# BioDIGS

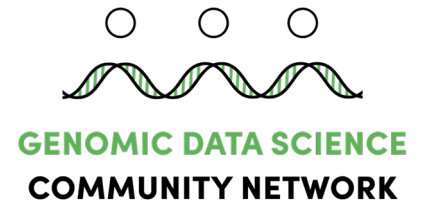## BioDiversity and Informatics for Genomics Scholars
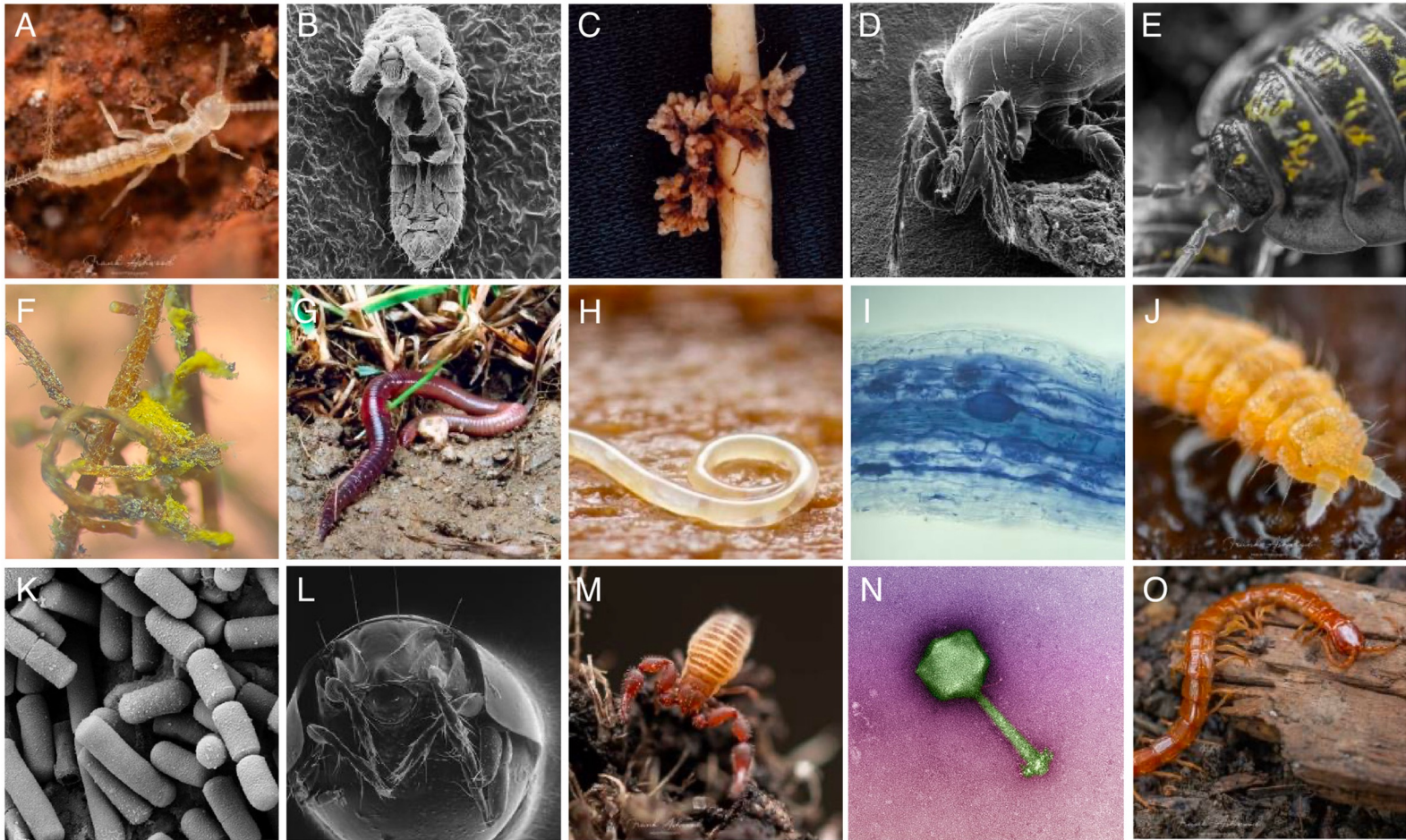
Michael Schatz
February 8, 2024
AGBT

# Disclosures

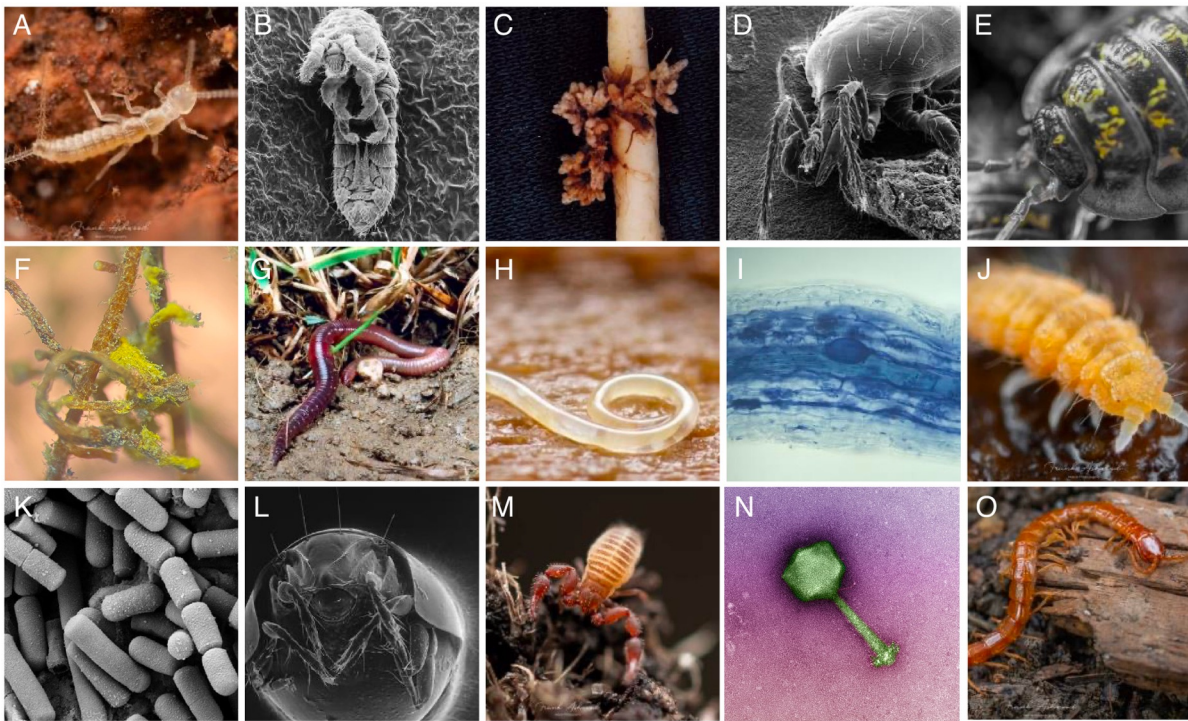I (Michael Schatz) am a Scientific Advisory Board Member for CosmosID

