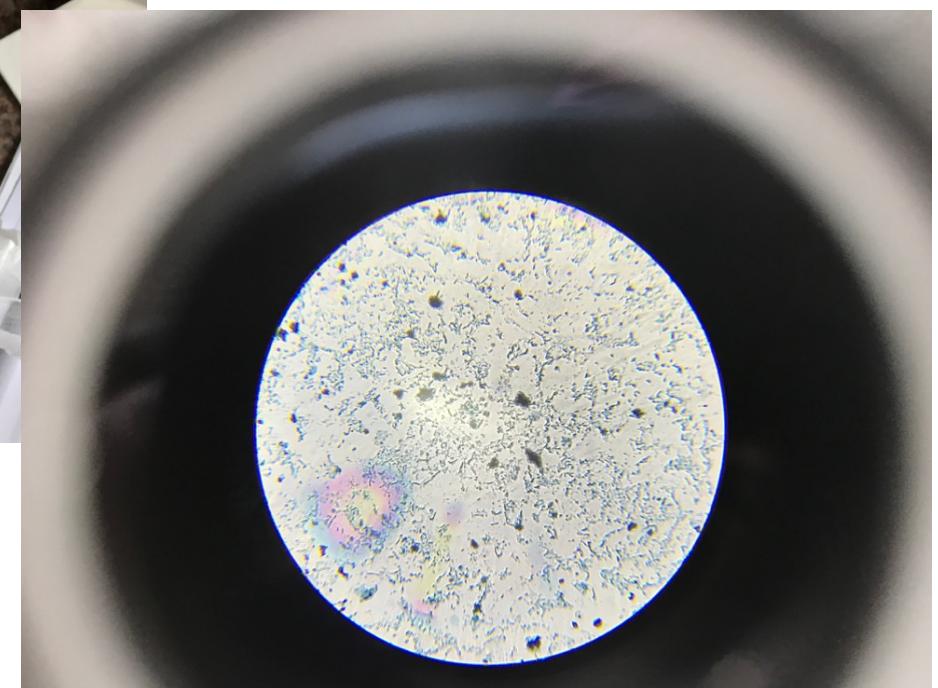
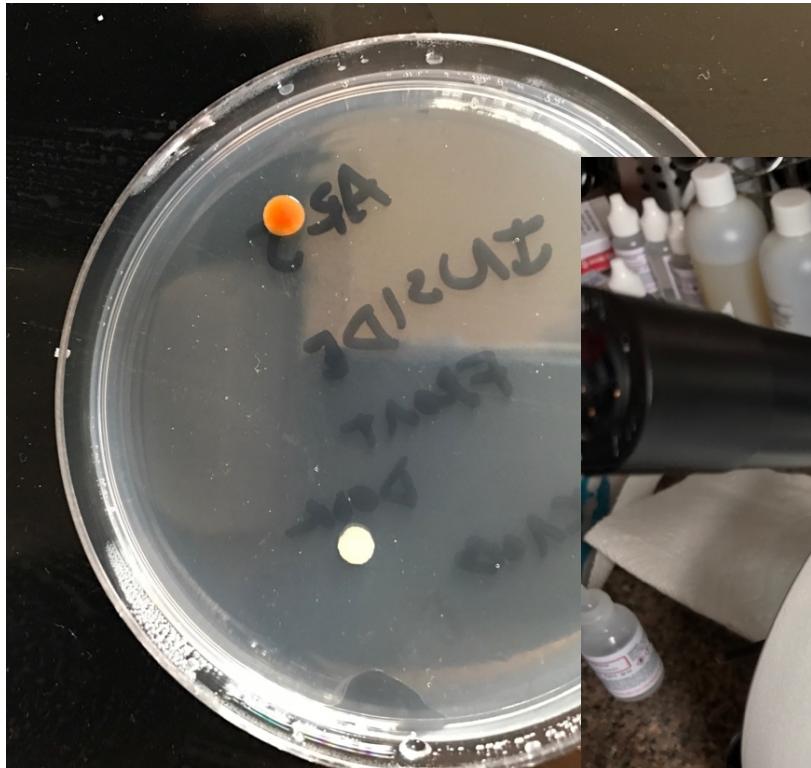
**Total mass:**  
**~3.3 lbs**

# Okay, maybe not 10x more cells but still a lot! 😊



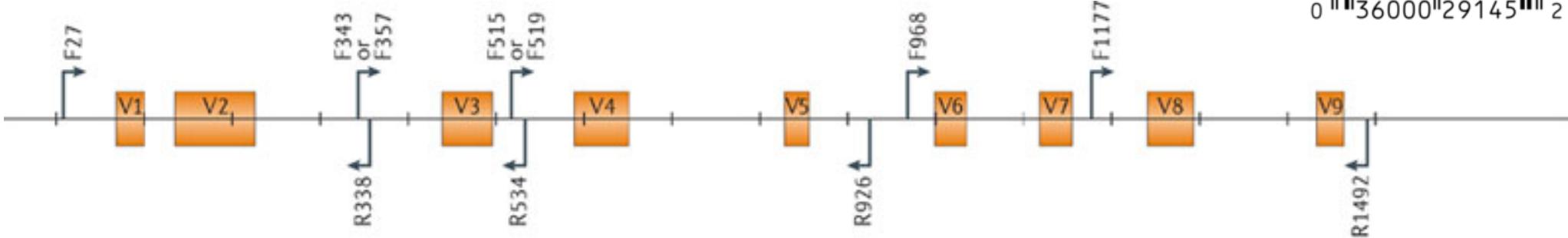
population segment	body weight [kg]	age [y]	blood volume [L]	RBC count [10 <sup>12</sup> /L]	colon content [g]	bac. conc. [10 <sup>11</sup> /g wet] <sup>(1)</sup>	total human cells [10 <sup>12</sup> ] <sup>(2)</sup>	total bacteria [10 <sup>12</sup> ]	B:H
ref. man	70	20–30	4.9	5.0	420	0.92	30	38	1.3
ref. woman	63		3.9	4.5	480	0.92	21	44	2.2
young infant	4.4	4 weeks	0.4	3.8	48	0.92	1.9	4.4	2.3
infant	9.6	1	0.8	4.5	80	0.92	4	7	1.7
elder	70	66	3.8 <sup>(3)</sup>	4.8	420	0.92	22	38	1.8
obese	140		6.7	5.0 <sup>(4)</sup>	610 <sup>(5)</sup>	0.92	40	56	1.4

# Pre-PCR: Gram-Staining



Gram staining differentiates bacteria by the chemical and physical properties of their cell walls by detecting peptidoglycan, which is present in the cell wall of Gram-positive bacteria

# 16S rRNA



**The 16S rRNA gene is a section of prokaryotic DNA found in all bacteria and archaea. This gene codes for an rRNA, and this rRNA in turn makes up part of the ribosome.**

**The 16S rRNA gene is a commonly used tool for identifying bacteria for several reasons.** First, traditional characterization depended upon phenotypic traits like gram positive or gram negative, bacillus or coccus, etc. Taxonomists today consider analysis of an organism's DNA more reliable than classification based solely on phenotypes. Secondly, researchers may, for a number of reasons, want to identify or classify only the bacteria within a given environmental or medical sample. Thirdly, the 16S rRNA gene is relatively short at 1.5 kb, making it faster and cheaper to sequence than many other unique bacterial genes.



# Rapid determination of 16S ribosomal RNA sequences for phylogenetic analyses

(reverse transcriptase/dideoxynucleotide)

DAVID J. LANE\*, BERNADETTE PACE\*, GARY J. OLSEN\*, DAVID A. STAHL†‡, MITCHELL L. SOGIN†,  
AND NORMAN R. PACE\*§

\*Department of Biology and Institute for Molecular and Cellular Biology, Indiana University, Bloomington, IN 47405; and †Department of Molecular and Cellular Biology, National Jewish Hospital and Research Center, Denver, CO 80206

Communicated by Ralph S. Wolfe, June 26, 1985

**ABSTRACT** Although the applicability of small subunit ribosomal RNA (16S rRNA) sequences for bacterial classification is now well accepted, the general use of these molecules has been hindered by the technical difficulty of obtaining their sequences. A protocol is described for rapidly generating large blocks of 16S rRNA sequence data without isolation of the 16S rRNA or cloning of its gene. The 16S rRNA in bulk cellular RNA preparations is selectively targeted for dideoxynucleotide-terminated sequencing by using reverse transcriptase and synthetic oligodeoxynucleotide primers complementary to universally conserved 16S rRNA sequences. Three particularly useful priming sites, which provide access to the three major 16S rRNA structural domains, routinely yield 800–1000 nucleotides of 16S rRNA sequence. The method is evaluated with respect to accuracy, sensitivity to modified nucleotides in the template RNA, and phylogenetic usefulness, by examination of several 16S rRNAs whose gene sequences are known. The relative simplicity of this approach should facilitate a rapid expansion of the 16S rRNA sequence collection available for phylogenetic analyses.

described here rapidly provides partial sequences of 16S rRNA that are useful for phylogenetic analysis.

## MATERIALS AND METHODS

**Purification of RNA Templates.** Bulk, cellular RNA was purified by phenol extraction of French pressure cell lysates as detailed by Pace *et al.* (6), except that ribosomes were not pelleted before extraction. High molecular weight RNA was then prepared by precipitation with 2 M NaCl (6). Although not essential, NaCl precipitation of the RNA generally increased the amount of legible sequence data and reduced backgrounds on gels, presumably by eliminating fragmented DNA from the reactions. RNA was stored at 2 mg/ml in 10 mM Tris·HCl (pH 7.4) at –20°C.

**Oligodeoxynucleotide Primers.** Oligodeoxynucleotide primers were synthesized manually by using the appropriate blocked and protected nucleoside diisopropylphosphoramidites and established coupling protocols (7). Deblocked products were purified by polyacrylamide gel electrophore-

## Box 1 | Species definitions and concepts in microbiology

### Definitions

Microbes are currently assigned to a common species if their reciprocal, pairwise DNA re-association values are  $\geq 70\%$  in DNA–DNA hybridization experiments under standardized conditions and their  $\Delta T_m$  (melting temperature) is  $\leq 5^\circ\text{C}$ <sup>79</sup>. In addition, all strains within a species must possess a certain degree of phenotypic consistency, and species descriptions should be based on more than one type strain<sup>11</sup>. A species name is only assigned if its members can be distinguished from other species by at least one diagnostic phenotypic trait<sup>79</sup>. Microbes with 16S ribosomal RNAs (rRNAs) that are  $\leq 98.7\%$  identical are always members of different species, because such strong differences in rRNA correlate with  $<70\%$  DNA–DNA similarity<sup>80</sup>. However, the opposite is not necessarily true, and distinct species have been occasionally described with 16S rRNAs that are  $>98.7\%$  identical. Most uncultured microbes cannot be assigned to a classical species because we do not know their phenotype. In some cases, uncultured microbes can be assigned a provisional ‘*Candidatus*’ designation if their 16S rRNA sequences are sufficiently different from those of recognized species, if experimental *in situ* hybridization can be used to specifically detect them and if a basic description of their morphology and biology has been provided<sup>81</sup>.

## Box 1 | Species definitions and concepts in microbiology

### Definitions

Microbes are currently assigned to a common species if their reciprocal, pairwise DNA re-association values are  $\geq 70\%$  in DNA–DNA hybridization experiments under standardized conditions and their  $\Delta T_m$  (melting temperature) is  $\leq 5^\circ\text{C}$ <sup>79</sup>. In addition, all strains within a species must possess a certain degree of phenotypic consistency, and species descriptions should be based on more than one type strain<sup>11</sup>. A species name is only assigned

diagnostic pH  
 $\leq 98.7\%$  identical differences in rRNAs that are classical species microbes can sequences are in situ hybridized their morphology

### Concepts

Various concepts have been suggested for microbial species, but none have been generally accepted<sup>9</sup>. The following quotes represent several published concepts that were chosen to illustrate the lack of consensus:

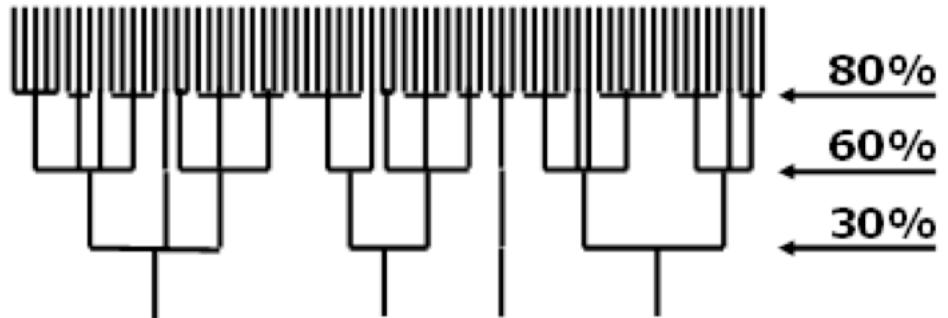
- “A species could be described as a monophyletic and genetically coherent cluster of individual organisms that show a high degree of overall similarity in many independent characteristics, and is diagnosable by a discriminative phenotypic property.” (REF. 9)
- “Species are considered to be an irreducible cluster of organisms diagnosably different from other such clusters and within which there is a parental pattern of ancestry and descent.” (REF. 82)
- “A species is a group of individuals where the observed lateral gene transfer within the group is much greater than the transfer between groups.” (REF. 83)
- “Microbes ... do not form natural clusters to which the term “species” can be universally and sensibly applied.” (REF. 84)
- “Species are (segments of) metapopulation lineages.” (REF. 7)

### Microbial diversity and the genetic nature of microbial species

Achtman & Wagner (2008) Nature Reviews Microbiology. doi:10.1038/nrmicro1872

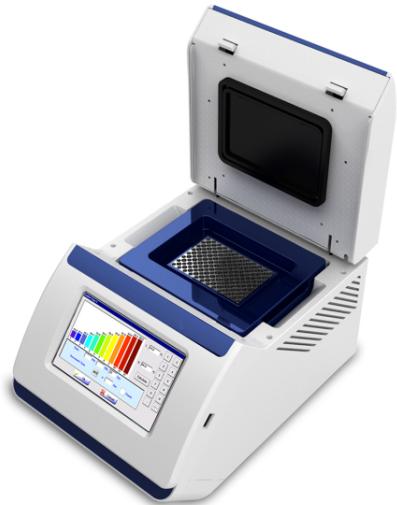
# Operational Taxonomic Units (OTUs)

***OTUs take the place of “species” in many microbiome diversity analyses because named species genomes are often unavailable for particular marker sequences.***



- Although much of the 16S rRNA gene is highly conserved, several of the sequenced regions are variable or hypervariable, so small numbers of base pairs can change in a very short period of evolutionary time.
- Because 16S regions are typically sequenced using only a single pass, there is a fair chance that they will thus contain at least one sequencing error. This means that requiring tags to be 100% identical will be extremely conservative and treat essentially clonal genomes as different organisms.
- Some degree of sequence divergence is typically allowed - 95%, 97%, or 99% are sequence similarity cutoffs often used in practice [18] - and the resulting cluster of nearly-identical tags (and thus assumedly identical genomes) is referred to as an Operational Taxonomic Unit (OTU) or sometimes phylotype.

