# Introduction to metagenomic analysis

Eric A. Franzosa, Ph.D.
Galeb Abu-Ali, Ph.D.

Harvard University CFAR Workshop on
Metagenomics and Transcriptomics

16 September 2014

Huttenhower Research Group
Harvard School of Public Health
Department of Biostatistics

MiRiBA

BROAD INSTITUTE

# Content

- http://huttenhower.sph.harvard.edu/content/cfar2014

# Plan

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
    - MetaPhlAn
- Tools for functional profiling
    - HUMAnN
    - ShortBRED
    - PICRUSt
- Tools for testing associations
    - LEfSe
    - MaAsLin
    - CCREPE
- Resources
- Research vignette (time permitting)

# What's metagenomics?

Total collection of **microorganisms** within a **community**

Also **microbial community** or **microbiota**

THE MICROFLORA AND THE PRODUCTIVITY OF LEACHED AND NON-LEACHED ALKALI SOIL

J. E. GREAVES[1]

*Utah Agricultural Experiment Station*

Received for publication July 2, 1926

Chemistry & Biology October 1998, 5:R245–249

**Molecular biological access to the chemistry of unknown soil microbes: a new frontier for natural products**

Jo Handelsman[1], Michelle R Rondon[1], Sean F Brady[2], Jon Clardy[2] and Robert M Goodman[1]

ness that they are thought to contain. The methodology has been made possible by advances in molecular biology and eukaryotic genomics, which have laid the groundwork for cloning and functional analysis of the collective genomes of soil microflora, which we term the metagenome of the soil.

Total **genomic potential** of a microbial community

www.sciencemag.org SCIENCE VOL 292 11 MAY 2001

**Commensal Host-Bacterial Relationships in the Gut**

Lora V. Hooper and Jeffrey I. Gordon*

ber our somatic and germ cells (3). The Nobel laureate Joshua Lederberg has suggested using the term "microbiome" to describe the collective genome of our indigenous microbes (microflora), the idea being that a comprehensive genetic view of *Homo sapiens* as a life-form should include the genes in our microbiome (4).

Study of **uncultured microorganisms** from the environment, which can include humans or other living hosts

Total **biomolecular repertoire** of a microbial community

4

# Examples of metagenomic studies: Global Ocean Sampling



J. Craig Venter INSTITUTE

2003/2004 - ongoing



The Sorcerer II Expedition
Global Ocean Sampling Route



The Biodiversity of Each New Region is Different



Proteorhodopsins Vary by Region



JTC Sequencer Lab
Capacity: 240,000 sequences/day or 80 million lanes/year at 24 runs per day

The NIH **H**uman **M**icrobiome **P**roject (**HMP**):
A comprehensive microbial survey

- ***What is a "normal" human microbiome?***
- 300 healthy human subjects
- Multiple body sites
  - 15 male, 18 female
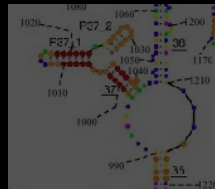- Multiple visits
- Clinical metadata



*www.hmpdacc.org*

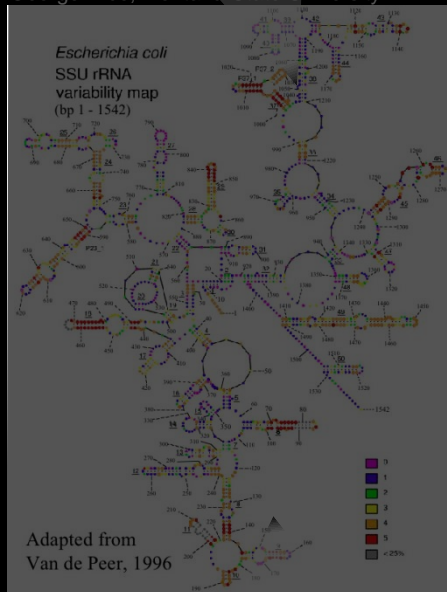# Sequencing as a tool for microbial community analysis

Lyse cells
Extract DNA (and/or RNA)

**Meta'omic**

**16S amplicons**

George Rice, Montana State University

*Escherichia coli*
SSU rRNA
variability map
(bp 1 - 1542)

Adapted from
Van de Peer, 1996
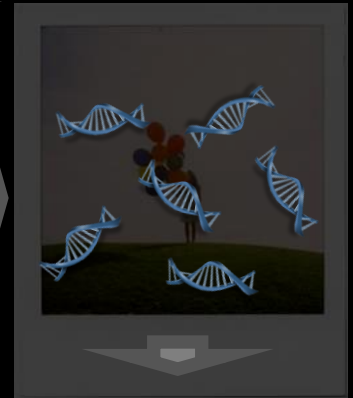
PCR to amplify the single
16S rRNA marker gene

Hello
my name is

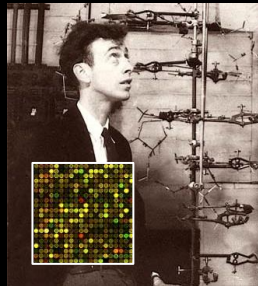Classify sequence
→ microbe

Samples

Microbes

Relative
abundances

Genes,
Genomes,
Metabolic profiling,
Relative abundances,
Genetic variants...

# What to do with your metagenome?

Reservoir of gene and protein functional information

Comprehensive snapshot of microbial ecology and evolution

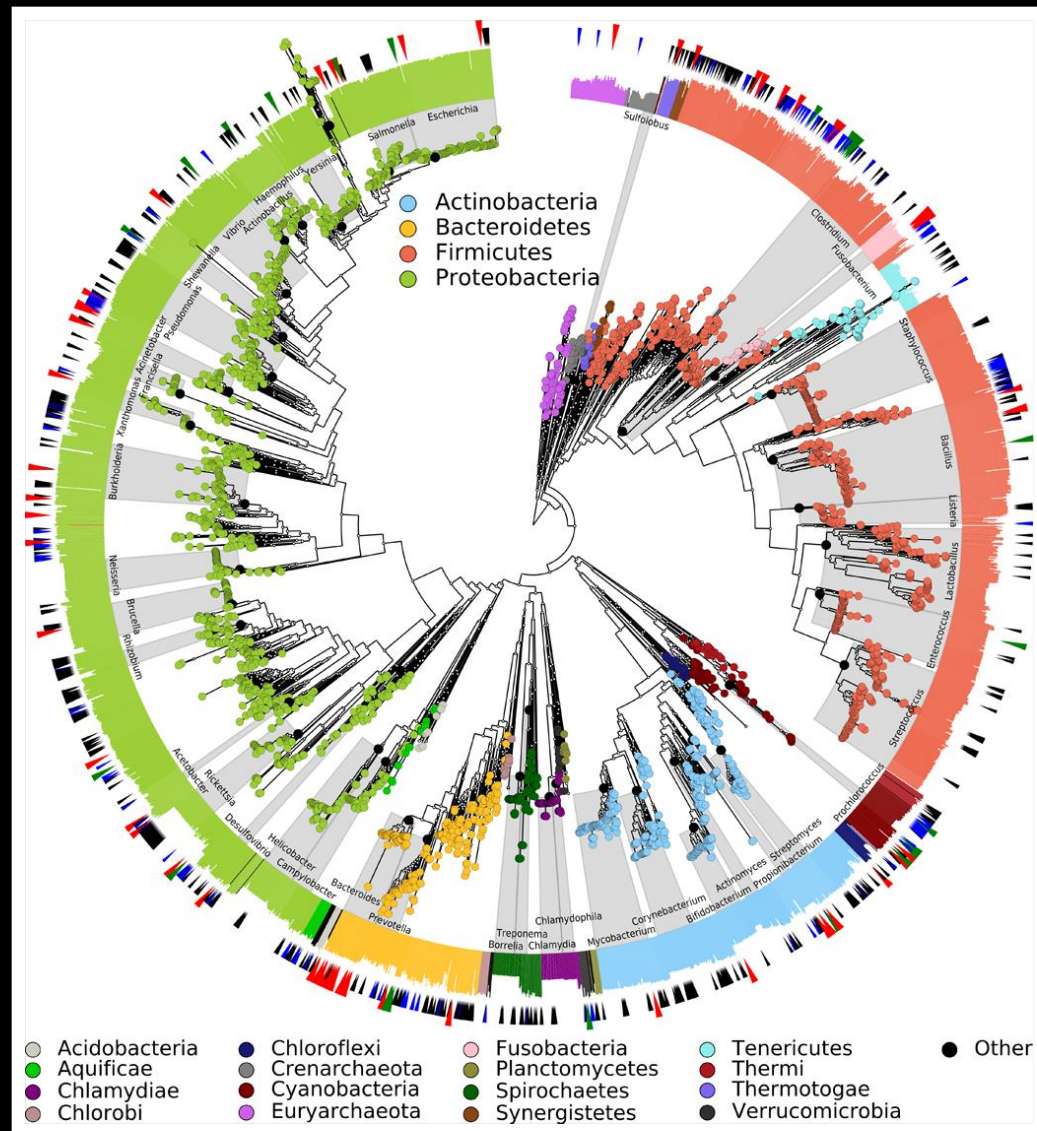## Basic science

## Translational

Public health tool monitoring population health and epidemiology

Diagnostic or prognostic biomarker for host disease
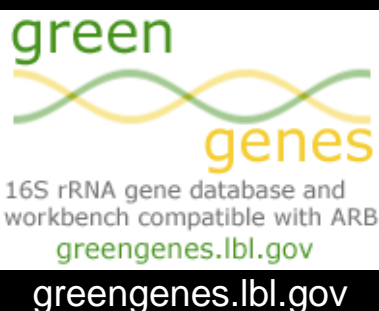
# Composition-based analyses

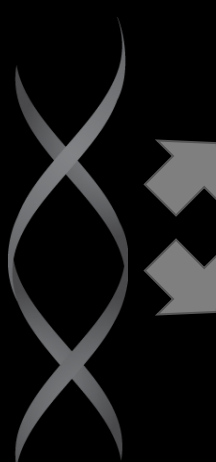# Microbiome composition analyses: **phylotypes** and **binning**


rdp.cme.msu.edu


greengenes.lbl.gov


www.arb-silva.de



- Actinobacteria
- Bacteroidetes
- Firmicutes
- Fusobacteria
- Other Proteobacteria
- Verrucomicrobia
- Alphaproteobacteria
- Betaproteobacteria
- Deltaproteobacteria
- Epsilonproteobacteria
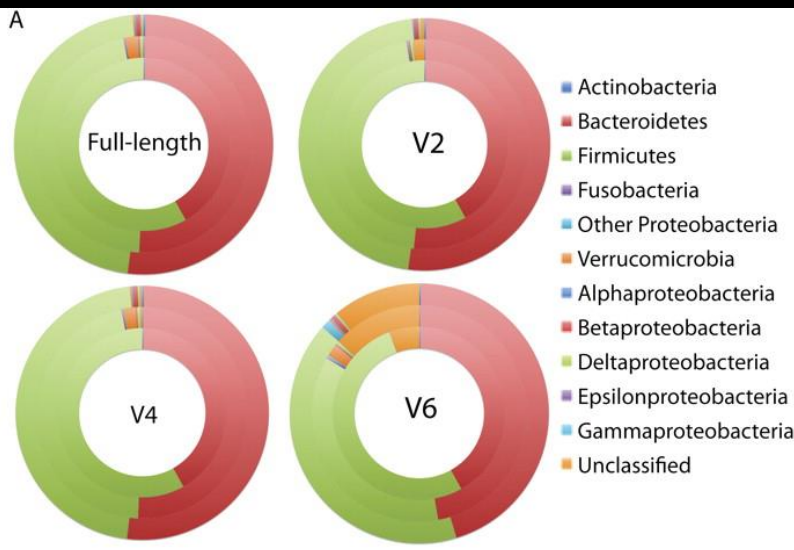- Gammaproteobacteria
- Unclassified

**Binning**: nontrivial assignment of reads to phylotypes

**Phylotype** or **operational taxonomic unit (OTU)**: organisms clonal to within some tolerance (e.g. 95%); "species"

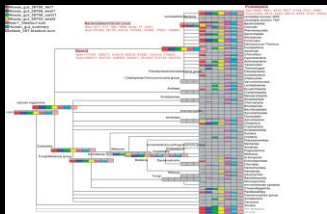Indirect binning: BLAST etc. Relies on high similarity, reference seq.

Direct binning: analyzes seq. characteristics (GC, codons, etc.) Relies on long reads
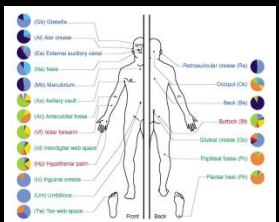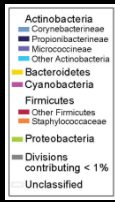
Hamady, 2009     PhyloPythia: McHardy, 2007     TETRA: Chan, 2008     Phymm: Brady, 2009     MetaPhlAn: Segata, 2012

# Microbiome composition analyses: **diversity**

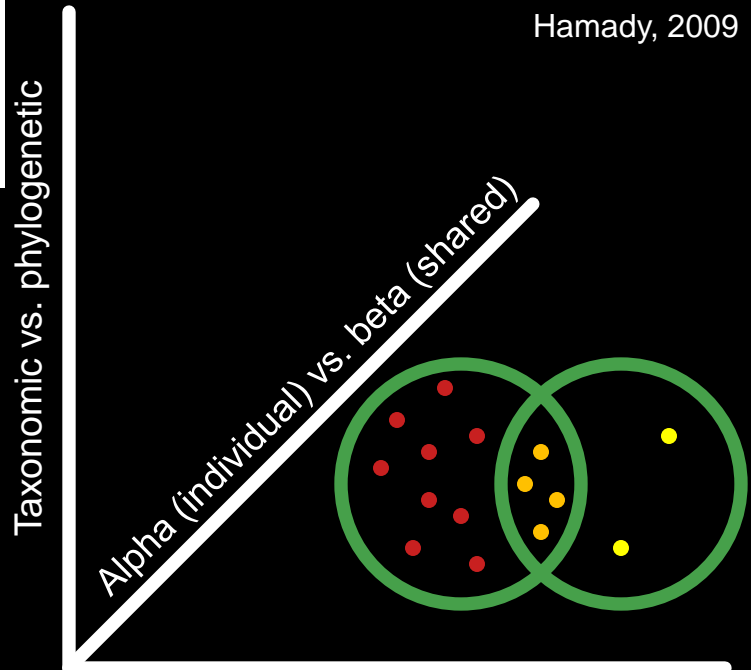**Diversity**: broadly, a community's number and distribution of organisms
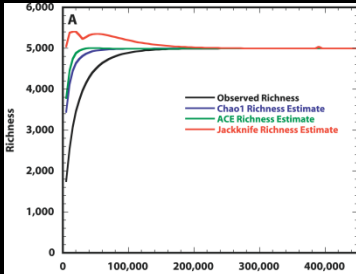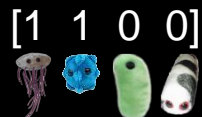
Also **community composition** or **structure**

Mitra, 2009

Hamady, 2009

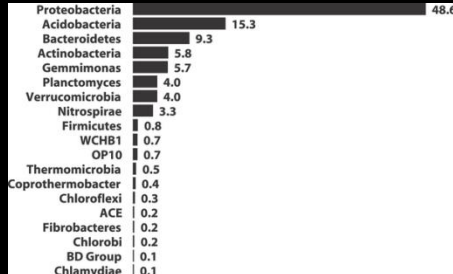Taxonomic vs. phylogenetic

Alpha (individual) vs. beta (shared)

Qualitative vs. quantitative

Susan Holmes, Stanford

Actinobacteria
- Corynebacterineae
- Propionibacterineae
- Micrococcineae
- Other Actinobacteria

Bacteroidetes
Cyanobacteria
Firmicutes
- Other Firmicutes
- Staphylococcaceae
Proteobacteria
Divisions contributing < 1%
Unclassified

[1  1  0  0]          [10  6  1  4]

Richness

- Observed Richness
- Chao1 Richness Estimate
- ACE Richness Estimate
- Jackknife Richness Estimate