Rita Colwell is the founder of CosmosID and
an Advocacy Board Member
for the Genomic Data Science Community Network (GDSCN)

# Genomic Diversity

# Genomic Diversity



"… soil is likely home to 59% of life … making it the singular most biodiverse habitat on Earth."
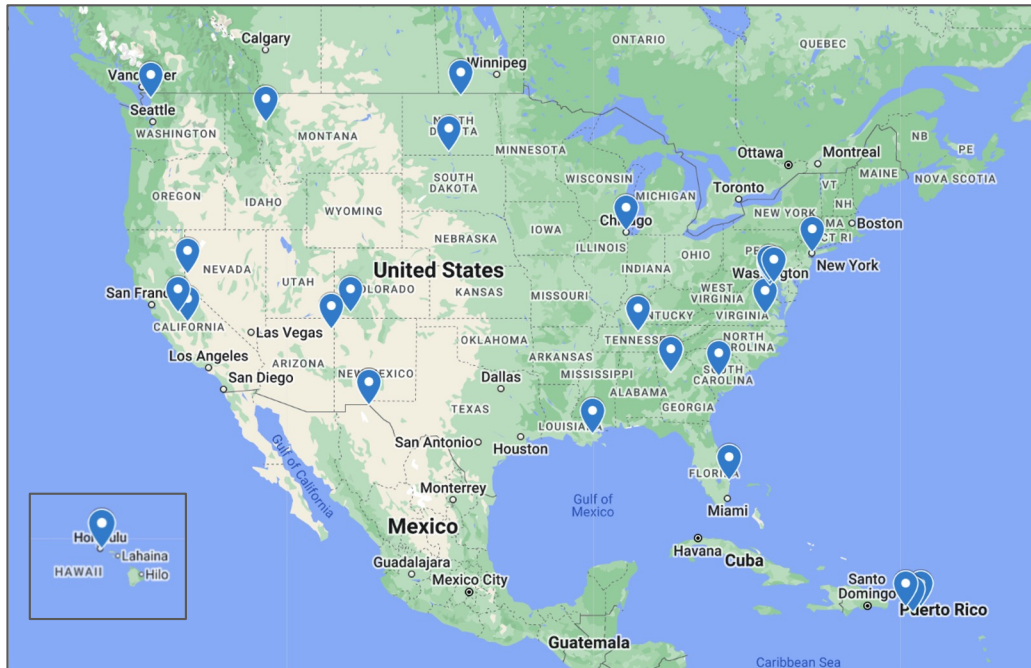
*Enumerating soil biodiversity*
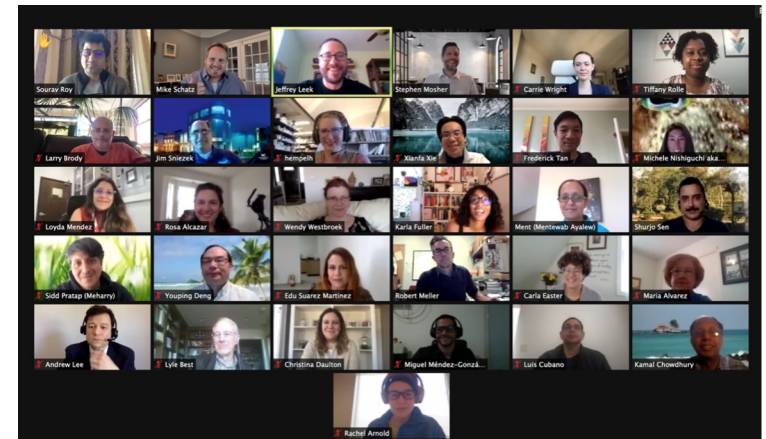Anthony et al (2023) PNAS.
doi:10.1073/pnas.2304663120

# Genomic Data Science Community Network

Promoting education and research in genomic data science at HBCUs, HSIs, TCUs, CCs, and other underserved



https://www.gdscn.org/

# GDSCN Needs & Resources



"Diversifying the Genomic Data Science Research Community." Genome Research (2022). doi:10.1101/gr.276496.121

**NORTH CAROLINA AGRICULTURAL AND TECHNICAL STATE UNIVERSITY**

Kristen Rhinehardt

GRADS-4C

# Foundations of GRADS-4C

*GRADS-4C: Genomic Research and Data Science Center for Computation and Cloud Computing*

Our Mission is to develop an educational and hands-on research training network with resources in computational genomics, data sciences, and cloud computing (CGDS) to investigate and improve human health, as well as develop the future CGDS workforce, particularly among underserved populations

Knowledge Transfer and Sustainability

Partnership

Growth Model Infrastructure Development

Trust

*Website: GRADS4C.ncat.edu*

AGGIES **DO**

Workshops

Course Integration Training

Symposia

Conferences

R Modules

Data Analytics

Access

Genomics Fundamentals

Computing Fundamentals and research ethics

Cloud interfacing

Platform training

HPC basics

Hands on workflow management

Scientific Communication principles

Scholarship and Infrastructure

**Diverse Genomic Workforce**

**New Cloud Based Tools**

Data Access and Seqr

Galaxy, Jupyter and Bioconductor

Terra

Dockstore and Outputs

DEI Training

Webinars

Career Pivoting

Social Media

# BioDIGS: What can we learn from the soil?



- *Microbiology & Metagenomics*
  - What's there? How does the genomic composition change in time and space?

- *Genomics & Bioinformatics*
  - Optimal approaches for metagenome assembly and classification? Merits of short- vs long-read sequencing?

- *Agriculture & Environment*
  - How do characteristics of the soil & soil microbiome modulate plant & animal development?

- *Public Health*
  - How does the soil microbiome influence the human microbiome & human health outcomes?