# 16S versus shotgun NGS



**16S**

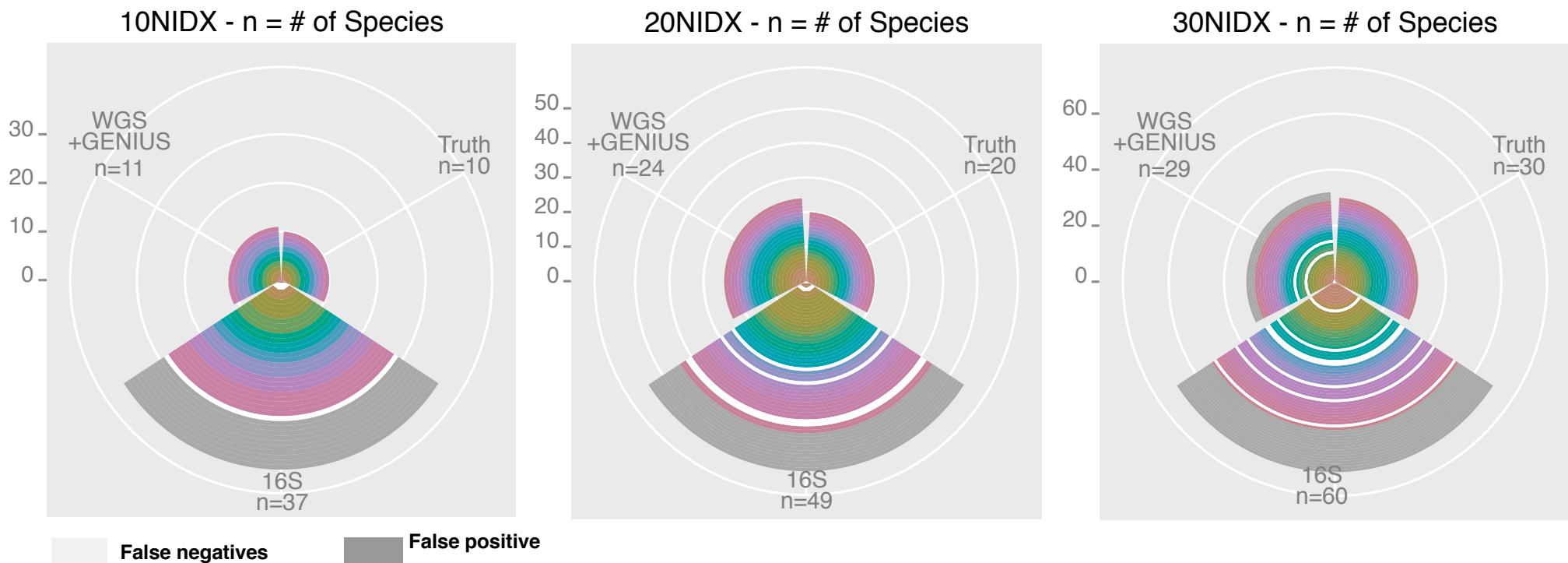
Fast (minutes – hours)  
Directed analysis  
Cheap per sample  
Family/Genus Identification

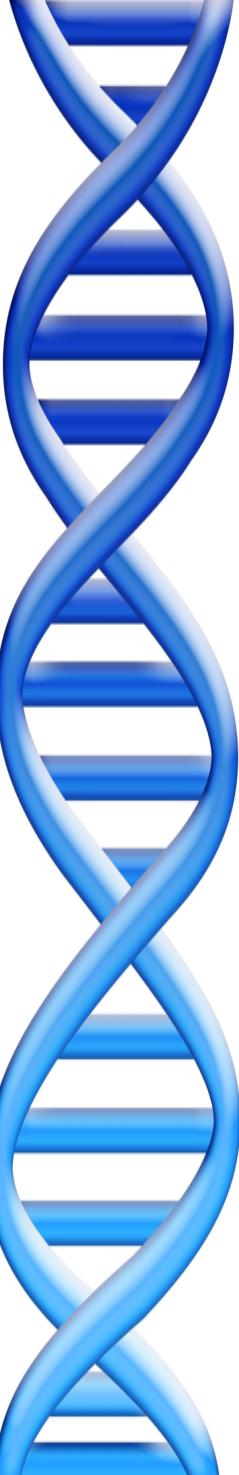


**NGS**

Slower (hours to days)  
Whole Metagenome  
More expensive per sample  
Species/Strain Identification  
Genes presence/absence  
Variant analysis  
Eukaryotic hosts  
Can ID fungi, viruses, etc.

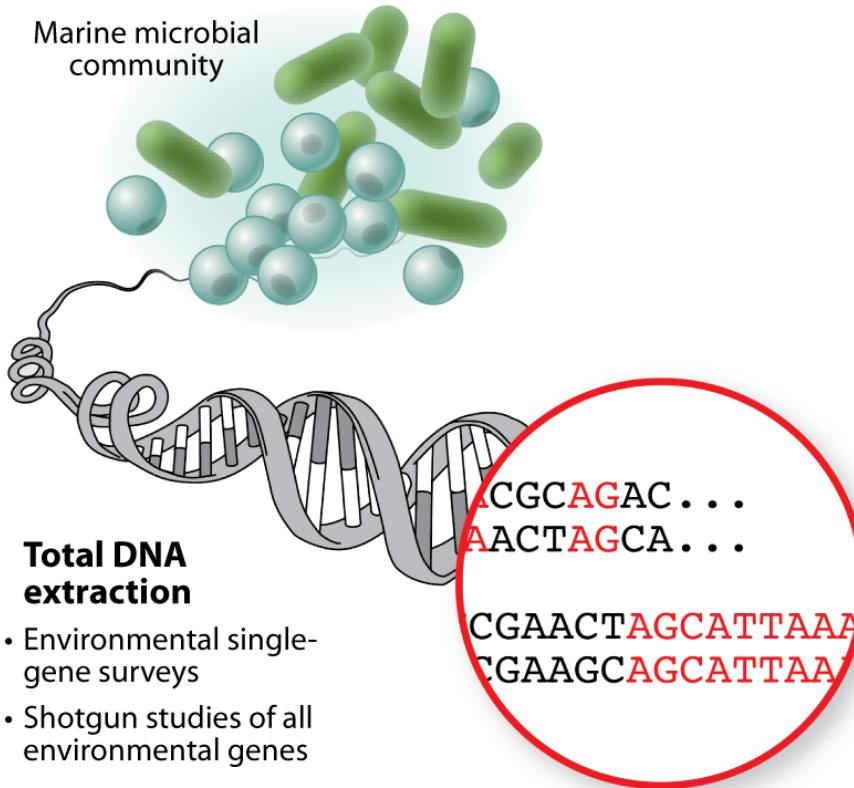
# 16S Overestimates Diversity





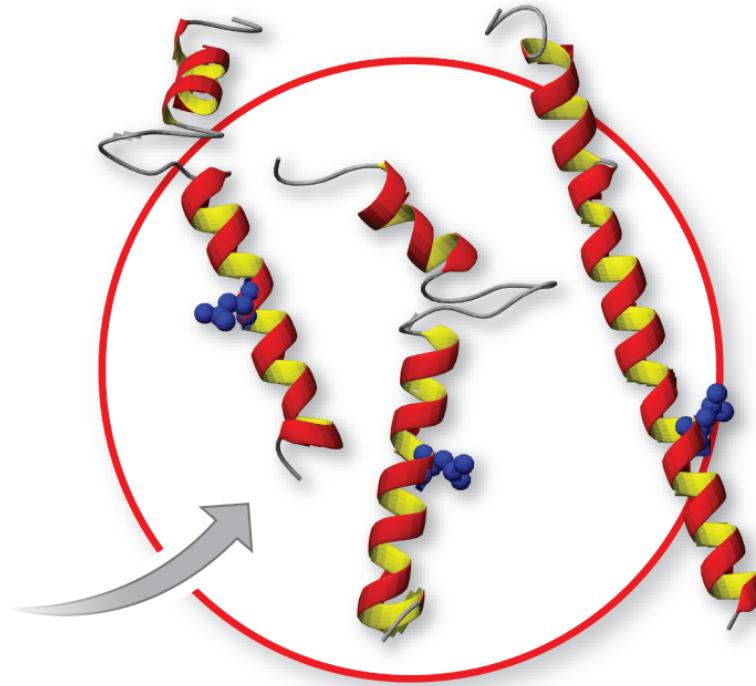
## Part II: Methods

# Sequencing Based Analysis



## DNA sequencing

- Identify common genes within a community
- Identify genome contents favored by current environmental conditions



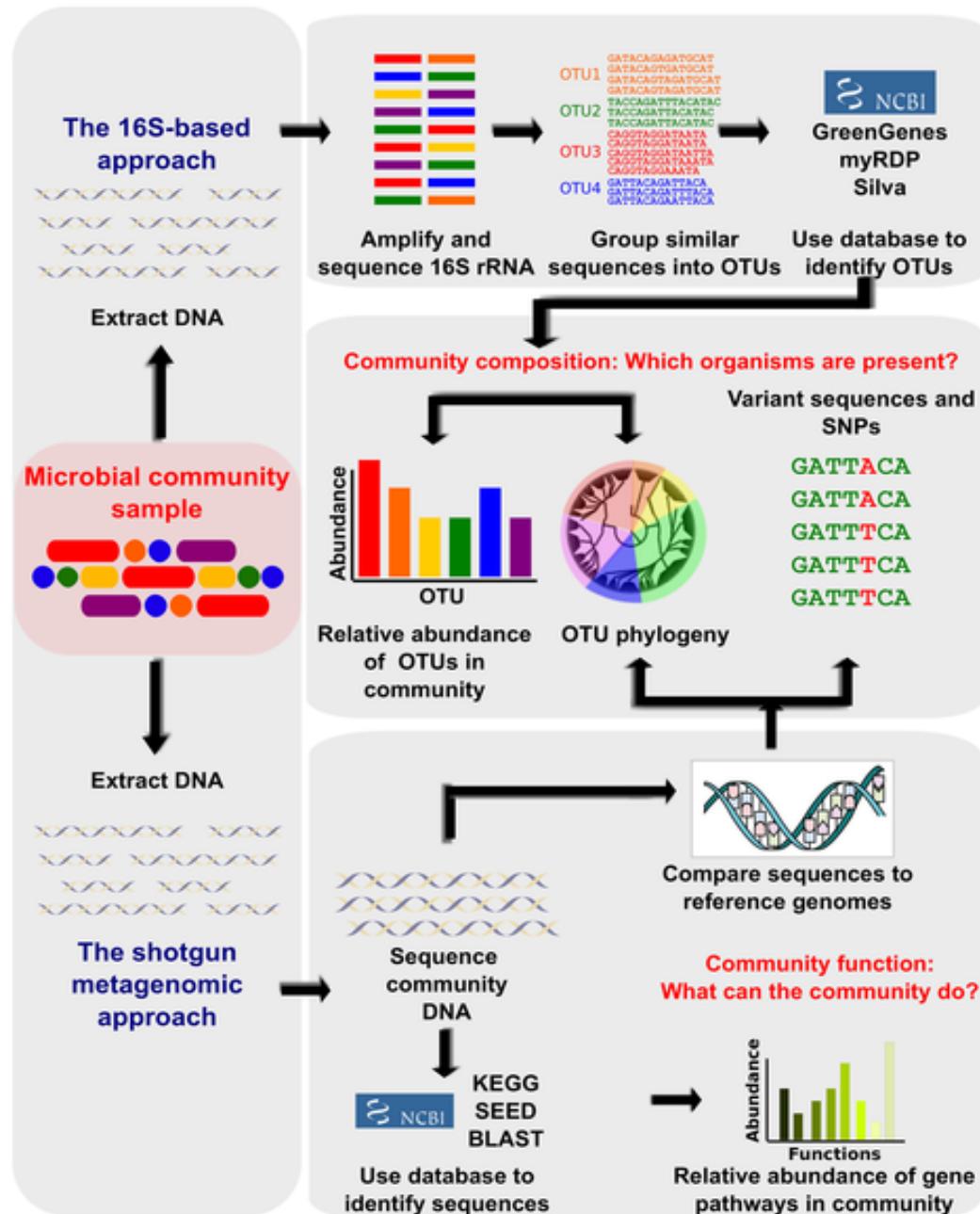
## Protein annotation

Use metagenomics studies as a tool to answer broader ecological or evolutionary questions

Also can do host DNA suppression/  
microbial enrichment!



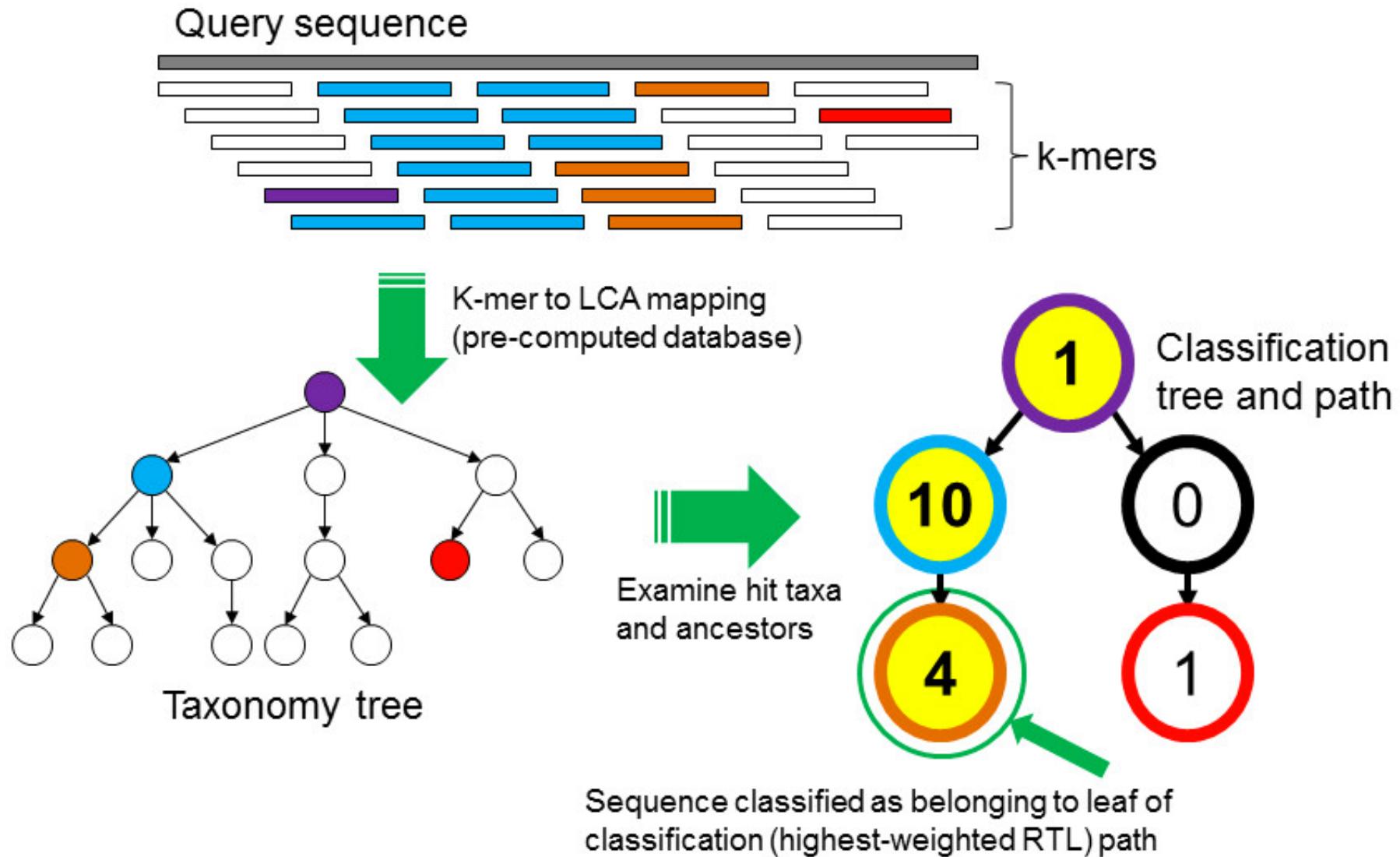
Gilbert JA, Dupont CL. 2011.  
Annu. Rev. Mar. Sci. 3:347–71



## Chapter 12: Human Microbiome Analysis

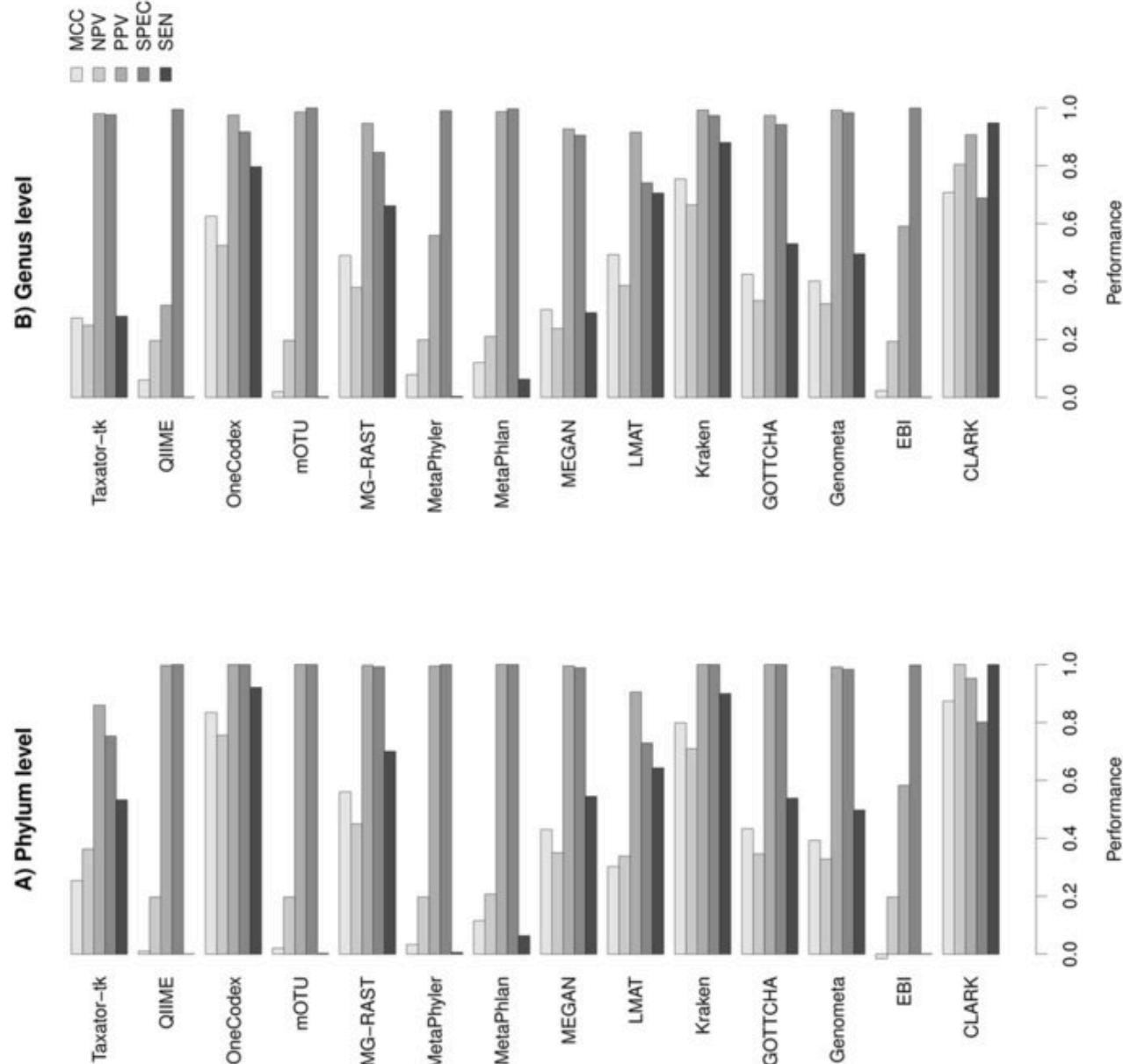
Morgan & Huttenhower (2012) PLOS Comp Bio. <https://doi.org/10.1371/journal.pcbi.1002808>

