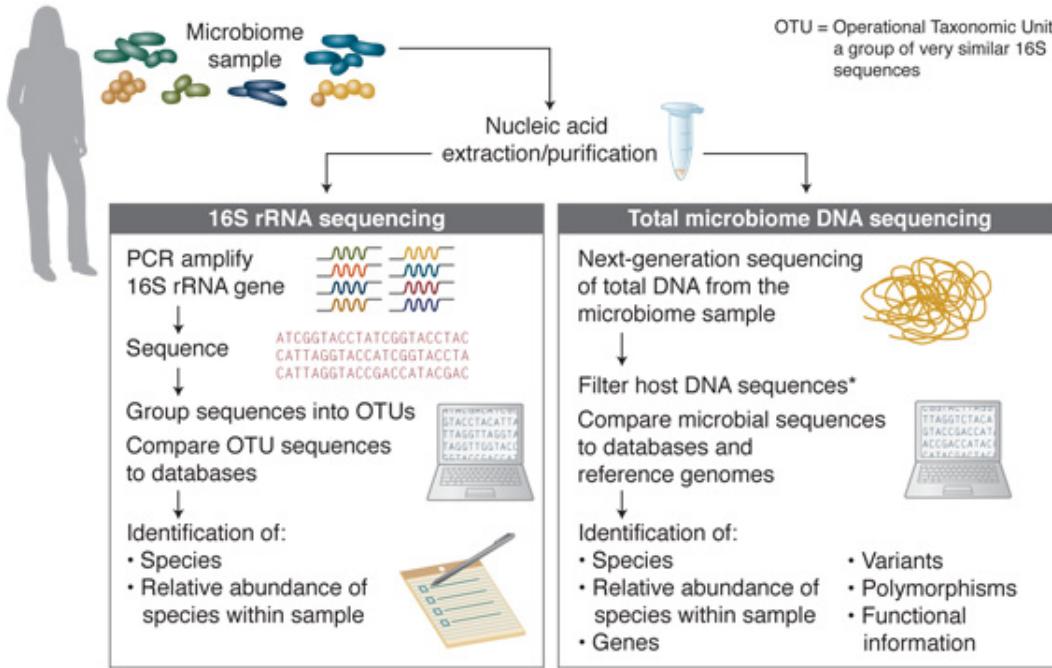




UNIVERSITY OF  
SOUTH FLORIDA

A PREEMINENT  
RESEARCH  
UNIVERSITY



## Best practices for microbiome sample-handling and nucleic acid-processing

Swamy Adapa

# Human Microbiome

**Microbiota** is the coevolved ecological community of commensal, symbiotic and pathogenic microorganisms that share our body space.

## General health and disease

Obesity (Energy use, Extraction/Absorption of calories)

Arthritis

IBD

Neuropsychiatric illnesses (Autism, Schizophrenia, OCD, ADHD, CFS)

## Disease severity

Autoimmune disorders

Infectious diseases

Asthma

## Response to medications

Metabolism, efficacy and toxicity

## The Importance of the **MICROBIOME** by the Numbers



**90%**

Up to 90% of all disease can be traced in some way back to the gut and health of the microbiome

**10-100 trillion**

Number of symbiotic microbial cells harbored by each person, primarily bacteria in the gut, that make up the human microbiota

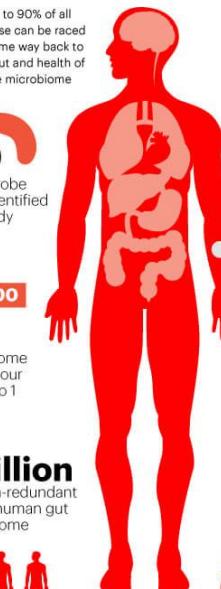


**>10,000**

Number of different microbe species researchers have identified living in the human body

**10X**

There are 10 times as many outside organisms as there are human cells in the human body



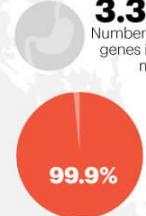
**100 to 1**

The genes in our microbiome outnumber the genes in our genome by about 100 to 1



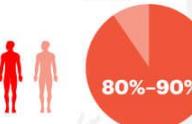
**22,000**

Approximate number genes in the human gene catalog



**99.9%**

Percentage individual humans are identical to one another in terms of host genome



**80%-90%**

Percentage individual humans are different from one another in terms of the microbiome

# Current Methods for Microbial Analysis

- Culture
- Gene cloning (Pan 16S rRNA) and Sanger sequencing
- Microarray
- Next-generation sequencing
  - 16S rRNA sequencing
  - Whole genome sequencing
- MALDI
- qPCR–Target dependent
  - 16S rRNA gene
  - Other relevant gene (antibiotic resistance gene, virulence factor gene)