# BioDIGS Analysis



GENOMIC DATA SCIENCE
COMMUNITY NETWORK

## Galaxy Workflow



1. Flye (assembly)
2. Porechop (trim & QC)
3. BWA-MEM (align reads)
4. Racon (polish contigs)
5. Kraken2 (classify contigs)
6. AntiSMASH (BGC annotation)
7. Kallisto (BGC abundance)

Galaxy PROJECT

Katherine Ulbricht

Nia Davis

Ayalew lab
Spelman College

# BioDIGS Analysis

## Genomic Diversity



## Environmental Associations



## Human Health & Disease

# BioDIGS Analysis

### Genomic Diversity

### Environmental Associations

### Human Health & Disease

# DC+Baltimore pilot study sampling



**Montgomery Co.
24 samples
(students from MC)**

**Baltimore
24 samples
(students from
NDMU, MC, and CSM)**

Interactive map at: http://biodigs.org

Stony Run (near JHU)

Gwynns Falls Trailhead

Lake Needwood

# DC+Baltimore Data

## Sampling & Sequencing

- Zymo DNA/RNA shield
- Qiagen DNeasy PowerSoil Pro kit
- 2x150bp @ NovaSeq SP
- ~20M reads / site



FastQC: Mean Quality Scores

## Soil Testing

- Routine soil testing: Soil pH, Organic matter, Mehlich 3 extractable nutrients (P, K, Ca, Mg, Mn, Zn, Cu, Fe, B, S, Na and Al), Phosphorus Saturation Ratio
- Heavy metals: Arsenic, Cadmium, Chromium, Copper, Nickel, Lead, Zinc

grcf
grcf.jhmi.edu

UNIVERSITY OF DELAWARE
COOPERATIVE EXTENSION

# Taxonomic classification

- Mostly common soil microbes; Several species associated with nitrogen fixation

**Improved metagenomic analysis with Kraken 2**
Wood, Lu, Langmead (2019) Genome Biology do:10.1186/s13059-019-1891-0

# Taxonomic classification



- Mostly common soil microbes; Several species associated with nitrogen fixation

- Several hits to the genus Streptomyces, which produce many naturally occurring antimicrobials

- Certain species of Verrucomicrobia are important human probiotics for GI metabolism (Akkermansia)

**Improved metagenomic analysis with Kraken 2**
Wood, Lu, Langmead (2019) Genome Biology do:10.1186/s13059-019-1891-0

# Species classification

| Name | Number of raw reads | Classified reads | Chordate reads | Artificial reads | Unclassified reads | Microbial reads | Bacterial reads | Viral reads | Fungal reads | Protozoan reads |
|---|---|---|---|---|---|---|---|---|---|---|
| HHGCVDRX3-1-ACTCGGCAAT-TTCAGTTGTC_S41_L002 | 5,779,462 | 15.5% | 0.00114% | 0% | 84.5% | 12.5% | 8.79% | 0.00147% | 0.0535% | 0.00019% |
| HHGCVDRX3-1-ACTCGGCAAT-TTCAGTTGTC_S41_L001 | 5,501,528 | 15.4% | 0.002% | 0% | 84.6% | 12.4% | 8.7% | 0.00129% | 0.0515% | 0.000309% |
| HHGCVDRX3-1-CAATCGGCTG-TTCCTACAGC_S39_L002 | 5,857,162 | 13.3% | 0.00099% | 0% | 86.7% | 11.3% | 9.08% | 0.00082% | 0.0453% | 0.000137% |
| HHGCVDRX3-1-CAATCGGCTG-TTCCTACAGC_S39_L001 | 5,617,529 | 13.2% | 0.00141% | 0.0000178% | 86.8% | 11.2% | 8.96% | 0.000819% | 0.0445% | 0.000196% |
| HHGCVDRX3-1-GAACTGAGCG-CGCTCCACGA_S1_L002 | 5,315,099 | 12.9% | 0.00122% | 0% | 87.1% | 10.9% | 8.6% | 0.000978% | 0.0396% | 0.000263% |
| HHGCVDRX3-1-GAACTGAGCG-CGCTCCACGA_S1_L001 | 5,087,303 | 12.7% | 0.001% | 0% | 87.3% | 10.8% | 8.5% | 0.00106% | 0.0392% | 0.000216% |
| HHGCVDRX3-1-GATCAAGGCA-ATTAACAAGG_S8_L001 | 7,260,595 | 11.8% | 0.00169% | 0% | 88.2% | 9.51% | 7.06% | 0.000689% | 0.0122% | 0.0000826% |
| HHGCVDRX3-1-GATCAAGGCA-ATTAACAAGG_S8_L002 | 7,676,877 | 11.8% | 0.00168% | 0% | 88.2% | 9.56% | 7.12% | 0.000391% | 0.0122% | 0.0000912% |
| HHGCVDRX3-1-ACCGGCCGTA-AATATTGCCA_S36_L002 | 7,379,576 | 11.5% | 0.000474% | 0.0000136% | 88.5% | 10.8% | 9.31% | 0.0011% | 0.0394% | 0.0000813% |
| HHGCVDRX3-1-CGTCTCATAT-AGCTACTATA_S3_L002 | 5,157,206 | 11.5% | 0.000756% | 0% | 88.5% | 10.7% | 8.93% | 0.0014% | 0.0358% | 0.000427% |
| HHGCVDRX3-1-ACCGGCCGTA-AATATTGCCA_S36_L001 | 6,918,535 | 11.4% | 0.000795% | 0% | 88.6% | 10.7% | 9.19% | 0.00116% | 0.041% | 0.0000434% |
| HHGCVDRX3-1-CGTCTCATAT-AGCTACTATA_S3_L001 | 4,917,122 | 11.4% | 0.000936% | 0% | 88.6% | 10.6% | 8.82% | 0.00122% | 0.0348% | 0.000285% |
| HHGCVDRX3-1-CTAGTGCTCT-TACTGTTCCA_S7_L002 | 5,701,275 | 11.4% | 0.0016% | 0% | 88.6% | 8.31% | 6.55% | 0.000368% | 0.0175% | 0.0000526% |
| HHGCVDRX3-1-GGTTGCGAGG-TTGCTCTATT_S26_L002 | 9,214,523 | 11.4% | 0.00158% | 0% | 88.6% | 11% | 9.57% | 0.0014% | 0.0435% | 0.000184% |
| HHGCVDRX3-1-CTAGTGCTCT-TACTGTTCCA_S7_L001 | 5,464,649 | 11.3% | 0.00135% | 0% | 88.7% | 8.22% | 6.46% | 0.000421% | 0.0161% | 0.000146% |

Showing 1 to 15 of 96 entries

Previous 1 2 3 4 5 6 7 Next

# SMAG catalogue: 40,039 soil MAGs

**A genomic catalogue of soil microbiomes boosts mining of biodiversity and genetic resources**
Ma *et al* (2023) Nat Communication https://doi.org/10.1038/s41467-023-43000-z

# Genomic Diversity



**Metagenome profiling and containment estimation through abundance-corrected k-mer sketching with sylph**
Shaw and Yu (2023) bioRxiv. doi:10.1101/2023.11.20.567879