# Kraken



**Kraken: ultrafast metagenomic sequence classification using exact alignments**  
Wood and Salzberg (2014) Genome Biology. DOI: 10.1186/gb-2014-15-3-r46

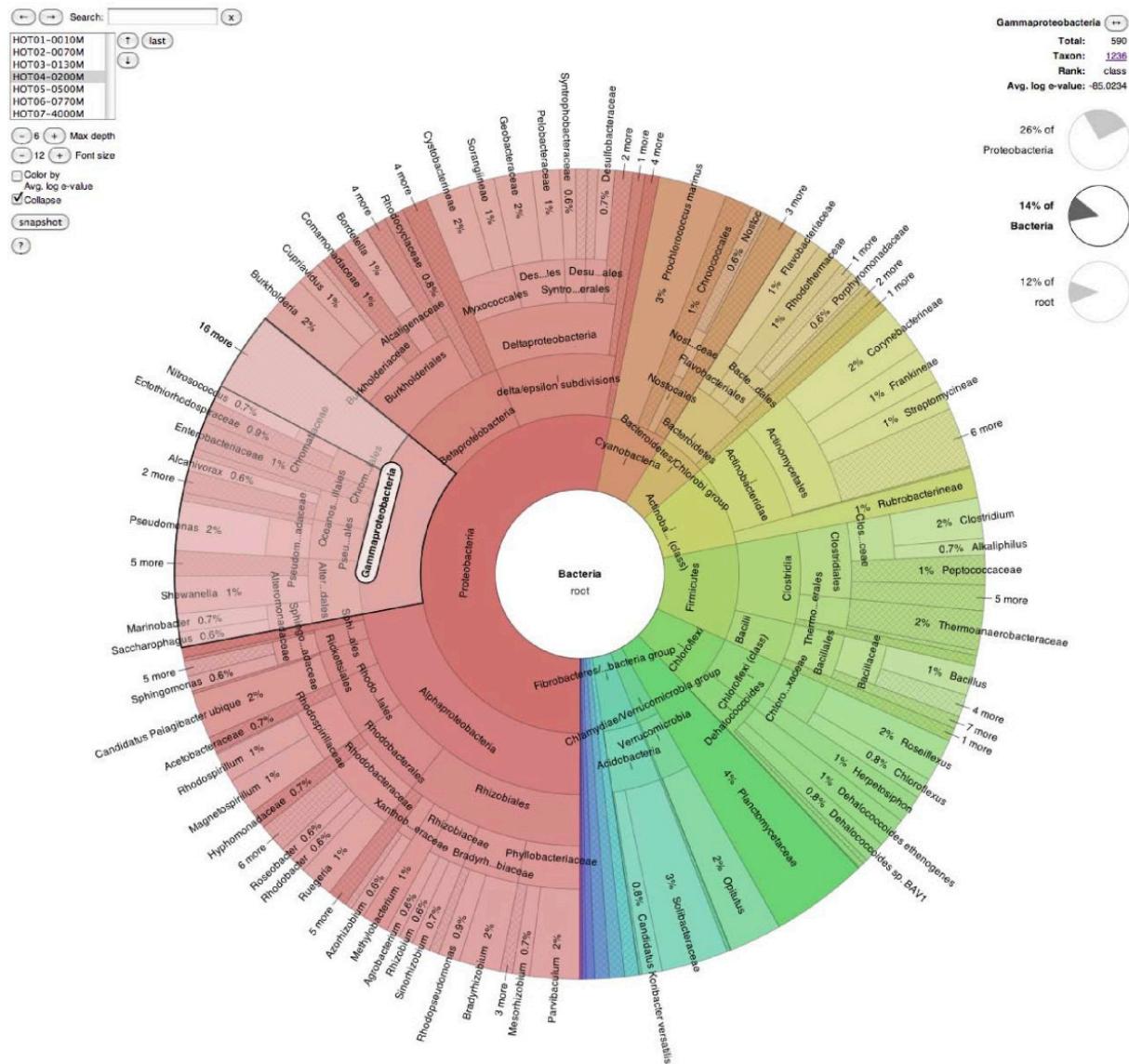
# Metagenomics Benchmarking



**An evaluation of the accuracy and speed of metagenome analysis tools**

Lindgreen et al (2016) Scientific Reports. doi:10.1038/srep19233

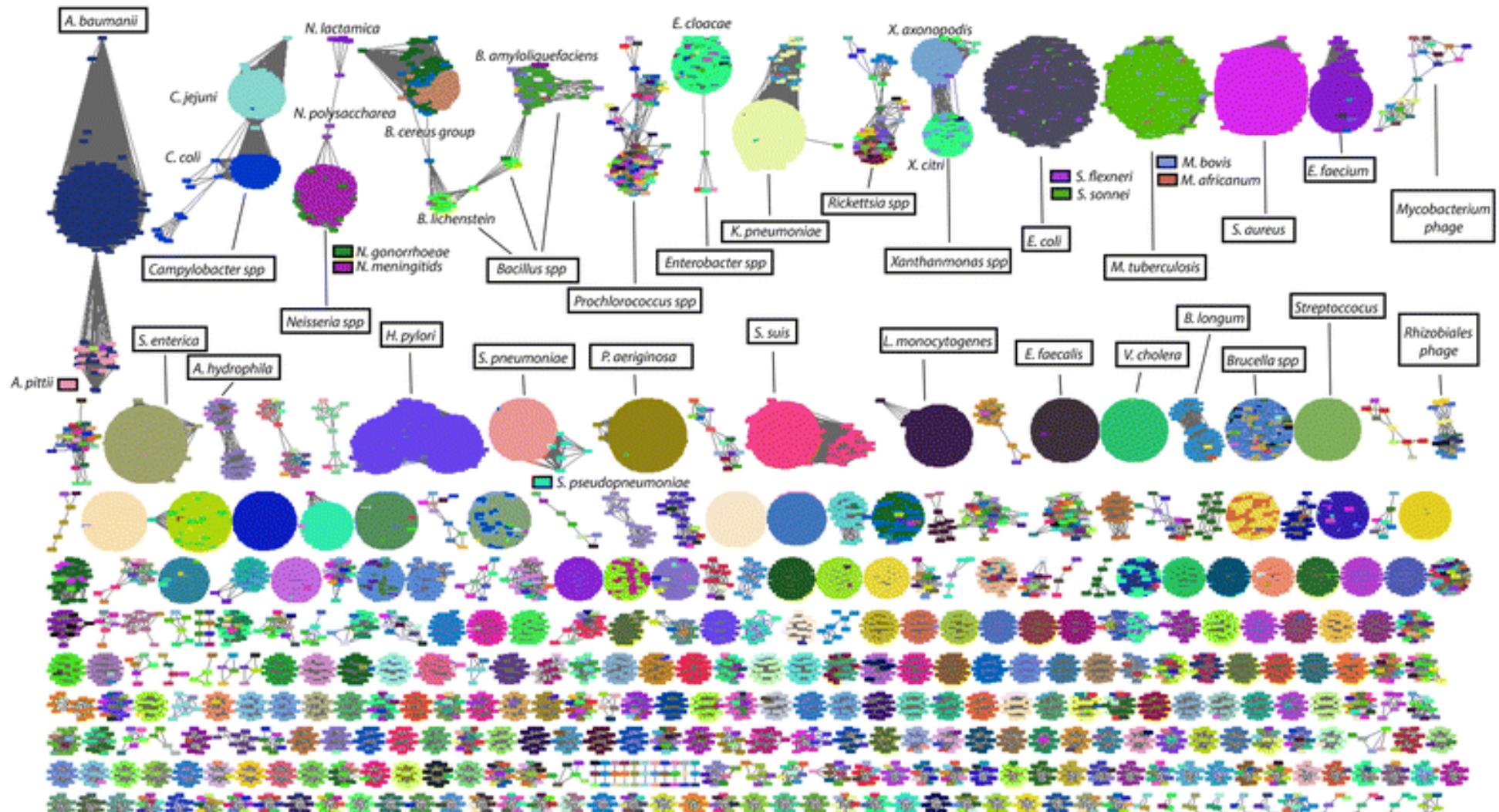
# Krona Plots



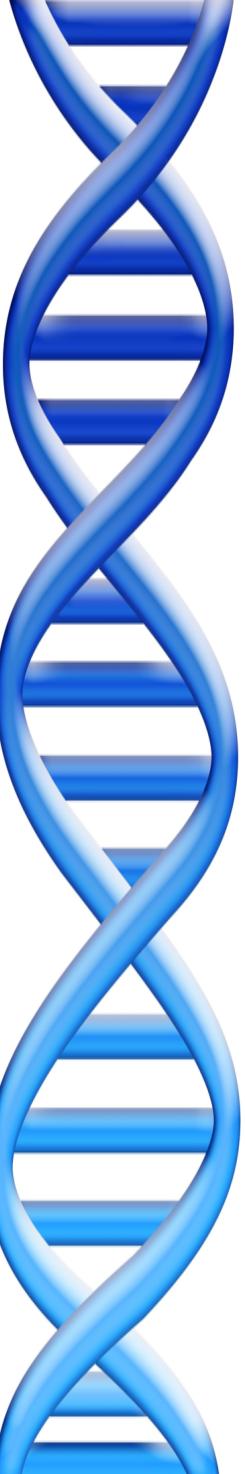
**Interactive metagenomic visualization in a Web browser**

Ondov et al (2011) BMC Bioinformatics. DOI: 10.1186/1471-2105-12-385

# Min-Hash: Comparing all 54,118 RefSeq genomes in 1 day on a laptop

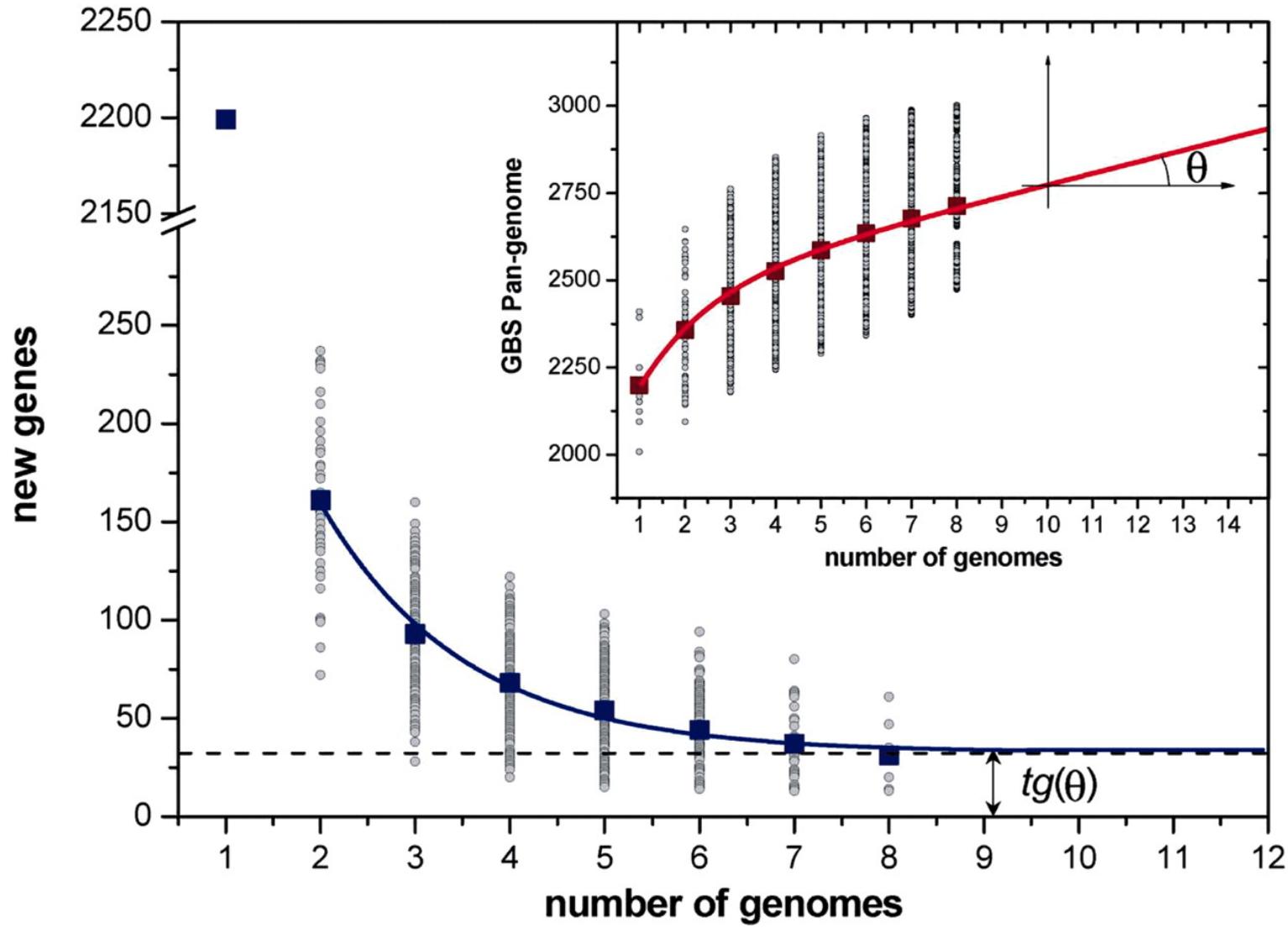


**Mash: fast genome and metagenome distance estimation using MinHash**  
Ondov et al. (2016) Genome Biology. DOI: 10.1186/s13059-016-0997-x