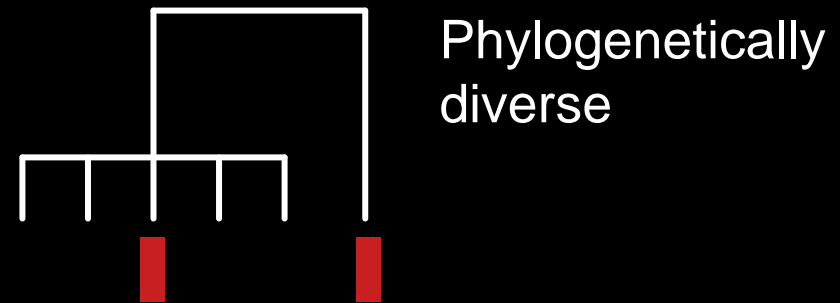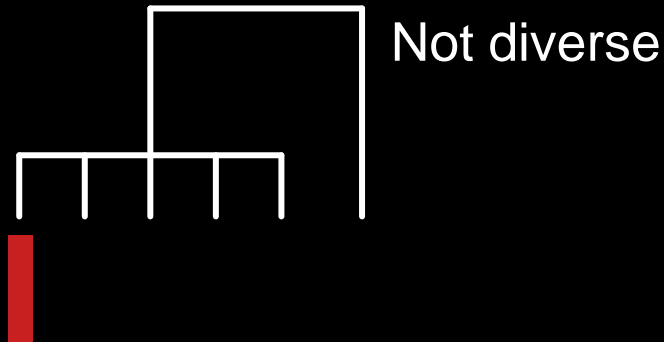| | % |
|---|---|
| Proteobacteria | 48.6 |
| Acidobacteria | 15.3 |
| Bacteroidetes | 9.3 |
| Actinobacteria | 5.8 |
| Gemmimonas | 5.7 |
| Planctomyces | 4.0 |
| Verrucomicrobia | 4.0 |
| Nitrospirae | 3.3 |
| Firmicutes | 0.8 |
| WCHB1 | 0.7 |
| OP10 | 0.7 |
| Thermomicrobia | 0.5 |
| Coprothermobacter | 0.4 |
| Chloroflexi | 0.3 |
| ACE | 0.2 |
| Fibrobacteres | 0.2 |
| Chlorobi | 0.2 |
| BD Group | 0.1 |
| Chlamydiae | 0.1 |

**# of sequences**      Schloss, 2006      **% of sequences**

# Microbiome composition analyses: **alpha diversity (1-sample) scenarios**
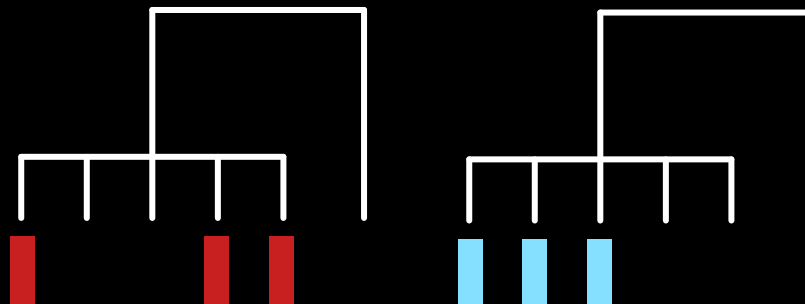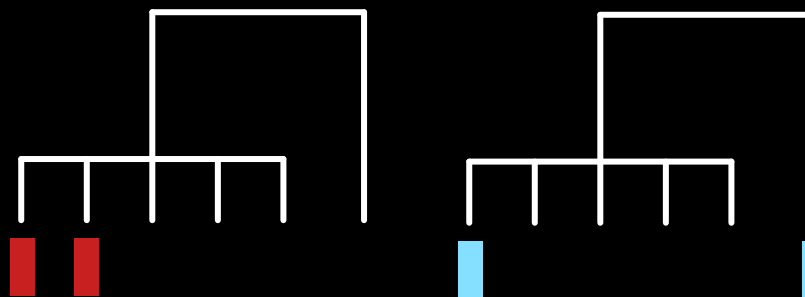
Not diverse

Qualitatively diverse
Taxonomically diverse

Phylogenetically diverse

Quantitatively diverse
Taxonomically diverse

# Microbiome composition analyses:
## beta diversity (2-sample) scenarios

# Which human body sites harbor the greatest microbial diversity per individual?



Within–Sample Alpha Diversity

Legend:
- Phylotypes (16S)
- OTUs (16S)
- Reference genomes (WGS)
- Metabolic modules (WGS)
- Gene index (WGS)

y-axis: $\log_2$ (Relative Alpha Diversity)

x-axis body sites: Anterior nares, L Antecubital fossa, R Antecubital fossa, L Retroauricular crease, R Retroauricular crease, Buccal mucosa, Keratinized gingiva, Hard palate, Palatine Tonsils, Saliva, Throat, Tongue dorsum, Subgingival plaque, Supragingival plaque, Stool, Mid vagina, Posterior fornix, Vaginal introitus

# Which human body sites share the greatest microbial diversity among individuals?



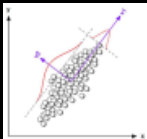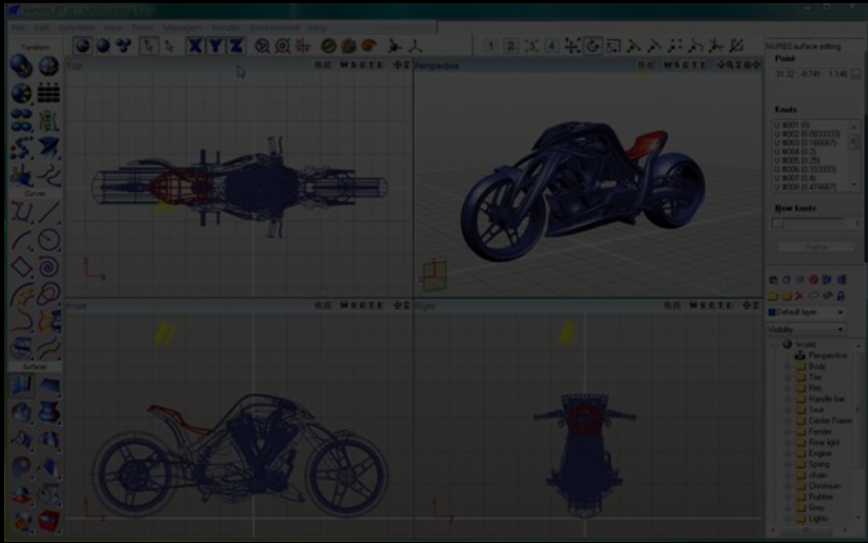Between–Sample Beta Diversity

# Microbiome composition analyses: ordination

**Ordination** is the constrained projection of high-dimensional data into fewer dimensions.

**PCA** or **Principal Component Analysis** guarantees the new dimensions maximize normal variation.
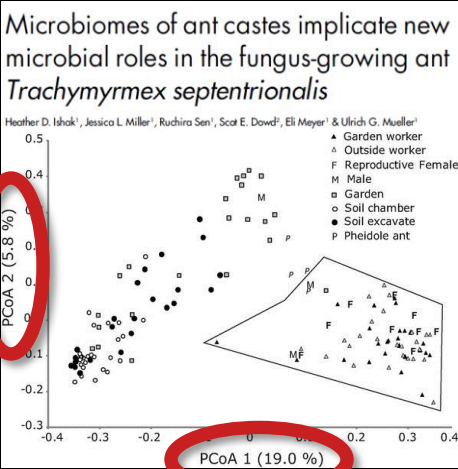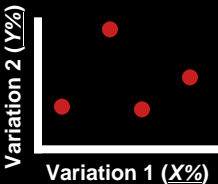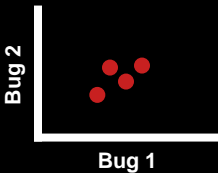
**NMDS** or **Nonmetric Multidimensional Scaling**, also called **PCoA** or **Principal Coordinates Analysis**, guarantees the new dimensions maximize an arbitrary similarity score (such as UniFrac beta-diversity).
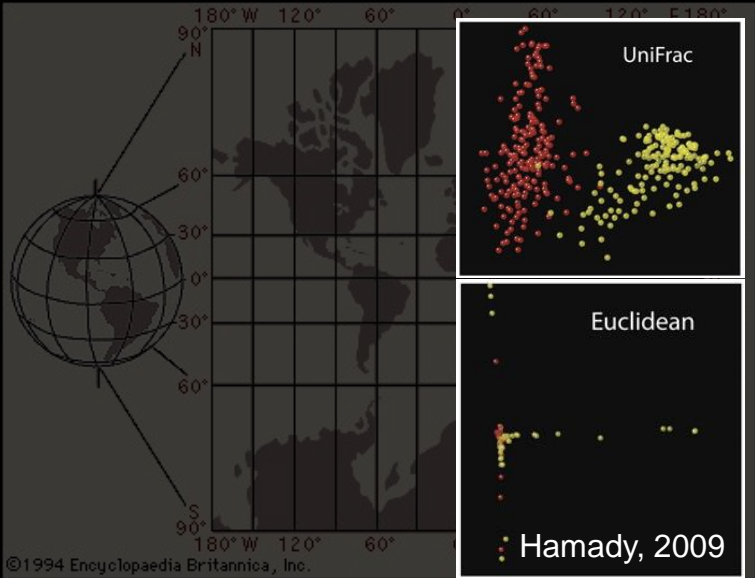
**Samples →**
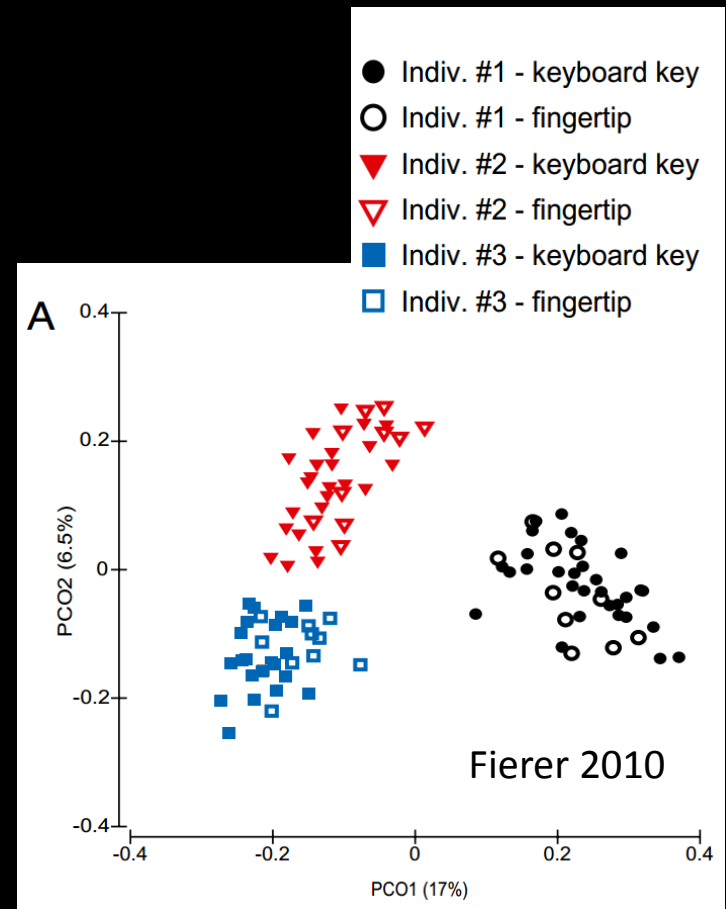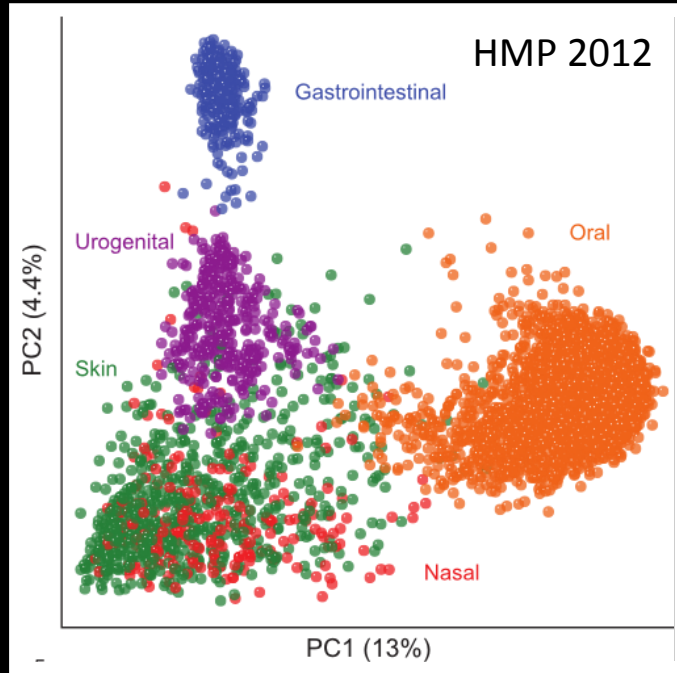
**← Microbes**

Distance between points is Euclidean

Bug 2 / Bug 1

Variation 2 (*Y%*) / Variation 1 (*X%*)

Distance between points is a *proportional function* of their *similarity*



Microbiomes of ant castes implicate new microbial roles in the fungus-growing ant *Trachymyrmex septentrionalis*

Heather D. Ishak[1], Jessica L. Miller[1], Ruchira Sen[1], Scot E. Dowd[2], Eli Meyer[1] & Ulrich G. Mueller[1]

▲ Garden worker
△ Outside worker
F Reproductive Female
M Male
□ Garden
○ Soil chamber
● Soil excavate
P Pheidole ant

PCoA 2 (5.8 %)
PCoA 1 (19.0 %)

UniFrac

Euclidean

Hamady, 2009
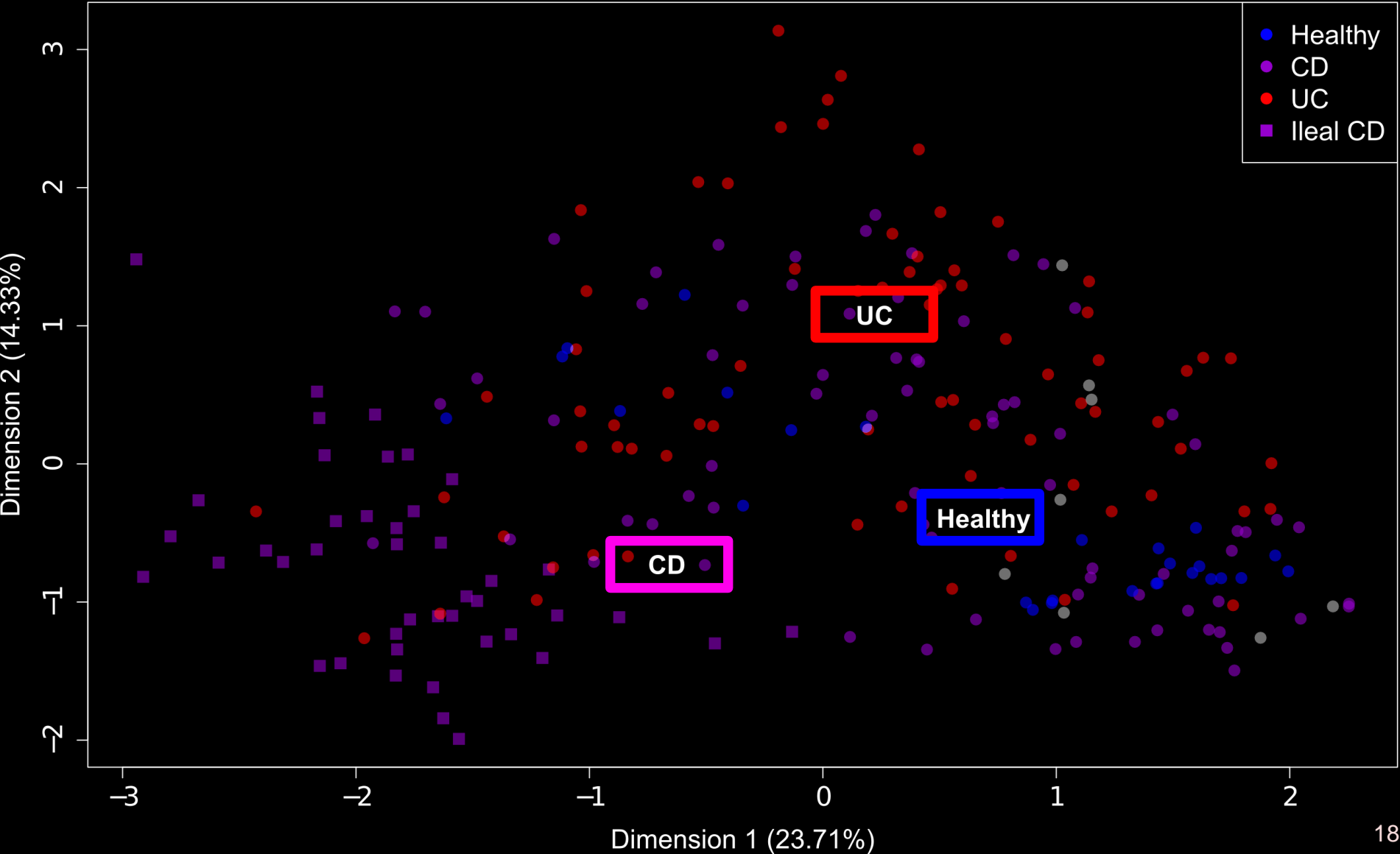
©1994 Encyclopaedia Britannica, Inc.