# BioDIGS Analysis
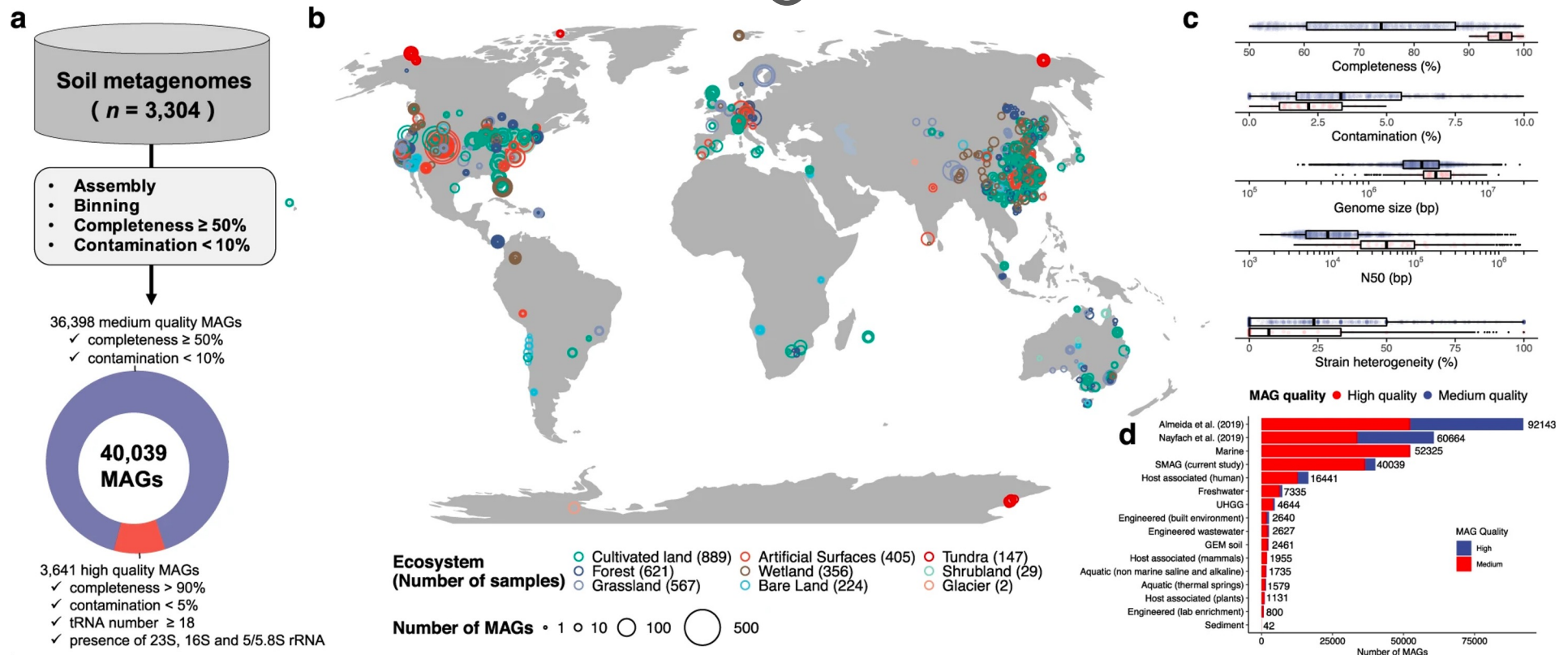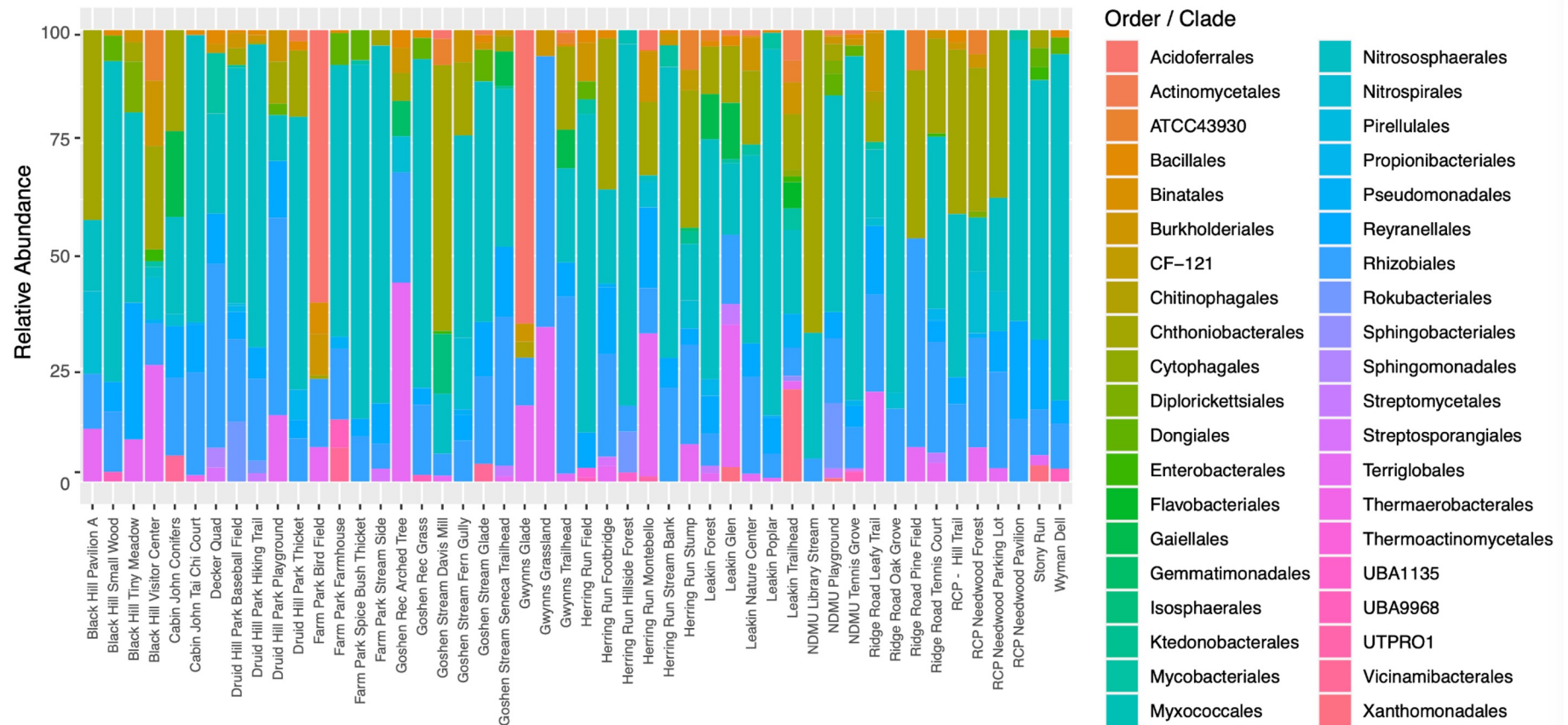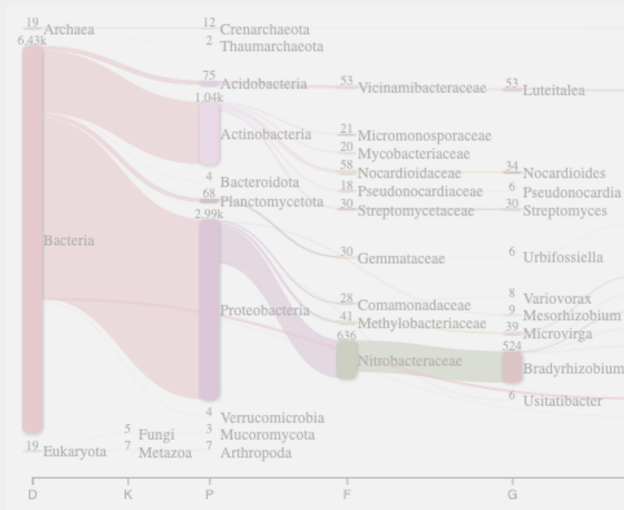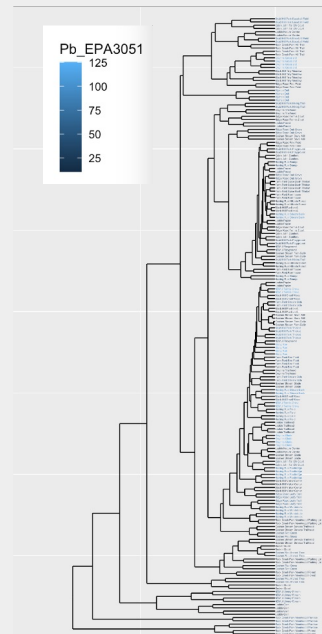