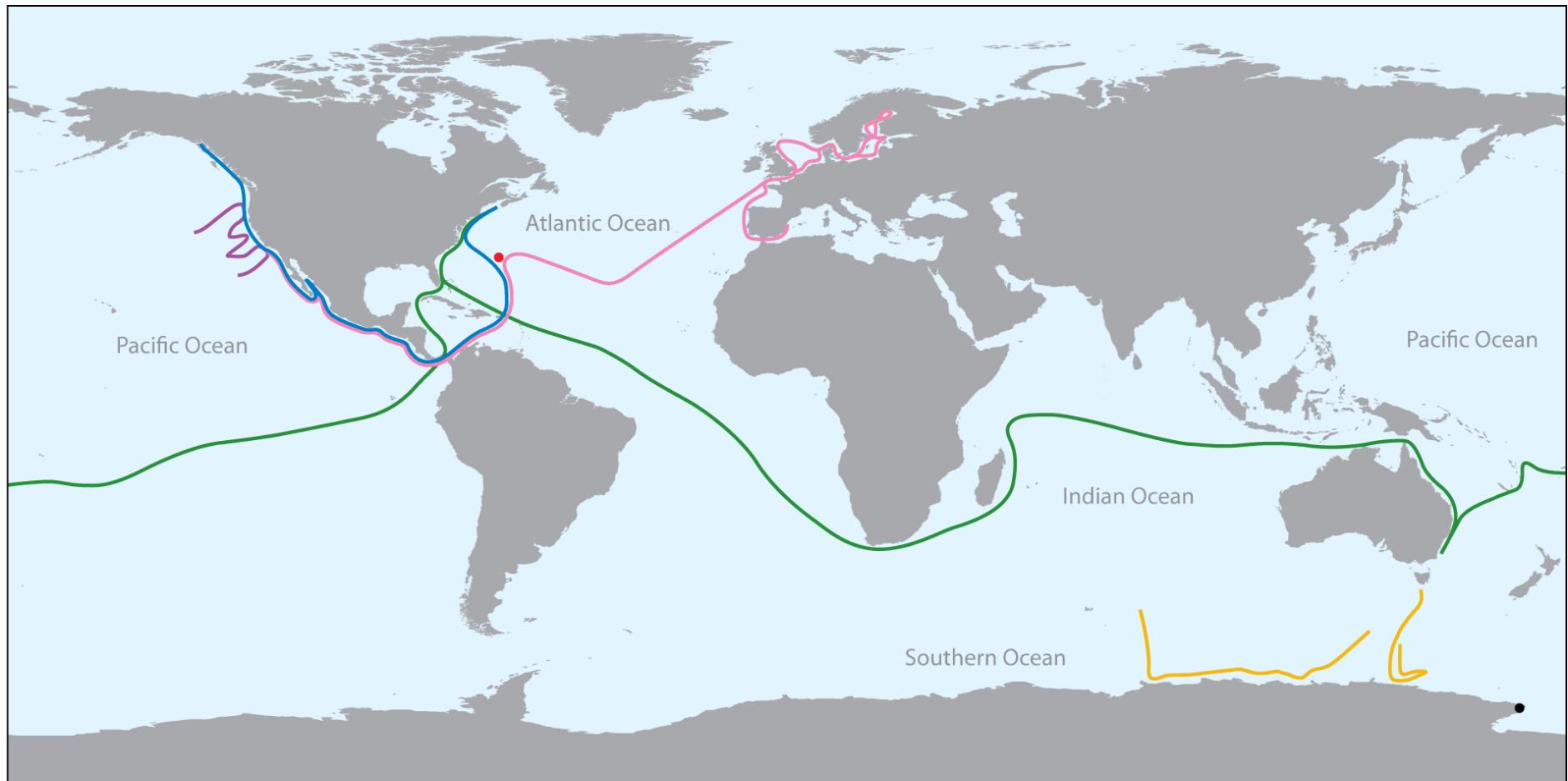
## Part III: Results

# Pan genome of *Streptococcus agalactiae*



Hervé Tettelin et al. PNAS 2005;102:13950-13955

# Global Ocean Survey



- 2003 Sargasso Sea pilot study
- 2003–2006 circumnavigation
- 2006–2007 Antarctica cruises
- 2007 east-to-west coast USA
- 2007 collaborative cruises
- 2009 Antarctica sea ice and water samples
- 2009–2010 Europe expedition

# Global Ocean Survey

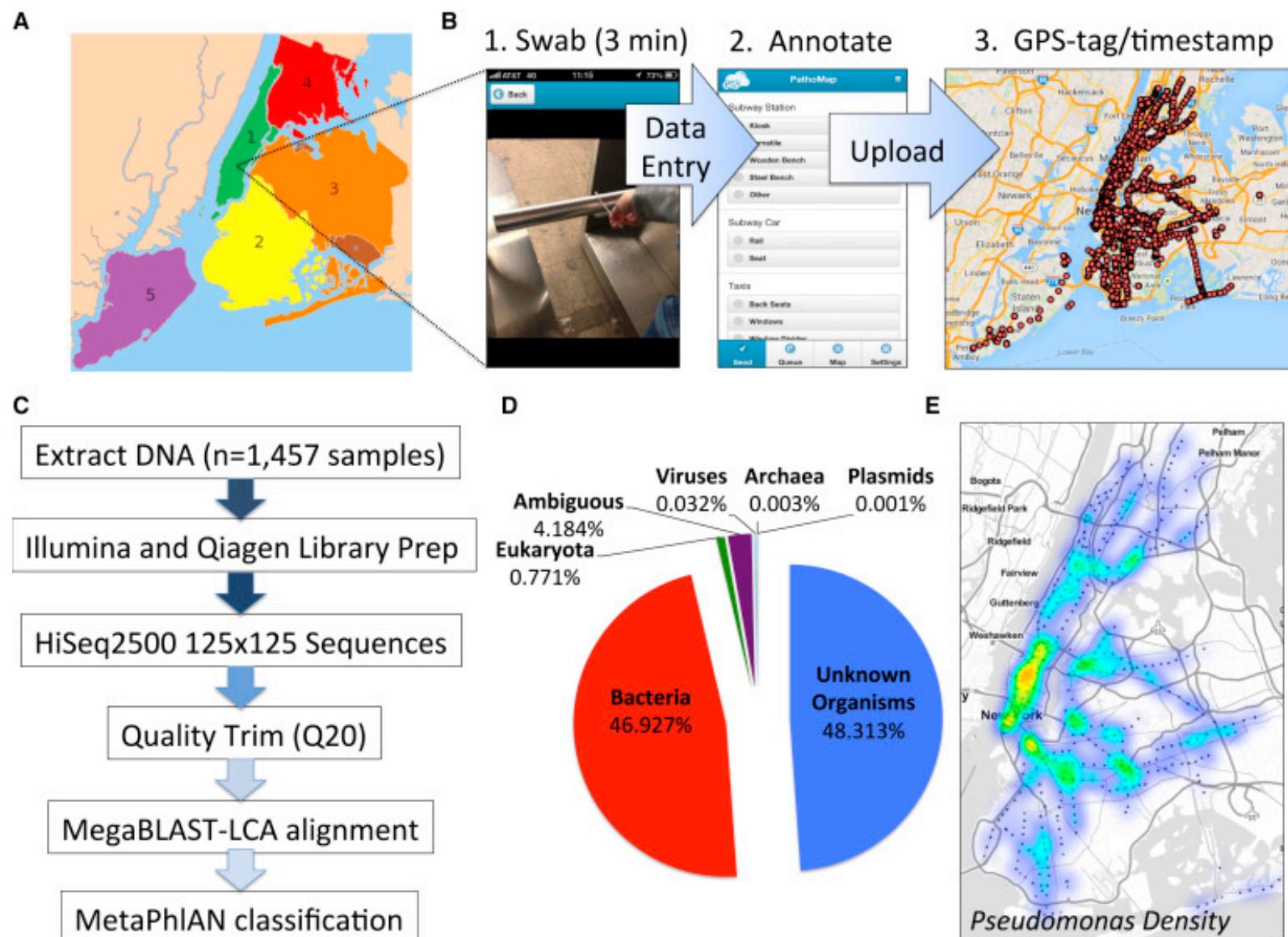


The combined set of predicted proteins in NCBI-nr, PG, TGI-EST, and ENS, as expected, has a lot of redundancy. For instance, most of the PG protein predictions are in NCBI-nr. Removing exact substrings of longer sequences (i.e., 100% identity) reduces this combined set to 3,167,979 predicted proteins. When we perform the same filtering on the GOS dataset, 5,654,638 predicted proteins remain.

***Thus, the GOS-predicted protein set is 1.8 times the size of the predicted protein set from current publicly available datasets.***

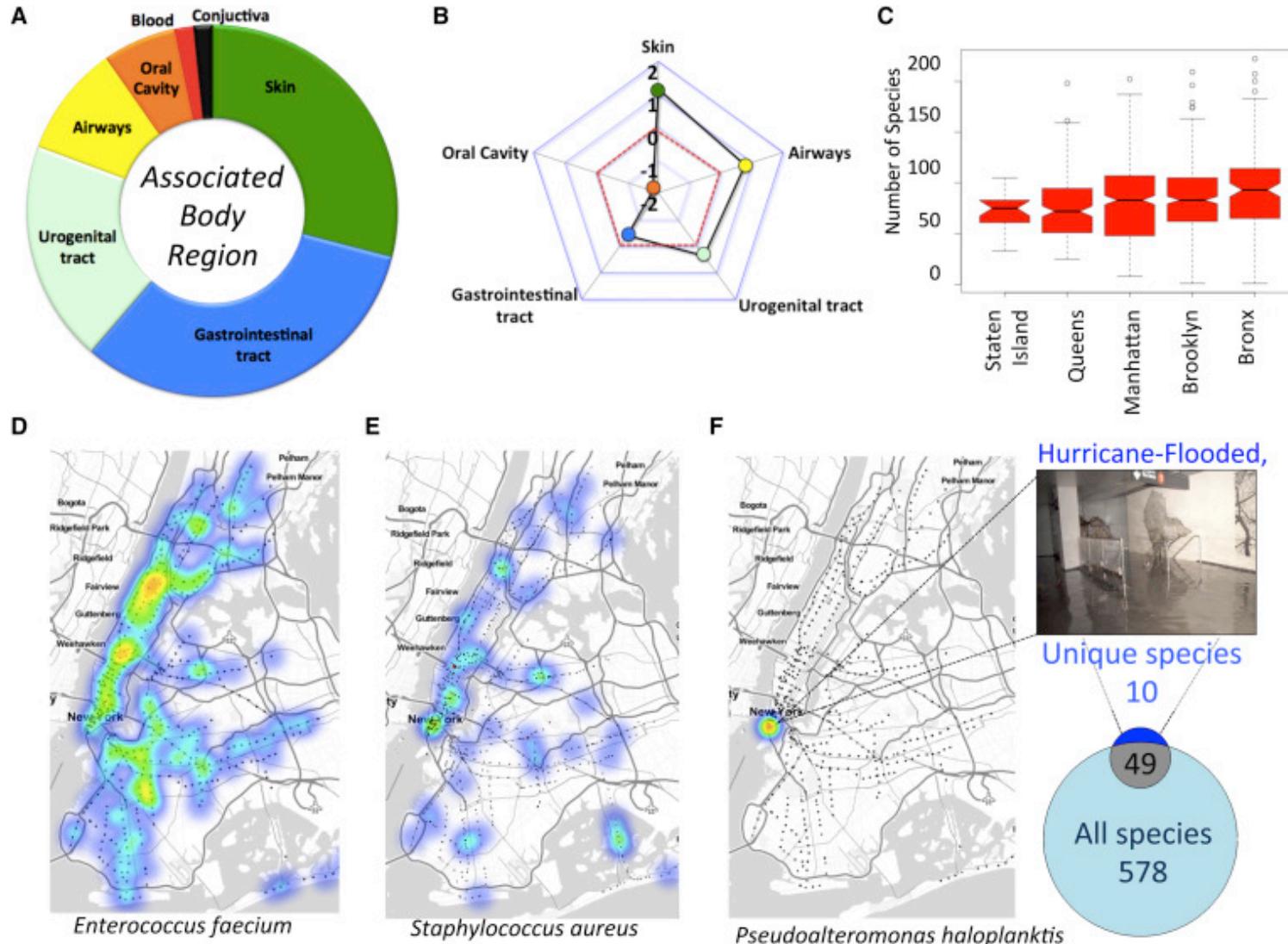
- 2003 Sargasso Sea pilot study
- 2003–2006 circumnavigation
- 2006–2007 Antarctica cruises
- 2007 east-to-west coast USA
- 2007 collaborative cruises
- 2009 Antarctica sea ice and water samples
- 2009–2010 Europe expedition

# Metasub



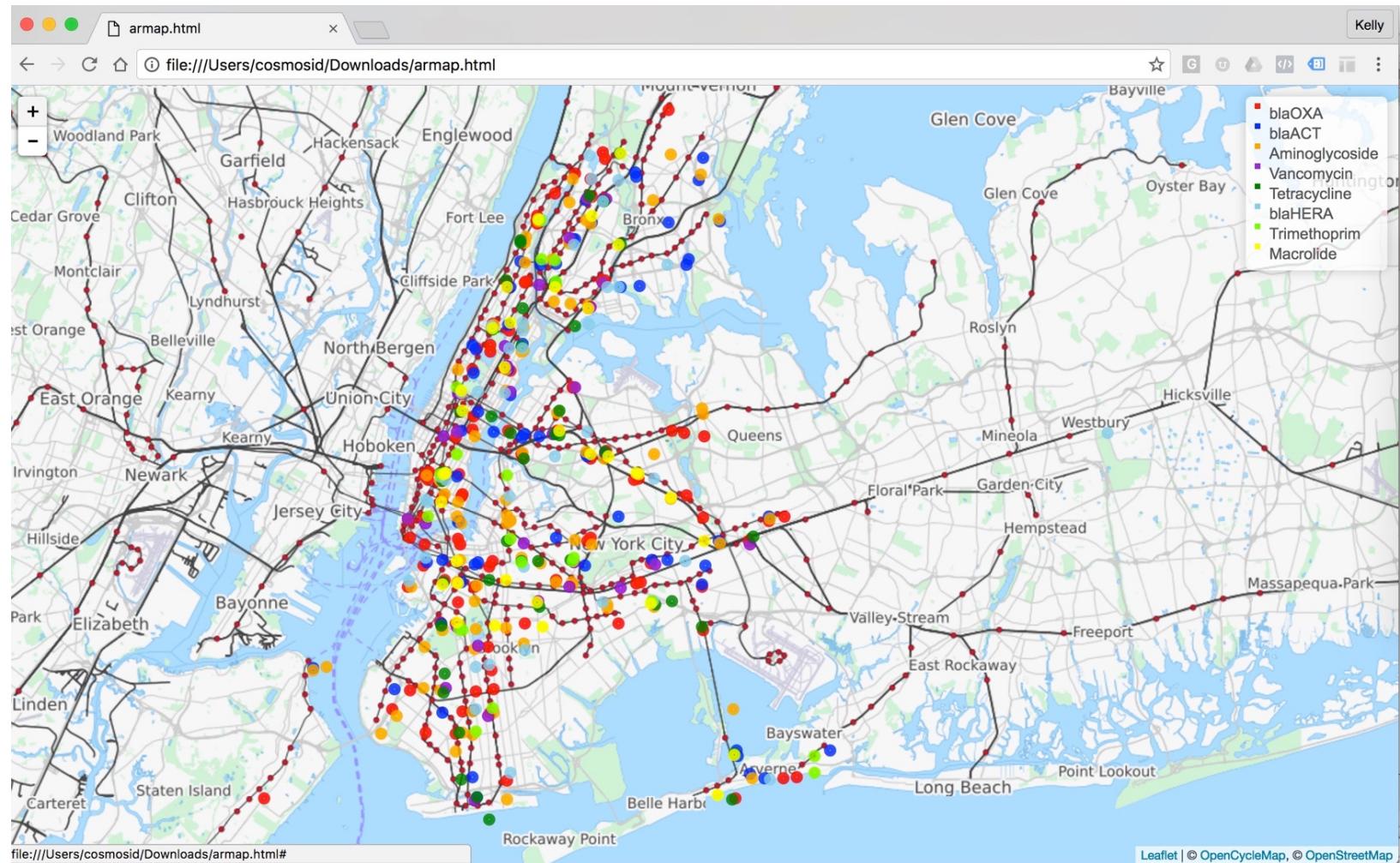
**Geospatial Resolution of Human and Bacterial Diversity with City-Scale Metagenomics**  
Afshinnekoo et al (2016) Cell Systems. <http://dx.doi.org/10.1016/j.cels.2015.01.001>

# Different subway stations resembled different body sites



**Geospatial Resolution of Human and Bacterial Diversity with City-Scale Metagenomics**  
Afshinnekoo et al (2016) Cell Systems. <http://dx.doi.org/10.1016/j.cels.2015.01.001>

# Mapping Antimicrobial Resistance Factors: PathoMap



Antibiotic resistance genes that were found most frequently in samples were plotted on the map of New York City based on their origin.

