# Meta'omic Analysis with MetaPhlAn & LEfSe

## Eric Franzosa

Postdoctoral Fellow / Huttenhower Lab

**Symposium and Workshop on New Methods**
**for Phylogenomics and Metagenomics**
The University of Texas at Austin
17 February 2013

Harvard School of Public Health
Department of Biostatistics

# Metagenomic Phylogenetic Analysis

Fast and accurate metagenomic profiling of microbial community composition using unique clade-specific marker genes

# LDA Effect Size

High-dimensional biomarker discovery and explanation

http://http://huttenhower.sph.harvard.edu/content/metaphlan-tutorial

# Tutorial Outline

▶ Introduction to MetaPhlAn

▶ MetaPhlAn Demo

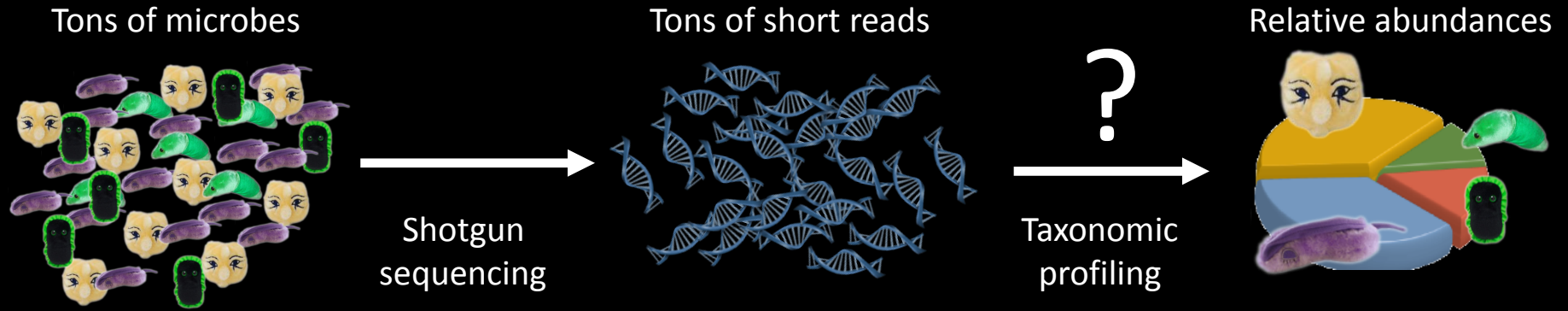▶ Introduction to LEfSe

▶ LEfSe Demo

▶ Links, other tools, and Q&A

# Tutorial Outline

▶ **Introduction to MetaPhlAn**

▶ MetaPhlAn Demo

▶ Introduction to LEfSe

▶ LEfSe Demo

▶ Links, other tools, and Q&A

# MetaPhlAn Motivation

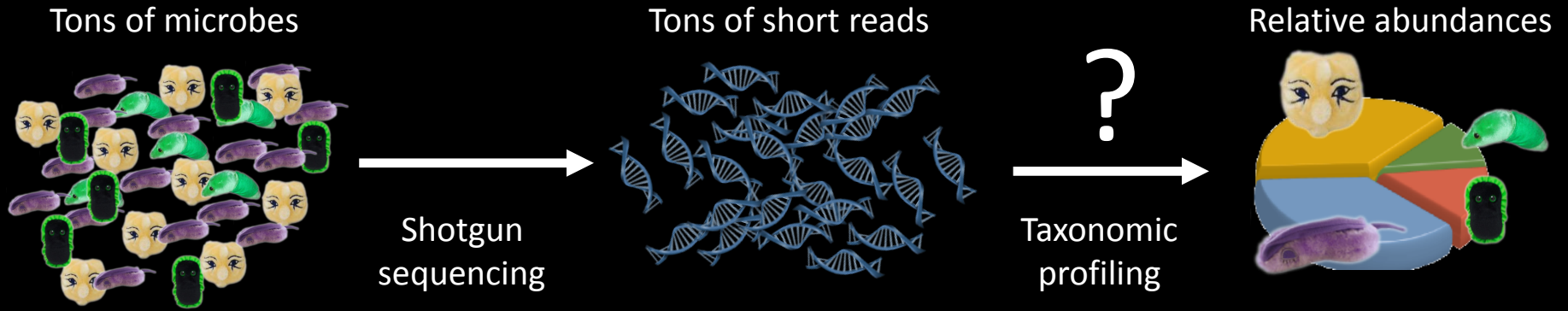Tons of microbes

Tons of short reads

Relative abundances

Shotgun
sequencing

?

