**University at Buffalo** *The State University of New York* | REACHING OTHERS

# "So you want to do microbiome research…"

Michael J. Buck, Ph.D.
mjbuck@buffalo.edu

Maria Tsompana, Ph.D.
tsompana@buffalo.edu

**University at Buffalo** *The State University of New York* | REACHING OTHERS

# Outline of Discussion

- Bacteria, archaea, fungi, or viruses???
- Sample collection, storage and processing
- 16S or shotgun
- Library construction and sequencing
- How the experiments are done?

# Microbiota defined

- We are born consisting not only of our own eukaryotic human cells, but over the first few days of our life, our skin surface, oral cavity and gut are colonized by a tremendous diversity of **bacteria**, **archaea**, **fungi**, and **viruses -** a new microbial ecosystem defined as the **human microbiota**.
- The human microbiota contains almost **ten times** as many cells as are in the rest of our bodies and accounts for several pounds of body weight.
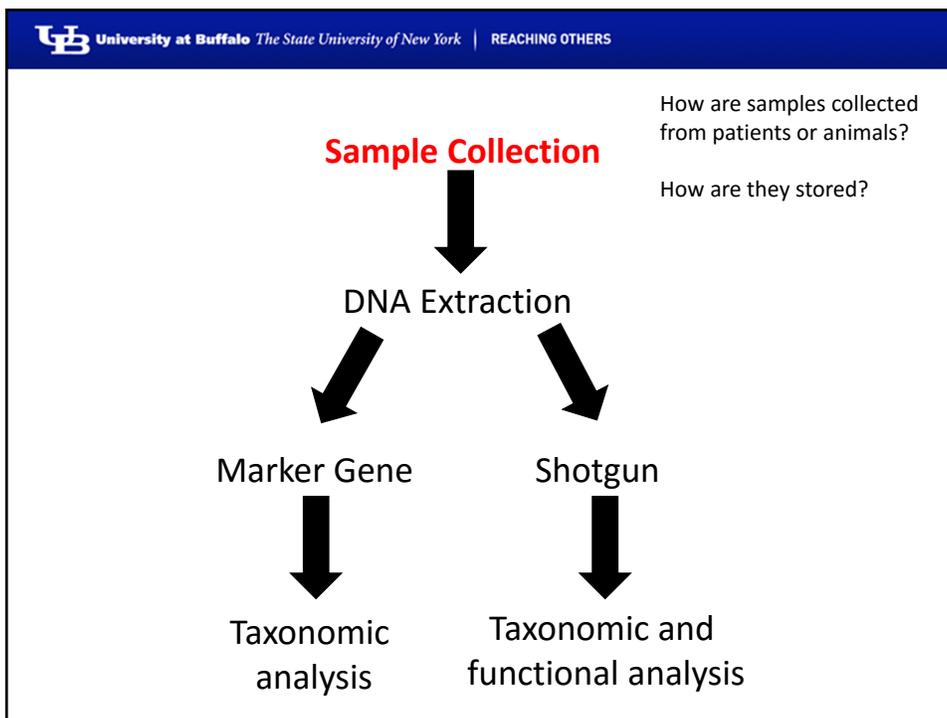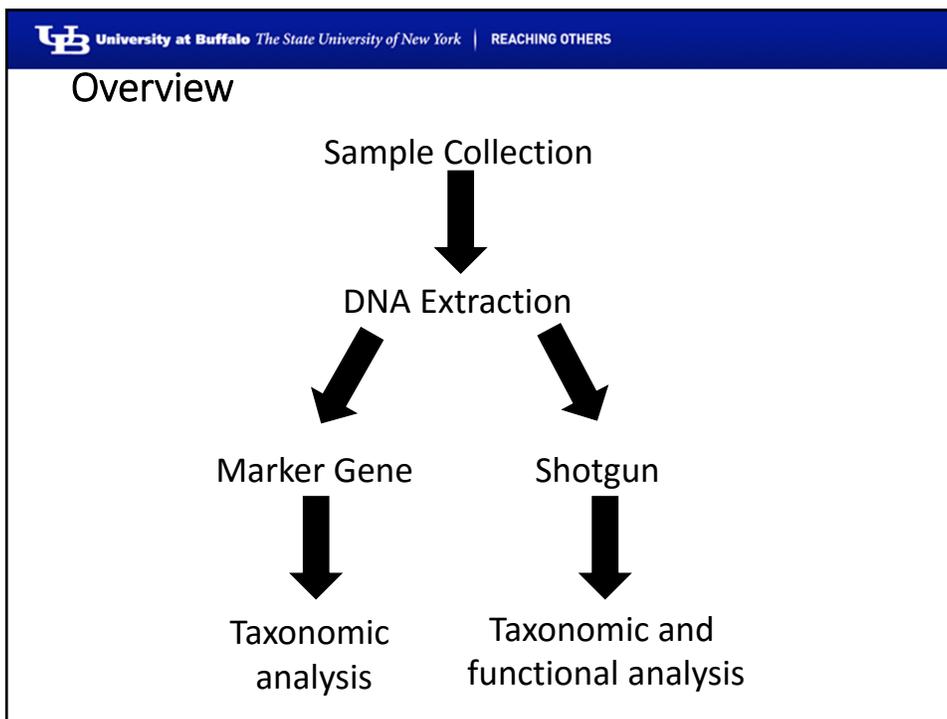
# Microbiota research