# Current Platforms

Next Generation Sequencing  
Amplified Single Molecule Sequencing

## Illumina



## Thermo Fisher Scientific



# Current Platforms

Third Generation Sequencing  
Single Molecule Sequencing

## PacBio



PacBio RSII



PacBio Sequel

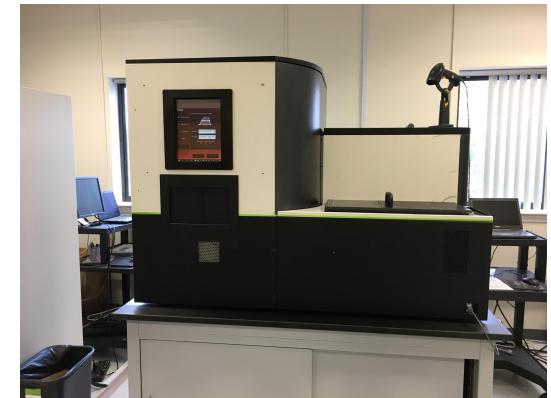
## Oxford Nanopore



GridION

PromethION

## SeqLL



tSMS sequencer

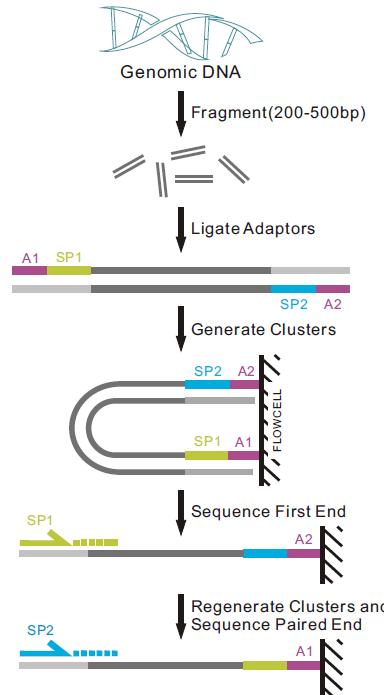
# Illumina

	iSeq 100	MiniSeq	MiSeq	NextSeq 550
<b>Read Length</b>	2 x 150 bp	2 x 150 bp	2 x 300 bp	2 x 150 bp
<b>Throughput</b>	1.2 Gb	7.5 Gb	15 Gb	120 Gb
<b>Reads per run</b>	4 million	50 million	50 million	800 million
<b>Accuracy</b>	99,9 % (>80%)	99,9 % (>80%)	99,9 % (>70%)	99,9 % (>80% of the bases)
<b>Run Time</b>	17.5 hours	24 hours	55 hours	29 hours

	HiSeq 2500 / 3000 / 4000	HiSeq X	NovaSeq 6000
<b>Read Length</b>	2 x 125 / 2 x 150 / 2 x 150 bp	2 x 150 bp	2 x 150 bp
<b>Throughput</b>	1000 / 750 / 1500 Gb	1800 Gb	850 – 3000 Gb
<b>Reads per run</b>	4 / 2,5 / 5 billion	6 billion	2.8 – 10 billion
<b>Accuracy</b>	99,9 % (>80% of the bases)	99,9 % (>75%)	99,9 % (>75% of the bases)
<b>Run Time</b>	6 / 3,5 / 3,5 days	< 3 days	36 – 44 hours

**Workflow**      Library preparation      ->      Bridge amplification      ->      Reversible termination sequencing

# Pair-End Sequencing



- The forward and reverse reads are stored in two fastq files.
- The order of pairs and naming is identical, except the designation of forward and reverse.

ID\_R1\_001.fasta

```
@HISEQ:100:C3MG8ACXX:5:1101:1160:2
197 1:N:0:ATCACG
CAGTTGCATGAGAGCGTTGAGAAGTATAATAGG
AGTTAACTGAGTAACAGGATAAGAAATAGTGAG
ATATGGAAACGTTGTGGTCTGAAAAGAAGATGT
+
B@CFFFFFFHHHHHGJJJJJJJJFHHIIIIJJ
JIHGIJJJJIJIJJJJJJJJJJJIEIHHIJ
HGHHHHHHDFFFEDDDDDCDCDDCDCDDDDDDCDC
```

ID\_R2\_001.fasta

```
@HISEQ:100:C3MG8ACXX:5:1101:1160:
2197 2:N:0:ATCACG
CTTCGTCCACTTCATTATTCTCTTACATAG
CTCTCCGGTTTAGGGTACTCTTGACCTGGCCTT
TTTCAAGACGTCCCTGACTTGATCTTGAAACG
+
CCCFFFFFFHHHHHJJJJJJJJJJJJJJJJJJJJ
JJJJJJJJIJIJGIJHBGHIIIIJJJJJJJJJJ
JJJHFFFFFFDDDDDDDDDDDDDDDDDEDCCDDDD
```

Figure 1-2-1 Pipeline of paired-end sequencing ([www.illumina.com](http://www.illumina.com))

# Sequencing depth

How many base pairs of sequence data will I get?