## Genomic Diversity

## Environmental Associations

## Human Health & Disease

# Heavy metal concentration

# Heavy metal associations

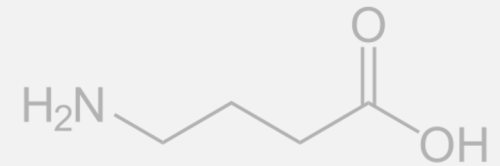# Heavy metal associations

# BioDIGS Analysis
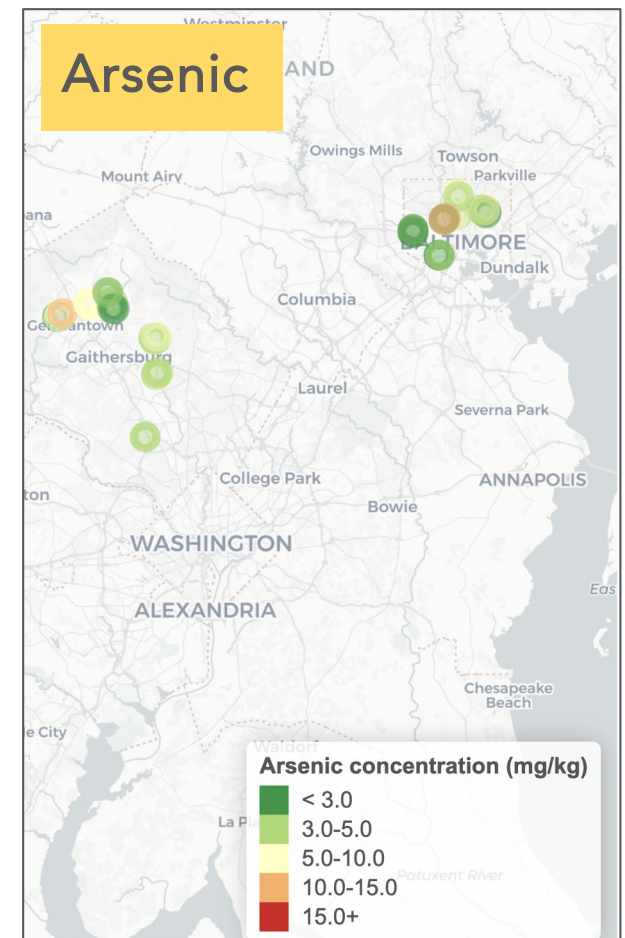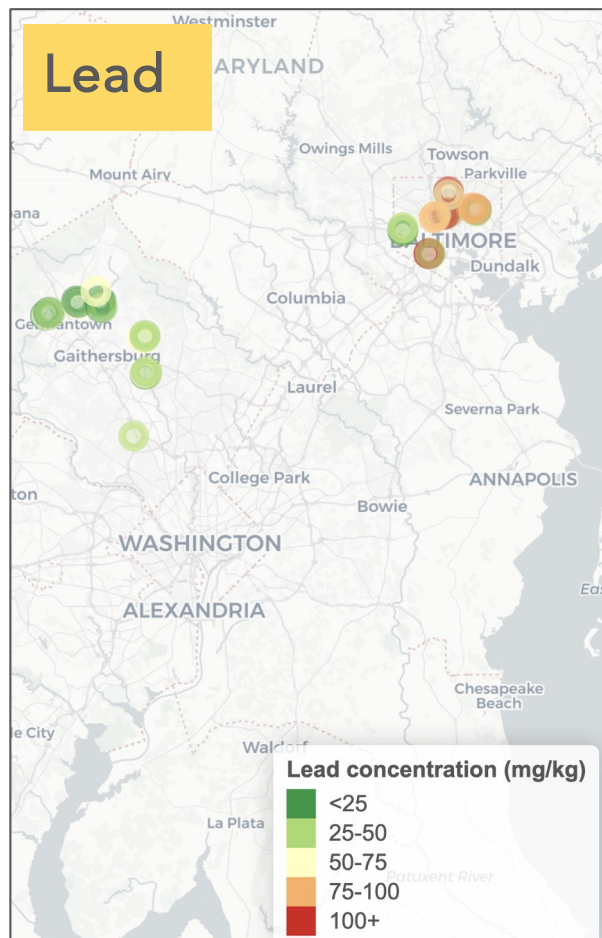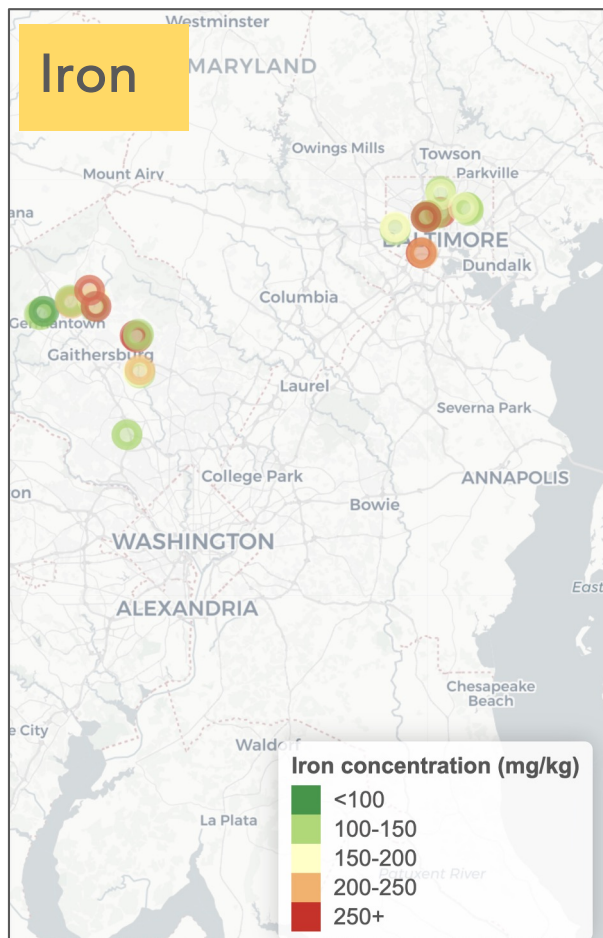