Taxonomic
profiling

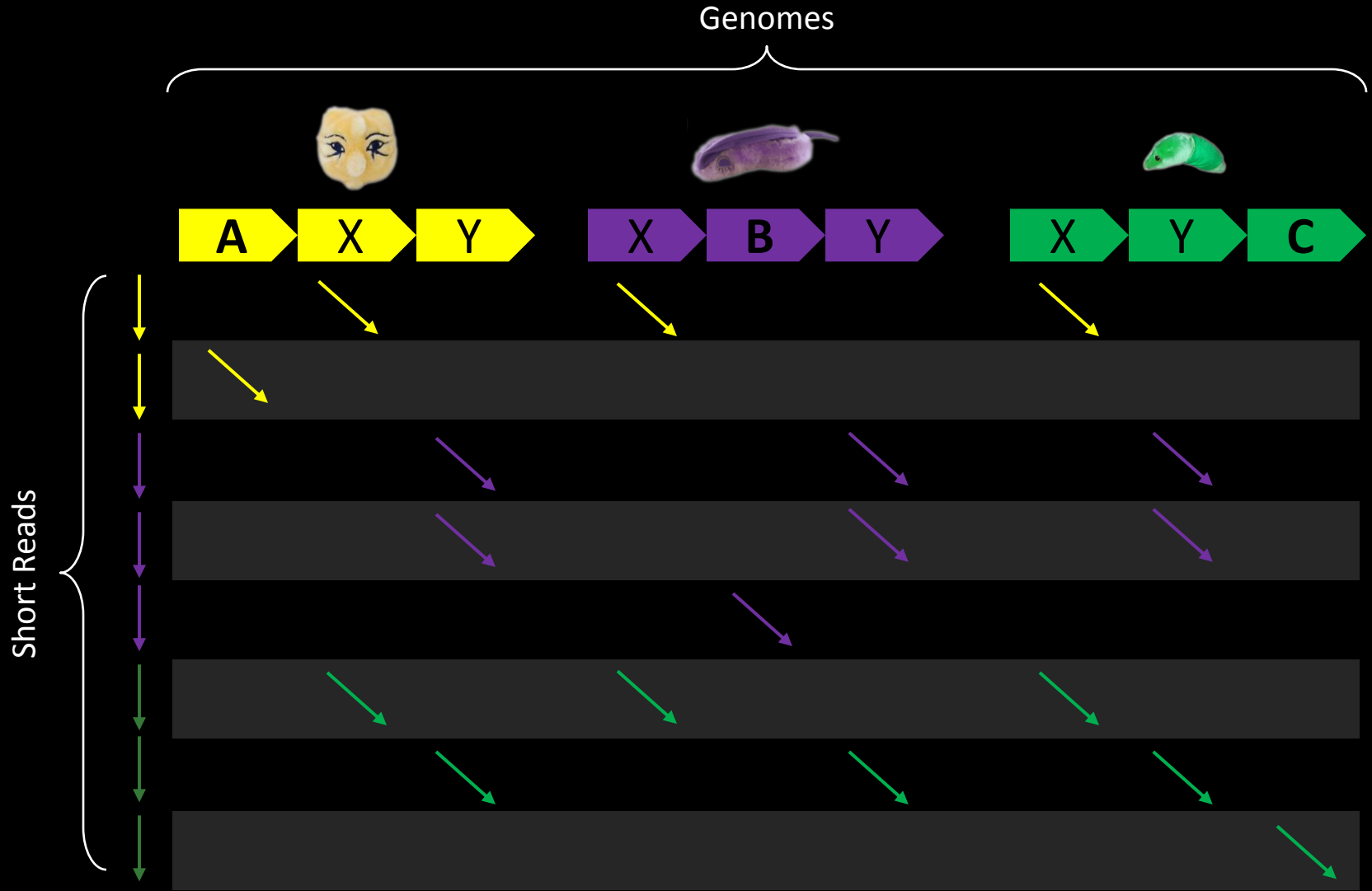**Profile the taxonomic composition of microbial communities from whole shotgun metagenomic data**

- Which clades (e.g. species, genera, classes) are there?
- What is the relative abundance of each clade in the community?

# MetaPhlAn Motivation

Tons of microbes        Tons of short reads        Relative abundances



Shotgun sequencing      **?**    Taxonomic profiling

**Profile the taxonomic composition of microbial communities from whole shotgun metagenomic data**

- Which clades (e.g. species, genera, classes) are there?
- What is the relative abundance of each clade in the community?

**Several challenges**

- Species-level resolution
- Computationally feasibility
- Organismal relative abundance rather than DNA concentrations
- Consistent detection confidence for all clades, including archaea
- High accuracies for very short reads (as short as ~50nt)
- Detection of organisms without sequenced genomes at higher taxonomic levels