- Number of reads being sequenced
- Read length (if paired consider them as individuals)
- Number of samples being sequenced
- Expected percentage of usable data

$$bpPerSample = \frac{readLength * readCount}{sampleCount} * 0.8$$

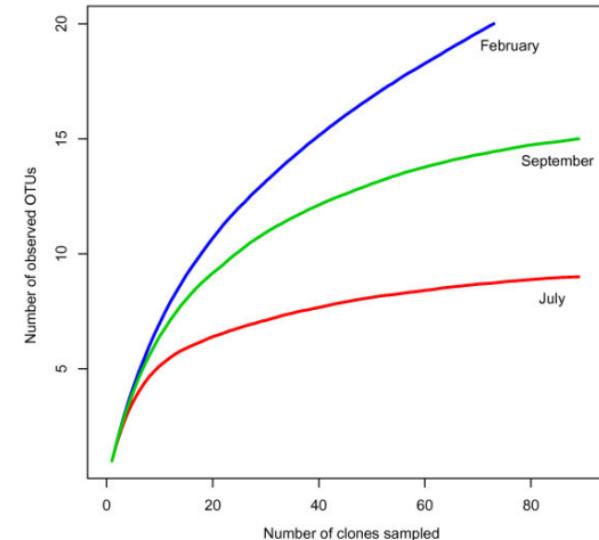
- Reads and read length data - manufacturer's website (search for specifications)
- Always use the lower end of the estimate.

# For Amplicon Sequencing

- Number of reads being sequenced
- Proportion that is diversity sample (e.g. PhiX)
- Number of samples being pooled in the run

$$\frac{\text{reads}}{\text{sample}} = \frac{\text{reads\_sequenced} * (1 - \text{diversity\_sample})}{\text{num\_samples}}$$

- Illumina ‘recommends’ 100K per sample
  - For 16S – 50k minimum.
- After Pilot Seq - Plot rarefaction curves to determine if saturation is achieved.



# Cost Estimation

- Sample Prep
  - DNA/RNA extraction and QA/QC (Bioanalyzer/Gels)
  - Library QA/QC (Bioanalyzer and Qubit/microplate reader)
  - Pooling
- Sequencing (Number of Lanes / runs)
- Bioinformatics (General rule is to estimate the same amount as data generation, i.e. double your budget)

# Sequencing Costs

- Metagenomics
  - Shotgun libraries intended to sequence random genomic sequences from the entire bacterial community.
  - Can be costly per sample (\$500 to multi thousands per sample)
  - Better resolution and sensitivity to characterize the sample
  - Due to cost, can only do relatively few samples
- Amplicon Sequencing
  - Sequence only one regions of one gene (e.g. 16s, ITS, LSU)
  - Cheap per sample (at scale, down to \$20/sample)
  - Due to cost, can do many hundreds of samples make more global inferences

# Bioinformatics Costs

Bioinformatics includes:

- Storage of data, metadata maintenance
- Access and use of computational resources and software
- System Administration time
- Bioinformatics Data Analysis time
- Back and forth consultation/analysis to extract biological meaning

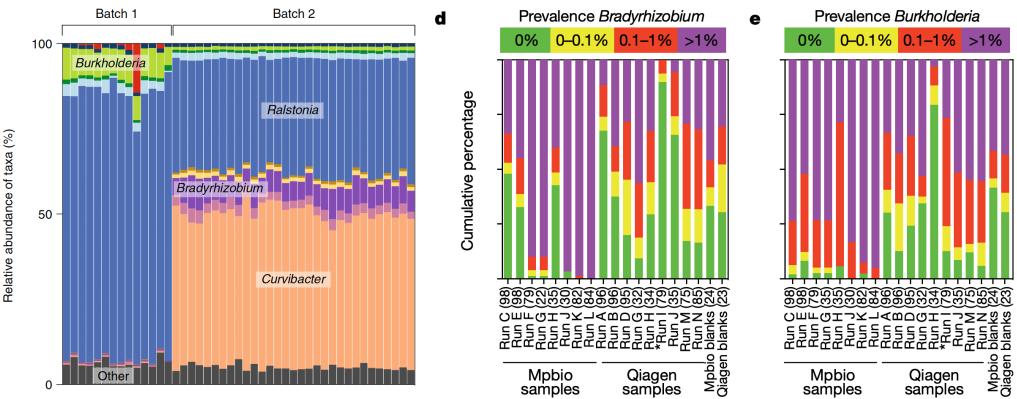
Rule of thumb:

Bioinformatics can and should cost as much (sometimes more) as the cost of data generation.