## Genomic Diversity

## Environmental Associations

## Human Health & Disease

# Soil metagenomes and human health

# Enter long reads...

## Oxford Nanopore PromethION

N: 14.82M
Total Yield: 70.2Gb
10kbp+: 8.3Gb
Read N50: 5,788bp

## PacBio HiFi Sequel IIe

N: 1.31M
Total Yield: 5.9Gb
10kbp+: 0.1Gb
Read N50: 5,180bp

## PacBio HiFi Revio

N: 9.97M
Total Yield: 85.2Gb
10kbp+: 39.5Gb
Read N50: 9,600bp

# Assembly Results



Revio HiFi

Total: 10.9 Gb
10kb+: 5.4 Gb
metaMDBG

Nanopore

Total: 4.3 Gb
10kb+: 3.5 Gb
metaFlye

Sequel HiFi

Total: 0.4 Gb
10kb+: 0.04 Gb
metaMDBG

Illumina

Total: 0.07 Gb
10kb+: 0.001 Gb
megahit

# Assembly Results



Complete Genomes

**Revio HiFi**
Total: 10.9 Gb
10kb+: 5.4 Gb
metaMDBG

**Nanopore**
Total: 4.3 Gb
10kb+: 3.5 Gb
metaFlye

**Sequel HiFi**
Total: 0.4 Gb
10kb+: 0.04 Gb
metaMDBG

**Illumina**
Total: 0.07 Gb
10kb+: 0.001 Gb
megahit

# Assembly Results



**Complete Genomes**

**Easy to bin**

Legend:
- **Revio HiFi**
  Total: 10.9 Gb
  10kb+: 5.4 Gb
  metaMDBG
- **Nanopore**
  Total: 4.3 Gb
  10kb+: 3.5 Gb
  metaFlye
- **Sequel HiFi**
  Total: 0.4 Gb
  10kb+: 0.04 Gb
  metaMDBG
- **Illumina**
  Total: 0.07 Gb
  10kb+: 0.001 Gb
  megahit

Axes:
- Y-axis: Contig length ($10^2$, $10^3$, $10^4$, $10^5$, $10^6$)
- X-axis: Cumulative sequence length (0.00e+00, 2.50e+09, 5.00e+09, 7.50e+09, 1.00e+10, 1.25e+10)

# Assembly Results

# PacBio HiFi - MAG Analysis



Analysis: Dan Portik, PacBio

https://github.com/PacificBiosciences/pb-metagenomics-tools

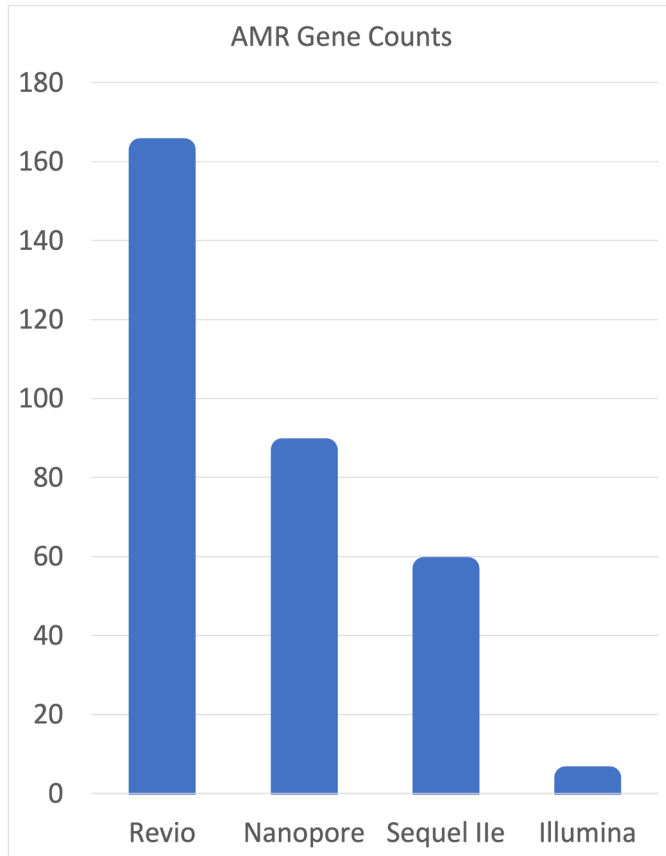Out of 158 MAGs (55 HQ), only 1 could be assigned to a known species
- 27 single contig genomes!
- ~25x more HQ MAGs / site than short reads
- Archaea are highly represented