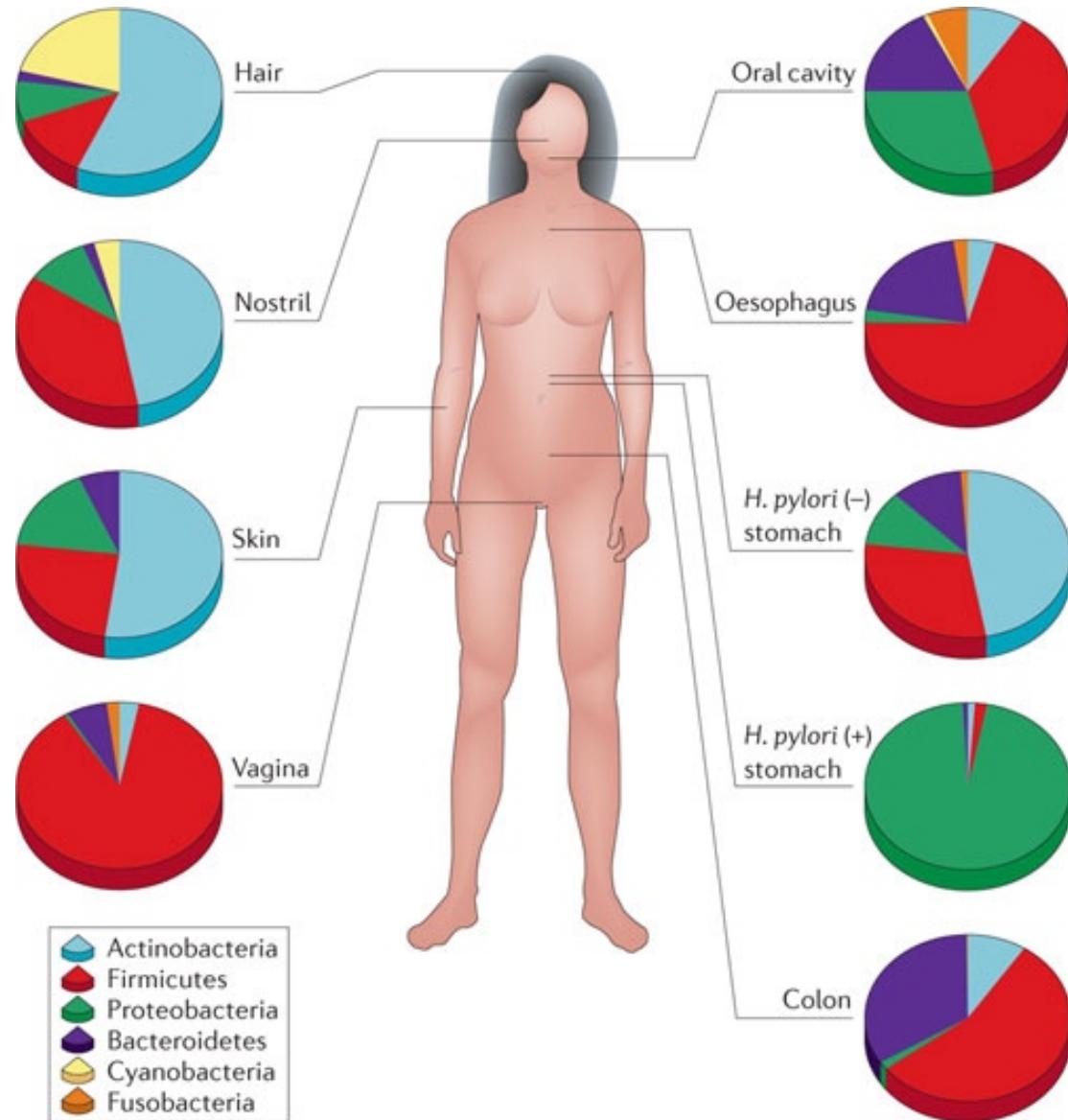
# Microbes and Human Health



**“MICROBE DIET** Mice fed microbes from obese people tend to gain fat. Microbes from lean people protect mice from excessive weight gain, even when animals eat a high-fat, low-fiber diet.”

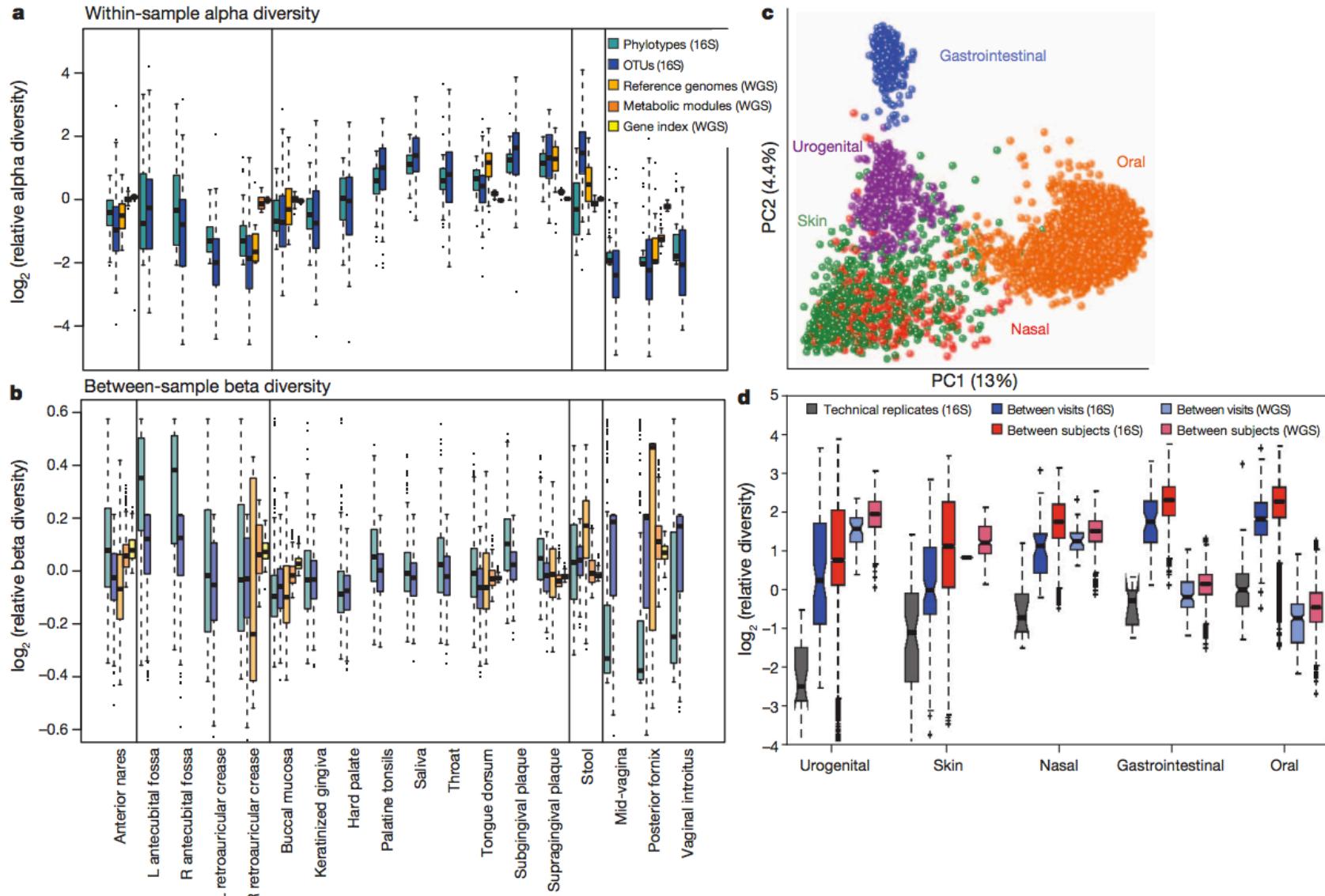
**Gut Microbiota from Twins Discordant for Obesity Modulate Metabolism in Mice**  
Ridaura et al (2013) Science. doi: 10.1126/science.1241214

# Microbes and Human Health



***The human microbiome: at the interface of health and disease***  
Cho & Blaser (2012) Nature Reviews Genetics. doi:10.1038/nrg3182

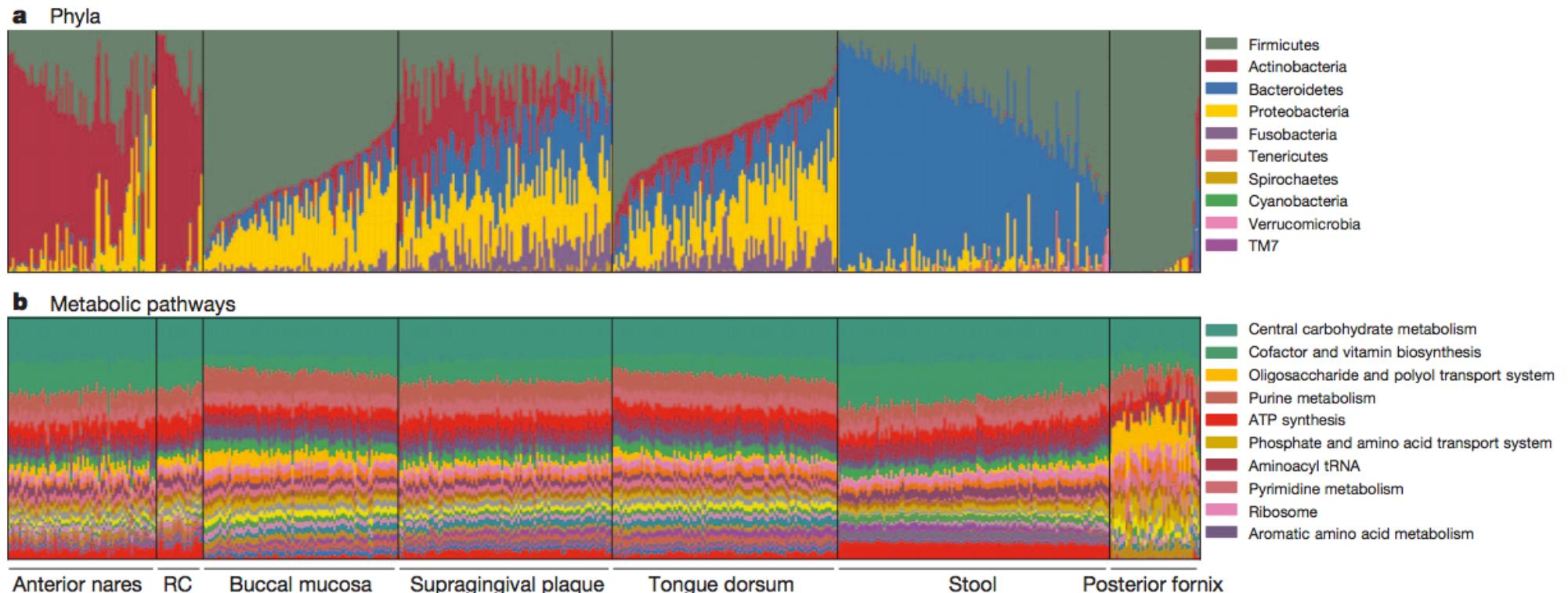
# Human Microbiome Project



**Structure, function and diversity of the healthy human microbiome**

The Human Microbiome Project Consortium (2012) Nature. doi:10.1038/nature11234

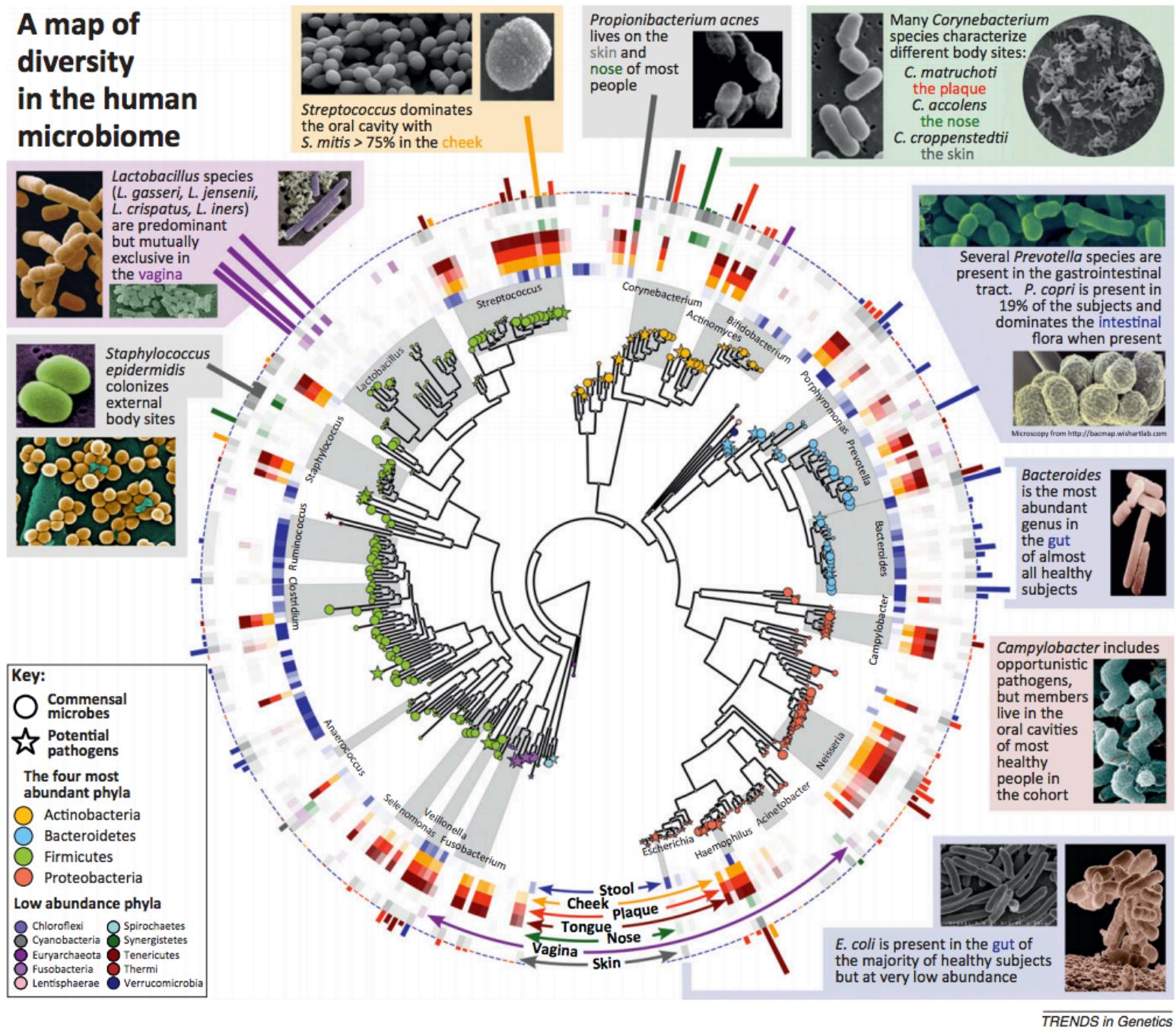
# Functional composition tends to be more stable than genome composition



