

# Metagenomics: Characterization of Microbial Communities using NGS

WGS

**UTSouthwestern**  
Medical Center | BICF

*Brandi Cantarel, PhD*  
10/19/2016

- Community Structure Based on WGS
  - 16S VS WGS
  - Least Common Ancestor
  - Reference Genome databases
  - Assignment Strategies
- Data Processing Workflow
  - QC
  - Functional Databases
  - Analysis Platforms
  - DIY Analysis
- Non-DNA based technologies: Metatranscriptomics, Metaproteomics and Metabolomics
  - Expression (RNASeq)
  - Translation (Proteomics)
  - Metabolites (Metabolomics)

- Community Structure Based on WGS
  - 16S VS WGS
  - Least Common Ancestor
  - Reference Genome databases
  - Assignment Strategies
- Data Processing Workflow
  - QC
  - Functional Databases
  - Analysis Platforms
  - DIY Analysis
- Non-DNA based technologies: Metatranscriptomics, Metaproteomics and Metabolomics
  - Expression (RNASeq)
  - Translation (Proteomics)
  - Metabolites (Metabolomics)

## Taxonomic Assessment using 16S

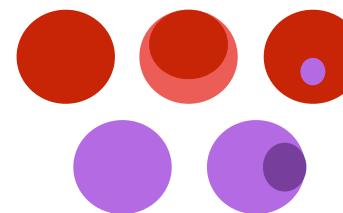
- 16S is targeted sequencing for a single gene which acts as a marker for organisms
- Pros
  - Well established
  - Relatively inexpensive \$50-\$100/sample
  - Amplifies only bacteria not host or environmental fungi, plants, etc
- Cons
  - Amplifies only bacteria not viruses, microbial fungi, archaea, etc
    - Although can be paired with 18S and archaeal specific 16S
  - Is based on a very well conserved gene, making it hard to resolve species and strains
  - V-region choice can bias results

# Taxonomic Assignment using WGS

- WGS (whole genome shotgun) aims to sequence the “whole” metagenome
- Pros
  - Not biased by amplicon primer set
  - Not limited to by conservation of the amplicon
  - Can also provide functional information
- Cons
  - Environmental contamination, including host
  - More expensive - \$1000+/sample
  - Complex data analysis
    - Requires high performance computing, high memory, high compute capacity

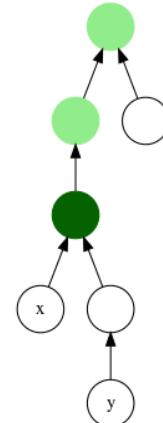
## Taxonomic Assignment: Complex Analysis

- All of the organism mixed together
  - It's hard to bin all of the reads from one organism (strain or species) for deconvolution
  - Reads are short
  - Reads can potentially share similarity to multiple taxa
- Lateral gene transfer
  - Not all of the genes in a genome “shares” the same evolutionary history



# Least Common Ancestor Taxonomic Assignment

- Reads can potentially share similarity to multiple taxa
- Least Common Ancestor allows for the taxonomic assignment when similarity is shared to multiple taxa
- Dependent on the taxonomic tree and similarity to genomes
  - Remember there are different versions of bacterial taxonomy



## Sources of Reference Genomes for Comparison

**JGI GOLD**  
GENOMES ONLINE DATABASE

Home Search Distribution Graphs Biogeographical Metadata Statistics References Team Help News

Welcome to the Genomes OnLine Database GOLD Release v.5

GOLD: Genomes Online Database, is a World Wide Web resource for comprehensive access to information regarding genome and metagenome sequencing projects, and their associated metadata, around the world.

1. Register      2. Annotate      3. Publish

Register your project information and Metadata in the Genomes Online Database

Annotate your microbial genome or metagenome with IMG/ER or IMG/MER

Publish your genome or metagenome in open access standards-supportive journal.

Download Excel Data file  
File last generated: 28 Sep, 2015

NCBI Resources How To

Genome

### Reference and representative genomes

**HMP**  
NIH HUMAN MICROBIOME PROJECT

Current News  
• January 2015 Metagenome Analysis Workshop March 3-6

OVERVIEW REFERENCE GENOMES MICROBIOME ANALYSIS HEALTH & ETHICS RESOURCES OUTREACH DATA BROWSER

home > reference genomes

Microbial Reference Genomes

The HMP plans to sequence, or collect from publicly available sources, a total of 3000 reference genomes isolated from human body sites. The information gained from the reference genomes will aid in taxonomic assignment and functional annotation of 16S rRNA and metagenomic wgs sequence, respectively, from microbiome samples. More information can be found below and on the NIH Common Fund Site.