- It has long been recognized that many microbes visualized by microscopy cannot by cultivated.
- Despite advances in cultivation technology, >99% of the microbial species residing in various habitats cannot be recovered by available techniques, a phenomenon termed the '**great plate count anomaly**'.
- Currently most scientists use a PCR- and sequence-based approach that exploits 16S ribosomal DNA (rDNA) to profile bacterial diversity.

University at Buffalo *The State University of New York* | REACHING OTHERS

## Bacteria, archaea, fungi, viruses or all of the above.

- Depending on the research question, different parts of the microbiome should be studied

- Bacteria and archaea – 16S
- Fungi and other eukaryotic – 18S
- Viruses – targeted or shotgun
- Everything - shotgun

University at Buffalo *The State University of New York* | REACHING OTHERS

## How do we study the microbiome?

- Marker gene
  - 16S / 18S
  - Amplify region and compare
  - Cheap ($50 per sample), biased but effective

- Shotgun
  - Extract all genomic DNA
  - Fragment, sequence and analyze
  - Expensive ($500 per sample), information rich, should be less biased

University at Buffalo *The State University of New York* | REACHING OTHERS

## Overview

Sample Collection

↓

DNA Extraction

↙          ↘

Marker Gene          Shotgun

↓                    ↓

Taxonomic analysis    Taxonomic and functional analysis

University at Buffalo *The State University of New York* | REACHING OTHERS

How are samples collected from patients or animals?

How are they stored?

**Sample Collection**

↓

DNA Extraction

↙          ↘

Marker Gene          Shotgun

↓                    ↓

Taxonomic analysis    Taxonomic and functional analysis

# Sample collection, storage and processing

- Bacteria like to grow, *E. coli* doubles every 20-30 minutes!
- Anaerobic versus aerobic bacteria will grow at different rates in sample collection tubes.
    - So if a person collects a sample at home and stores it in in the fridge or even the freezer the population of bacteria will change over time.

- Samples need to be chemically preserved or flash frozen at -80

# Sample storage

Journal of Microbiological Methods
Volume 95, Issue 3, December 2013, Pages 381–383

ELSEVIER

Note

Differential recovery of bacterial and archaeal 16S rRNA genes from ruminal digesta in response to glycerol as cryoprotectant

Nest McKain[a], Buğra Genc[b], Timothy J. Snelling[a], R. John Wallace[a]

Show more

doi:10.1016/j.mimet.2013.10.009

Get rights and content

"Samples frozen with and without glycerol as cryoprotectant indicated a major loss of *Bacteroidetes* in unprotected samples"

Sample Collection

**DNA Extraction**

What protocol are you using to extract DNA?