# MetaPhlAn Overview
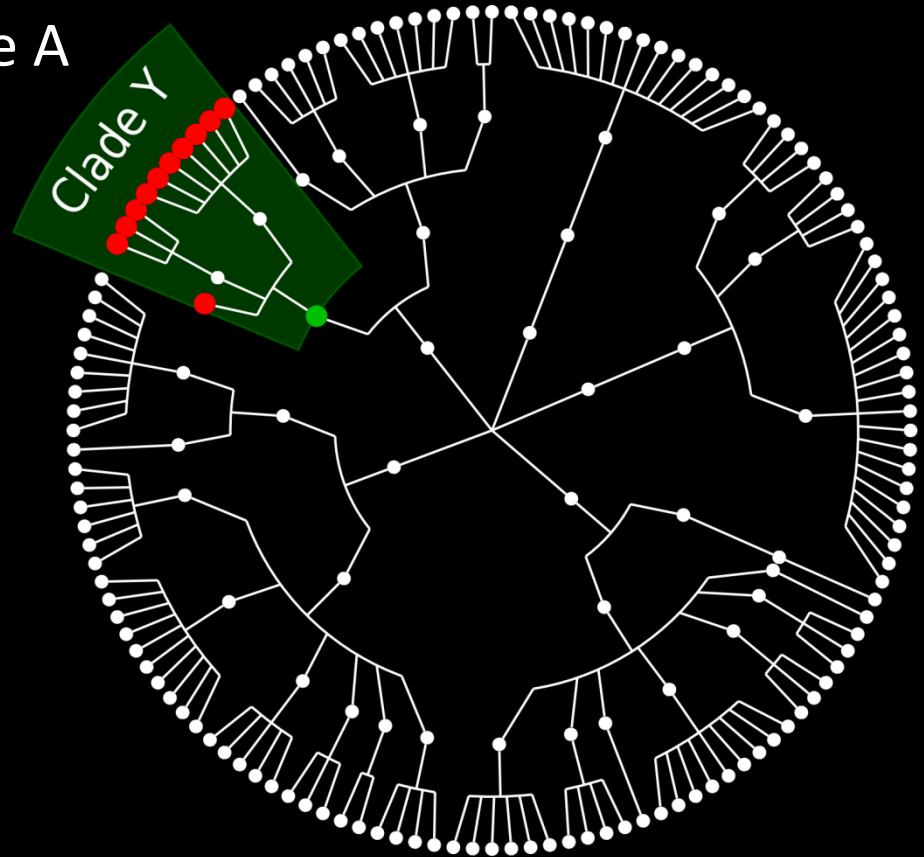
Genomes

Short Reads

A Y X B Y Y C

A X X B X Y C

# MetaPhlAn Overview

A is a core gene for clade Y
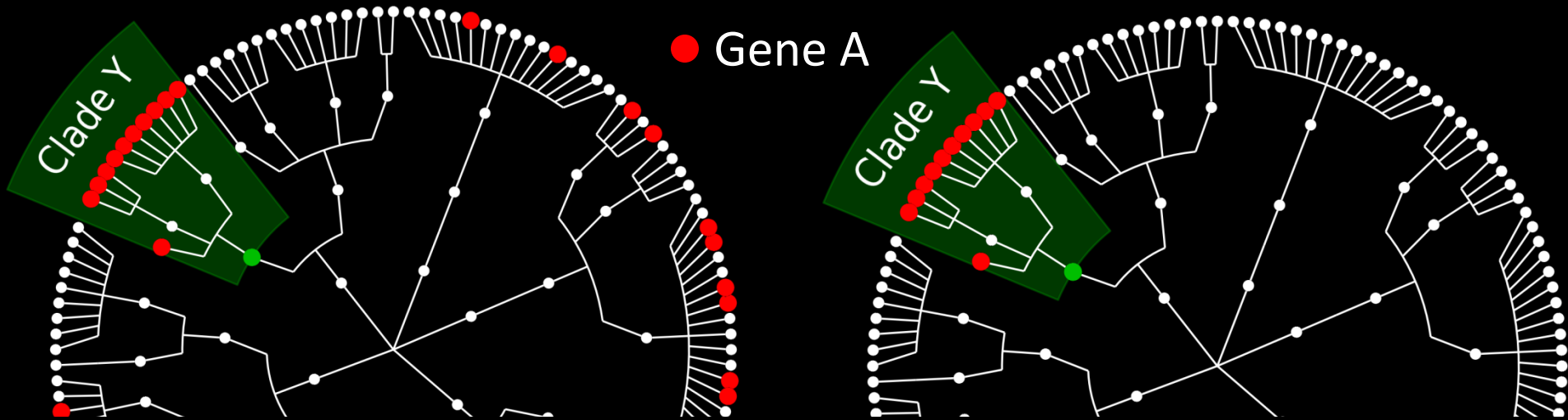
A is a unique marker gene for clade Y

● Gene A

# MetaPhlAn Overview

A is a core gene for clade Y

A is a unique marker gene for clade Y

● Gene A



ChocoPhlAn (offline pipeline)

Unique marker genes DB

- Identify all **core genes** for all clades
- Screen core genes for unique **marker genes**
- Select most representative marker genes

Available reference genomes

MetaPhlAn

Metagenome

- Blast reads agains the marker genes
- Assign, count, normalize reads

# Tutorial Outline

▶ Introduction to MetaPhlAn

▶ **MetaPhlAn Demo**

▶ Introduction to LEfSe

▶ LEfSe Demo

▶ Links, other tools, and Q&A

# MetaPhlAn Demo: Setup

▶ Download program and marker database from:
http://huttenhower.sph.harvard.edu/metaphlan

▶ Requires python with **numpy** installed

▶ Requires **BLAST** or **Bowtie2** for alignment
(**Bowtie2** recommended)

▶ Today's sample data available at:
http://huttenhower.sph.harvard.edu/content/metaphlan-tutorial

# MetaPhlAn (command line)

▶ Show minimal MetaPhlAn setup

▶ Run MetaPhlAn on downsampled HMP FASTA

▶ Examine marker output file

▶ Example abundance output file

▶ Discuss taxonomic organization

▶ Discuss "unclassified" species

# MetaPhlAn (command line)

# MetaPhlAn (command line)

# MetaPhlAn (command line)

# MetaPhlAn (command line)

# MetaPhlAn (command line)

# Tutorial Outline

▶ Introduction to MetaPhlAn

▶ MetaPhlAn Demo

▶ **Introduction to LEfSe**

▶ LEfSe Demo

▶ Links, other tools, and Q&A

# Turning anecdotes into biology

| Site | Oral | Gut |
|---:|:---:|:---:|
| Clade1 | 0.40 | 0.87 |
| Clade1\|Bug1 | 0.40 | 0.56 |
| Clade1\|Bug2 | 0.00 | 0.30 |
| Clade2 | 0.60 | 0.13 |
| Clade2\|Bug3 | 0.11 | 0.00 |
| Clade2\|Bug4 | 0.49 | 0.13 |