# Cost Estimation

- Amplicons
  - 384 Samples
  - Amplicon generation (~\$20/sample)= \$7,680
  - Sequencing PE300, target 50K reads per sample
  - Bioinformatics
- Metagenome
  - 12 samples (DNA) = \$400/sample = \$4,800
  - Expectations: Host Proportion < 50%, use average genome size of EColi, Target the 1% and coverage of 20
  - Sequencing PE150
  - Bioinformatics

# Kit'ome - Reagent microbiome



## ARTICLE

Corrected: Author Correction

<https://doi.org/10.1038/s41586-019-1451-5>

## Human placenta has no microbiome but can contain potential pathogens

Marcus C. de Goffau<sup>1,2,8</sup>, Susanne Lager<sup>3,4,5,8</sup>, Ulla Sovio<sup>3,4</sup>, Francesca Gaccioli<sup>3,4</sup>, Emma Cook<sup>3</sup>, Sharon J. Peacock<sup>1,6,7</sup>, Julian Parkhill<sup>1,2\*</sup>, D. Stephen Charnock-Jones<sup>3,4,9</sup> & Gordon C. S. Smith<sup>3,4,9</sup>

### Recognizing the reagent microbiome

A noticeable part of the microbiome literature, especially that working with low-biomass samples, is plagued by reagent contamination. Here we describe visual, statistical, methodical and ecological techniques to facilitate recognition of signals that represent contamination.

Marcus C. de Goffau, Susanne Lager, Susannah J. Salter, Josef Wagner, Andreas Kronbichler, D. Stephen Charnock-Jones, Sharon J. Peacock, Gordon C. S. Smith and Julian Parkhill

APPLIED AND ENVIRONMENTAL MICROBIOLOGY, Aug. 1998, p. 3110–3113  
0099-2240/98/\$04.00+0  
Copyright © 1998, American Society for Microbiology. All Rights Reserved.

Vol. 64, No. 8

### Specific Ribosomal DNA Sequences from Diverse Environmental Settings Correlate with Experimental Contaminants

MICHAEL A. TANNER, BRETT M. GOEBEL,<sup>†</sup> MICHAEL A. DOJKA, AND NORMAN R. PACE\*

Departments of Plant and Microbial Biology and Molecular and Cell Biology,  
University of California, Berkeley, California 94720-3102

Received 23 April 1998/Accepted 10 June 1998

# Kit'ome - Reagent microbiome

- Breast tissue from cancer patients and normal controls
  - Different abundances of *Methylobacterium radiotolerans* and *Sphingomonas yanoikuyaе*.
  - *Methylobacteriaceae* and *Sphingomonadaceae* popped up in some of the DNA extraction kits.
- *Bradyrhizobium enterica*, in patients with a colitis syndrome.
  - *Bradyrhizobium* were also detected in DNA extraction kits

The NEW ENGLAND JOURNAL of MEDICINE

OPEN  ACCESS Freely available online



## Microbial Dysbiosis Is Associated with Human Breast Cancer

Caiyun Xuan<sup>1</sup>, Jaime M. Shamonki<sup>2</sup>, Alice Chung<sup>3</sup>, Maggie L. DiNome<sup>4</sup>, Maureen Chung<sup>4</sup>, Peter A. Sieling<sup>1</sup>, Delphine J. Lee<sup>1\*</sup>

**1** Dirks/Dougherty Laboratory for Cancer Research, Department of Translational Immunology, John Wayne Cancer Institute, Santa Monica, California, United States of America, **2** Pathology Department, Saint John's Health Center, Santa Monica, California, United States of America, **3** Department of Medicine, Cedars-Sinai Medical Center, Los Angeles, California, United States of America, **4** Margie Petersen Breast Center, Saint John's Health Center, Santa Monica, California, United States of America