# Antimicrobial Resistance

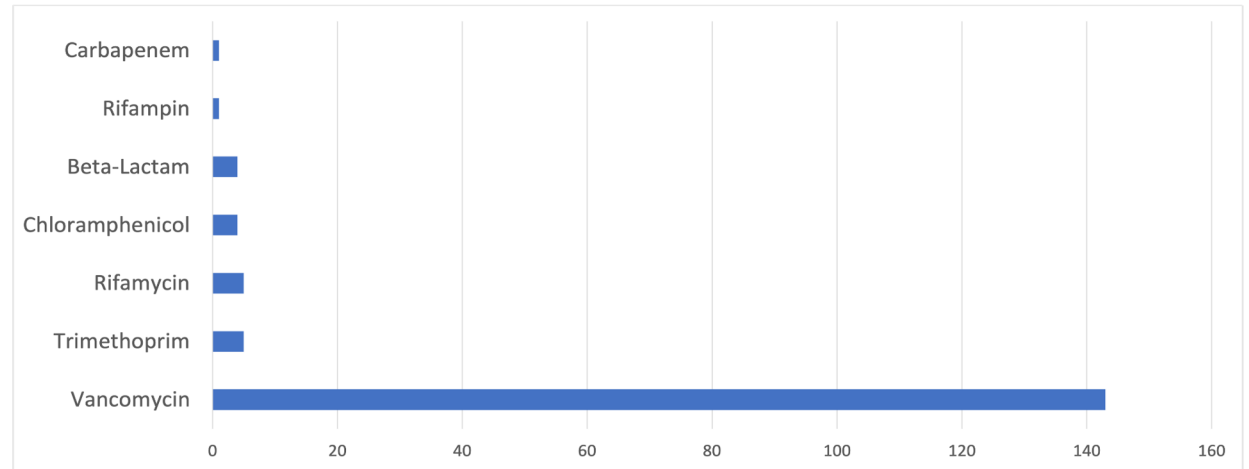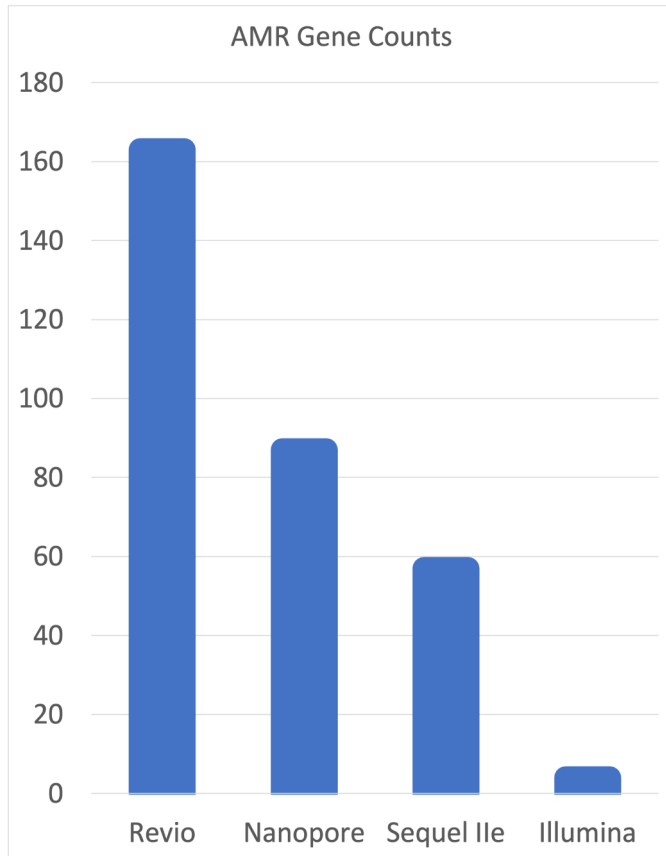

AMR Gene Counts

https://github.com/tseemann/abricate

# Antimicrobial Resistance



https://github.com/tseemann/abricate

# Antimicrobial Resistance



https://github.com/tseemann/abricate

# Metabolites in Humans & Microbes

**The neuroactive potential of the human gut microbiota in quality of life and depression**
Valles-Colomer et al. (2019) Nature Microbiology. doi:10.1038/s41564-018-0337-x

# Metabolites in Humans & Microbes
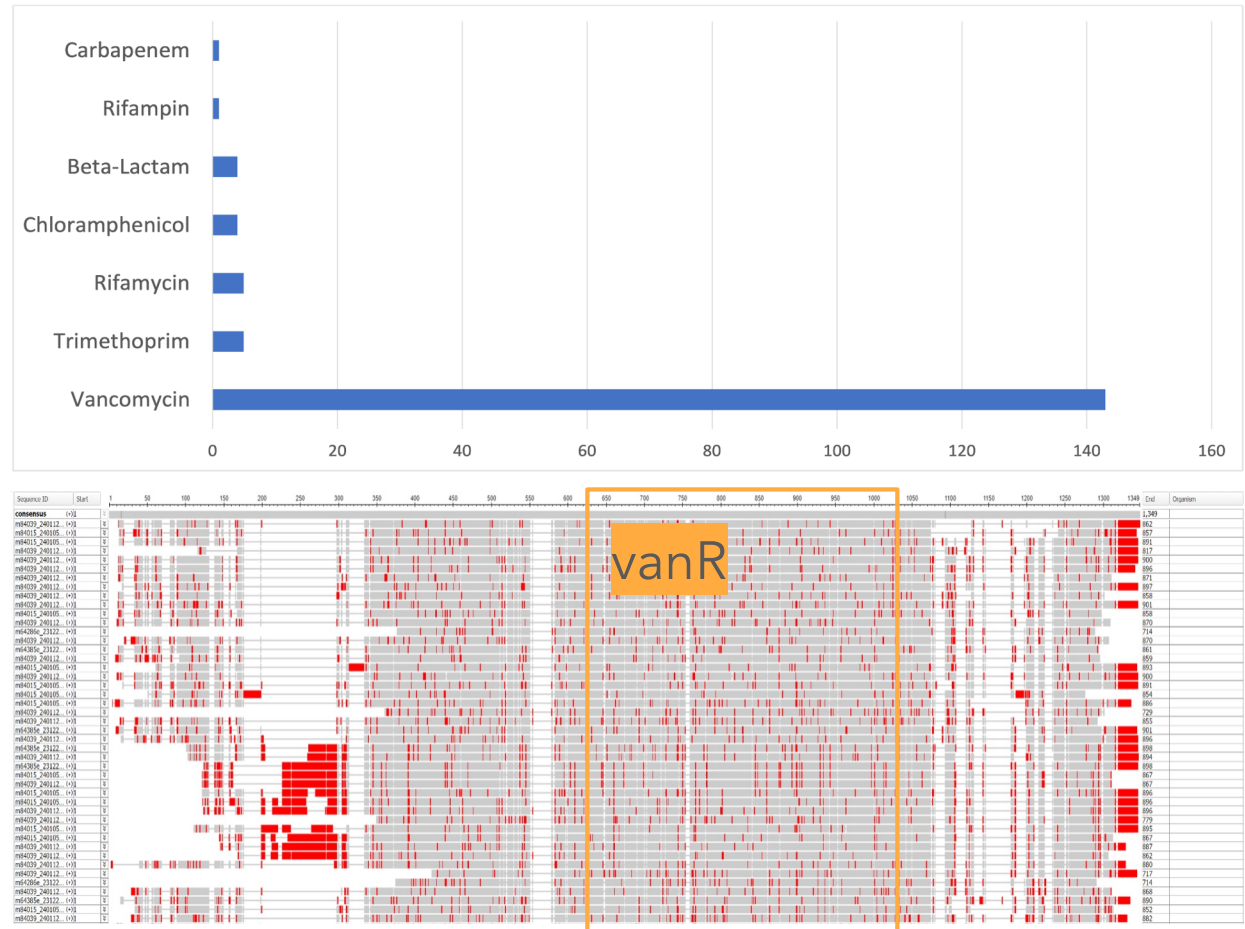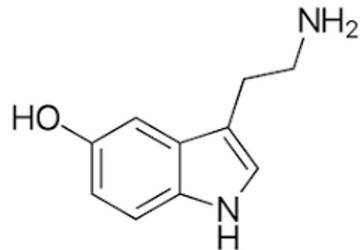
GENOMIC DATA SCIENCE
COMMUNITY NETWORK