# Turning anecdotes into biology

- More samples are a start...
- But the statistics are still non-trivial
  - Data are noisy
  - Compositional nature ($\Sigma = 1$)
  - High dynamic range
  - Hierarchical organization

| Site | Oral | Gut | Oral | Gut | Oral | Gut |
|---|---|---|---|---|---|---|
| Clade1 | 0.40 | 0.87 | 0.43 | 0.68 | 0.47 | 0.32 |
| Clade1\|Bug1 | 0.40 | 0.56 | 0.07 | 0.31 | 0.42 | 0.27 |
| Clade1\|Bug2 | 0.00 | 0.30 | 0.36 | 0.37 | 0.04 | 0.05 |
| Clade2 | 0.60 | 0.13 | 0.57 | 0.32 | 0.53 | 0.68 |
| Clade2\|Bug3 | 0.11 | 0.00 | 0.10 | 0.32 | 0.15 | 0.23 |
| Clade2\|Bug4 | 0.49 | 0.13 | 0.47 | 0.00 | 0.39 | 0.45 |

# Turning anecdotes into biology

- We may also want to leverage (or control for) additional metadata

| Sample # | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| **Profession** | Student | Postdoc | Postdoc | Professor | Student | Student |
| **Gender** | Male | Female | Female | Male | Male | Female |
| **Site** | Oral | Gut | Oral | Gut | Oral | Gut |
| **Clade1** | 0.40 | 0.87 | 0.43 | 0.68 | 0.47 | 0.32 |
| **Clade1\|Bug1** | 0.40 | 0.56 | 0.07 | 0.31 | 0.42 | 0.27 |
| **Clade1\|Bug2** | 0.00 | 0.30 | 0.36 | 0.37 | 0.04 | 0.05 |
| **Clade2** | 0.60 | 0.13 | 0.57 | 0.32 | 0.53 | 0.68 |
| **Clade2\|Bug3** | 0.11 | 0.00 | 0.10 | 0.32 | 0.15 | 0.23 |
| **Clade2\|Bug4** | 0.49 | 0.13 | 0.47 | 0.00 | 0.39 | 0.45 |

# **LEfSe**: *finding metagenomic biomarkers*

# Tutorial Outline

▶ Introduction to MetaPhlAn

▶ MetaPhlAn Demo

▶ Introduction to LEfSe

▶ **LEfSe Demo**

▶ Links, other tools, and Q&A

# LEfSe Demo: Galaxy Version



Available at: http://huttenhower.sph.harvard.edu/galaxy/

# Load your data table



Tab delimited text, consisting of a class (e.g. oral vs. gut), an optional subclass (e.g. male vs. female), an optional sample ID, and then your features.

# My input to Galaxy...

| Site | Oral | Gut | Oral | Gut | Oral | Gut |
|---|---|---|---|---|---|---|
| Clade1 | 0.40 | 0.87 | 0.43 | 0.68 | 0.47 | 0.32 |
| Clade1\|Bug1 | 0.40 | 0.56 | 0.07 | 0.31 | 0.42 | 0.27 |
| Clade1\|Bug2 | 0.00 | 0.30 | 0.36 | 0.37 | 0.04 | 0.05 |
| Clade2 | 0.60 | 0.13 | 0.57 | 0.32 | 0.53 | 0.68 |
| Clade2\|Bug3 | 0.11 | 0.00 | 0.10 | 0.32 | 0.15 | 0.23 |
| Clade2\|Bug4 | 0.49 | 0.13 | 0.47 | 0.00 | 0.39 | 0.45 |

- MetaPhlAn output for many HMP subjects
- Buccal_mucosa *is a proxy for* oral
- Stool *is a proxy for* gut
- Input file included with the sample data

# My input to Galaxy…

| Site | Oral | Gut | Oral | Gut | Oral | Gut |
|---|---|---|---|---|---|---|
| **Clade1** | 0.40 | 0.87 | 0.43 | 0.68 | 0.47 | 0.32 |
| **Clade1\|Bug1** | 0.40 | 0.56 | 0.07 | 0.31 | 0.42 | 0.27 |
| **Clade1\|Bug2** | 0.00 | 0.30 | 0.36 | 0.37 | 0.04 | 0.05 |
| **Clade2** | 0.60 | 0.13 | 0.57 | 0.32 | 0.53 | 0.68 |
| **Clade2\|Bug3** | 0.11 | 0.00 | 0.10 | 0.32 | 0.15 | 0.23 |
| **Clade2\|Bug4** | 0.49 | 0.13 | 0.47 | 0.00 | 0.39 | 0.45 |

- MetaPhlAn output for many HMP subjects
- Buccal_mucosa *is a proxy for* oral
- Stool *is a proxy for* gut
- Input file included with the sample data

- Features don't have to come from MetaPhlAn
- Any tab-delimited continuous data will work!

# Tell LEfSe which rows are class/subclass/ID



We're just using class (STSite, Stool versus Buccal_mucosa)
Set the other boxes to "no subclass" and "no subject"

# Run the algorithm



Because the oral and gut communities are quite different, I adjusted the threshold from the default of 2 to 4.5 (logs) to only show the most extreme differences.