**Structure, function and diversity of the healthy human microbiome**

The Human Microbiome Project Consortium (2012) Nature. doi:10.1038/nature11234

# A map of diversity in the human microbiome

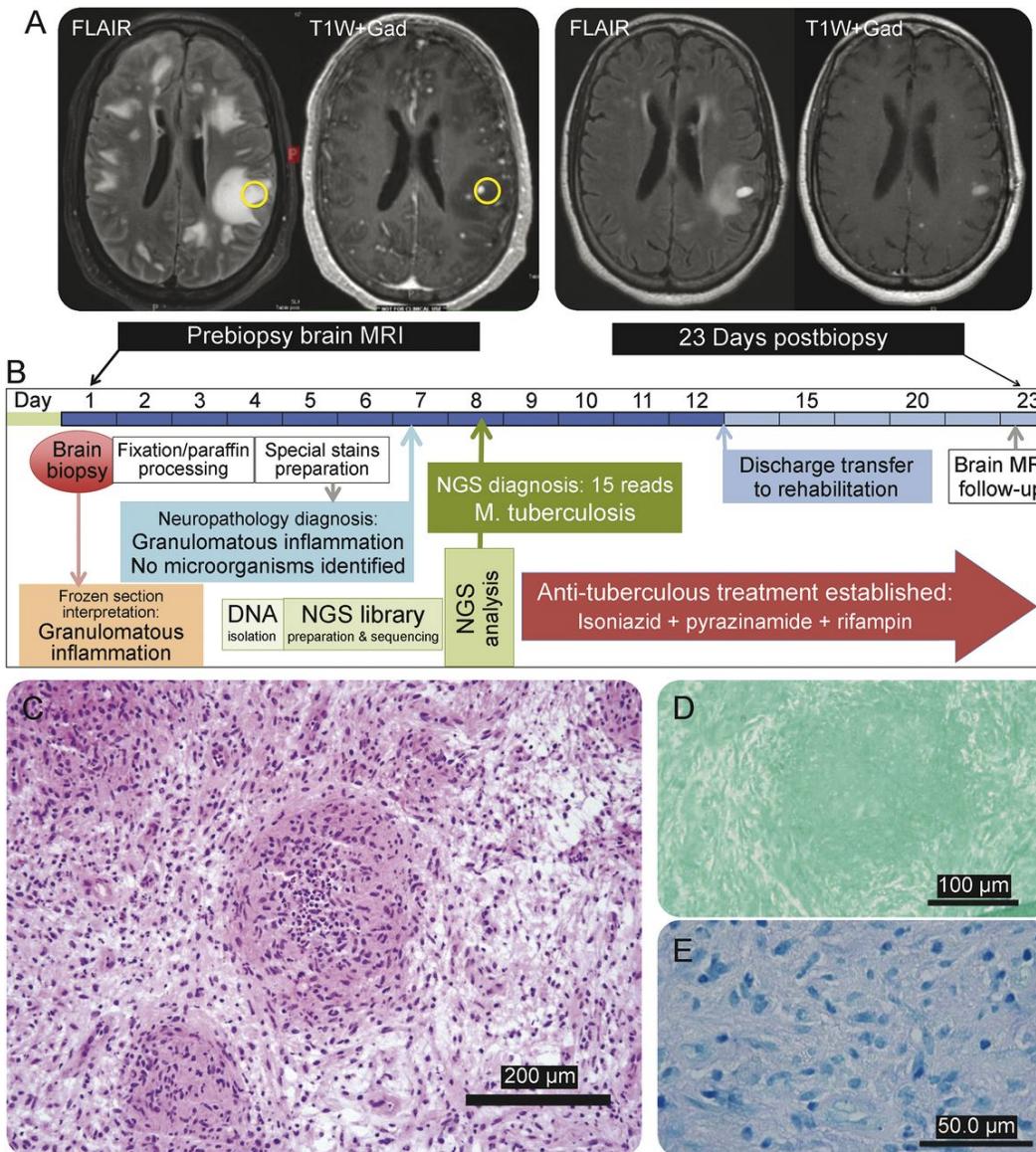


TRENDS in Genetics

## Biodiversity and functional genomics in the human microbiome

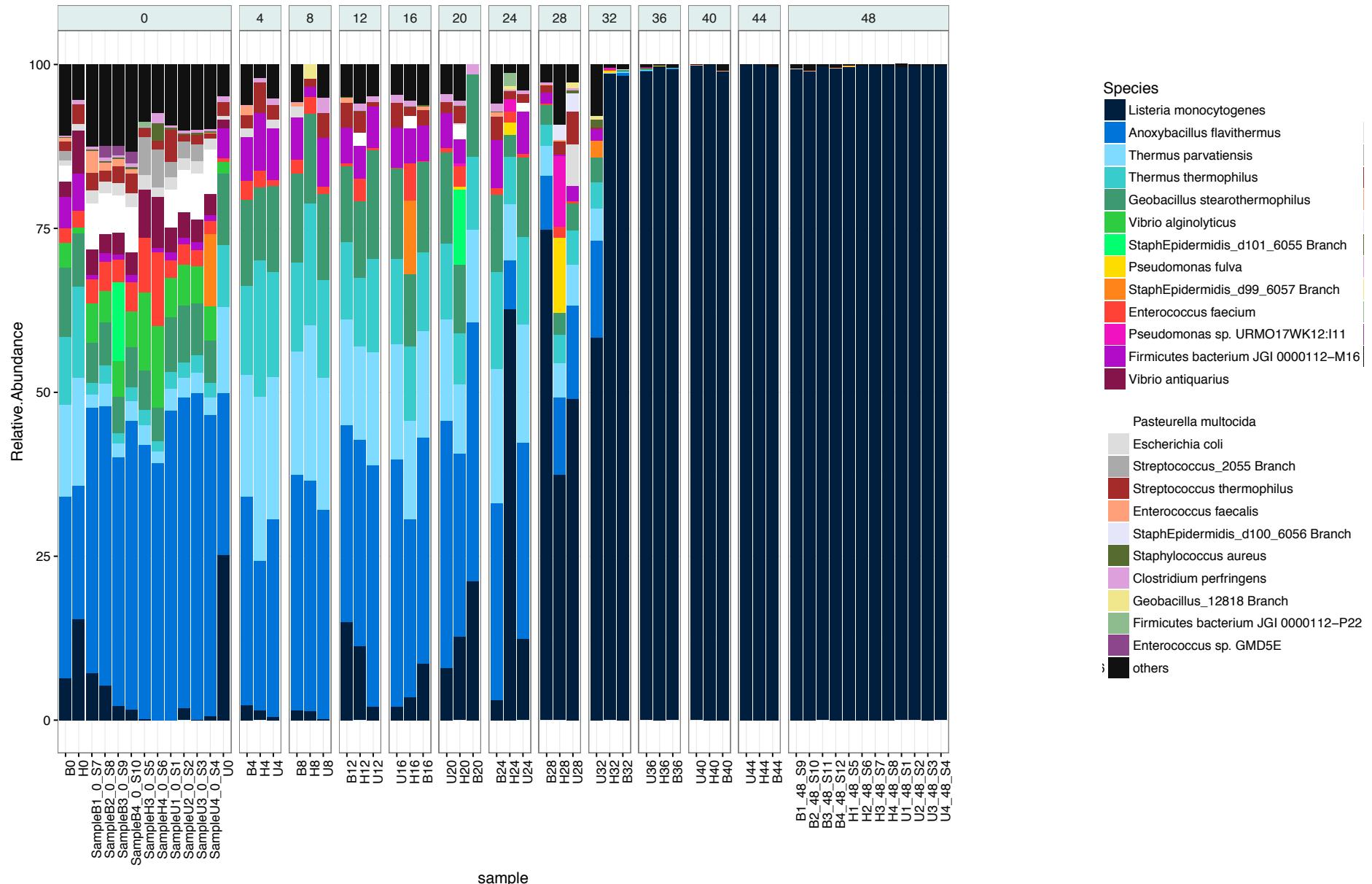
Morgan et al (2013) Trends in Genetics. <http://doi.org/10.1016/j.tig.2012.09.005>

# Diagnosing Brain Infections with NGS

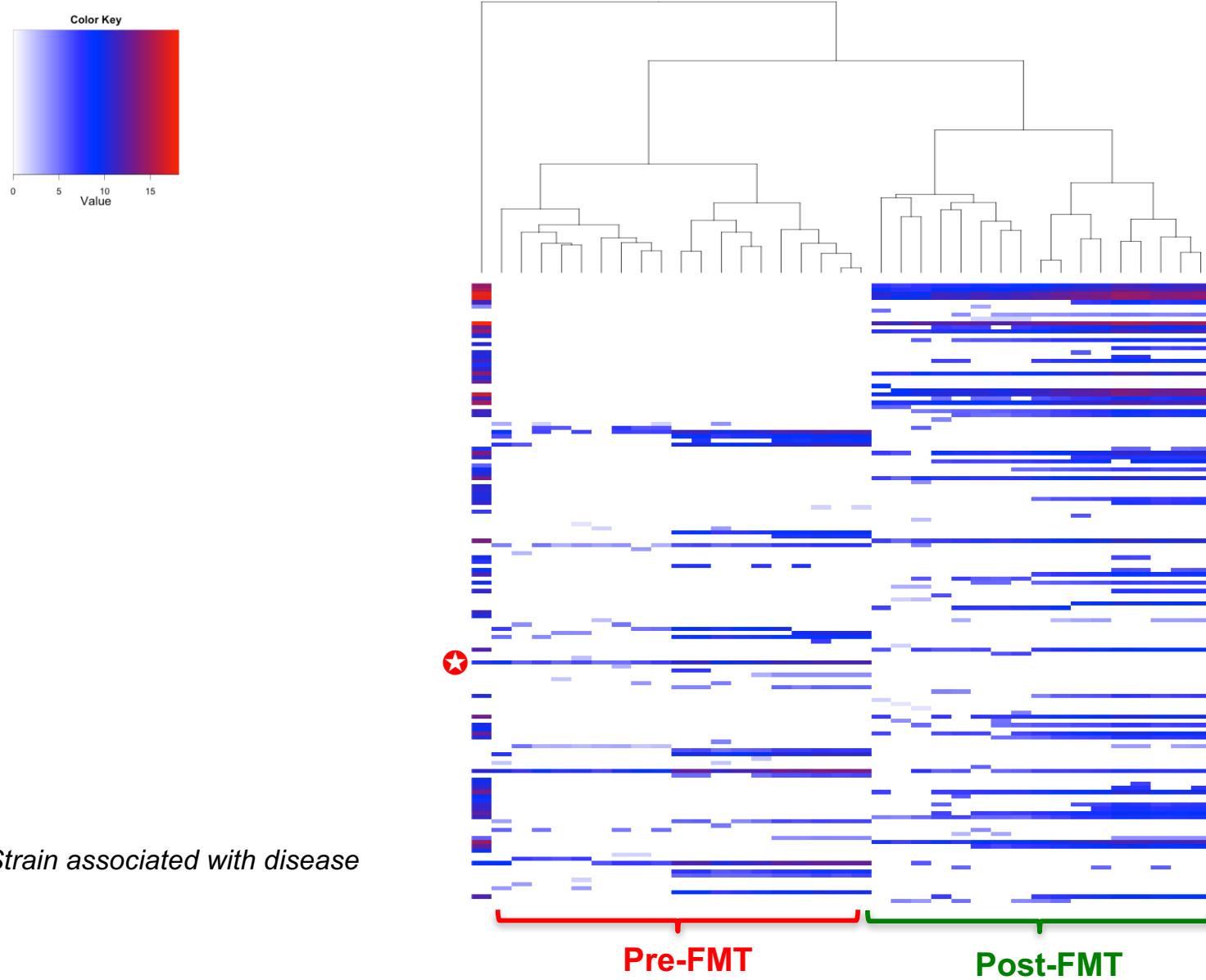


**Next-generation sequencing in neuropathologic diagnosis of infections of the nervous system**  
Salzberg et al (2016) Neurol Neuroimmunol Neuroinflamm dx. doi.org/10.1212/NXI.0000000000

# Listeria in ice cream



# Pre and Post Fecal Medical Transplant Microbial Shift



# National Institute of Cholera and Enteric Disease Kolkata, India



[www.niced.org.ind](http://www.niced.org.ind)

# Microbiome of Diarrheal Patients Compared with Healthy Individuals

Total # of NICED samples:	74
Indian Healthy Control ( <b>HC</b> ):	20
Sick with Unknown Etiology ( <b>UE</b> ):	28
Sick with Known Etiology ( <b>KE</b> ):	26
Healthy Human Microbiome Project ( <b>HMP</b> ):	20

## Enteric Pathogens Monitored By NICED

### Bacteria

***Vibrio cholerae***

*Vibrio parahaemolyticus*  
*Vibrio fluvialis*  
*Aeromonas spp.*  
*Campylobacter jejuni*  
*Campylobacter coli*  
*Shigella*  
*Salmonella*  
*Escherichia coli*

### Parasites

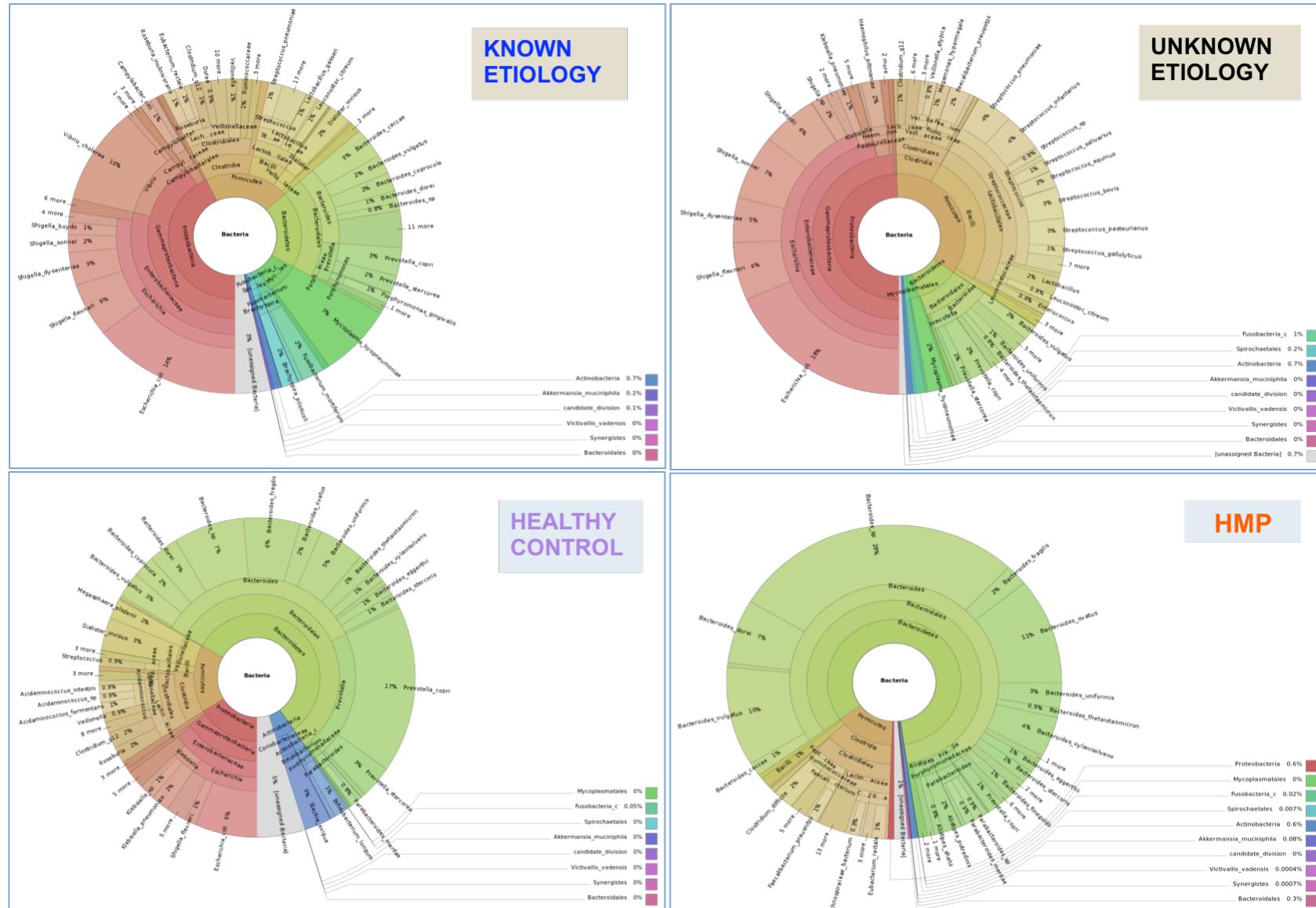
*Giardia lamblia*  
*Cryptosporidium parvum*  
*Entamoeba histolytica*  
*Blastocystis hominis*

### Viruses

Rotavirus  
Adenovirus  
Norovirus  
Sapovirus  
Astrovirus

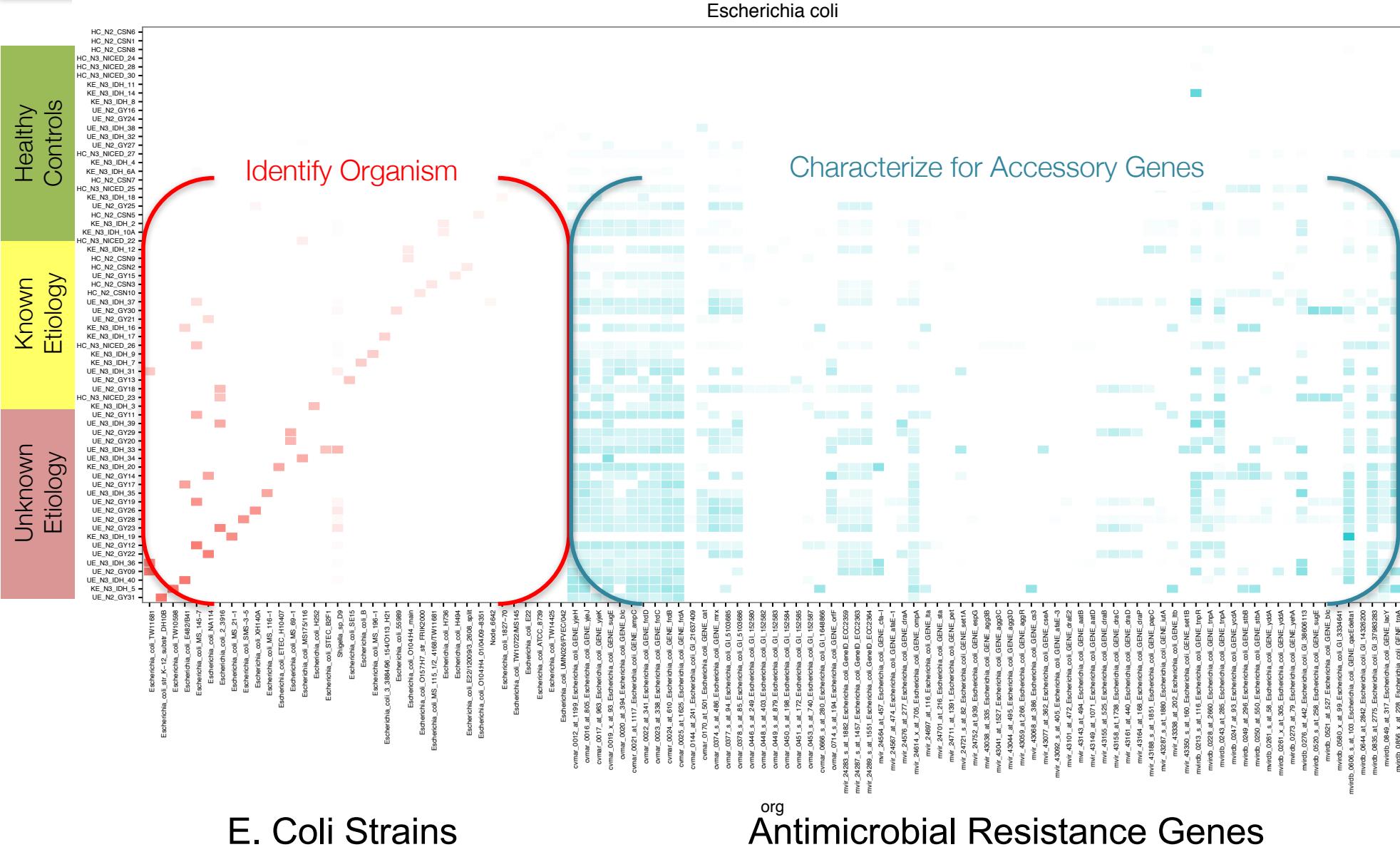
In collaboration with the **National Institute of Cholera and Enteric Disease** (NICED), Calcutta, India