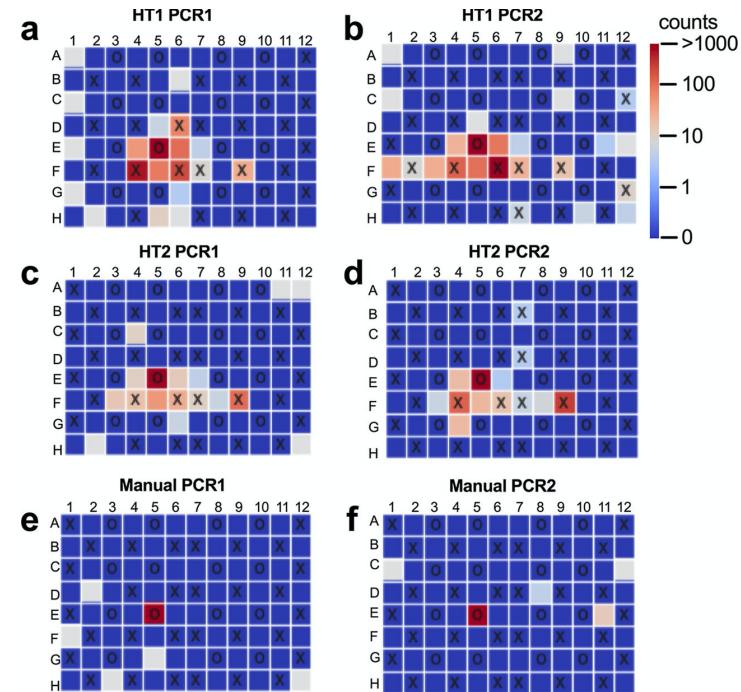
ORIGINAL ARTICLE

## Sequence-Based Discovery of *Bradyrhizobium enterica* in Cord Colitis Syndrome

Ami S. Bhatt, M.D., Ph.D., Samuel S. Freeman, B.S.E., Alex F. Herrera, M.D., Chandra Sekhar Pedamallu, Ph.D., Dirk Gevers, Ph.D., Fujiko Duke, B.S., Joonil Jung, Ph.D., Monia Michaud, M.Sc., Bruce J. Walker, B.S., Sarah Young, Ph.D., Ashlee M. Earl, Ph.D., Aleksander D. Kostic, Ph.D., Akinyemi I. Ojesina, M.D., Ph.D., Robert Hasserjian, M.D., Karen K. Ballen, M.D., Yi-Bin Chen, M.D., Gabriela Hobbs, M.D., Joseph H. Antin, M.D., Robert J. Soiffer, M.D., Lindsey R. Baden, M.D., Wendy S. Garrett, M.D., Ph.D., Jason L. Hornick, M.D., Ph.D., Francisco M. Marty, M.D., and Matthew Meyerson, M.D., Ph.D.

# Well-to-Well Contamination

- Decreased cost of 16S rRNA -> Automation
- Outsourced samples are always processed using High-throughput platforms.
- Between-sample contamination is often neglected.
- Study found
  - DNA extraction primary contamination.
  - Lesser extent during library preparation
  - Barcode leakage was negligible.
- Recommendations:
  - Samples should be randomized across plates
  - Samples of similar biomasses should be processed together
  - Manual single-tube extractions or hybrid plate-based cleanups should be employed.



# Need for Controls – Blanks & Mock

- Positive or Mock Community control
  - To assess sequencing error rate
- Negative or water extraction sample.
  - Essential for low biomass samples
  - Probably a good idea for high biomass samples
- Airborne contaminants.
- Autoclave does not kill DNA.

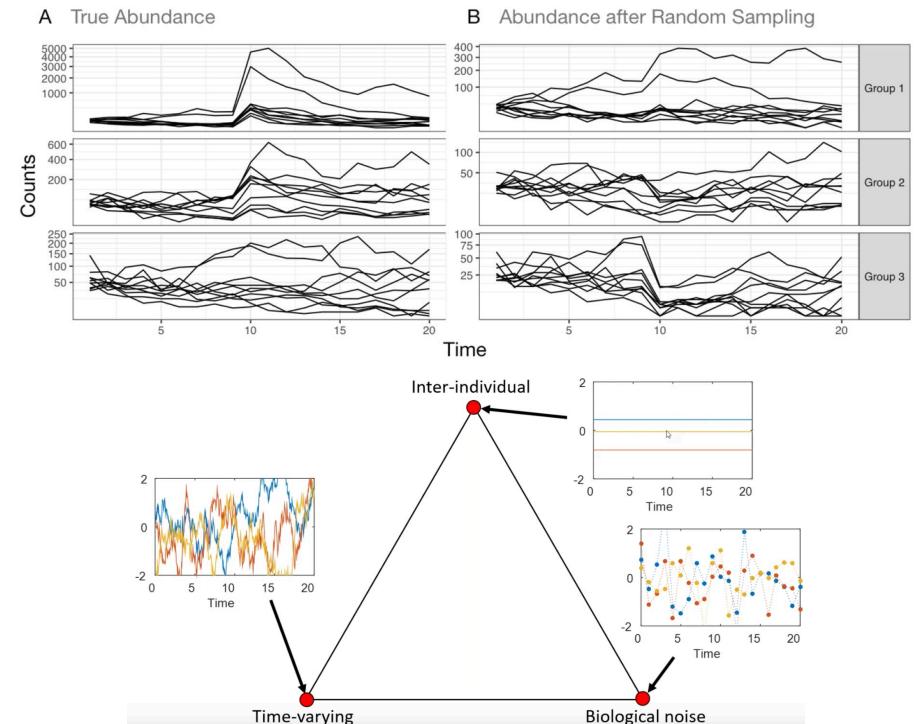
Avoid simplistic removals of taxa or operational taxonomic units (OTUs) appearing in negative controls – The contaminant could be **real**.