# Visualize the results

# LEfSe Output: *Raw Numbers*

k__Bacteria.p__Bacteroidetes
class: Buccal_mucosa
class: Stool

# LEfSe at the command line

▶ Source code available at:
https://bitbucket.org/nsegata/lefse

▶ Requires python and R with several additional packages installed (see README)

▶ Today's sample data available at:
http://huttenhower.sph.harvard.edu/content/metaphlan-tutorial

▶ We'll look at a second example included with the LEfSe distribution involving subclasses

# LEfSe (command line) demo

▶ Show LEfSe setup

▶ Show example folder

▶ Highlight LEfSe commands

▶ Show output

▶ Highlight more complicated subclass treatment

# LEfSe (command line) demo

# Tutorial Outline

▶ Introduction to MetaPhlAn

▶ MetaPhlAn Demo

▶ Introduction to LEfSe

▶ LEfSe Demo

▶ **Links, other tools, and Q&A**

# If you're thirsty for more…

▶ Tutorial materials available here:
http://huttenhower.sph.harvard.edu/content/metaphlan-tutorial

▶ My email (Eric Franzosa):
franzosa@hsph.harvard.edu

# If you're thirsty for more…

▶ For more about our tools, check out:
http://huttenhower.sph.harvard.edu/research

▶ Individual tool pages contain links to source code and documentation…

▶ …or try them out on Galaxy:
http://huttenhower.sph.harvard.edu/galaxy/

# If you're thirsty for more…

▶ `MetaPhlAn 2.0` coming soon with support for viruses, eukaryotes, and more bacteria/archaea

▶ `PhyloPhlAn` uses MetaPhlAn's **core gene** identification pipeline to assist in tree building
http://huttenhower.sph.harvard.edu/phylophlan

▶ `MaAsLin` is an evolution of LEfSe that can consider a wider number and variety of metadata
http://huttenhower.sph.harvard.edu/maaslin

# Thank you!

**Curtis Huttenhower**
Levi Waldron
Xochi Morgan
Tim Tickle

**Dirk Gevers**
Kat Huang

## Human Microbiome Project

| | |
|---|---|
| Owen White | Sahar Abubucker |
| Joe Petrosino | Brandi Cantarel |
| George Weinstock | Alyx Schubert |
| Karen Nelson | Mathangi Thiagarajan |
| Lita Proctor | Beltran Rodriguez-Mueller |
| Erica Sodergren | Makedonka Mitreva |
| Anthony Fodor | Yuzhen Ye |
| Marty Blaser | Mihai Pop |
| Jacques Ravel | Larry Forney |
| Pat Schloss | Barbara Methe |