# Overrepresentation of Subpopulation in Diarrheal Patients



γ - Proteobacteria – Firmicutes - Bacteroidetes

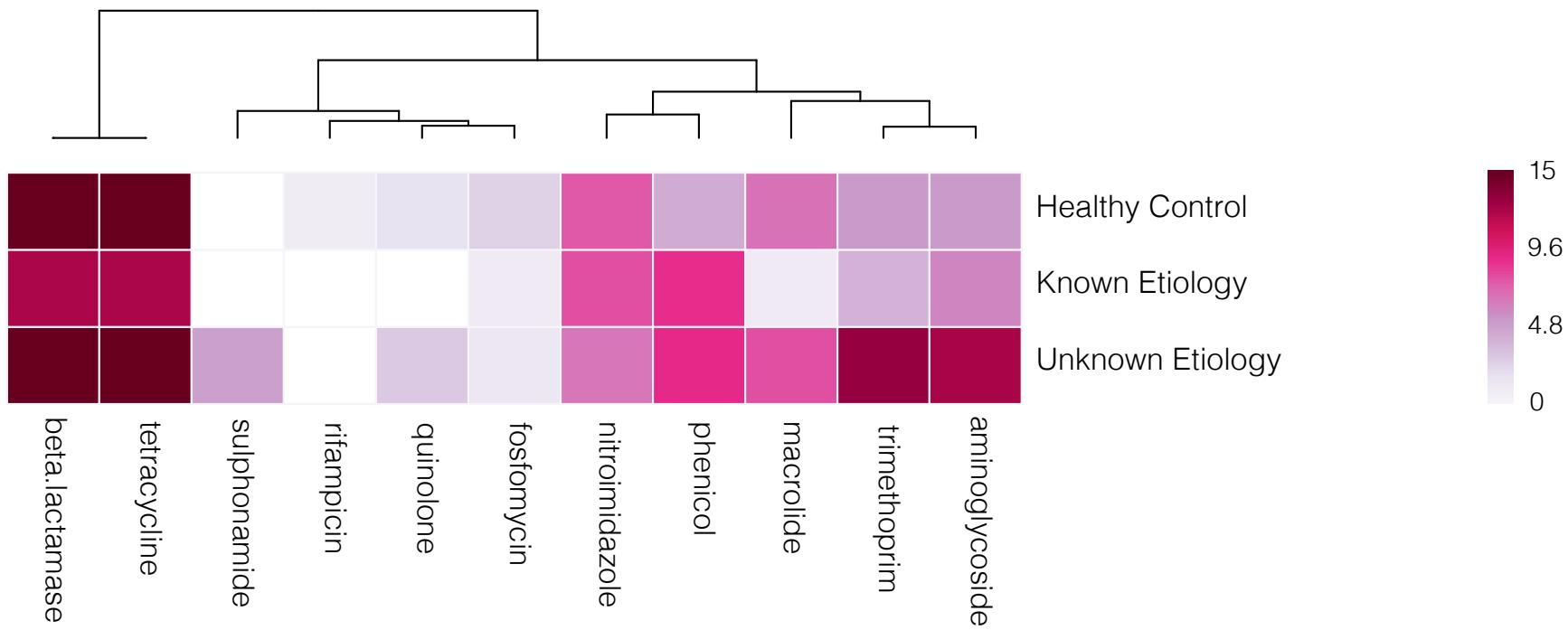
# E. Coli Superfamily in Unknown Etiology Samples



# AMR Genes Present in Microbiome

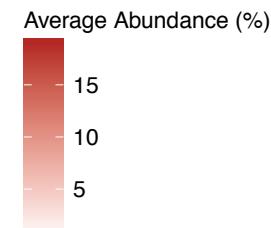
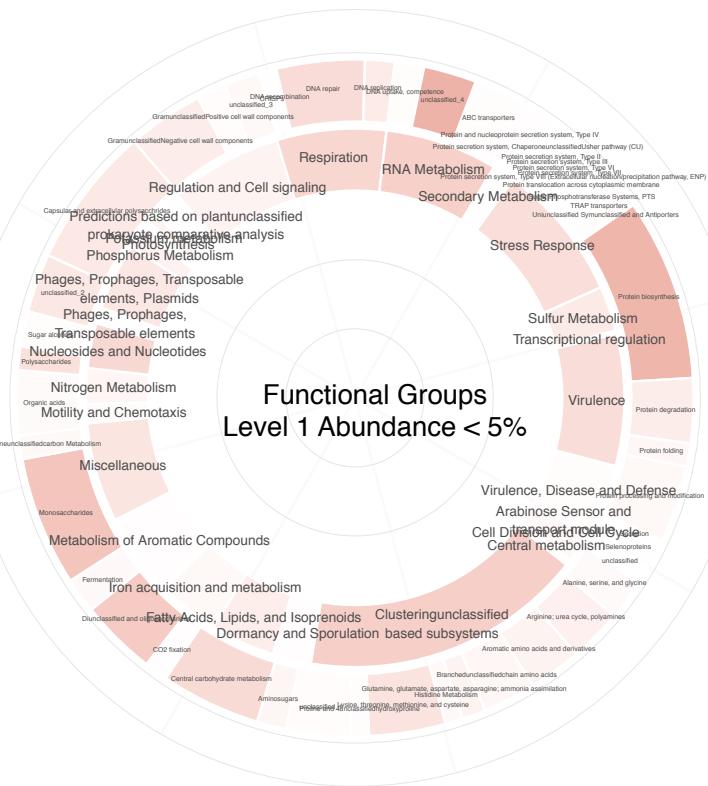
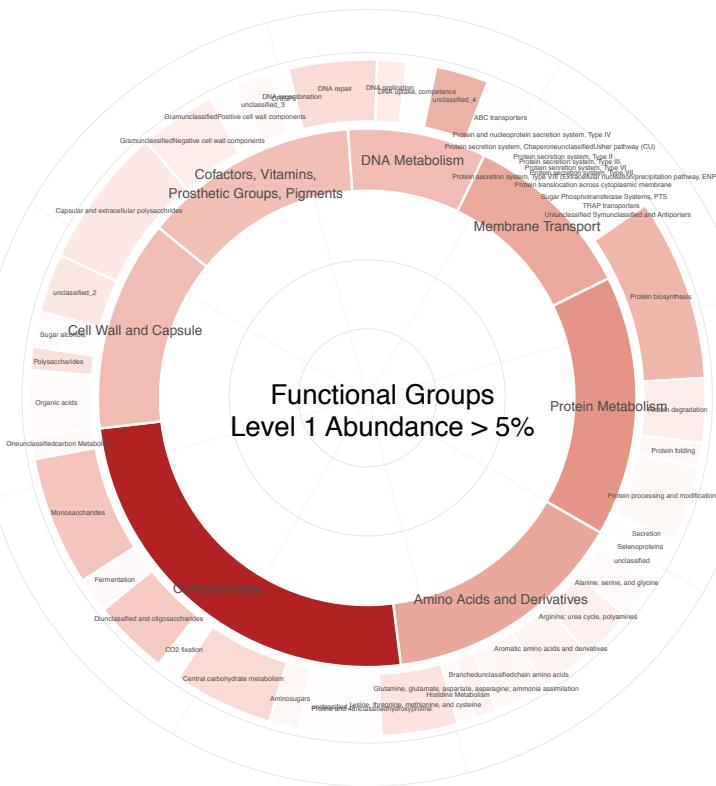
Genes which match at > 50% coverage

HMP samples had no genes present which matched at this level of coverage



# Functional Analysis

## Predominance of genes related to carbohydrate metabolism

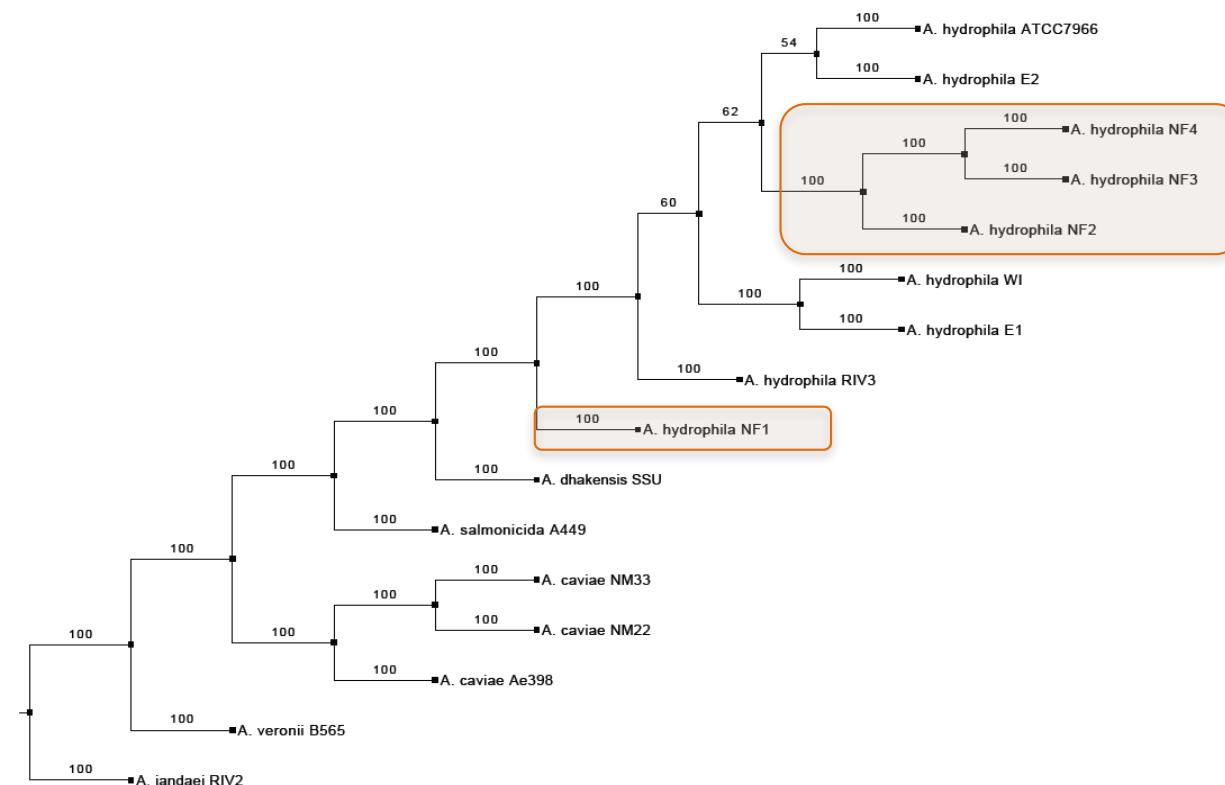


# Necrotizing Fasciitis

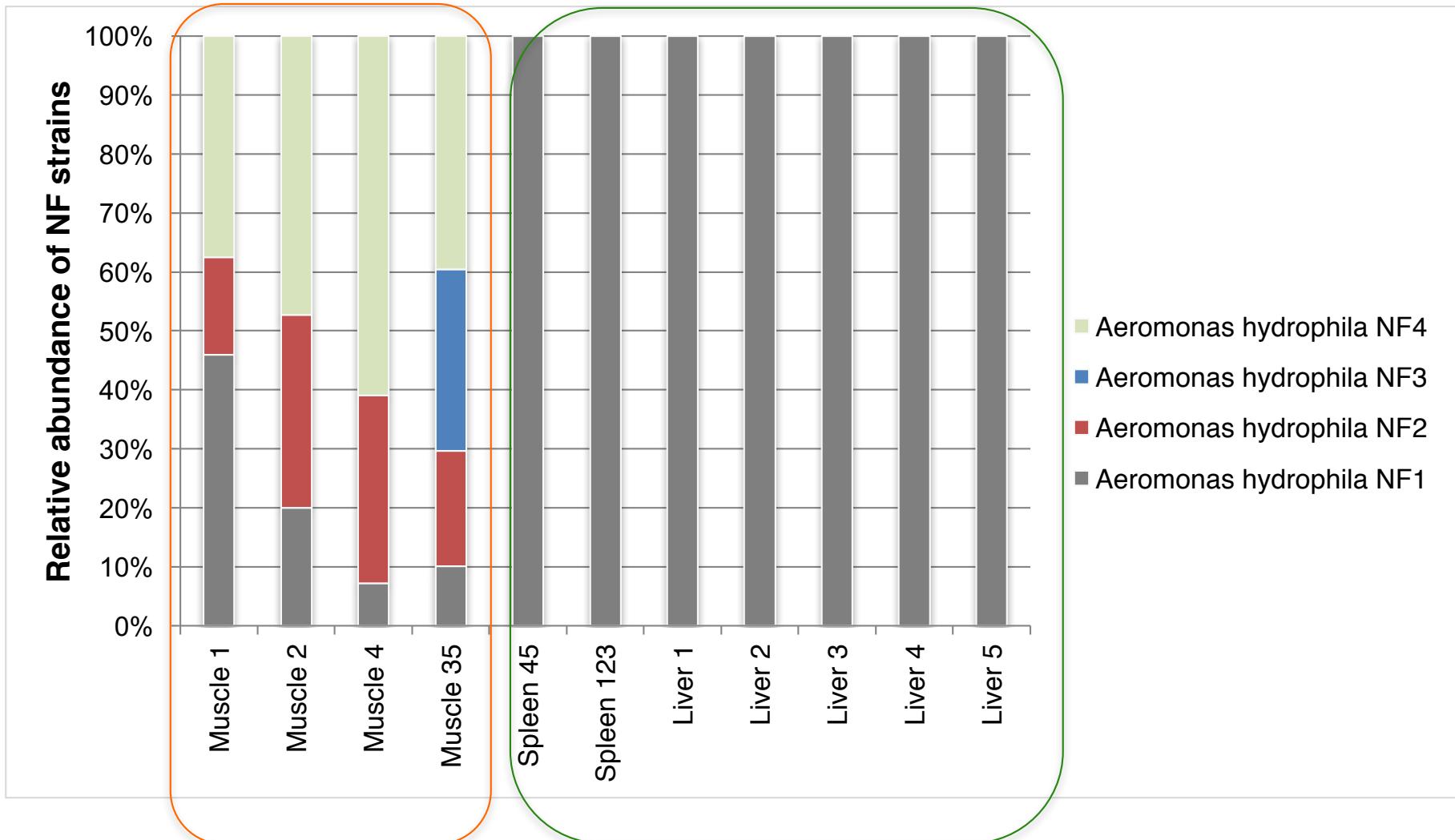
## Cross-talk among flesh-eating *Aeromonas hydrophila* strains in mixed infection leading to necrotizing fasciitis

Duraisamy Ponnusamy<sup>a,1</sup>, Elena V. Kozlova<sup>a,1</sup>, Jian Sha<sup>a</sup>, Tatiana E. Erova<sup>a</sup>, Sasha R. Azar<sup>a</sup>, Eric C. Fitts<sup>a</sup>, Michelle L. Kirtley<sup>a</sup>, Bethany L. Tiner<sup>a</sup>, Jourdan A. Andersson<sup>a</sup>, Christopher J. Grim<sup>b</sup>, Richard P. Isom<sup>c</sup>, Nur A. Hasan<sup>c,d</sup>, Rita R. Colwell<sup>c,d,e,2</sup>, and Ashok K. Chopra<sup>a,2</sup>