Does extraction method isolate all organisms equally?

How is contamination being controlled for?

Marker Gene

Shotgun

Taxonomic analysis

Taxonomic and functional analysis

---



# DNA extraction techniques can introduce bias

Wesolowska-Andersen et al. Microbiome 2014, 2:19
http://www.microbiomejournal.com/content/2/1/19

**Microbiome**

**RESEARCH**                    **Open Access**

## Choice of bacterial DNA extraction method from fecal material influences community structure as evaluated by metagenomic analysis

Agata Wesolowska-Andersen[1], Martin Iain Bahl[2], Vera Carvalho[2], Karsten Kristiansen[3], Thomas Sicheritz-Pontén[1], Ramneek Gupta[1*] and Tine Rask Licht[2*]

"We observed significant differences in distribution of bacterial taxa depending on the method"
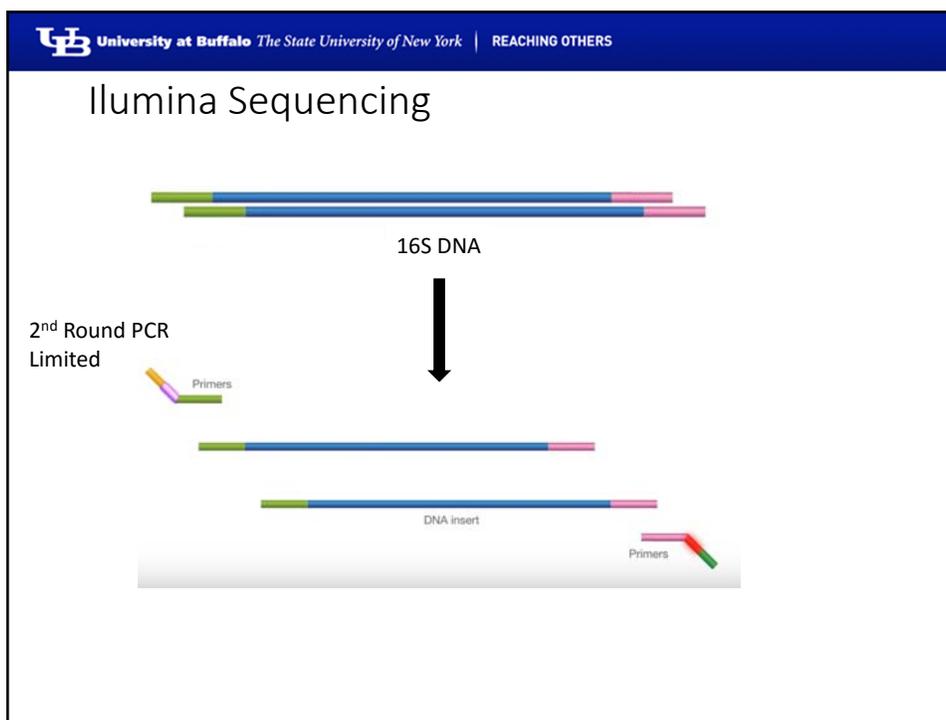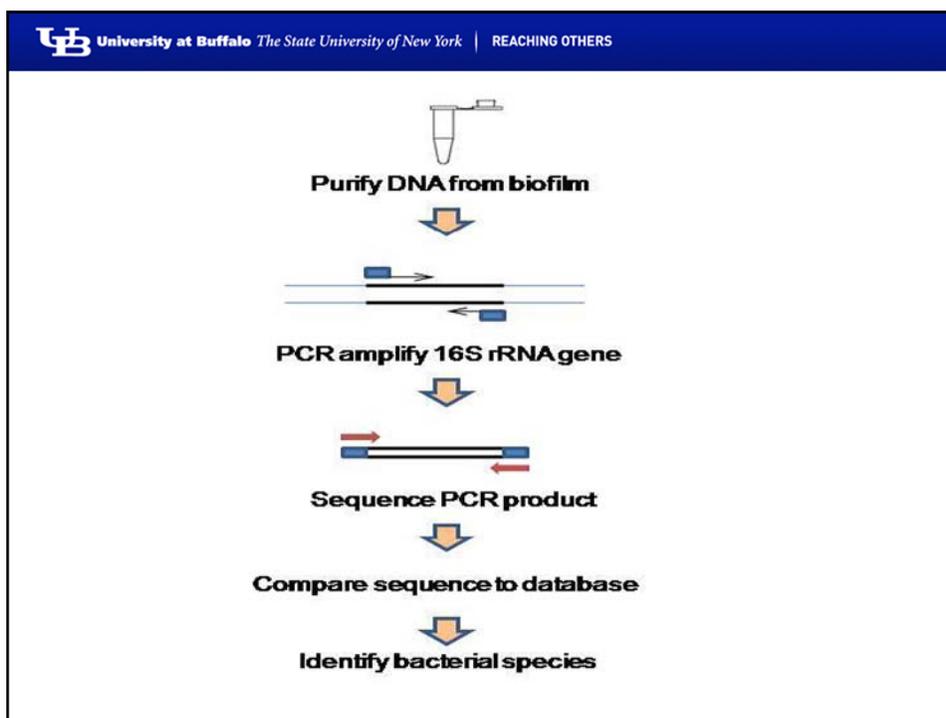
University at Buffalo *The State University of New York* | REACHING OTHERS