# Batch Effect – Run Preparation

- Do not just sequence samples as they come.
- Human Microbiome Project
  - Largest signal in that study was where the DNA extractions were done - St. Louis or Houston (Regional effect on the microbiome?)
  - Geographic region was perfectly confounded with who did the sampling, extractions, amplifications, and sequencing.
- Salter et. al, found that the differences they observed were due to differences in the contaminants that were present in the DNA extraction kits that they used.
  - Two batches of extractions used different lots of extraction kits
  - Treatment groups were confounded with the DNA extraction kit lots
- Wait until you have all of your samples.
- Either randomize the samples or thoughtfully block your samples.
  - So that treatment groups are evenly represented across your extraction sets and sequencing runs.

# Random sampling can lead to spurious conclusions in longitudinal studies.

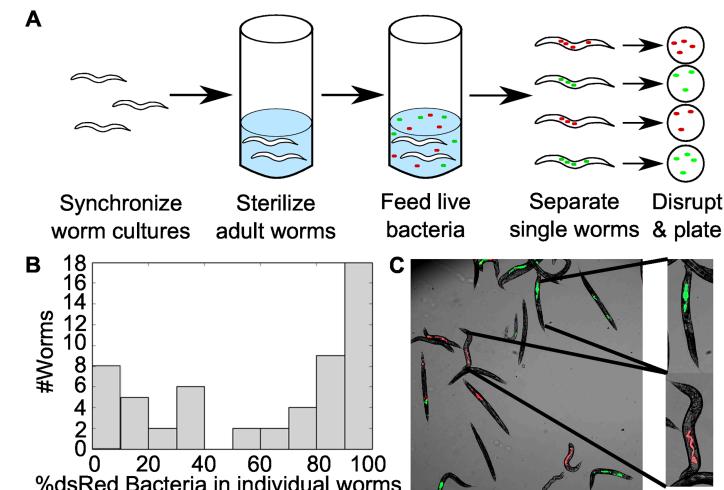
- Measurements of microbiota can differ substantially from the underlying true community structure.
- DNA of samples are randomly selected, processed, and then included in the subsequent step.
- The randomness of this data generation process introduces uncertainty into microbiome measurements.
- Technical variation from sample processing can influence ex vivo microbiota dynamics, accounting for 76% of community variation on hourly timescales.



<https://drive.google.com/file/d/1U36NyMinKAm1R2rwgRb0fJUFVhbeYS9J/view>

# Microbiome diversity is influenced by chance encounters

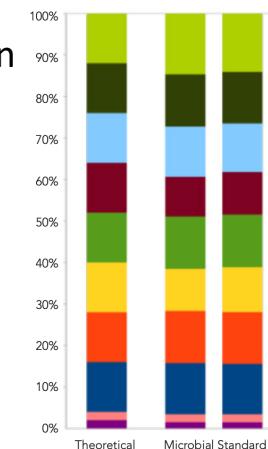
- Genetically identical worms were put into identical environments and fed them the same diet.
- The worms developed very different populations of bacteria in their gut, depending on which bacteria happened to make it there first.
- When the experiment began, the worms had no bacteria in their digestive tracts, but after being exposed to the identical bacterial diet, the worms each generated very different microbe populations in their guts.



- Studies involving multiple individuals over time should design the experiment in a way that will subsequently be utilized by the the data modeling scheme.
- Seasonality effects can also be important to consider in experimental design.
- While perhaps most obvious in environmental studies, dietary changes surrounding holidays or weekends, or even natural circadian variation can also be important factors to consider.