# Microbiome composition analyses: **ordination examples**

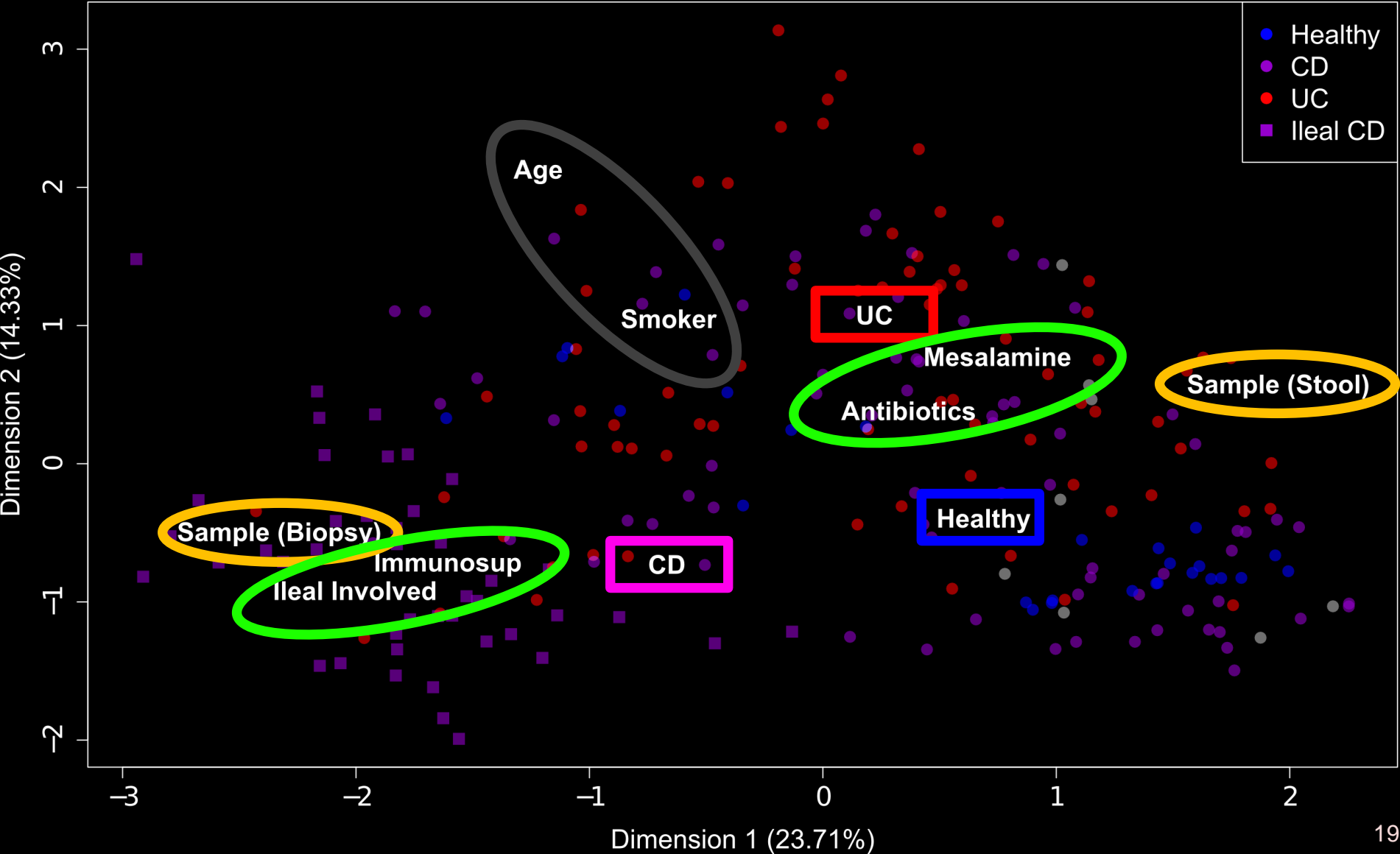# How is the gut microbiome disrupted during IBD and its treatment?
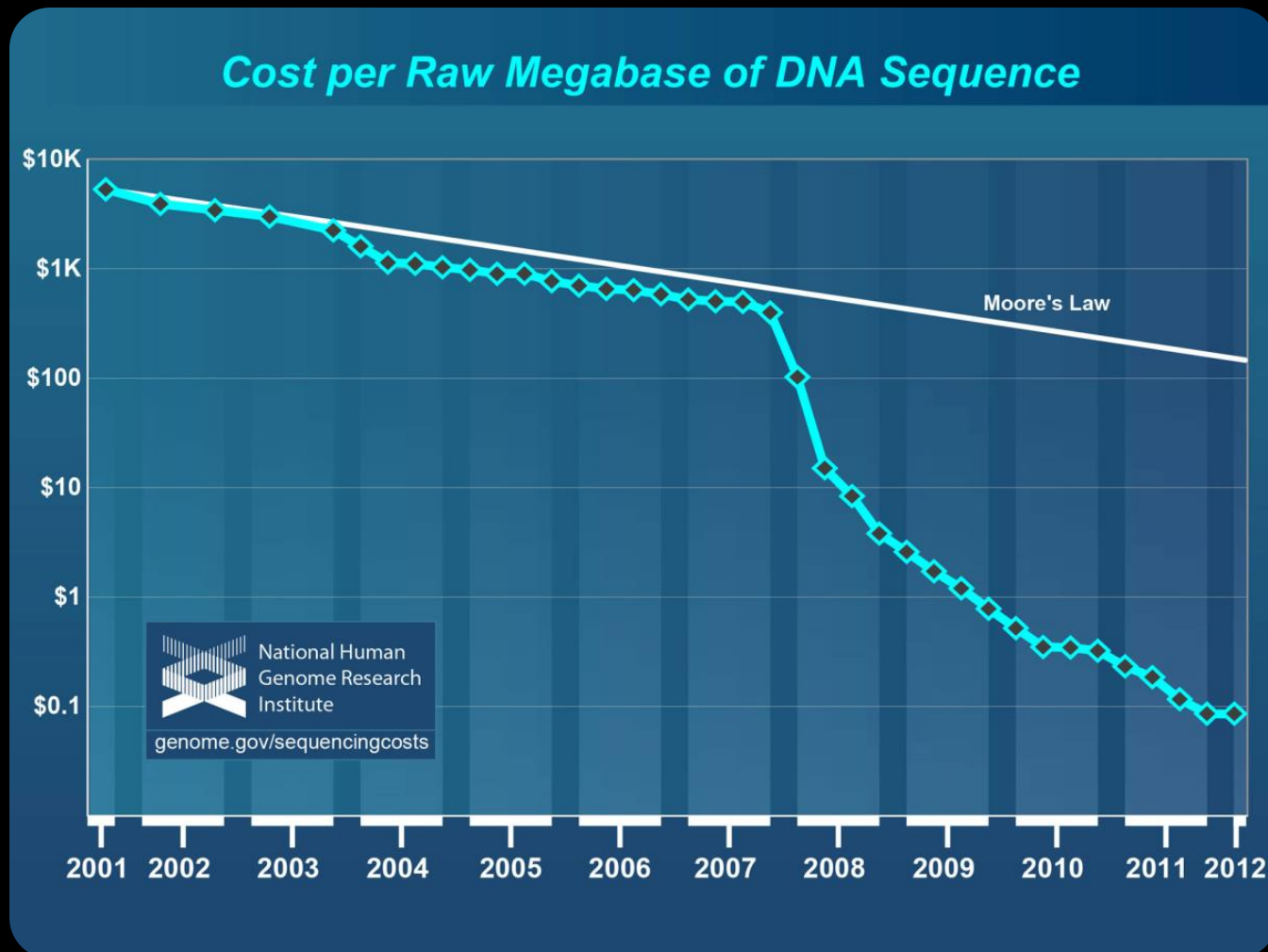
**With Ramnik Xavier, Bruce Sands**

# How is the gut microbiome disrupted during IBD and its treatment?
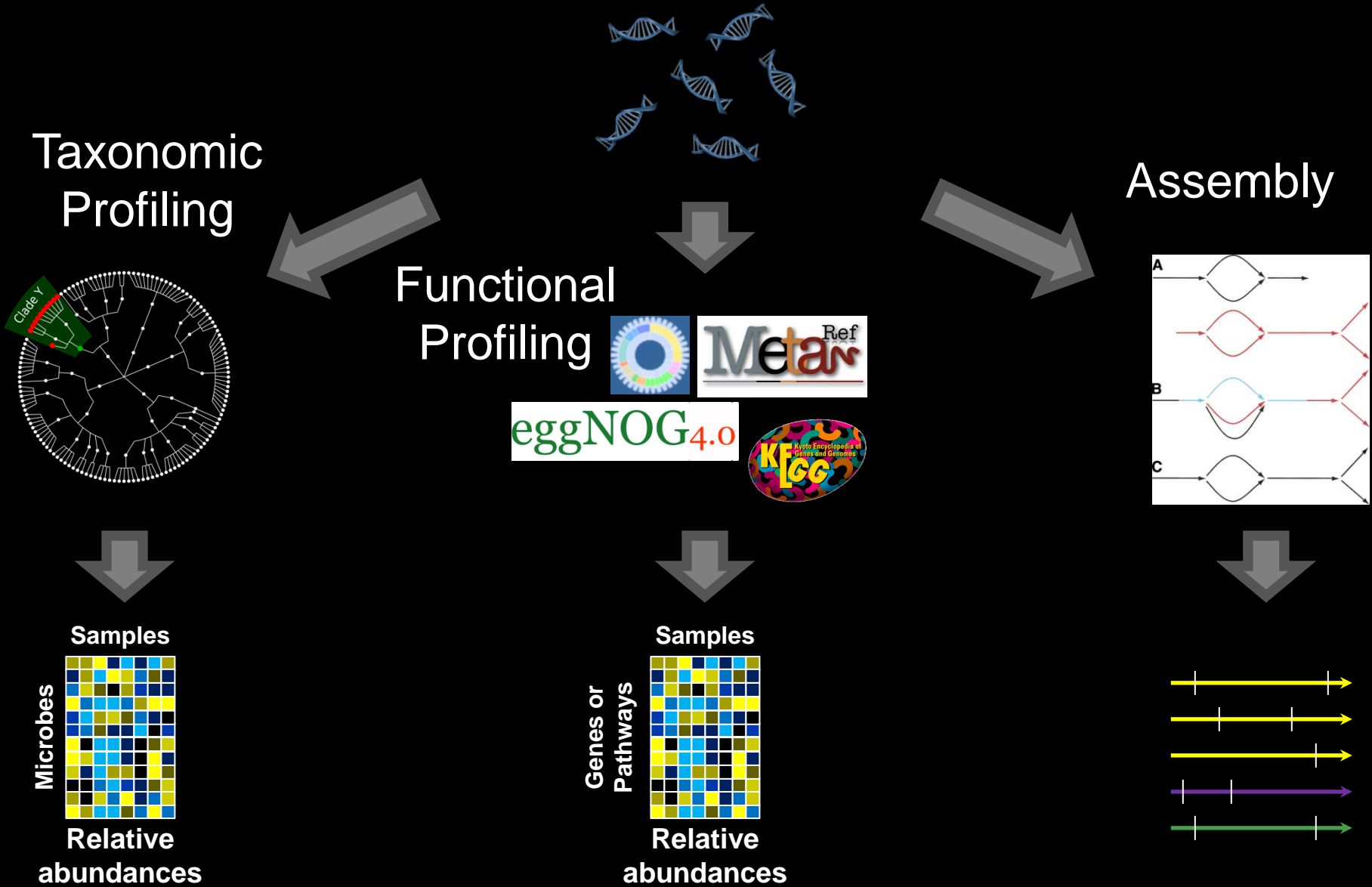
**With Ramnik Xavier, Bruce Sands**

# Meta'omic analyses



Cost per Raw Megabase of DNA Sequence

# Typical shotgun metagenome and metatranscriptome analyses

**Taxonomic Profiling**

**Assembly**

**Functional Profiling**

eggNOG4.0

**Samples**

Microbes

**Relative abundances**

**Samples**

Genes or Pathways

**Relative abundances**

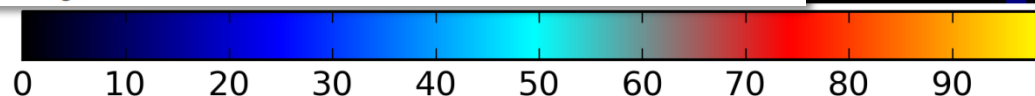# Profiling microbial communities and ecology at species-level resolution