Sample Collection

↓

DNA Extraction

↓ ↓

Marker Gene    **Shotgun**    How are you going to analyze the data?

Is this overkill for your project?

↓ ↓

Taxonomic analysis    Taxonomic and functional analysis

---

University at Buffalo *The State University of New York* | REACHING OTHERS

Shotgun sequencing

- Isolate all DNA from a sample and sequence it
- Should be less biased compared to "marker studies"

**University at Buffalo** *The State University of New York* | REACHING OTHERS

**Overview of the pipeline used by EBI metagenomics to process raw sequence files and predict the functions and taxa present in a given sample.**

Sarah Hunter et al. Nucl. Acids Res. 2014;42:D600-D606

© The Author(s) 2013. Published by Oxford University Press.

**Nucleic Acids Research**



**University at Buffalo** *The State University of New York* | REACHING OTHERS

Sample Collection

DNA Extraction

Which marker gene are you going to use?

**Marker Gene**　　　Shotgun

How can you minimize amplification bias?

Taxonomic analysis　　　Taxonomic and functional analysis

**University at Buffalo** *The State University of New York* | REACHING OTHERS

# 16S rDNA

- 16S rDNA: is a component of the 30S small subunit of prokaryotic ribosomes.
- 16S rDNA: it satisfies the criteria of a marker by containing both highly conserved, ubiquitous sequences and regions that vary with greater or lesser frequency over evolutionary time.
- The products of the rRNA genes can fold into a complex, stable secondary structure, consisting of stems and loops. The sequences of some of the loops are conserved across nearly all bacterial species because of the essential functions involved, whereas the features of the structural parts are largely variant and specific to one or more classes.