**Serotonin**
mood, digestion, and sleep

**GABA**
chief inhibitory neurotransmitter

**Dopamine**
neurotransmitter for "pleasure"

**Norepinephrine**
Fight or flight response

**Acetylcholine**
Motor neuron neurotransmitter

**Butyrate**
host immune homeostasis

**The neuroactive potential of the human gut microbiota in quality of life and depression**
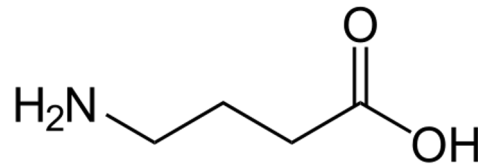Valles-Colomer et al. (2019) Nature Microbiology. doi:10.1038/s41564-018-0337-x

The neuroactive potential of the human gut microbiota in quality of life and depression
Valles-Colomer et al. (2019) Nature Microbiology. doi:10.1038/s41564-018-0337-x

# GABA pathway abundances



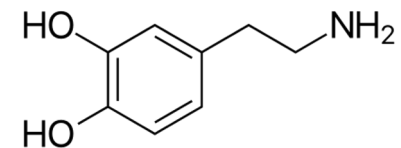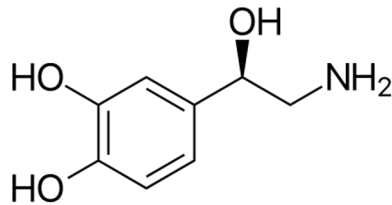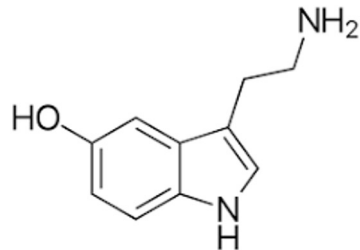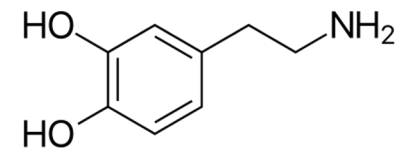Analysis: Kelly Moffat, CosmosID

**Species-level functional profiling of metagenomes and metatranscriptomes (HUMAnN 2.0)**
Franzosa et al (2018) Nature Methods. doi: doi.org/10.1038/s41592-018-0176-y

# Summary

- Exploring the interaction between soil, metagenomic diversity, and human health
    - Discover new genomes and genes
    - Discover new antimicrobial resistance genes, neurotransmitters present in the soil
    - Orders of magnitude improvements from long reads

- Empower the next generation of scientists
    - Teach state-of-the-art genomics & data science
    - Provide open access data & compute for community research

# Summary

- Exploring the interaction between soil, metagenomic diversity, and human health
    - Discover new genomes and genes
    - Discover new antimicrobial resistance genes, neurotransmitters present in the soil
    - Orders of magnitude improvements from long reads

- Empower the next generation of scientists
    - Teach state-of-the-art genomics & data science
    - Provide open access data & compute for community research

- Next Steps: Dynamics across space & time!
    - More institutions, longitudinal analysis
    - Exposures, Climate, Health data, etc
    - Training materials



We need your help!

linktr.ee/biodigs

# Student Acknowledgements

## Montgomery County
Madeline Graham, Daniel Chin, David Soussana

## Baltimore
Loraye Smith, Tyler Smith, Madeline Graham

## Seattle
Matheus Fernandes, Carl Pontino, Randon Serikawa, Joelle Taganna

## Clovis Community College
Malachi Whitford, Grace Freymiller, Domonique Advincula, Troy Burgess, Janet Castillo, Jennifer Elziade, Dorthy Esparza, Nicholas Foreman, Ana Hernandez, Glenda Medina, Christina Munoz, Nicole Potter, Quince Quintana, Nickie Ruiz, Ryan Wilder, Orion York

## Meharry Medical College
Lincoln Liburd II, Sydney Jamison, Destiney Ball, Claude Albritton, Arjun Pratap

## Virginia State University
Michael Marone

## Northern Virginia Community College
Rachel Marie Ametin, Joceph Duncan, Noha Elnawam, Sarah-Leila Kaci

## University of Texas - El Paso
Frida Delgadillo, Armando Jimenez, Keyan Ozuna
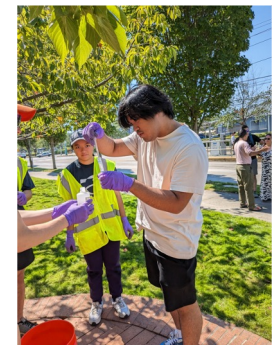
## El Paso Community College
Efren Barragan, Faith Chanhuhwa, Tania Da Silveira, Marco Ferrel, Josh Samuel Ikechi-Konkwo, Olivia Kelly, America Pinela, Ryley Stewart

## Spelman College
Natajha Graham, Nia Davis, Katherine Ulbricht

## University of Hawaii - Mānoa
Sudhir Kumar Rai, Yujia Qin, Ba Thong Nguyen, Mohammadamin Mahmanzar, Yu Chen, Isam Mohd Ibrahim, Donna Lee Kuehu, Asmita Pandey

# Acknowledgements

GENOMIC DATA SCIENCE
COMMUNITY NETWORK

National Human Genome Research Institute

COSMOSID Unlocking The Microbiome

PacBio

Oxford NANOPORE Technologies

tiny earth

Bloomberg Professors

**Mentewab (Ment) Ayalew, Ph.D.**
Associate Professor of Biology
Spelman College. Atlanta, GA

**Kristen L Rhinehardt, Ph.D.**
Assistant Professor
Computational Science & Engineering
North Carolina Agricultural and Technical
State University, Winston-Salem, NC

**Emily Biggane, Ph.D.**
Research Faculty
Environmental Toxicology
United Tribes Technical College
Bismarck, ND

**Siddharth (Sidd) Pratap, Ph.D.**
Director of Bioinformatics
Meharry Medical College
Nashville, TN

**Rosa Alcazar, Ph.D.**
Biology Faculty
Clovis Community College
Fresno, CA

**Ava Hoffman, Ph.D.**
Senior Staff Scientist
Fred Hutch Cancer Center
Seattle, WA

**Natalie Kucher**
Program Manager
Johns Hopkins University
Baltimore, MD

GDSCN
project

# Thank you!

http://biodigs.org