ARTICLE

## Enterotypes of the human gut microbiome

The New York Times
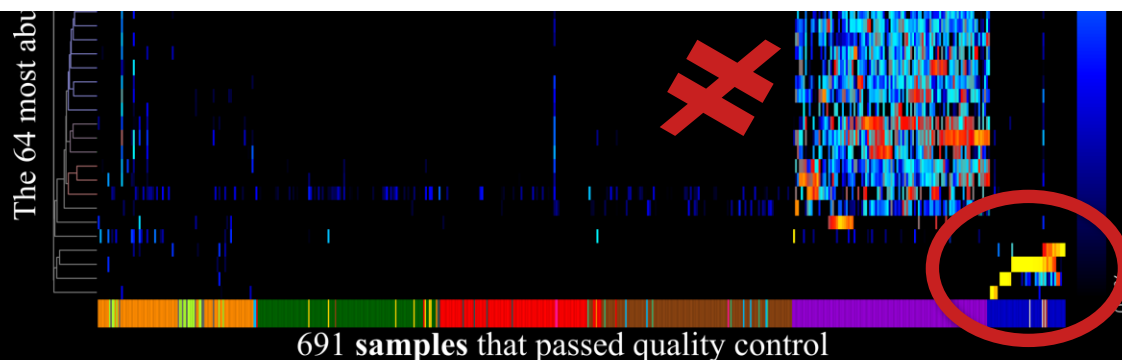
**Science**

| WORLD | U.S. | N.Y. / REGION | BUSINESS | TECHNOLOGY | SCIENCE | HEALTH | SPORTS |

ENVIRONMENT    SPACE & COS

Bacterial Ecosystems Divide People Into 3 Groups, Scientists Say

Lactobacillus crispatus
Atopobium vaginae
Prevotella amnii
Gardnerella vaginalis
Bifidobacterium breve
Lactobacillus gasseri
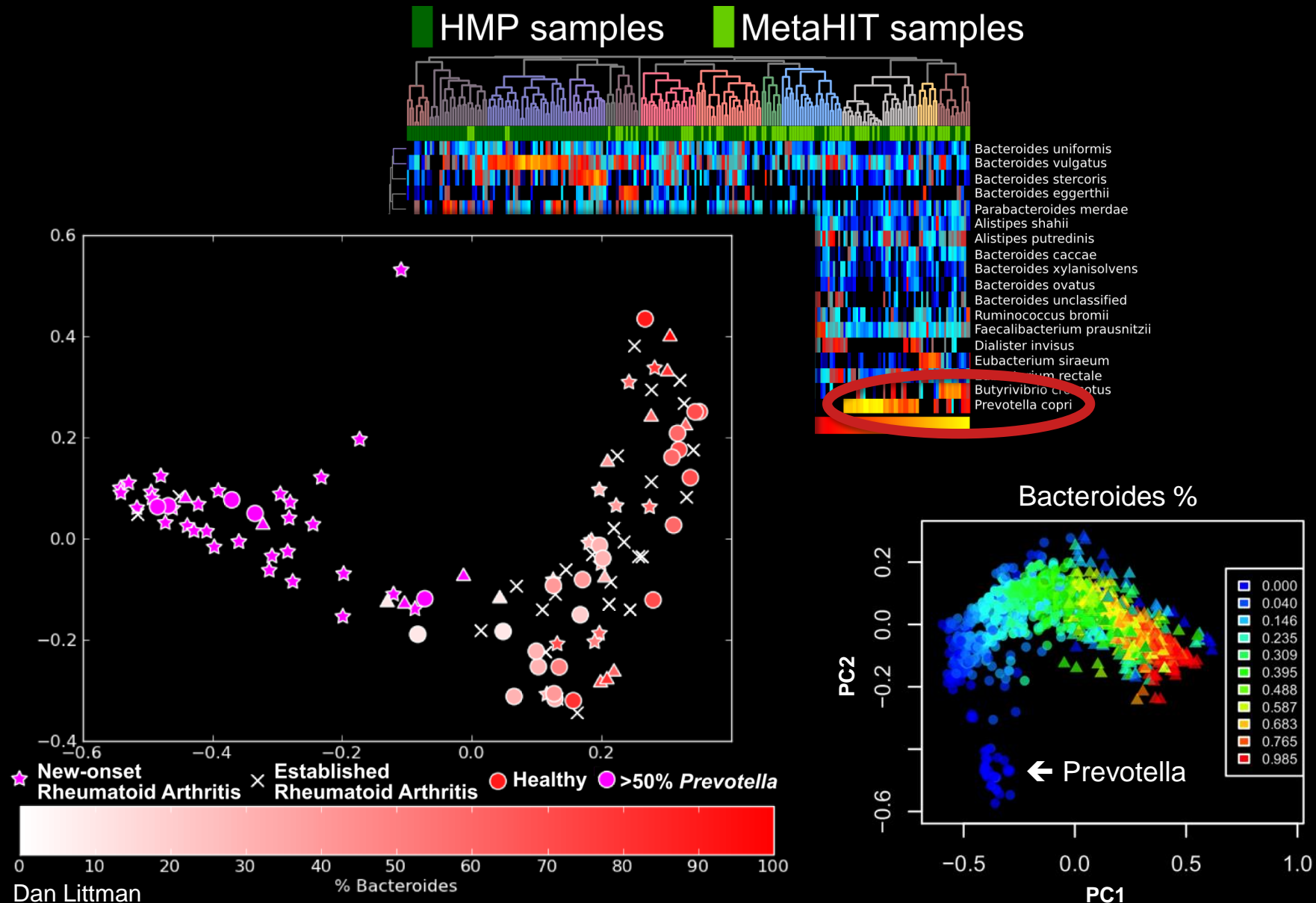Prevotella multiformis
Bifidobacterium dentium

0    10    20    30    40    50    60    70    80    90

**Gut**

**Vaginal**

691 **samples** that passed quality control

anterior nares — throat — subgingival plaque — saliva — posterior fornix
right retroauricular crease — buccal mucosa — tongue dorsum — palatine tonsils — vaginal introitus
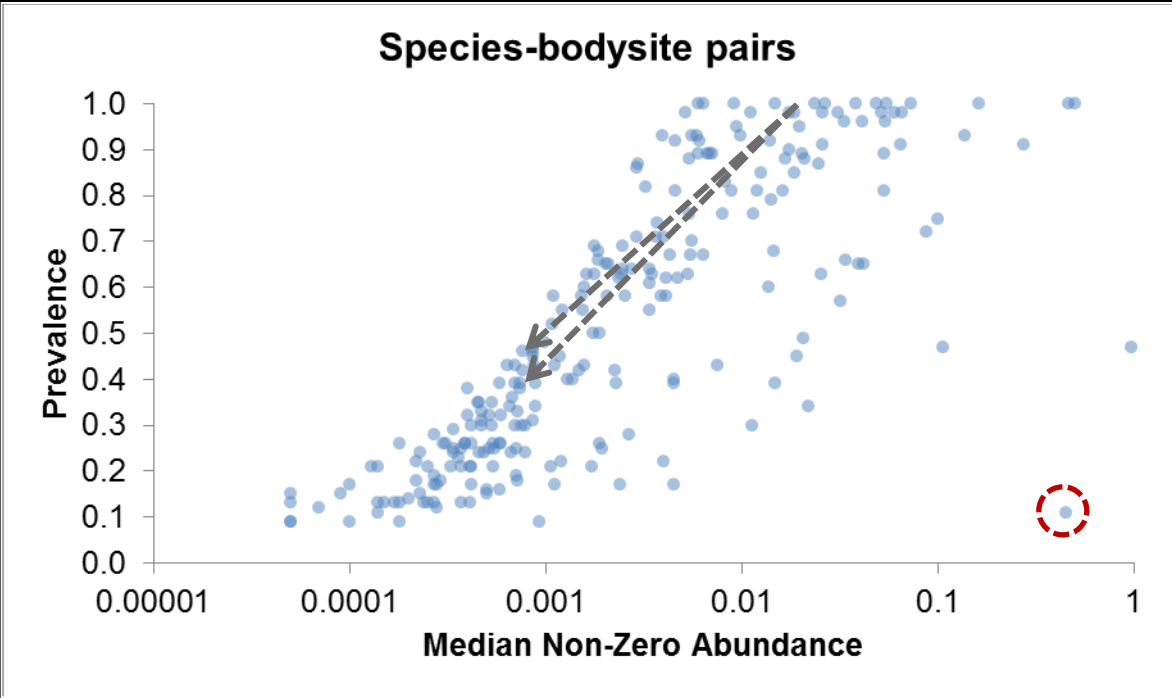left retroauricular crease — supragingival plaque — attached keratinized gingiva — stool — mid vagina

22

# Are there discrete "types" of typical human microbiomes?



Dan Littman

# Species prevalence vs. species abundance



| | site | bug | log10 med non0 abund | prevalence | tp 1-2 stability |
|---|---|---|---|---|---|
| | Stool | s__Prevotella_copri | -0.3 | 0.11 | 1.00 |

| | site | bug | log10 med non0 abund | prevalence | tp 1-2 stability |
|---|---|---|---|---|---|
| | Tongue_dorsum | s__Actinomyces_odontolyticus | -1.8 | 1.00 | 1.00 |
| | Buccal_mucosa | s__Actinomyces_odontolyticus | -3.1 | 0.47 | 0.78 |
| | Supragingival_plaque | s__Actinomyces_odontolyticus | -3.1 | 0.38 | 0.50 |

# Typical shotgun metagenome and metatranscriptome analyses

Taxonomic Profiling

Functional Profiling

Assembly

eggNOG4.0

**Samples**

**Microbes**

**Relative abundances**

**Samples**

**Genes or Pathways**
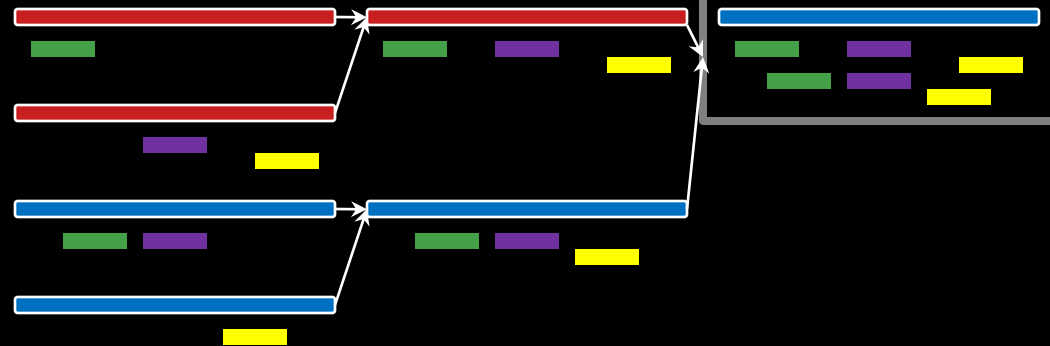
**Relative abundances**

# Metagenomic analyses: gene calling and proxygenes
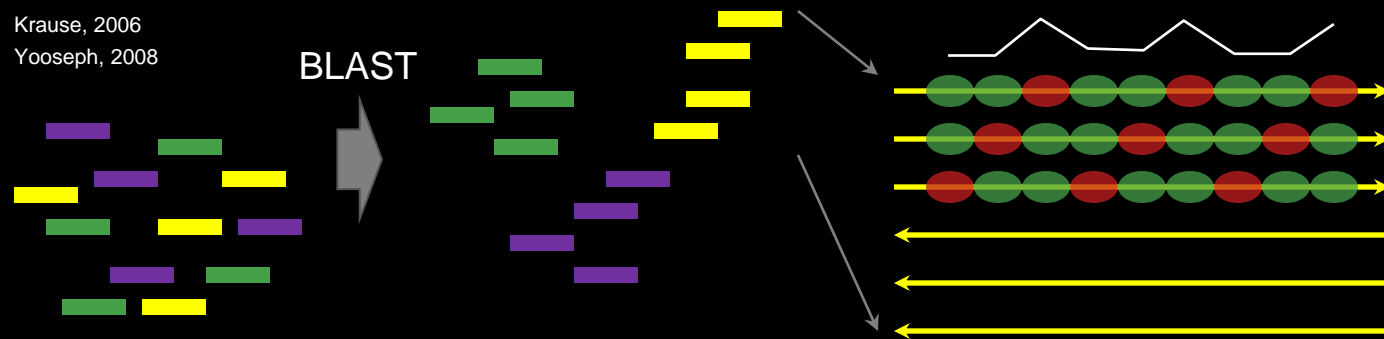
Dalevi, 2009

Extrinsic gene calling: BLAST etc. (proxygenes)
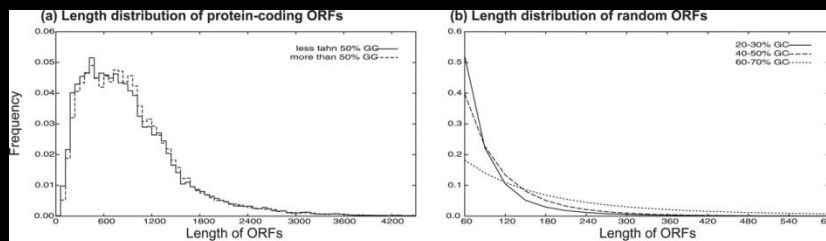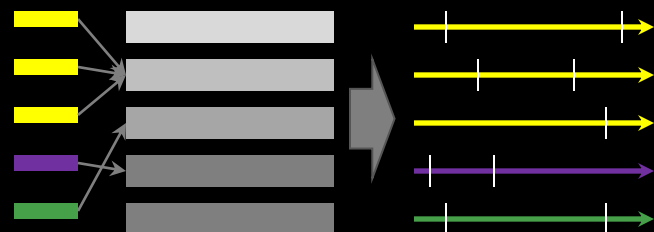
Intrinsic gene calling: ORF detection from seq.

Krause, 2006
Yooseph, 2008

BLAST

Orphelia: Hoff, 2009
MetaGene: Noguchi, 2006

HMM models

(a) Length distribution of protein-coding ORFs
less tahn 50% GG
more than 50% GG

Frequency

Length of ORFs

(b) Length distribution of random ORFs
20-30% GC
40-50% GC
60-70% GC

Length of ORFs

# Metagenomic analyses: molecular functions and biological roles

**Orthology**: Grouping genes by conserved sequence features
COG, KO, FIGfam…

**Structure**: Grouping genes by similar protein domains
Pfam, TIGRfam, SMART, EC…

**Biological roles**: Grouping genes by pathway and process involvement
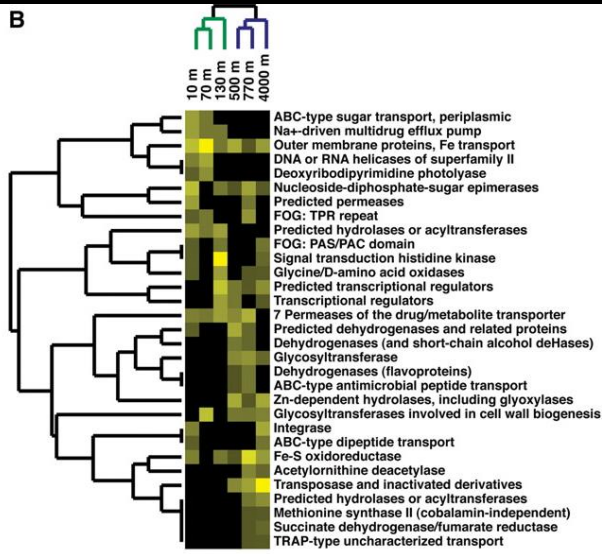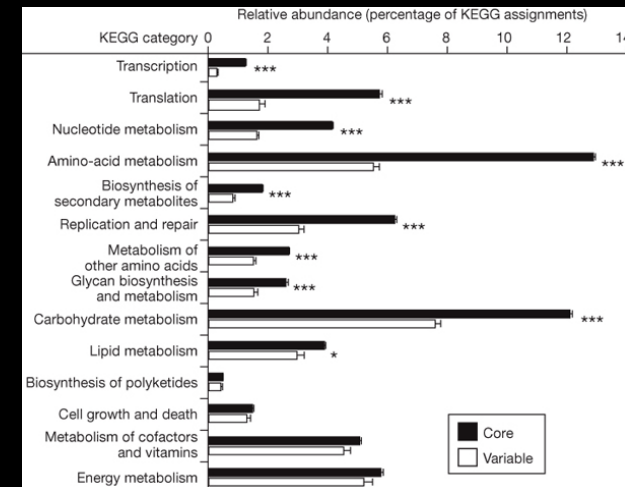GO, KEGG, MetaCyc, SEED…



DeLong, 2006



**Table 1 | Glycoside hydrolases and carbohydrate-binding modules**

| CAZy family* | Pfam HMM name† | Known activities‡ | Termite gut community§ |
|---|---|---|---|
| Glycoside hydrolase catalytic domains**†† | | | |
| GH1 | Glyco_hydro_1 | β-Glucosidase, β-galactosidase, β-mannosidase, others | 22 |
| GH2 | Glyco_hydro_2_C | β-Galactosidase, β-mannosidase, others | 23 |
| GH3 | Glyco_hydro_3 | β-1,4-Glucosidase, β-1,4-xylosidase, β-1,3-glucosidase, α-L-arabinofuranosidase, others | 69 |
| GH4 | Glyco_hydro_4 | α-Glucosidase, α-galactosidase, α-glucuronidase, others | 14 |
| GH5 | Cellulase | Cellulase, β-1,4-endoglucanase, β-1,3-glucosidase, β-1,4-endoxylanase, β-1,4-endomannanase, others | 56 |
| GH8 | Glyco_hydro_8 | Cellulase, β-1,3-glucosidase, β-1,4-endoxylanase, β-1,4-endomannanase, others | 5 |
| GH9 | Glyco_hydro_9 | Endoglucanase, cellobiohydrolase, β-glucosidase | 9 |
| GH10 | Glyco_hydro_10 | Xylanase, β-1,3-endoxylanase | 46 |
| GH11 | Glyco_hydro_11 | Xylanase | 14 |
| GH13 | Alpha-amylase | α-Amylase, catalytic domain, and related enzymes | 48 |
| GH16 | Glyco_hydro_16 | β-1,3(4)-Endoglucanase, others | 1 |
| GH18 | Glyco_hydro_18 | Chitinase, endo-β-N-acetylglucosaminidase, non-catalytic proteins | 17 |
| GH20 | Glyco_hydro_20 | β-Hexosaminidase, lacto-N-biosidase | 15 |