---

**University at Buffalo** *The State University of New York* | REACHING OTHERS

# 16S rDNA

0  100  200  300  400  500  600  700  800  900  1000  1100  1200  1300  1400  1500 bp

| V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 |

**CONSERVED REGIONS:** unspecific applications
**VARIABLE REGIONS:** group or species-specific applications

University at Buffalo *The State University of New York* | REACHING OTHERS

Purify DNA from biofilm

PCR amplify 16S rRNA gene

Sequence PCR product

Compare sequence to database

Identify bacterial species



University at Buffalo *The State University of New York* | REACHING OTHERS

## Ilumina Sequencing

16S DNA

2nd Round PCR Limited

Primers

DNA insert

Primers

## Slide 1



Region complementary to flow cell oligo | Region same as flow cell oligo

Sequencing primer binding site 1 | Sequencing primer binding site 2

Index 2 | Index 1

16S DNA

Cluster on the flowcell

## Slide 2 — Illumina Sequencing



adapter
DNA fragment
dense lawn of primers
adapter

Bind single-stranded fragments randomly to the inside surface of the flow cell channels

1. Prepare 16S DNA
2. Attach DNA to surface
3. Bridge amplification
4. Fragments become double stranded
5. Denature the double-stranded molecules
6. Complete amplification

## Illumina Sequencing

1. Prepare genomic DNA

2. Attach DNA to surface

3. Bridge amplification

4. Fragments become double stranded

5. Denature the double-stranded molecules

6. Complete amplification

Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification

## Illumina Sequencing

Attached terminus    free terminus    Attached terminus

1. Prepare genomic DNA

2. Attach DNA to surface

3. Bridge amplification

4. Fragments become double stranded

5. Denature the double-stranded molecules

6. Complete amplification

The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate

**Illumina Sequencing**

1. Prepare genomic DNA

2. Attach DNA to surface

3. Bridge amplification

4. Fragments become double stranded

5. Denature the double-stranded molecules

6. Complete amplification

Attached

Attached

Denaturation leaves single-stranded templates anchored to the substrate



**Illumina Sequencing**

1. Prepare genomic DNA

2. Attach DNA to surface

3. Bridge amplification

4. Fragments become double stranded

5. Denature the double-stranded molecules

6. Complete amplification

Clusters

Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell

Illumina Sequencing

7. Determine first base

8. Image first base

9. Determine second base

10. Image second chemistry cycle

11. Sequencing over multiple chemistry cycles

12. Align data

The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase



Illumina Sequencing

7. Determine first base

8. Image first base

9. Determine second base

10. Image second chemistry cycle

11. Sequencing over multiple chemistry cycles

12. Align data

After laser excitation, the emitted fluorescence from each cluster is captured and the first base is identified

## UB NSG Core

**Illumina HiSeq 2500**

**MiSeq**



8 flow cell lanes with 2 flow cells
500-1000 Gb
2 billion reads per flow cell
2 x 250 bp max read length

1 flow cell lane
0.3 – 15 Gb
25 million reads
2 x 300 max read length
- 1 billion 35-100 bp

## Sample analysis

- Quality filter -- Was the sequencing good
- Paired-end sequence joining
- OTU calling
  - Reference-based
  - Non-reference based

- **Reference-based analysis will change over time as databases are updated.**

## OTU Table

Counts need to converted into relative frequency
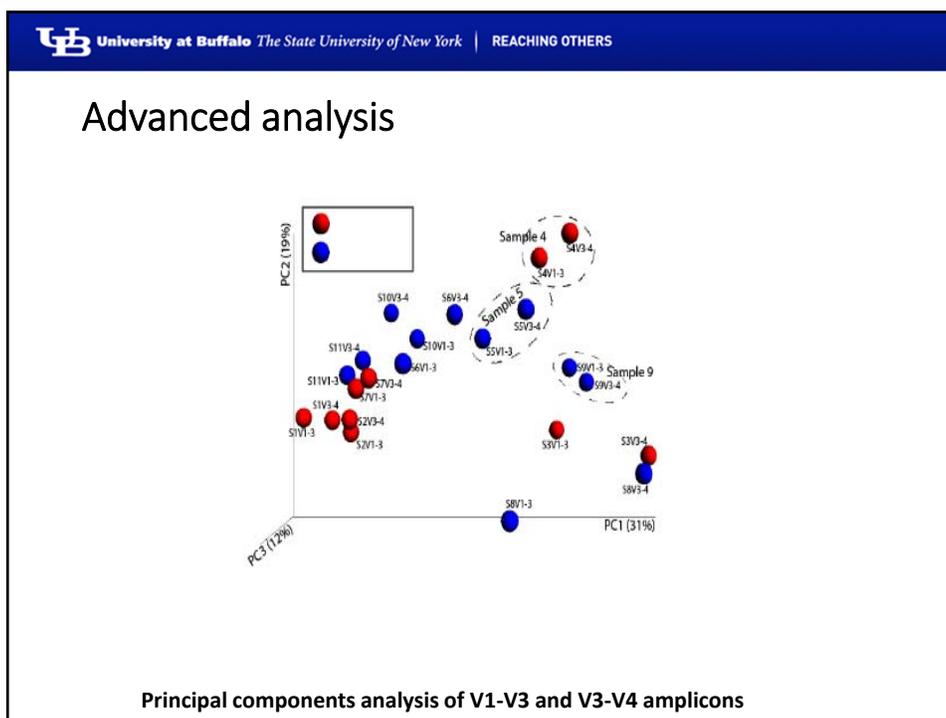
Frequency = (OTU count / reads per sample) *100