Vagheesh Narasimhan
Emma Schwager
**Nicola Segata**
Daniela Boernigen

**Ramnik Xavier**
**Harry Sokol**
**Dan Knights**
Moran Yassour

Bruce Birren   Mark Daly
Doyle Ward   Ashlee Earl
**BROAD INSTITUTE**

Joseph Moon
Jim Kaminski
Craig Bielski

Brian Palmer
**Ren Lu**
Hufeng Zhou

Wendy Garrett
Michelle Rooks

**Rob Beiko**
*Morgan Langille*

**Rob Knight**
*Greg Caporaso*
*Jesse Zaneveld*

**Mark Silverberg**
*Boyko Kabakchiev*
*Andrea Tyler*

Ruth Ley
Omry Koren

Jacques Izard
Katherine Lemon

**Bruce Sands**

**http://huttenhower.sph.harvard.edu**

43

# If you're thirsty for more...
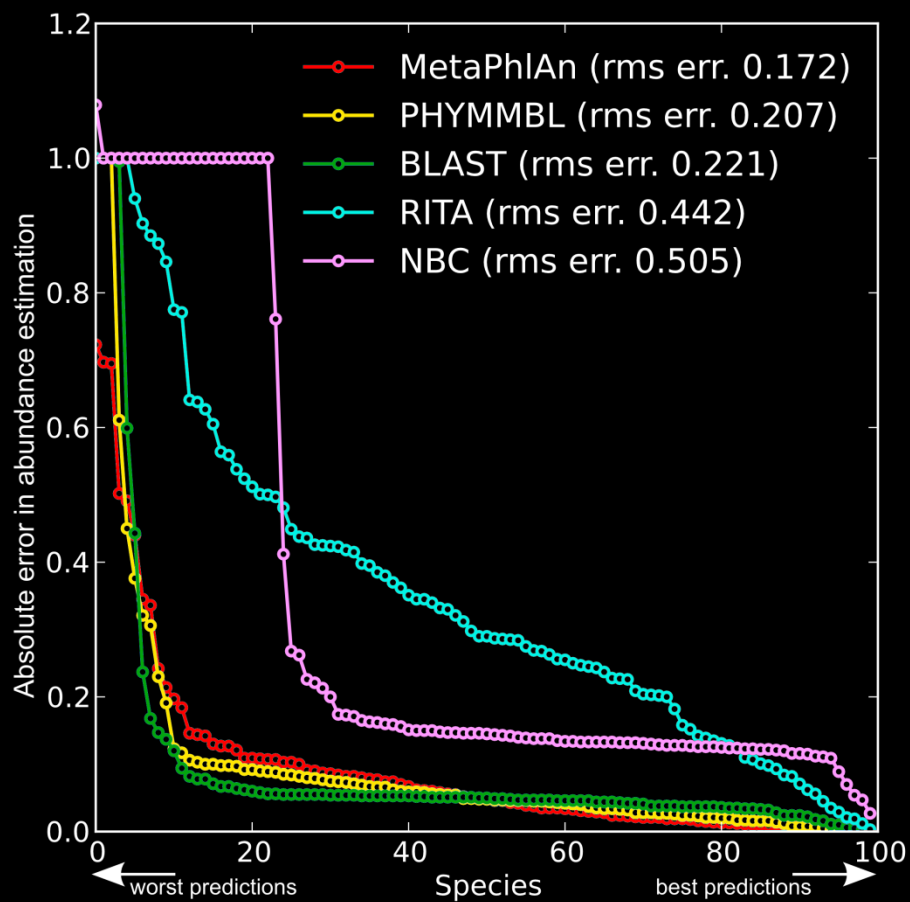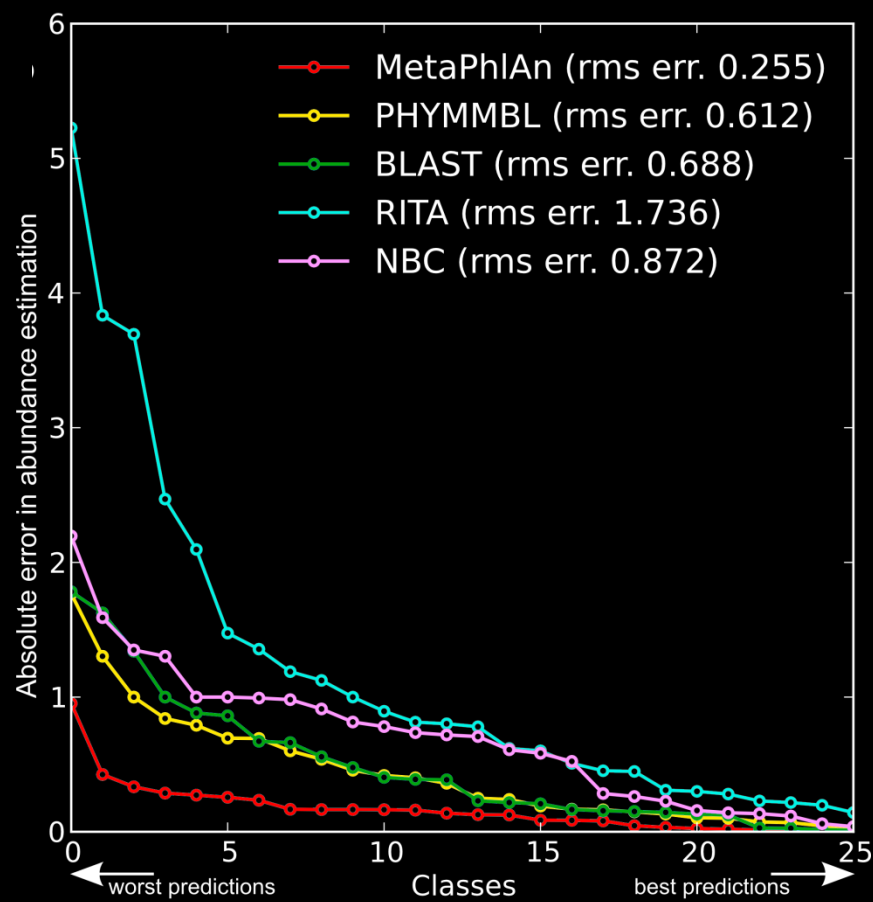
▶ Questions?

# MetaPhlAn Statistics

- Representing 2,887 genomes (107 Archaea)

- 1,222 species, 652 genera, 278 families, 130 orders, 66 classed, 33 phyla, 2,383 total clades

- ≈2M total unique marker genes

- ≈400k "most representative" unique marker genes

- 231±107 markers per species (350 fixed max)