<sup>a</sup>Department of Microbiology and Immunology, University of Texas Medical Branch, Galveston, TX 77555; <sup>b</sup>Center for Food Safety and Applied Nutrition, Office of Applied Research and Safety Assessment, Food and Drug Administration, Laurel, MD 20708; <sup>c</sup>CosmosID Inc., Rockville, MD 20850; <sup>d</sup>Center for Bioinformatics and Computational Biology, University of Maryland Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742; and <sup>e</sup>Bloomberg School of Public Health, The Johns Hopkins University, Baltimore, MD 21205

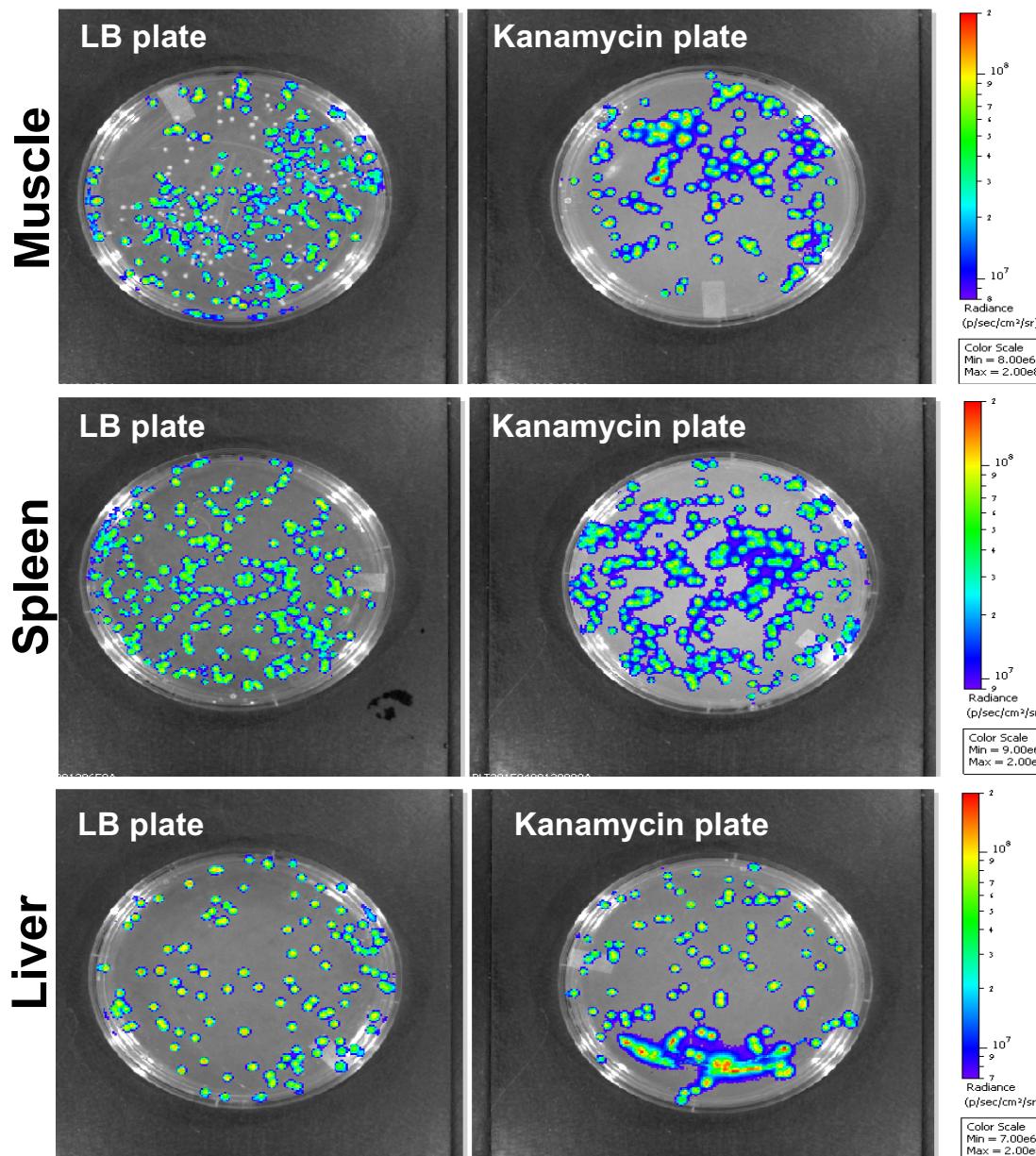


# Strain Level ID shows selective dissemination



Relative distribution of four *Aeromonas hydrophila* strains - NF1, 2, 3, and 4 into different metagenomic datasets derived from muscle, spleen and liver samples.

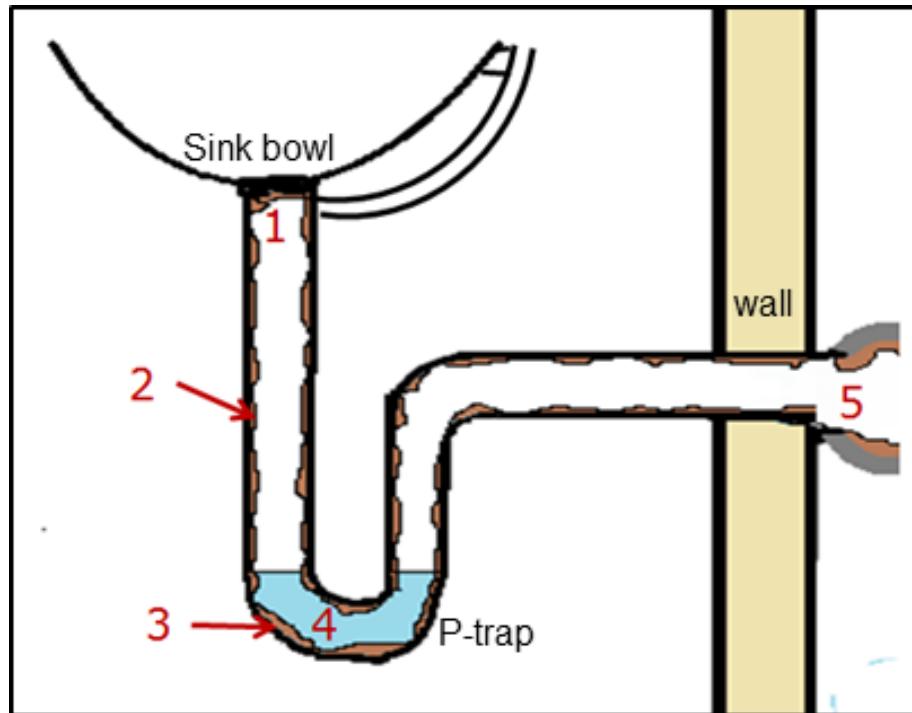
# Mixed infection of *Aeromonas* isolates (NF1, NF2) in mice



Ability to provide strain level ID facilitated better understanding about the dynamics of mixed infections:

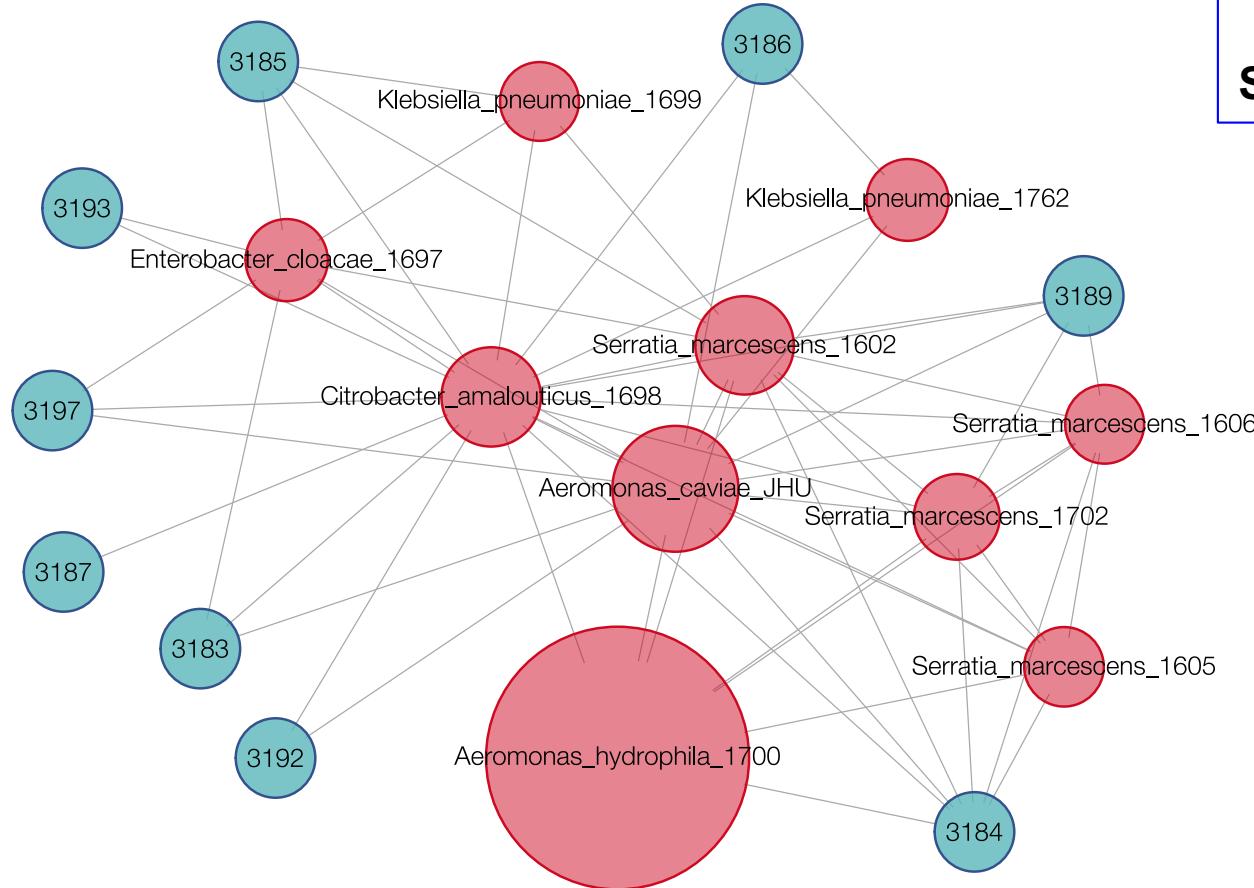
- exoA helps in better dissemination of NF1
- NF2 may facilitate NF1 dissemination in mixed infection
- Likely, unknown factor(s) from NF1 would prevent NF2 dissemination or block bacterial proliferation

# Hospital Biofilms: Source Tracking



Collaboration with Amy J. Mathers, UVA

# Patient Isolates Present in Sink Trap Biofilms



**Patient isolate (red)**

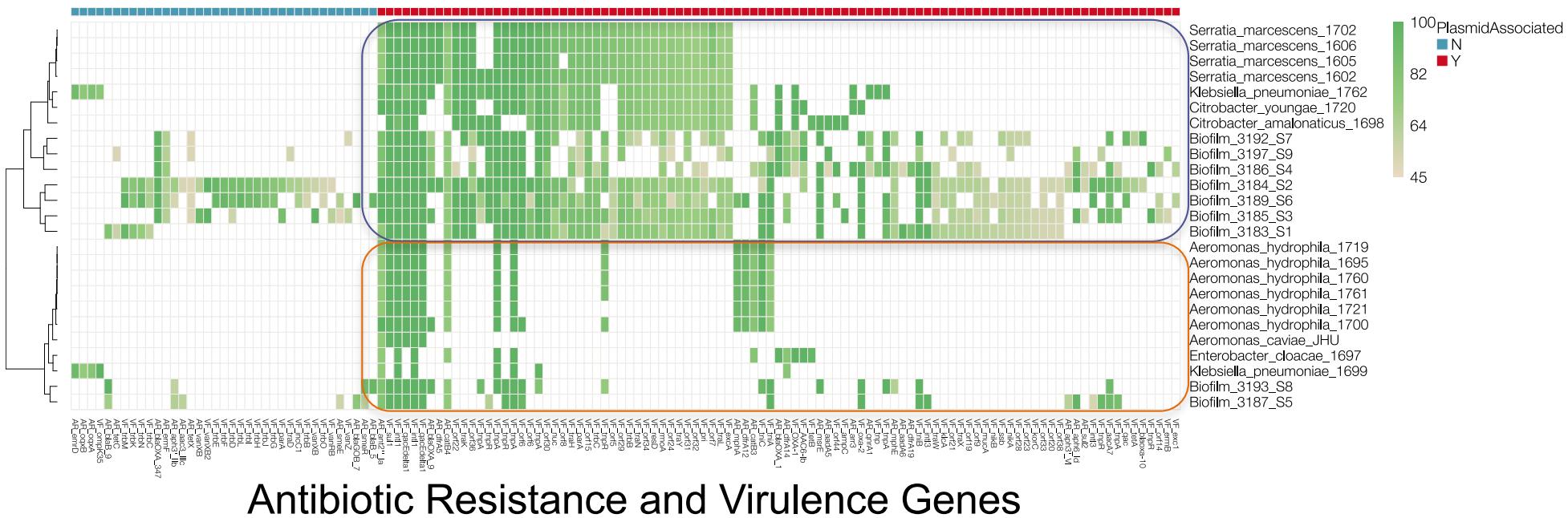
**Sink trap biofilm (blue)**

- An edge connects an isolate (in red) to a sink trap biofilm (in blue), if the isolate is present in the biofilm.
- An edge connects 2 isolates if they were both found in the same sink trap sample.
- The size of the isolate vertices represents the average abundance in the biofilm samples.

# Accessory Gene Profiling

Not Plasmid Associated

Plasmid Associated



# Your second genome?



**Human body:**  
**~10 trillion cells**

**Microbiome**  
**~100 trillion cells**

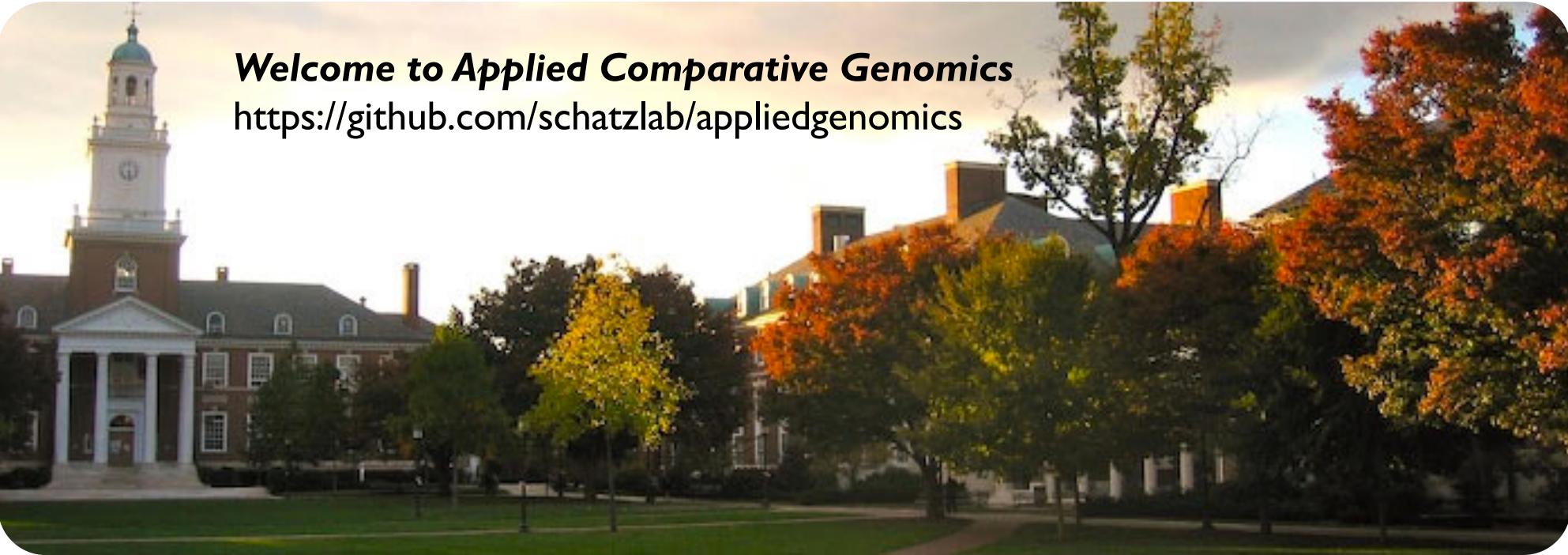
**Human brain:**  
**~3.3 lbs**

**Total mass:**  
**~3.3 lbs**

**Are We Really Vastly Outnumbered? Revisiting the Ratio of Bacterial to Host Cells in Humans**  
Sender et al (2016) Cell. <http://doi.org/10.1016/j.cell.2016.01.013>

# Next Steps

1. Questions on project?
2. Check out the course webpage



**Welcome to Applied Comparative Genomics**

<https://github.com/schatzlab/appliedgenomics>

**Questions?**