Total # of Reads per sample

Number of reads per sample will vary due to sequencing
>>>> Each sample needs to be normalized to each other <<<<



Statistical tests can then be performed

2/22/2016



Advanced analysis

Principal components analysis of V1-V3 and V3-V4 amplicons



How we do it in the NGS core

Sample Types

↓

DNA Extraction

↓

16S V1-V3 or V3-V4

↓

Taxonomic
analysis

18

Sample handling

**Sample Types**

DNA Extraction

16S V1-V3 or V3-V4

Taxonomic analysis

Saliva
Plaque
Fecal (Human, rat, mice)
Lavage

Samples are processed in our BSL-2 lab



Robotic DNA Extraction

Sample Types

**DNA Extraction**

16S V1-V3 or V3-V4
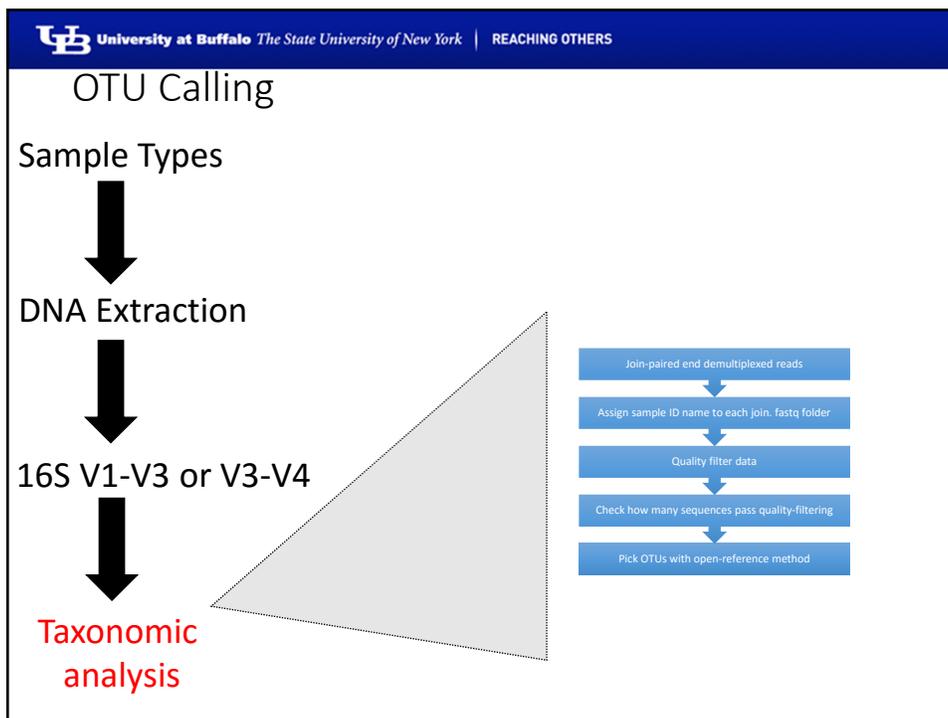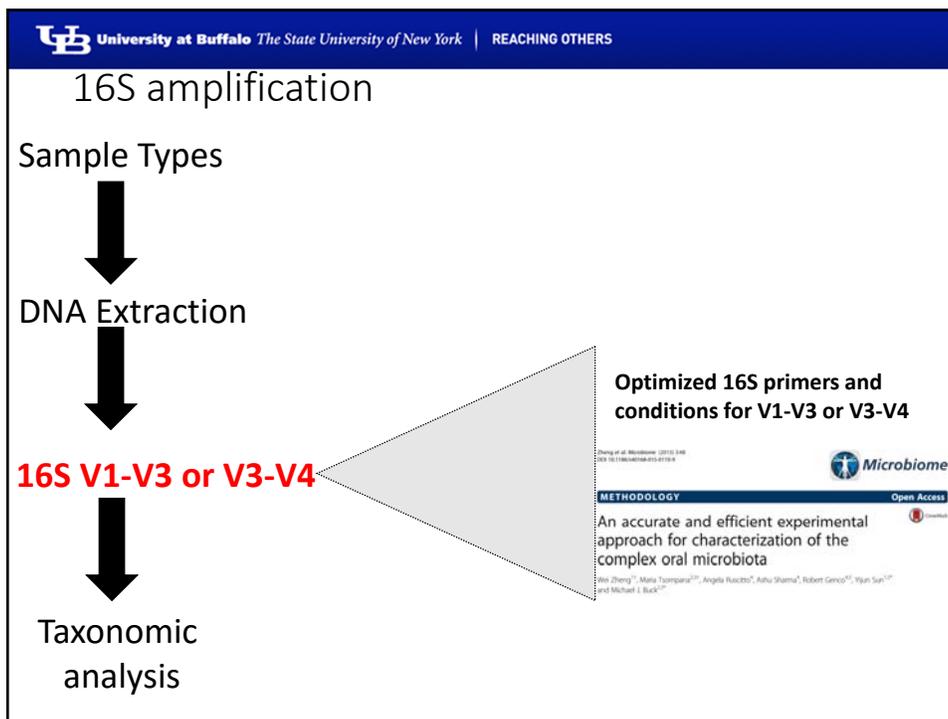
Taxonomic analysis

DNA is extracted 96 samples at time using a QIAGEN QIASYMPHONY

Samples are tracked by barcodes

**16S amplification**

Sample Types → DNA Extraction → **16S V1-V3 or V3-V4** → Taxonomic analysis

**Optimized 16S primers and conditions for V1-V3 or V3-V4**

An accurate and efficient experimental approach for characterization of the complex oral microbiota



**OTU Calling**

Sample Types → DNA Extraction → 16S V1-V3 or V3-V4 → Taxonomic analysis

Join-paired end demultiplexed reads
Assign sample ID name to each join. fastq folder
Quality filter data
Check how many sequences pass quality-filtering
Pick OTUs with open-reference method

University at Buffalo *The State University of New York* | REACHING OTHERS

## Limitations of 16S

- However, our ability to taxonomically characterize the microbiota using sequencing data is still restricted by the lack of universally accepted similarity thresholds, and the differential discriminatory power of the nine 16S rRNA hypervariable regions (V1-V9).
- Not all primer pairs work well for all genus/species. Amplification of non-representative genomic targets can heavily bias microbiome phylogenetic and diversity studies leading to inconclusive or inaccurate results.
- Requires PCR amplification, which can compress differences
- Does not capture viruses and eukaryotes (fungi)

University at Buffalo *The State University of New York* | REACHING OTHERS

## Bacterial DNA is everywhere!

- Sample collection tubes, collection liquids, processing liquids will all likely have low amount of bacterial DNA.
- Even **sterile** solutions have bacterial DNA!