Warnecke, 2007



Turnbaugh, 2009
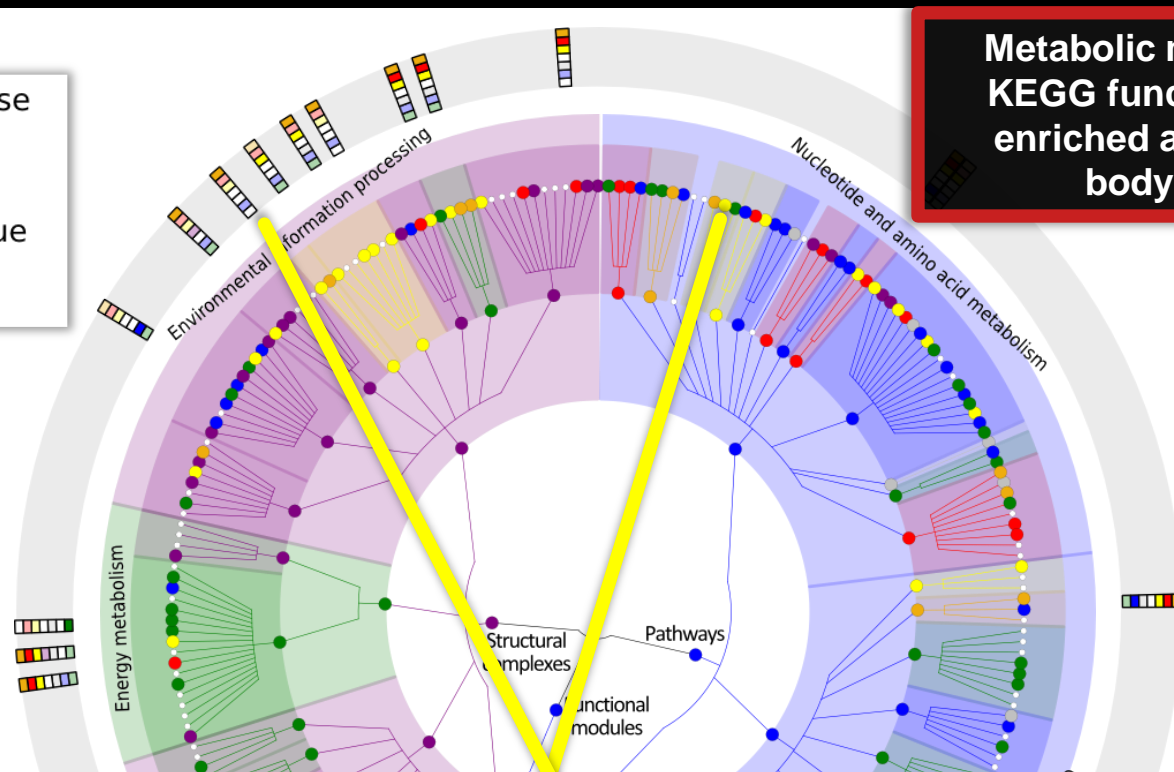
# Niche specialization in human microbiome function

**Metabolic modules in the KEGG functional catalog enriched at one or more body habitats**

Legend:
- Retroauricular crease
- Stool
- Anterior nares
- Posterior fornix
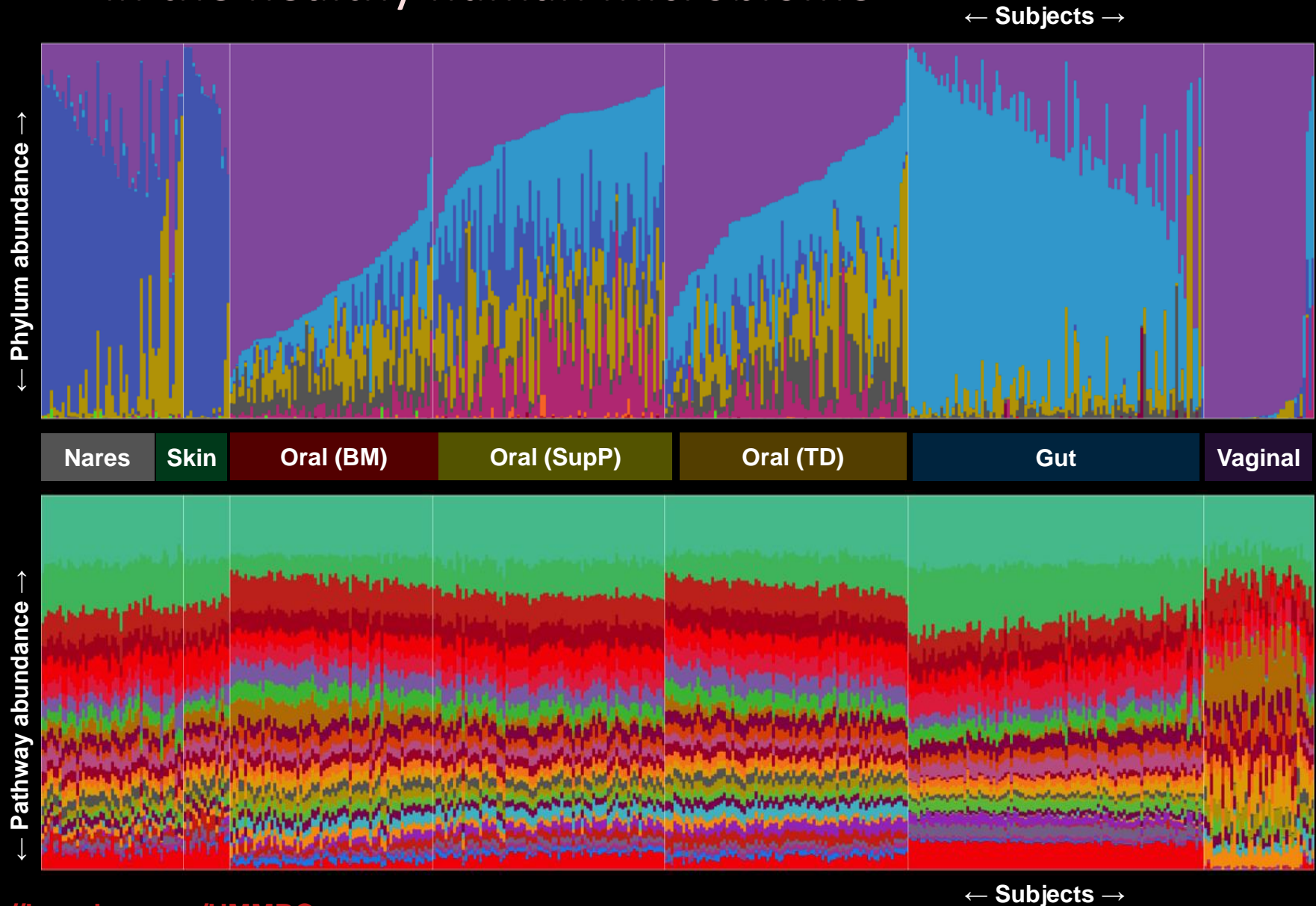- Supragingival plaque
- Buccal mucosa
- Tongue dorsum

M00142: Complex I (NADH dehydrogenase), NADH dehydrogenase I

- Most **_processes are "core"_**: <10% are *differentially present/absent* even by body site
  - Contrast **_zero_** microbes meeting this threshold!
- Most **_processes are habitat-adapted_**: >66% are *differentially abundant* by body site

28

# "Who's there," versus, "What they're doing," in the healthy human microbiome



← Subjects →

Phylum abundance →

Nares | Skin | Oral (BM) | Oral (SupP) | Oral (TD) | Gut | Vaginal

Pathway abundance →

← Subjects →

# Which *functions* of the gut microbiome are disrupted by IBD?

- Over <u>*six times*</u> as many microbial metabolic processes disrupted in IBD as microbes.
  - If there's a transit strike, everyone working for the MBTA is disrupted, not everyone named Smith or Jones.
  - Phylogenetic distribution of function is *consistent* but *diffuse*

- During IBD, microbes...

| **Stop** | **Start** |
| --- | --- |
| • Creating most amino acids<br>• Degrading complex carbs.<br>• Producing short-chain fatty acids | • Taking up more host products<br>• Dodging the immune system<br>• Adhering to and invading host cells |

# Plan

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
  - MetaPhlAn
- Tools for functional profiling
  - HUMAnN
  - ShortBRED
  - PICRUSt
- Tools for testing associations
  - LEfSe
  - MaAsLin
  - CCREPE
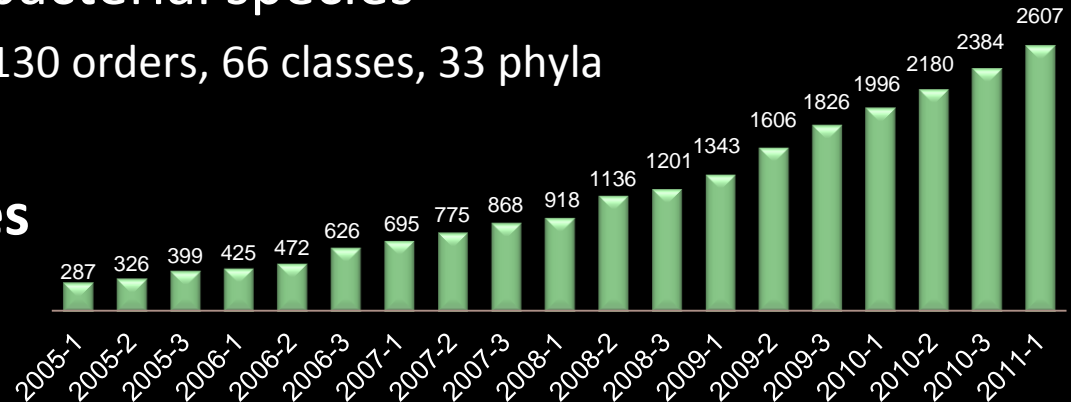- Resources
- Research vignette (time permitting)

# Who is there?
## (taxonomic profiling)

# What are they doing?
## (functional profiling)

# Reference genomes
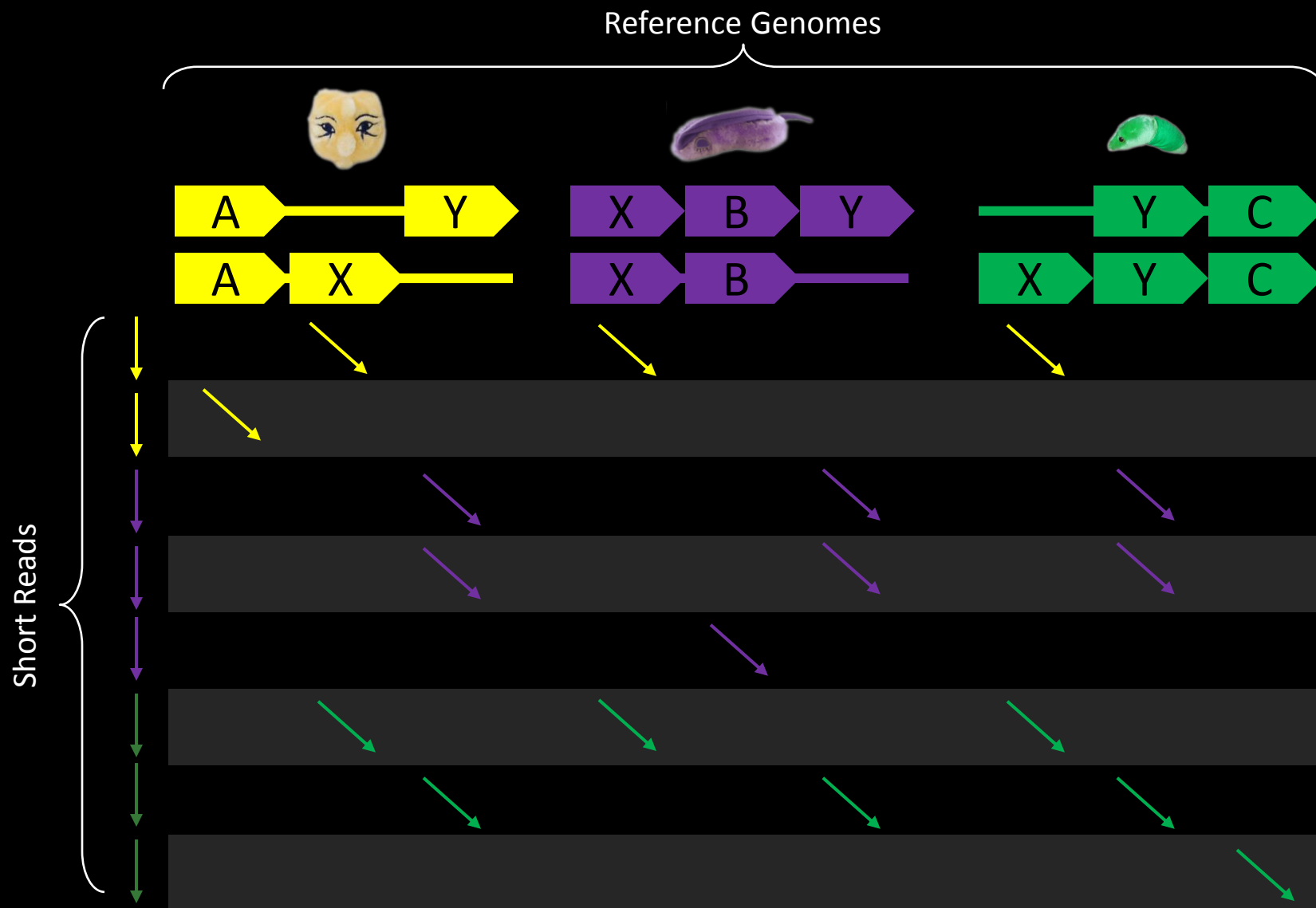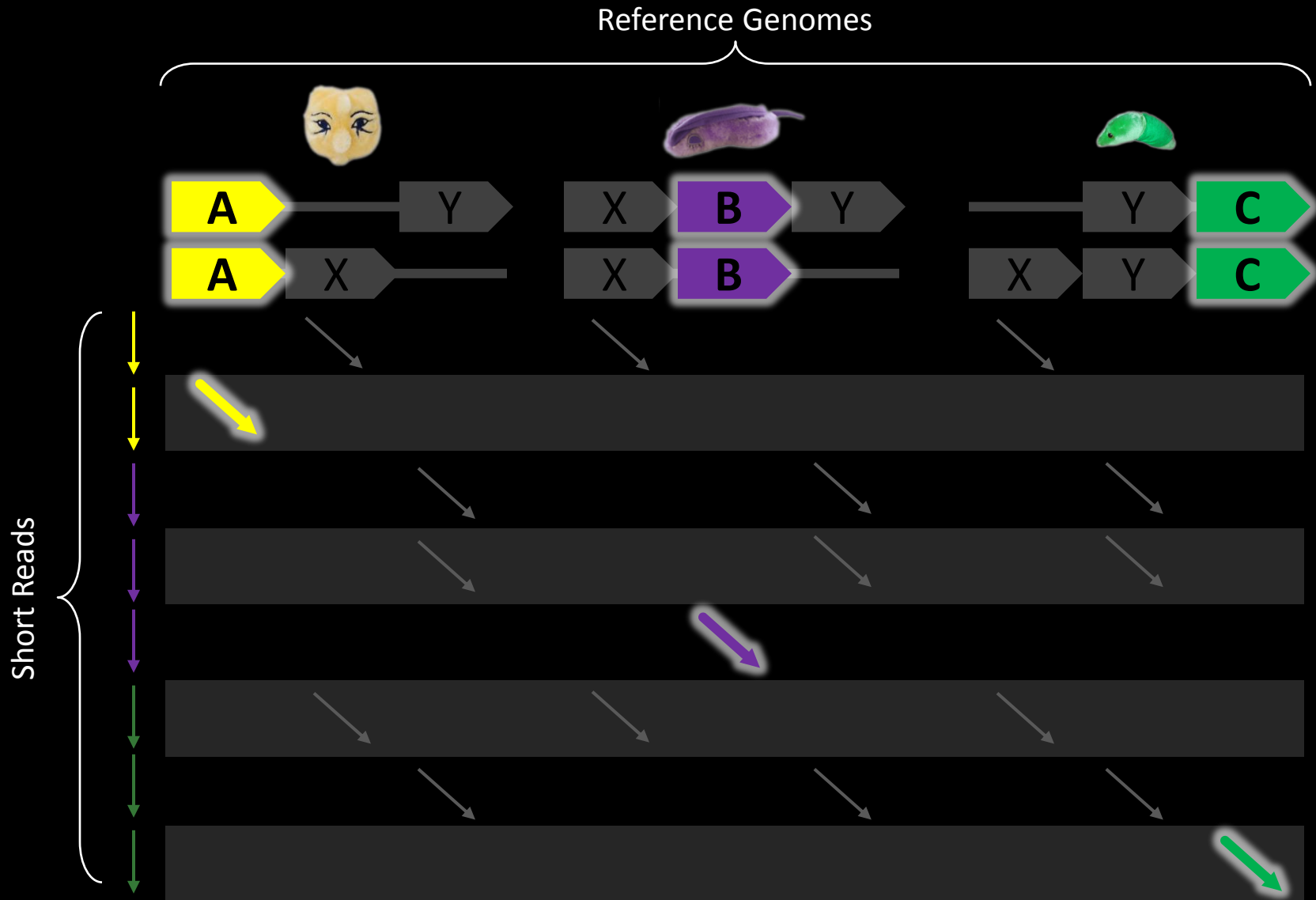
- IMG alone now contains ~3,100 bacterial genomes
  - Plus ~100 archaeal, ~100 eukaryotic, and a few thousand viruses
  - About half final and half draft
- These comprise 1,222 bacterial species
  - 652 genera, 278 families, 130 orders, 66 classes, 33 phyla
  - 2,383 total clades
- **And roughly 10M genes**