- The 400k database represent ≈4% of the total microbial sequence data available

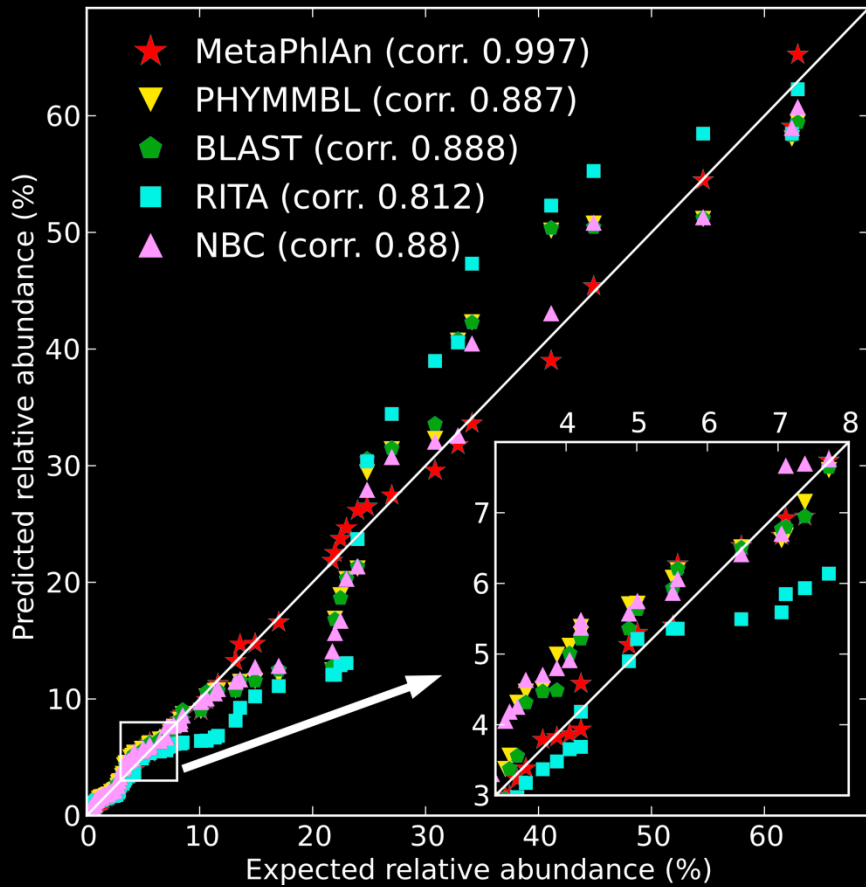# Evaluation of accuracy (1)
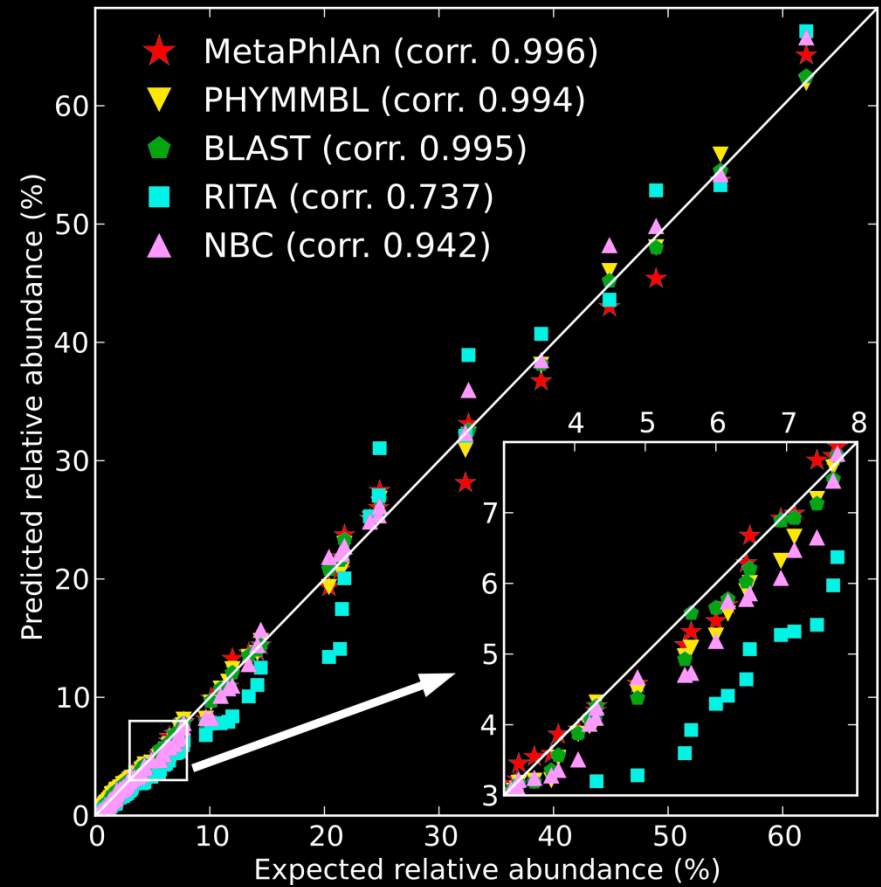


Species

Classes

(Validation on high-complexity uniformly distributed synthetic metagenomes.)
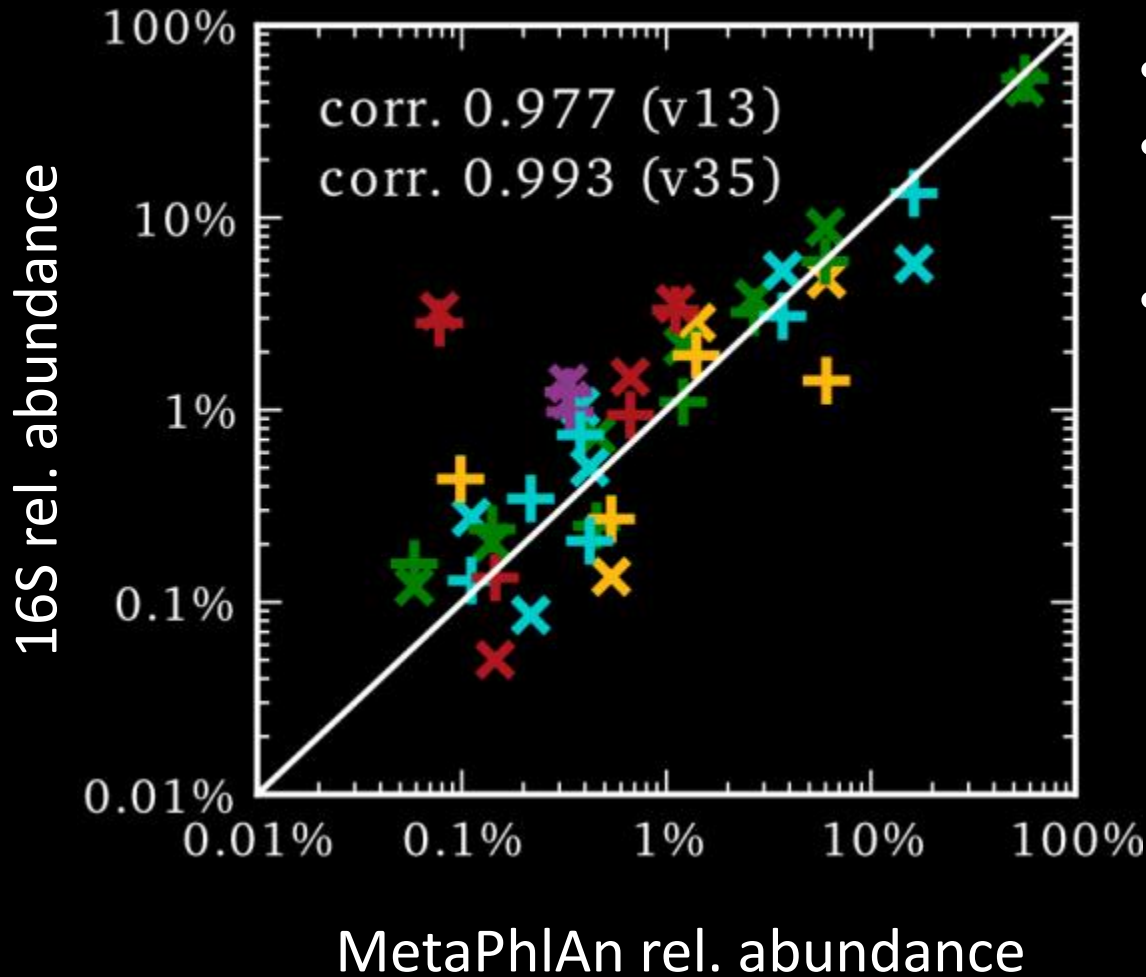
# **Evaluation of accuracy (2)**