- Good experimental design is essential

UB University at Buffalo *The State University of New York* | REACHING OTHERS
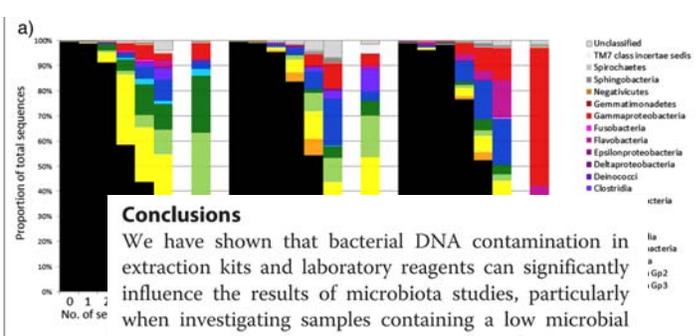
**BMC Biology**

**RESEARCH ARTICLE**      **Open Access**

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

Susannah J Salter[1*], Michael J Cox[2], Elena M Turek[2], Szymon T Calus[3], William O Cookson[2], Miriam F Moffatt[2], Paul Turner[4,5], Julian Parkhill[1], Nicholas J Loman[3] and Alan W Walker[1,6*]

- Sequence a pure culture of *Salmonella bongori*
- Extracted DNA using different kits
- Did serial dilutions of the pure culture to assess impact of contaminating species

---

UB University at Buffalo *The State University of New York* | REACHING OTHERS



**Conclusions**

We have shown that bacterial DNA contamination in extraction kits and laboratory reagents can significantly influence the results of microbiota studies, particularly when investigating samples containing a low microbial biomass. Such contamination is a concern for both 16S rRNA gene sequencing projects, which require targeted PCR amplification and enrichment, and also for shotgun metagenomic projects which do not. Awareness of this issue by the microbiota research community is important to ensure that studies are adequately controlled and erroneous conclusions are not drawn from culture-independent investigations.

**University at Buffalo** *The State University of New York* | REACHING OTHERS

# Acknowledgments

**UB Genomics & Bioinformatics Facility**

    Norma Nowak

    **Maria Tsompana**

    Sujith Valiyaparambil

    Natalie Waldron

    Jonathan Bard

    Brandon Marzullo

**WHI Microbiome Team**

    Jean Wactawski-Wende

    Robert Genco

    Mike Lamonte

    Amy Millen

    Chris Andrews

    Jo Freudenheim

    Yijun Sun

    Karen Falkner

    Kathy Hovey

    Wei Zheng

    Xiaodan Mai

---

**University at Buffalo** *The State University of New York* | REACHING OTHERS

Walker *et al. Microbiome*
DOI 10.1186/s40168-015-0087-4

**Microbiome**

**RESEARCH**      **Open Access**

# 16S rRNA gene-based profiling of the human infant gut microbiota is strongly influenced by sample processing and PCR primer choice

Alan W. Walker[1,2], Jennifer C. Martin[1], Paul Scott[2], Julian Parkhill[2], Harry J. Flint[1] and Karen P. Scott[1*]

**Abstract**

**Background:** Characterisation of the bacterial composition of the gut microbiota is increasingly carried out with a view to establish the role of different bacterial species in causation or prevention of disease. It is thus essential that the methods used to determine the microbial composition are robust. Here, several widely used molecular techniques were compared to establish the optimal methods to assess the bacterial composition in faecal samples from babies, before weaning.

**Results:** The bacterial community profile detected in the faeces of infants is highly dependent on the methodology used. Bifidobacteria were the most abundant bacteria detected at 6 weeks in faeces from two initially breast-fed babies using fluorescent in situ hybridisation (FISH), in agreement with data from previous culture-based studies. Using the 16S rRNA gene sequencing approach, however, we found that the detection of bifidobacteria in