Bar chart values: 287 (2005-1), 326 (2005-2), 399 (2005-3), 425 (2006-1), 472 (2006-2), 626 (2006-3), 695 (2007-1), 775 (2007-2), 868 (2007-3), 918 (2008-1), 1136 (2008-2), 1201 (2008-3), 1343 (2009-1), 1606 (2009-2), 1826 (2009-3), 1996 (2010-1), 2180 (2010-2), 2384 (2010-3), 2607 (2011-1)

- These genes and genomes are a tremendous resource to:
  - Identify conserved markers that can be used to infer phylogeny
  - Identify unique markers that can be used to infer taxonomy
  - Relate the microbial members of a community to their metagenomic functional potential
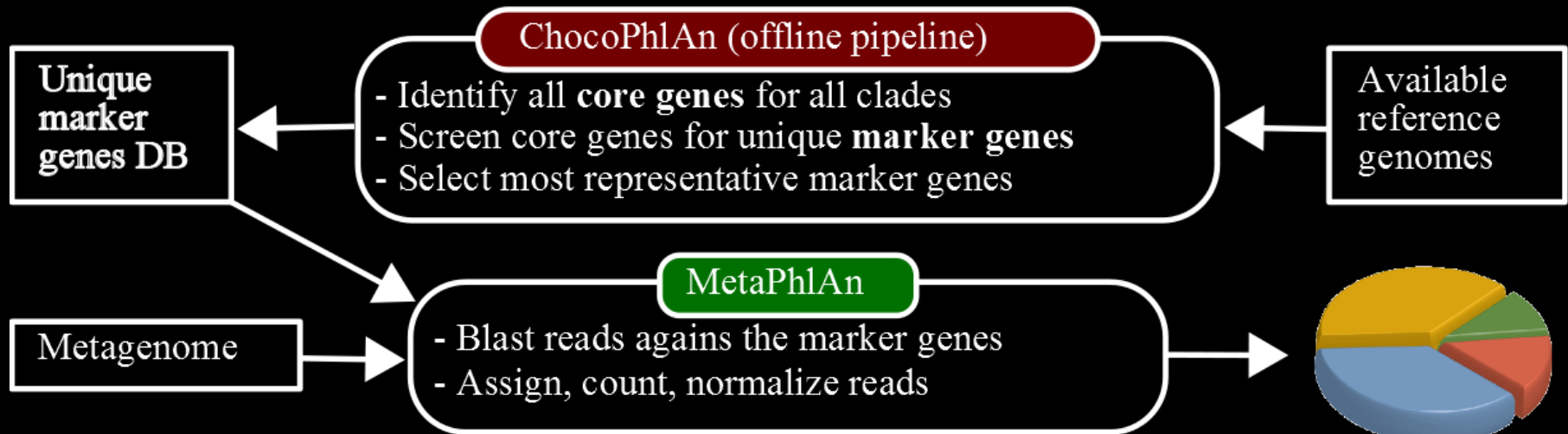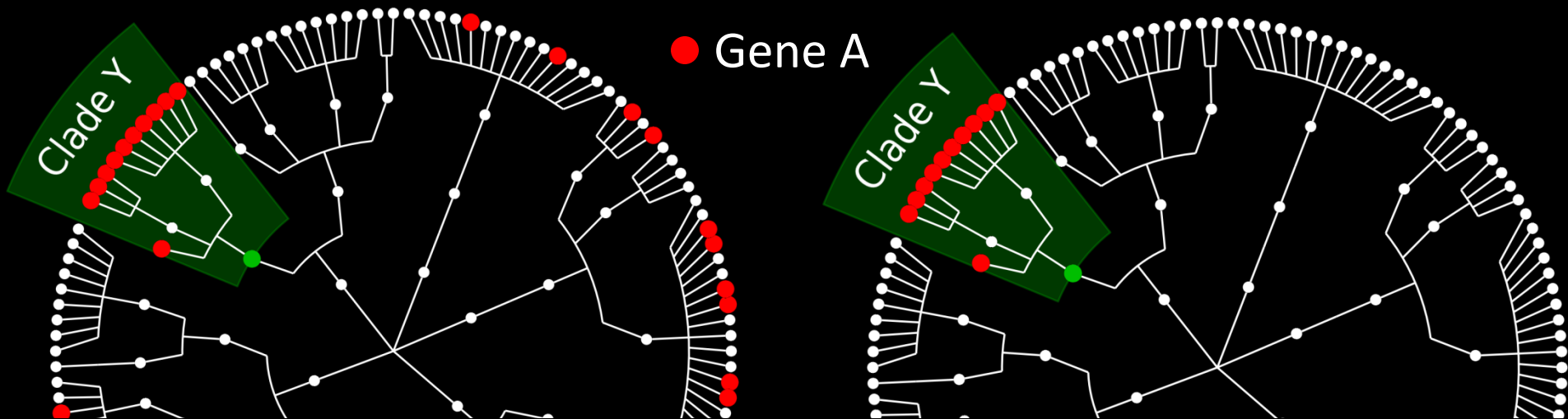
# MetaPhlAn
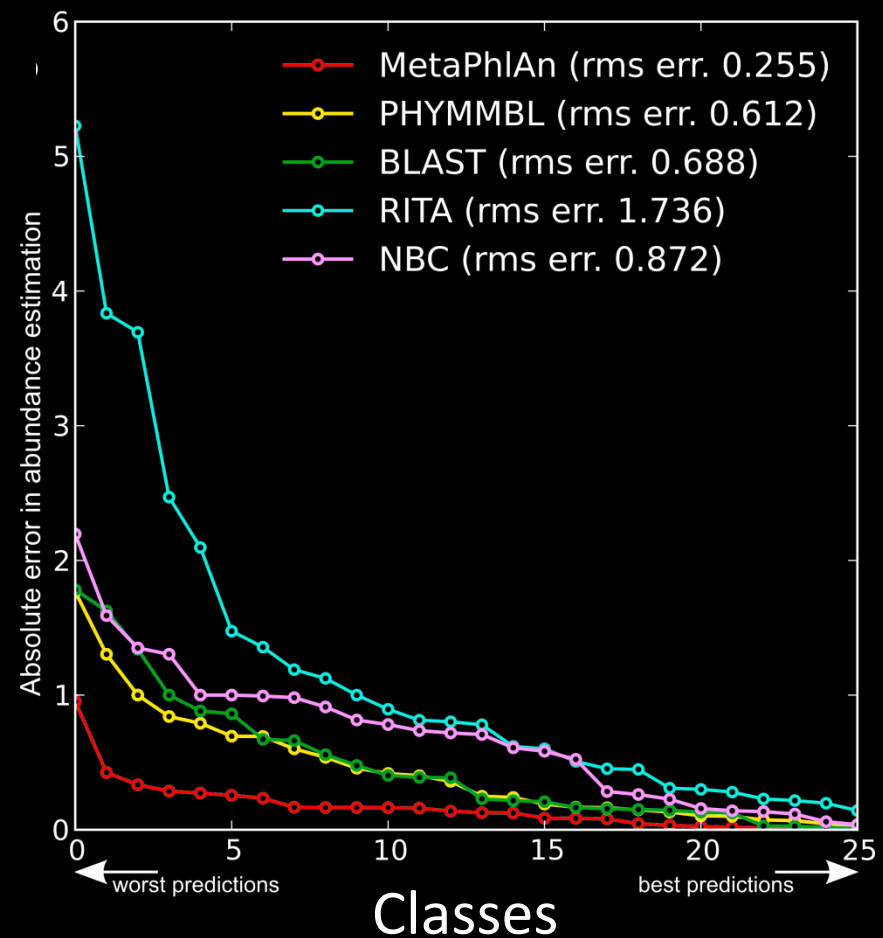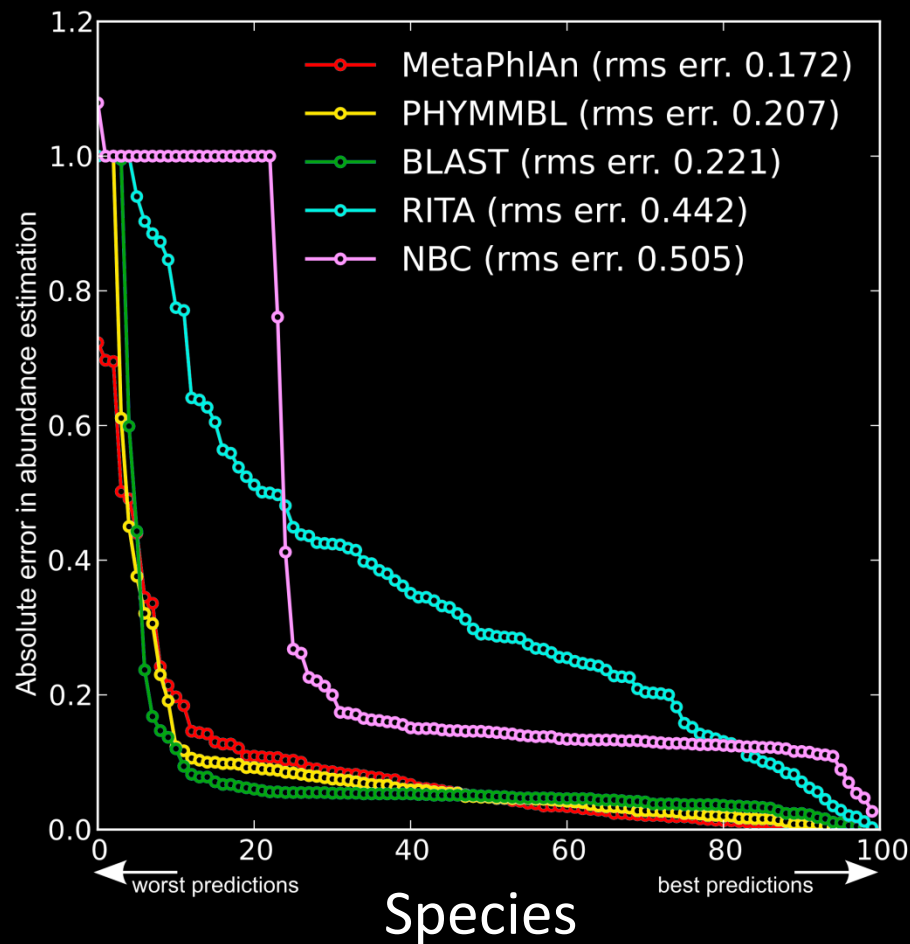## Metagenomic Phylogenic Analysis

*Nicola Segata*

# MetaPhlAn overview

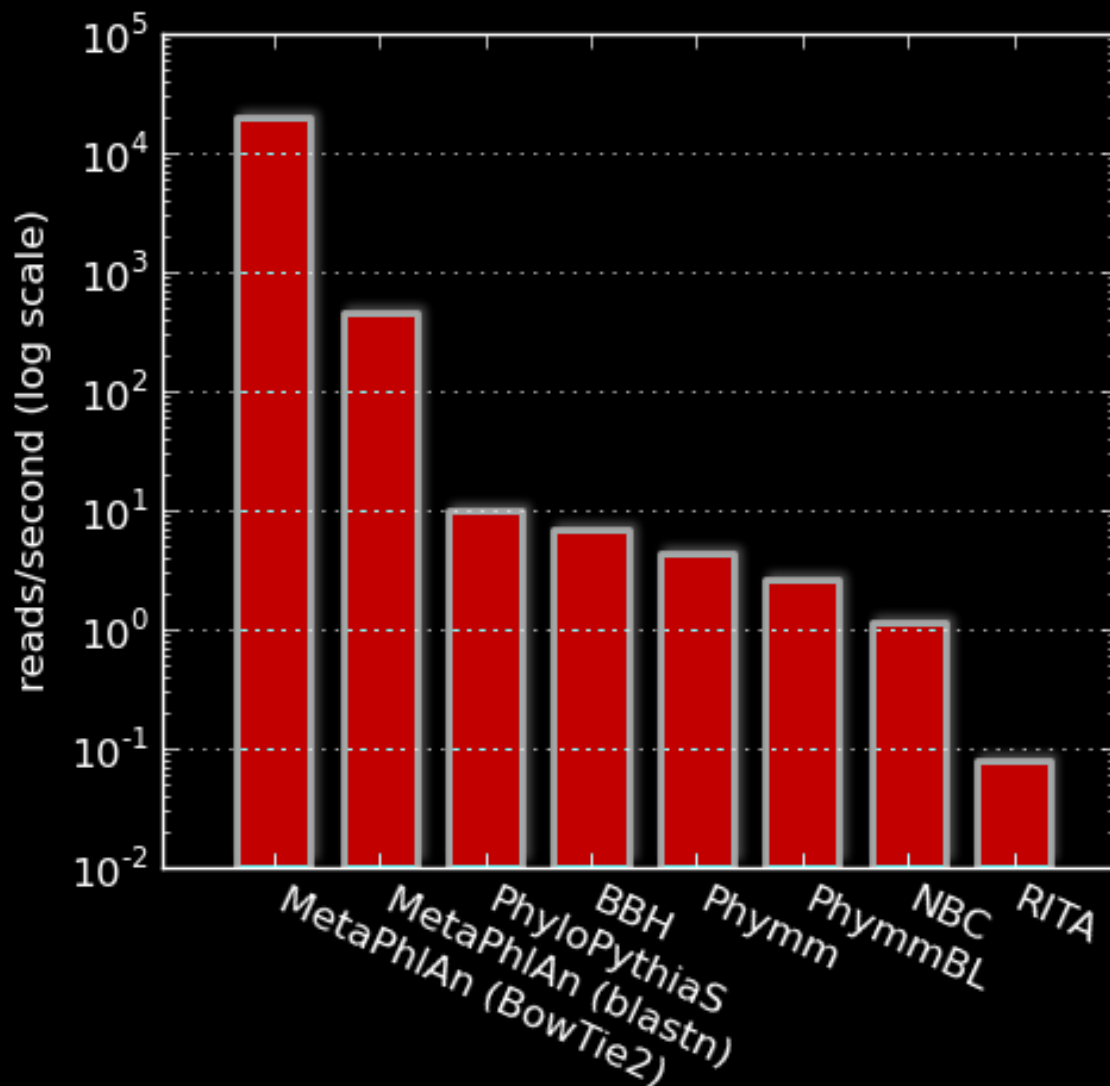A is a core gene for clade Y

A is a unique marker gene for clade Y

● Gene A



Unique marker genes DB

**ChocoPhlAn (offline pipeline)**
- Identify all **core genes** for all clades
- Screen core genes for unique **marker genes**
- Select most representative marker genes

Available reference genomes

Metagenome

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

# Evaluation of MetaPhlAn accuracy



(Validation on high-complexity uniformly distributed synthetic metagenomes.)