# Microbiome Standards - Mock Communities

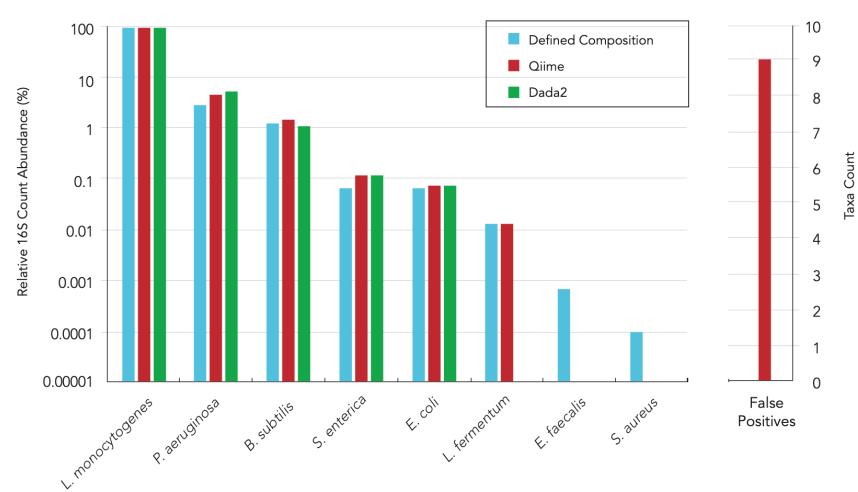
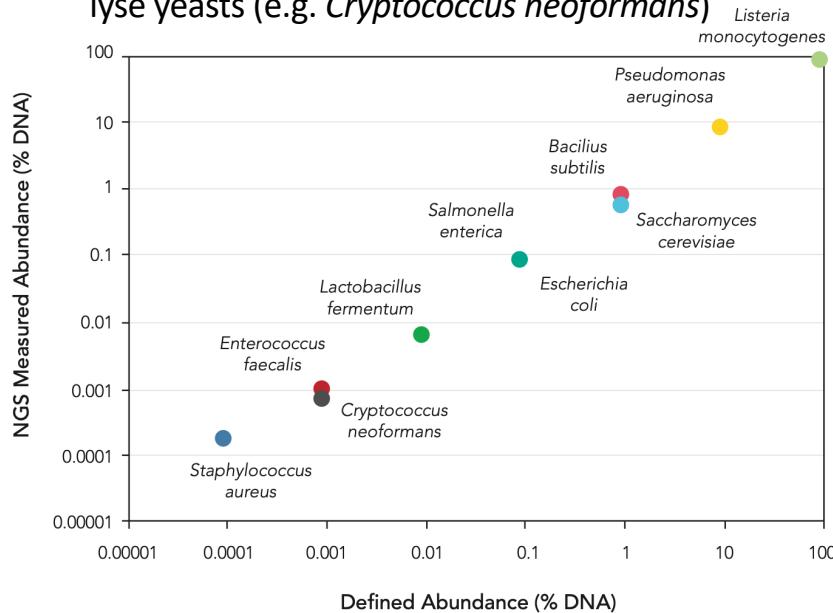
- Mock community control composed of predetermined ratios of Microbes/DNA from a mixture of bacterial species.
  - Quantification of sequencing error
  - Identify bias introduced during the sampling and library preparation processes.
- Pre-prepared bacterial communities are available in two different formats: DNA mock communities and whole-cell mock communities.
- Whole-cell mock communities are useful for establishing the efficiency of the DNA extraction step.
- DNA mock communities will only assess the efficiency of PCR, clean-up, sequencing, and analysis steps.
- Mock communities are available from the American Type Culture Collection (ATCC - BIE) and Zymo Research.



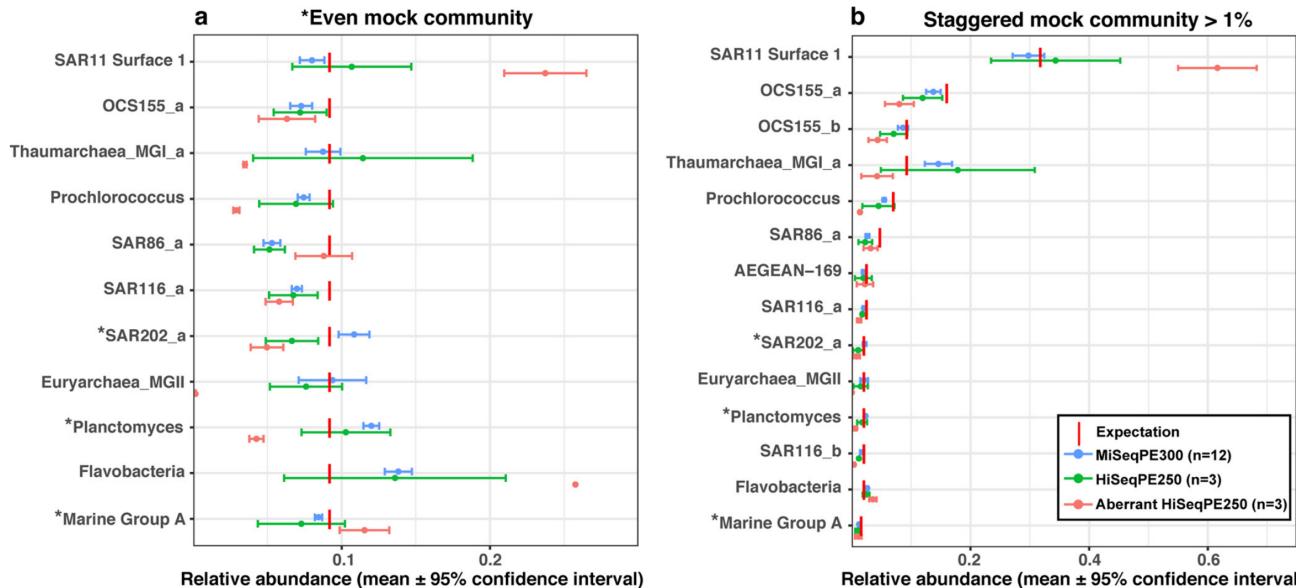
Species	Avg. GC (%)	Gram Stain	gDNA Abun. (%)
<i>Pseudomonas aeruginosa</i>	66.2	-	12
<i>Escherichia coli</i>	56.8	-	12
<i>Salmonella enterica</i>	52.2	-	12
<i>Lactobacillus fermentum</i>	52.8	+	12
<i>Enterococcus faecalis</i>	37.5	+	12
<i>Staphylococcus aureus</i>	32.7	+	12
<i>Listeria monocytogenes</i>	38.0	+	12
<i>Bacillus subtilis</i>	43.8	+	12
<i>Saccharomyces cerevisiae</i>	38.4	Yeast	2
<i>Cryptococcus neoformans</i>	48.2	Yeast	2