Species

Classes

(Validation on low-complexity log-normally distributed synthetic metagenomes.)
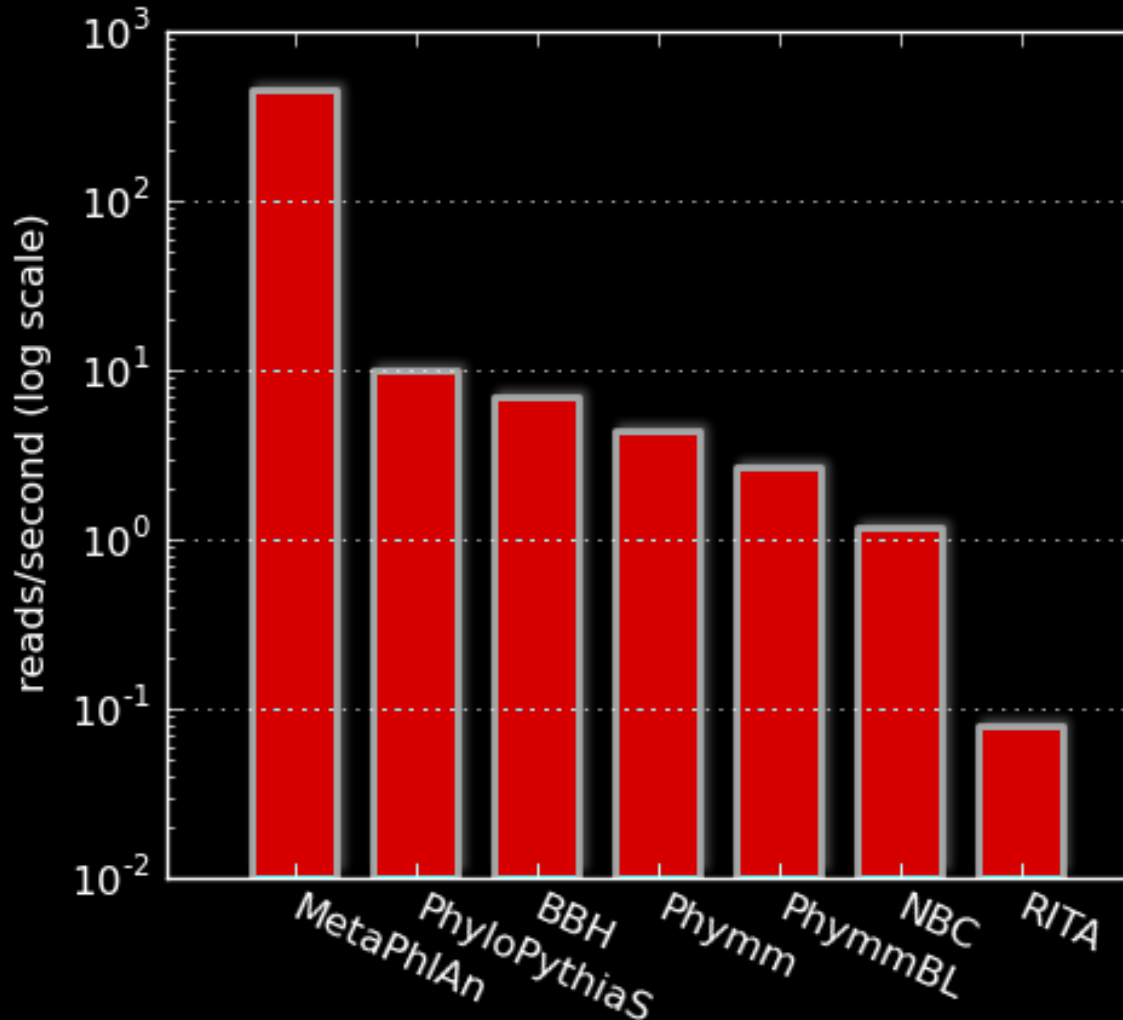
# Evaluation of accuracy (3)



- *Buccal mucosa example*
- Genus-level comparison with 16S-based estimation
- v13 ( + ) & v35 ( × ) regions

# Evaluation of performance



>50 times faster than existing methods

450 reads/sec (BLAST)

Up to 25,000 reads/sec (bowtie2)

Multi-threaded

Easily parallelizable