# Evaluation of MetaPhlAn performance



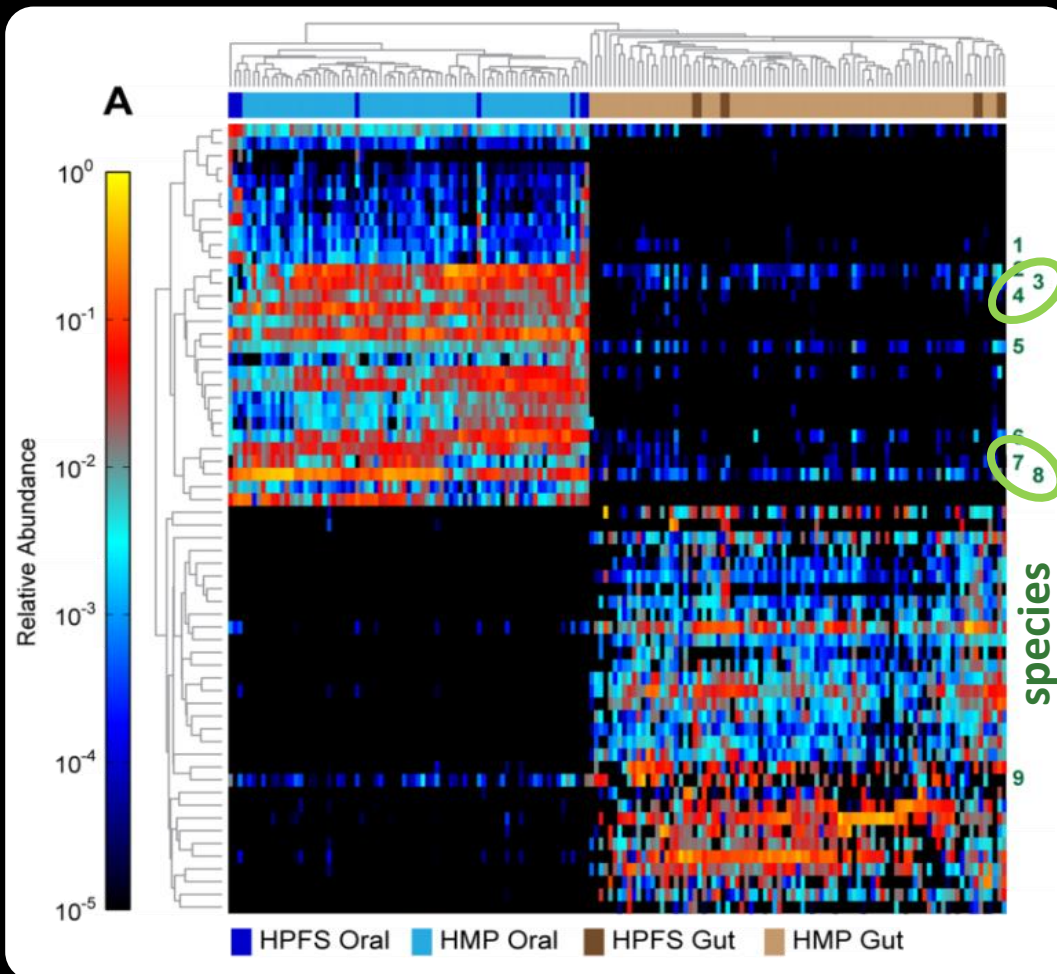>50 times faster than earlier methods

450 reads/sec (BLAST)

Up to 25,000 reads/sec (bowtie2)

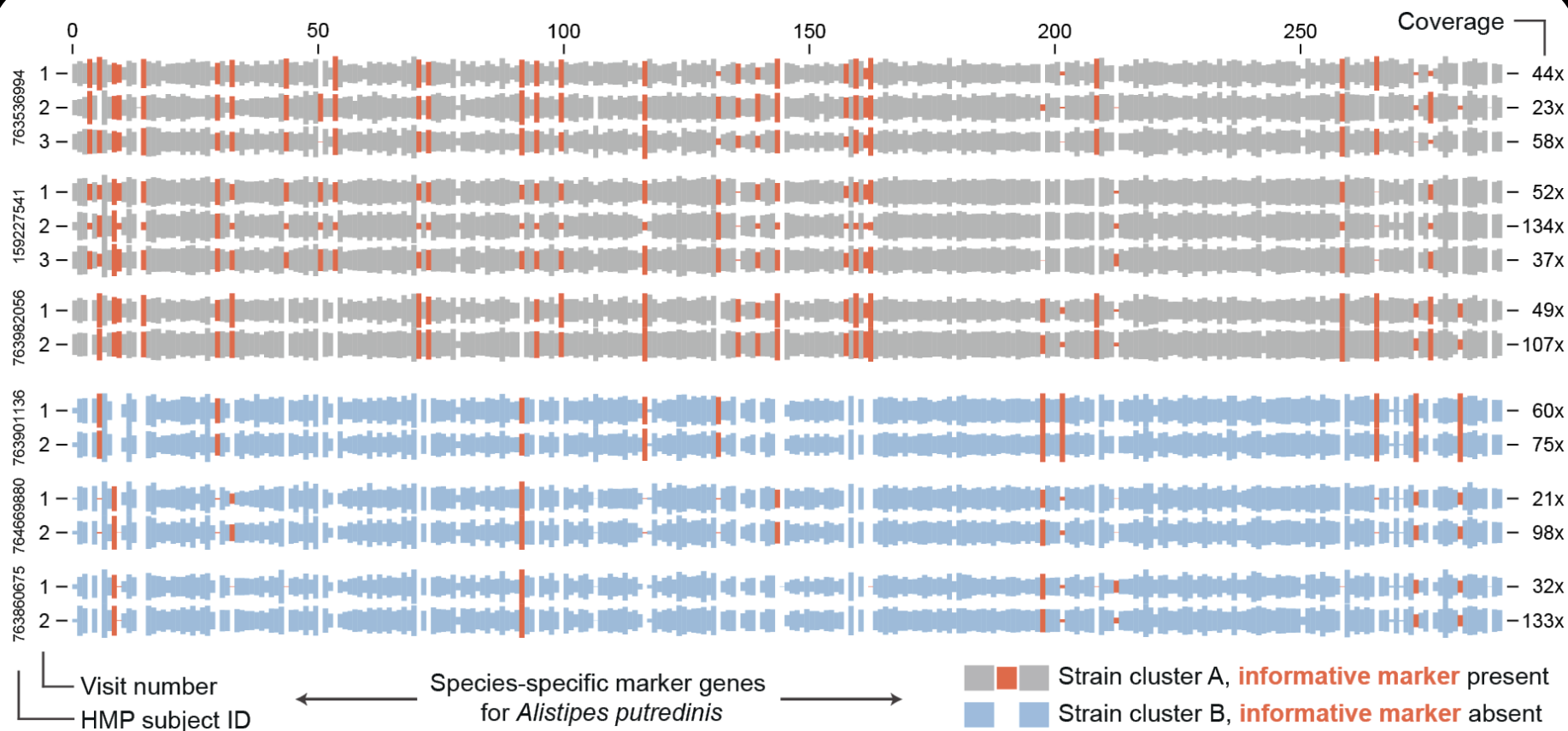Multi-threaded

Easily parallelizable

# MetaPhlAn in action

# MetaPhlAn in action: *strain profiling*



- In practice, not all markers are present
- Individual-specific marker "barcodes"
- Often very stable over time

# Plan

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
  - MetaPhlAn
- **Tools for functional profiling**
  - **HUMAnN**
  - **ShortBRED**
  - **PICRUSt**
- **Tools for testing associations**
  - **LEfSe**
  - **MaAsLin**
  - **CCREPE**
- **Resources**
- **Research vignette (time permitting)**

# Who is there?
(taxonomic profiling)

# What are they doing?
(functional profiling)

# (What we mean by "function")



INOSITOL PHOSPHATE METABOLISM

# HUMAnN
## *HMP Unified Metabolic Analysis Network*



Short reads + protein families

Translated BLAST search

$$c(g) = \frac{1}{|g|} \sum_r \frac{\sum_{a(r)} (1 - p_a) \Delta(a = g)}{\sum_{a(r)} 1 - p_a}$$

Weight hits by significance

Sum over families

Adjust for sequence length

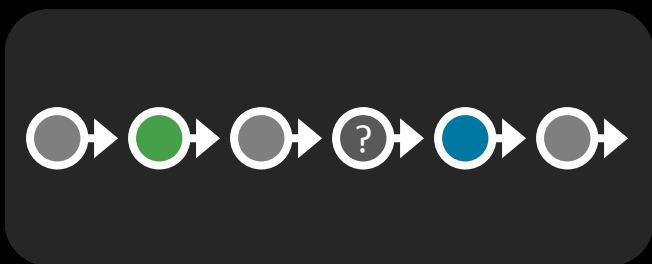Repeat for each metagenomic or metatranscriptomic sample
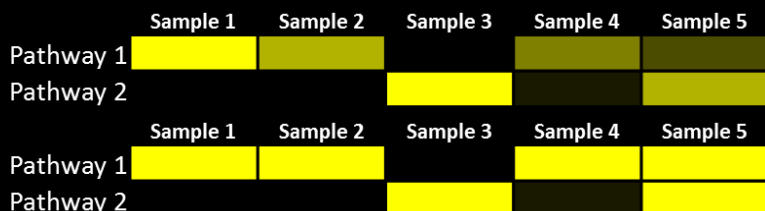
# HUMAnN
## *HMP Unified Metabolic Analysis Network*



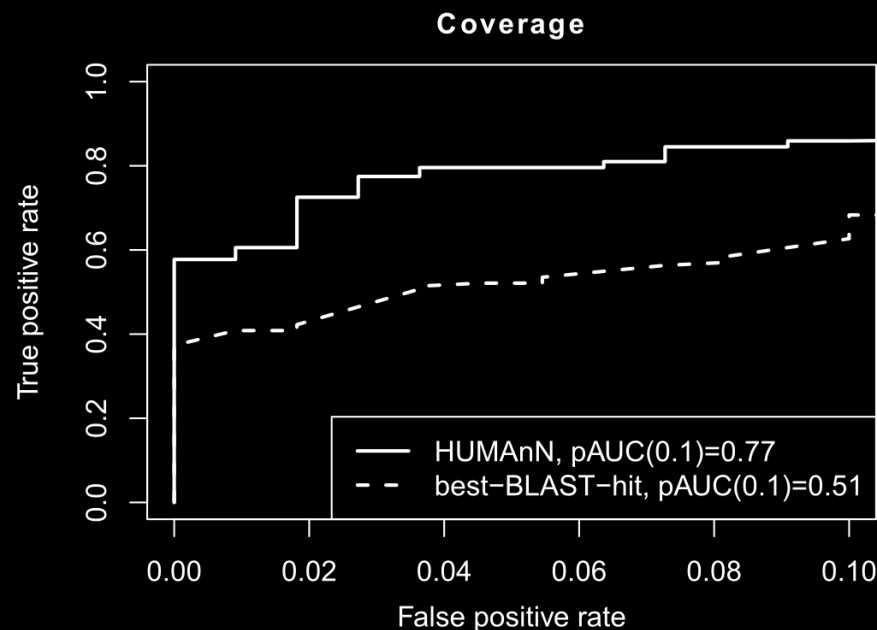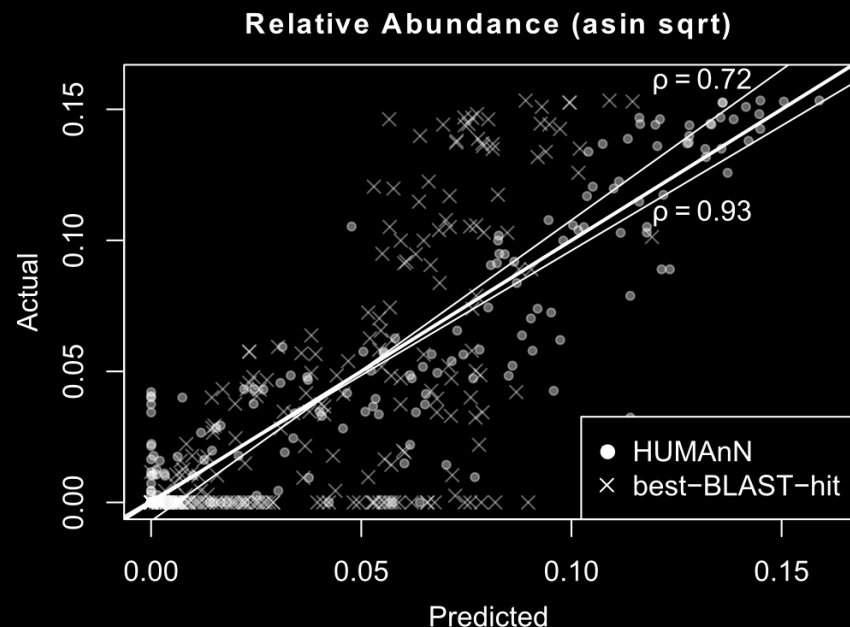Millions of hits are collapsed into thousands of gene families (KOs) (*still a large number*)



- Map genes to KEGG pathways

- Use MinPath (Ye 2009) to find simplest pathway explanation for observed genes

- Remove pathways unlikely to be present due to low organismal abundance

- Smooth/fill gaps



Collapsing KO abundance into KEGG pathway abundance (or presence/absence) yields a smaller, more tractable feature set
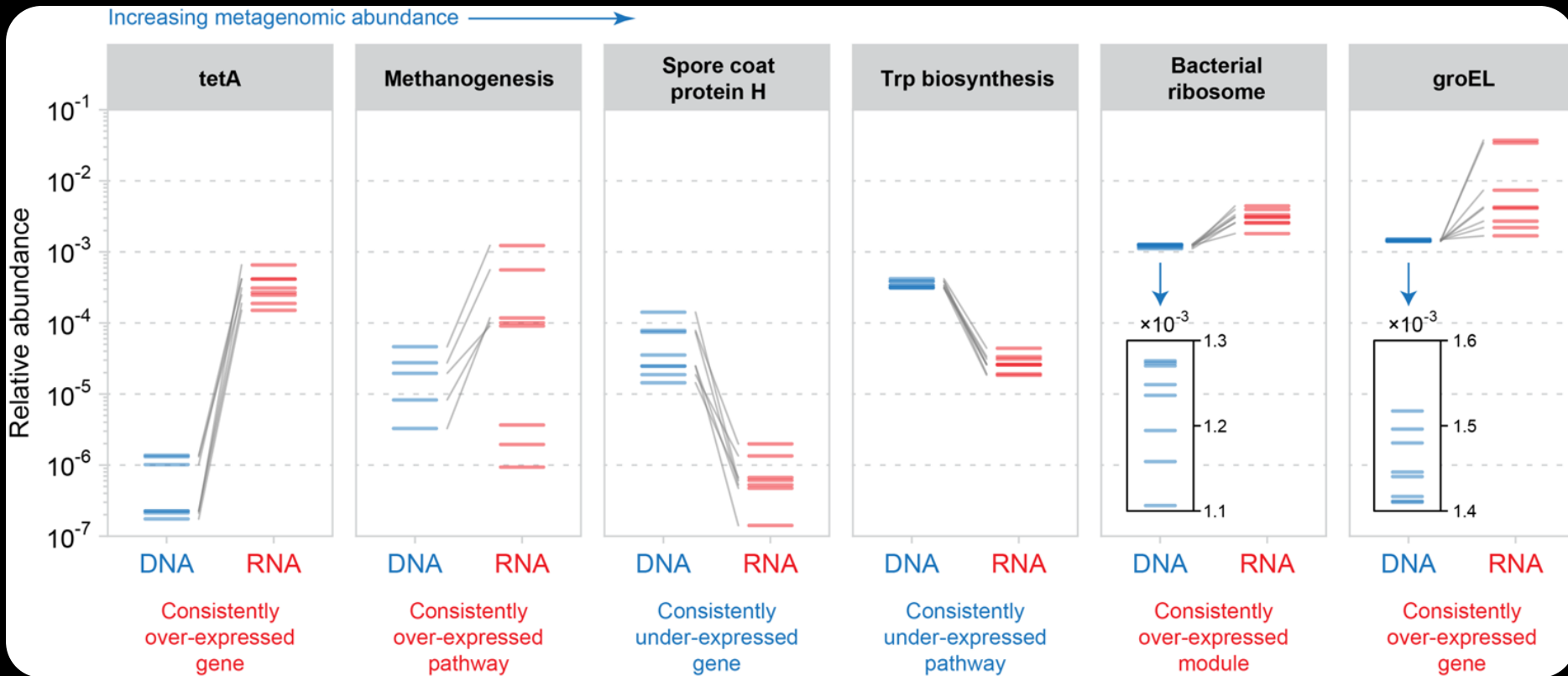
# HUMAnN accuracy