# Staggered Mock Communities

Zymo Research's Microbial Community Standard includes three easy-to-lyse Gram-negative bacteria (e.g. *Escherichia coli*), five tough-to-lyse Gram-positive bacteria (e.g. *Listeria monocytogenes*), and two tough-to-lyse yeasts (e.g. *Cryptococcus neoformans*)



# Mock communities help identify aberrant sequencing run



**FIG 1** Comparisons of “even” mock communities (a) and “staggered” mock communities (b) sequenced by MiSeqPE300 and HiSeqPE250. Values that are significantly different for a clone by MiSeqPE300 versus HiSeqPE250 are indicated with an asterisk before the clone name ( $P < 0.05$  by Wilcoxon rank sum test). Significant differences in the whole-community composition by MiSeqPE300 and HiSeqPE250 were found only in the even mock community ( $P < 0.05$  by ANOSIM test).

# Microbiome Standards - Spike-in Controls

- Unlike mock communities, these standards are added directly to samples.
- Quality control can be performed on a per-sample basis.
- Care must be taken to select bacteria which are highly unlikely to occur in the samples of interest.

Example:

*Imtechella halotolerans* (Gram-negative) and *Allobacillus halotolerans* (Gram-positive) – which are alien to human microbiome.

# Spike-in controls

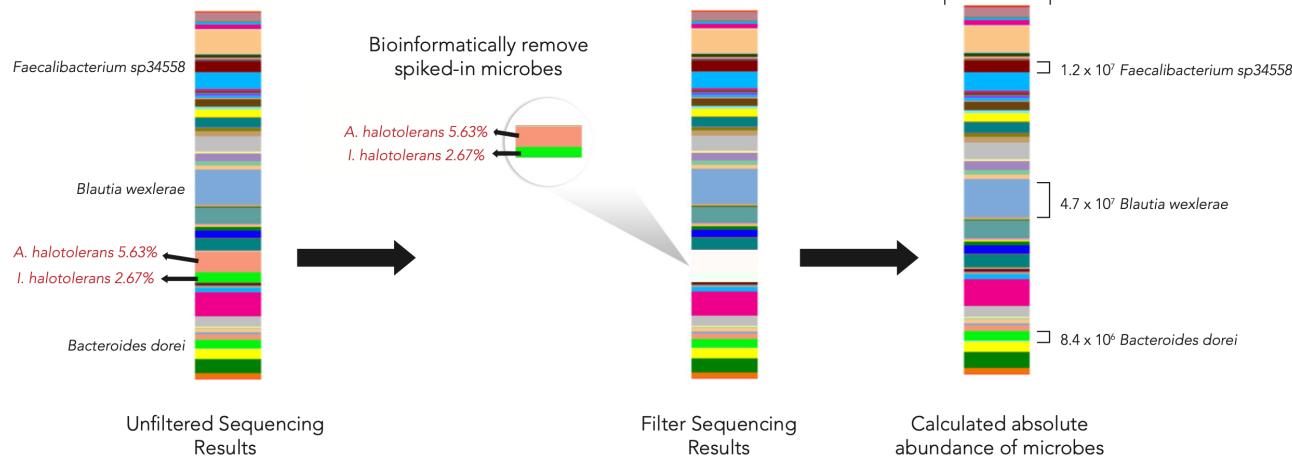
- Microbial Composition

Species	Per Prep (20 µl)		
	Cells	16S Copies <sup>1</sup>	Total DNA <sup>2</sup> (ng)
<i>Imtechella halotolerans</i>	$2 \times 10^7$	$6.0 \times 10^7$	67.2
<i>Allobacillus halotolerans</i>	$2 \times 10^7$	$1.4 \times 10^8$	58.2

<sup>1</sup> 16S copies = cells × 16S copy number per cell/genome.

<sup>2</sup> Total genomic DNA (ng) = cells × genome size (bp/genome) × DNA unit conversion constant (ng/bp). DNA unit conversion constant (ng/bp) =  $1.079 \times 10^{-12}$ .

- Differences in abundance ratio indicate potential workflow bias.
- If the abundance of *I. halotolerans* is much higher than that of *A. halotolerans*.
  - Bias during DNA extraction
  - A. halotolerans* is Gram-positive, tougher to lyse than *I. halotolerans*.





# UNIVERSITY OF SOUTH FLORIDA

A PREEMINENT RESEARCH UNIVERSITY