Validated against synthetic metagenome samples
(similar to MetaPhlAn validation)

Gene family abundance and pathway presence/absence
calls beat naïve best-BLAST-hit strategy

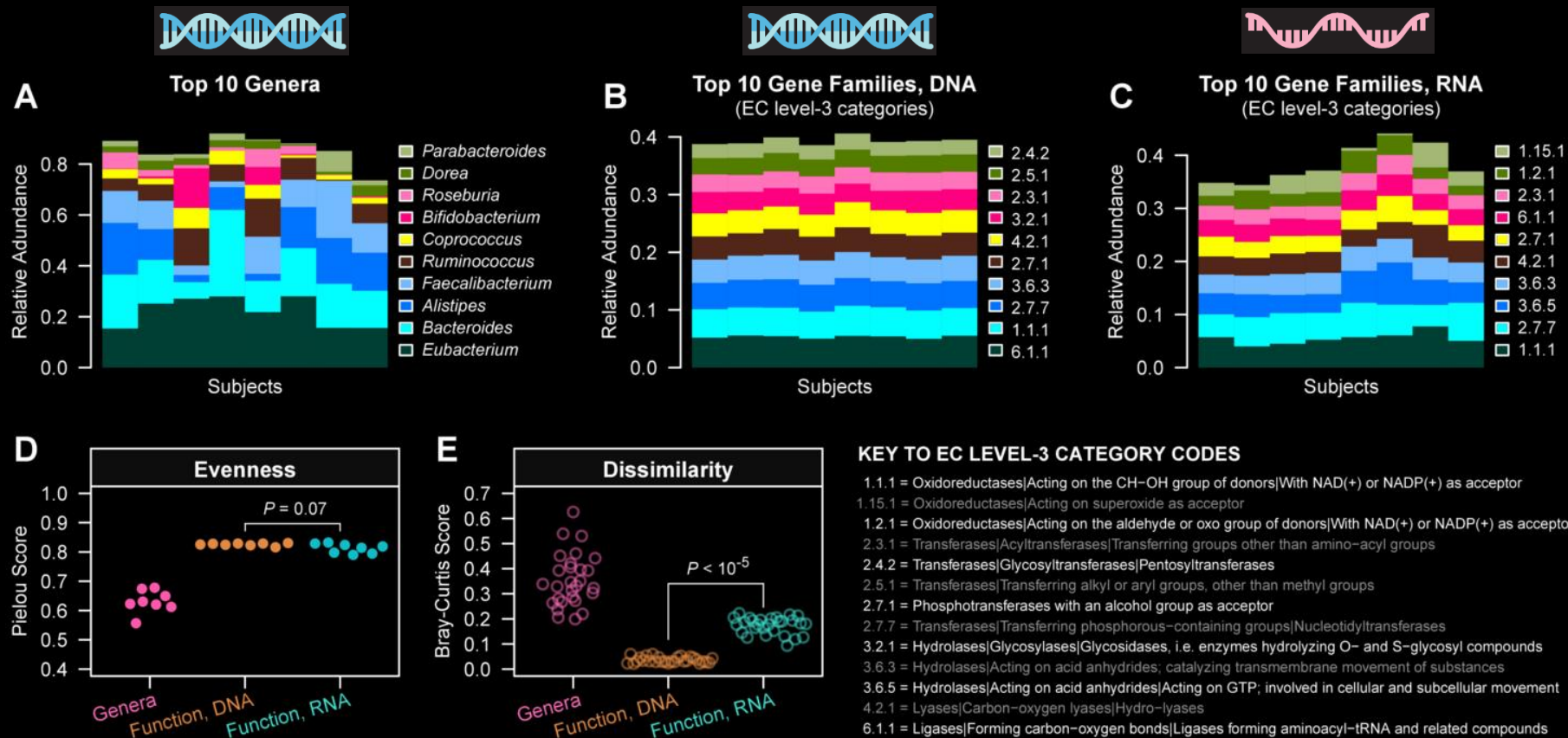# HUMAnN in action
## Functional potential (DNA) vs activity (RNA)



Functional metagenomics & metatranscriptomics
of 8 heathy human stool samples

Franzosa et al. *PNAS* 11:E2329-38 (2014)

# HUMAnN in action
## Conserved potential & variable activity



**A** Top 10 Genera

Relative Abundance (0.0–0.8) vs Subjects

Legend:
- Parabacteroides
- Dorea
- Roseburia
- Bifidobacterium
- Coprococcus
- Ruminococcus
- Faecalibacterium
- Alistipes
- Bacteroides
- Eubacterium

**B** Top 10 Gene Families, DNA (EC level-3 categories)

Relative Abundance (0.0–0.4) vs Subjects

Legend: 2.4.2, 2.5.1, 2.3.1, 3.2.1, 4.2.1, 2.7.1, 3.6.3, 2.7.7, 1.1.1, 6.1.1

**C** Top 10 Gene Families, RNA (EC level-3 categories)

Relative Abundance (0.0–0.4) vs Subjects

Legend: 1.15.1, 1.2.1, 2.3.1, 6.1.1, 2.7.1, 4.2.1, 3.6.3, 3.6.5, 2.7.7, 1.1.1

**D** Evenness

Pielou Score (0.4–1.0) vs Genera, Function, DNA, Function, RNA; $P = 0.07$

**E** Dissimilarity

Bray-Curtis Score (0.0–0.7) vs Genera, Function, DNA, Function, RNA; $P < 10^{-5}$

**KEY TO EC LEVEL-3 CATEGORY CODES**

1.1.1 = Oxidoreductases|Acting on the CH−OH group of donors|With NAD(+) or NADP(+) as acceptor
1.15.1 = Oxidoreductases|Acting on superoxide as acceptor
1.2.1 = Oxidoreductases|Acting on the aldehyde or oxo group of donors|With NAD(+) or NADP(+) as acceptor
2.3.1 = Transferases|Acyltransferases|Transferring groups other than amino−acyl groups
2.4.2 = Transferases|Glycosyltransferases|Pentosyltransferases
2.5.1 = Transferases|Transferring alkyl or aryl groups, other than methyl groups
2.7.1 = Phosphotransferases with an alcohol group as acceptor
2.7.7 = Transferases|Transferring phosphorous−containing groups|Nucleotidyltransferases
3.2.1 = Hydrolases|Glycosylases|Glycosidases, i.e. enzymes hydrolyzing O− and S−glycosyl compounds
3.6.3 = Hydrolases|Acting on acid anhydrides; catalyzing transmembrane movement of substances
3.6.5 = Hydrolases|Acting on acid anhydrides|Acting on GTP; involved in cellular and subcellular movement
4.2.1 = Lyases|Carbon−oxygen lyases|Hydro−lyases
6.1.1 = Ligases|Forming carbon−oxygen bonds|Ligases forming aminoacyl−tRNA and related compounds

# What's there: ShortBRED

Jim
Kaminski

- **ShortBRED** is a tool for <u>quantifying protein families in metagenomes</u>
  - Short Better REad Dataset

- Inputs:
  - FASTA file of proteins of interest
  - Large reference database of protein sequences (FASTA or blastdb)
  - Metagenomes (FASTA/FASTQ nucleotide files)
- Outputs:
  - Short, unique markers for protein families of interest (FASTA)
  - Relative abundances of protein families of interest in each metagenome (text file, RPKM)

- Compared to BLAST (or HUMAnN), this is:
  - Faster
  - More specific

# What's there: ShortBRED algorithm

- Cluster proteins of interest into families
  - Record consensus sequences

- Identify unique and common areas among proteins
  - Compared against each other
  - Compared against reference database
  - Remove all of these

- Remaining subseqs. uniquely ID a family
  - Record these as markers for that family

# What's there: ShortBRED marker identification

**Prots of interest**   **Reference database**   **Cluster into families**   **Identify short, common regions**

**True Marker**   **Junction Marker**   **Quasi Marker**

# What's there: ShortBRED family quantification



**Metagenome reads**

**ShortBRED markers**

**Translated search for high ID hits**

**Normalize relative abundances**

# What's there: ShortBRED is accurate



B. Antibiotic Resistance Genes Database
Correlation − 10% of Metagenome, 500 genes

ShortBRED: 0.95
Centroids: 0.815

Method
+ Centroids
O ShortBRED

Six synthetic metagenomes
from GemSim, spiked with
known proteins of interest:
ARDB = Antibiotic Resistance
VFDB = Virulence Factors

# What's there: ShortBRED is fast



Six synthetic metagenomes
from GemSim, spiked with
known proteins of interest:
ARDB = Antibiotic Resistance
VFDB = Virulence Factors

# Can we infer anything about function from 16S data?

Lyse cells
Extract DNA (and/or RNA)

**16S amplicons**

George Rice, Montana State University

**Meta'omic**

*Escherichia coli*
SSU rRNA
variability map
(bp 1 - 1542)

PCR to amplify the single
16S rRNA marker gene

Hello
my name is

Classify sequence
➔ microbe

Adapted from
Van de Peer, 1996

**Samples**

**Microbes**

**Relative abundances**

Genes,
Genomes,
Metabolic profiling,
Relative abundances,
Genetic variants...

# PICRUSt: Inferring community metagenomic potential from marker gene sequencing

**With Rob Knight, Rob Beiko**

One can recover **_general_** community function with reasonable accuracy from 16S profiles.

**http://picrust.github.com**

Pathways and modules     Orthologous gene families     Taxon abundances

**HUMAnN**



*y-axis:* PICRUSt Accuracy (Spearman, r)

*x-axis:* Average 16S distance to nearest reference genome (NSTI)

Legend:
- Hypersaline
- Soil
- Mammal
- Human

Relative abundance

# Plan

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
  - MetaPhlAn
- Tools for functional profiling
  - HUMAnN
  - ShortBRED
  - PICRUSt
- Tools for testing associations
  - LEfSe
  - MaAsLin
  - CCREPE
- Resources
- Research vignette (time permitting)

# Who is there?

# What are they doing?

| Sample # | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 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 |

# Who is there?

# What are they doing?

# What does it all mean?

| 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 |

# Properties of microbiome data

- Compositional nature (Σ = 1)
  - Abundance is relative, not absolute
- High dynamic range
- Often sparse (sample dominated by a few species)
- Noisy
- 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 |

# Properties of microbiome data

- General problem: correlate microbiome features with metadata (potentially controlling for other features)
- Intuitively summarize the results

| 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 |

# Recall that <u>ordination</u> is exploratory (no *p*-values for a trend, for example)

**LEfSe**: LDA Effect Size
*Finding metagenomic biomarkers*

*Nicola Segata*

**Data**

**Step 1** — Kruskal Wallis

**Step 2** — Wilcoxon on subclasses

**Step 3** — Signed LDA log-score

Class 1   Class 2

Features 1 through n

Samples 1 through m

Subclasses

0.13 ✗ A
0.01 ✓ B
0.19 ✗ C
0.00 ✓ D
0.03 ✓ E

0.03 ✓ Y
0.76 ✗ Z

Class 1

Class 2

**A** — Statistical consistency

**B** — Biological consistency

**C** — Biological Effect size

63

# Example LEfSe application:
## Find $O_2$-loving bugs (controlling for body site)

# Superimpose enrichments on the tree of life using GraPhlAn



LEfSe Associations



Metadata Rings

# MaAsLin
## Multivariate Association with Linear Models



- A more general solution for finding significant metagenomic associations in metadata-rich studies

**Tim Tickle**

# Microbiome downstream analyses: interaction network reconstruction



*It's a jungle in there* – microbial interactions follow patterns from classical macro-ecology.

## Mutualism



## Predation



## Competition



Given microbial relative abundance measurements over many samples, can we detect *co-occurrence and co-exclusion relationships?*

# Relative abundance data poses a problem for correlating metagenomic features

|      | Sample1 | Sample2 | Sample3 | sample4 | sample5 |
|------|---------|---------|---------|---------|---------|
| Bug1 | 100     | 100     | 100     | 100     | 100     |
| Bug2 | 1       | 10      | 100     | 1000    | 10000   |

|      | Sample1 | Sample2 | Sample3 | sample4 | sample5 |
|------|---------|---------|---------|---------|---------|
| Bug1 | 0.99    | 0.91    | 0.50    | 0.09    | 0.01    |
| Bug2 | 0.01    | 0.09    | 0.50    | 0.91    | 0.99    |



Absolute (cell) counts
*No bug1-bug2 correlation*

Relative abundance
*Spurious bug1-bug2 correlation*
(sequencing yields rel. ab.)

# CCREPE: Compositionality Corrected by REnormalization and PErmutation

Estimating a confidence interval



Estimating the null distribution

Emma
Schwager

# CCREPE: Compositionality Corrected by REnormalization and PErmutation



- Synthetic evaluation
- Random sample feature/tables
- No built-in correlation structure

# CCREPE: Compositionality Corrected by REnormalization and PErmutation



"Microbial co-occurrence relationships in the human microbiome."
Faust, et al. *PLoS Comp Biol*, 8:e1002606 (2012).

# Who is there?

# What are they doing?

# What does it all mean?

# Plan

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
  - MetaPhlAn
- Tools for functional profiling
  - HUMAnN
  - ShortBRED
  - PICRUSt
- Tools for testing associations
  - LEfSe
  - MaAsLin
  - CCREPE
- **Resources**
- **Research vignette (time permitting)**

# Using tools through Galaxy



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

# Tutorials available online

(click on your tool-of-interest)

# All tools are open source



http://bitbucket.org/biobakery/biobakery

# The bioBakery Virtual Machine

https://bitbucket.org/biobakery/biobakery/wiki/biobakery_wiki



Ubuntu base image preloaded and configured to run all
Huttenhower lab tools; one click up-and-running via Vagrant

# Thank you!

Curtis Huttenhower
Xochitl Morgan
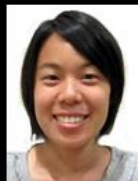Afrah Shafquat
Keith Bayer
George Weingart

Regina Joice
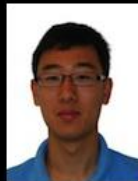Aleksandar Kostic
Chengwei Luo
Tiffany Hsu
Emma Schwager

Koji Yasuda
Kevin Oh
Boyu Ren
Andy Shi
Jim Kaminski

Joseph Moon
Randall Schwager
Levi Waldron
Nicola Segata

Wendy Garrett
Michelle Rooks

**Dirk Gevers**
Kat Huang

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

**Rob Beiko**
*Morgan Langille*

Jacques Izard

Katherine Lemon

Ruth Ley
Omry Koren

## 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 |

Bruce Birren   Mark Daly
Doyle Ward    Ashlee Earl

**BROAD** INSTITUTE

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

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

**Bruce Sands**

DANONE   CCFA   ALFRED SLOAN FOUNDATION

JDRF   NSF   NATIONAL INSTITUTES OF HEALTH

**http://huttenhower.sph.harvard.edu** 78

# Tutorial

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
  - **MetaPhlAn**
- Tools for functional profiling
  - HUMAnN
  - **ShortBRED**
  - PICRUSt
- Tools for testing associations
  - **LEfSe**
  - MaAsLin
  - CCREPE
- Resources
- Research vignette (time permitting)

# Plan

- Informal survey
- Metagenomics concepts & examples
- Tools for taxonomic profiling
  - MetaPhlAn
- Tools for functional profiling
  - HUMAnN
  - ShortBRED
  - PICRUSt
- Tools for testing associations
  - LEfSe
  - MaAsLin
  - CCREPE
- Resources
- **Research vignette (time permitting)**