GET DATA GET TOOLS

**e!EnsemblBacteria** Sequence Search | BLAST | Tools | Documentation

Search for a gene      Search for a genome

Search all species... Go Start typing the name of a genome...  
e.g. ftsZ or uridine\*

### Access to over 20,000 Bacterial Genomes

- Search for a gene - type the name of a gene or other identifier into the search box above
- Find a genome - click in the 'browse a genome' box above and start typing your genome name to find matching genomes
- View full list of all Ensembl Bacteria species
- Access Ensembl Bacteria programmatically

# Strategies for Taxonomic Assignment of WGS

- Compositional Based Taxonomic Assignment
  - This is assignment based on “base content”
- Sequence Alignment Based Taxonomic Assignment
  - This assignment is based on an alignment
- Maker Gene Based Taxonomic Assignment
  - This assignment is based similarity on a subset of the reads to conserved genes.

## Composition Based Taxonomic Assignment

- GC content (TETRA)
- K-mer based (naïve Bayes classifier)
- Pros
  - Speed
  - Require less compute power compared to alignment-based methods.
- Cons
  - Requires query sequences of sufficient length
  - Genomes in the same clade (genera, family, etc) can be quite heterogenous in some regions

ATTGCC	17
AGTGCC	10
CCGTGA	25
TTGTGA	57
CCGTGA	12

# Sequence Alignment Based Taxonomic Assignment

- BLAST/Megablast
- Malt/Diamond
- Kraken
- Pros
  - Higher assignment accuracy and specificity
- Cons
  - These methods are computationally intense because they either:
  - Require a high memory machine to generate the database and complete the searches
  - Require high number of cpus to complete the searches

# Marker Gene Based Composition

- MetaPhlAn2
  - Relies on ~1M unique clade-specific marker genes
- PhyloSift
  - Uses a reference database of protein and RNA sequences
- Taxonomer
  - Based on 16S rRNA
- Pros
  - Less computationally intensive
  - Accurate for the marker gene composition
- Cons
  - Only assigns a subset of the data ie can't determine taxonomy of certain function.

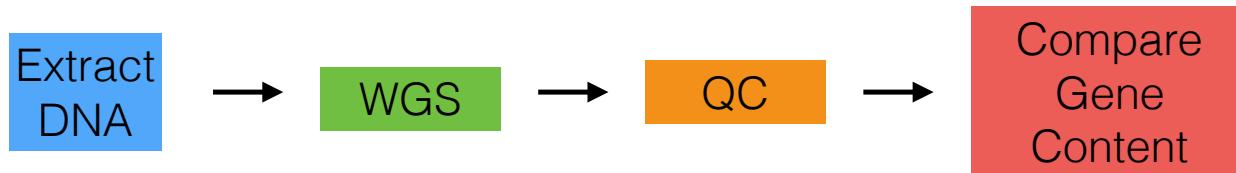
# WGS Taxonomy Assignment and Visualization

- Taxonomer
  - <http://taxonomer.iobio.io>
- Megan
  - <http://ab.inf.uni-tuebingen.de/software/megan5/>
  - Tool with WGS taxonomic assignment (based on BLAST) and functional assignment
- MG-RAST
  - <http://metagenomics.anl.gov/>
  - Online tools with WGS taxonomic assignment and functional assignment
- MetaCRAM

- Community Structure Based on WGS
  - 16S VS WGS
  - Least Common Ancestor
  - Reference Genome databases
  - Assignment Strategies
- Data Processing Workflow
  - QC
  - Functional Databases
  - Analysis Platforms
  - DIY Analysis
- Non-DNA based technologies: Metatranscriptomics, Metaproteomics and Metabolomics
  - Expression (RNASeq)
  - Translation (Proteomics)
  - Metabolites (Metabolomics)

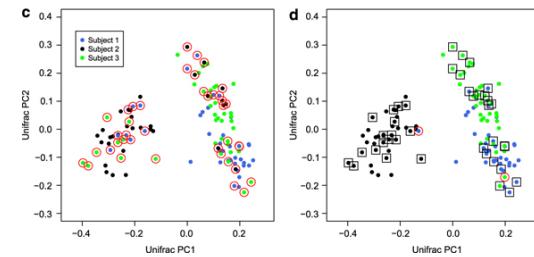
# What is the Functional Capability

- Identify organisms present — if closely related to organisms with sequenced genomes
- Identify gene families present — if homologs have been functionally characterized
- Identify functional pathways present — if homologs have been annotated to gene pathways
- Identify new species/strains — if assemblies are of sufficient depth



## Quality Control

- Negative Controls are the best way to identify microbial lab contamination
- Sequencing Errors
  - Low Quality Bases
  - Homopolymer Strings
  - Too short trimmed reads
- Biological and Technical Replicates
  - Helps to ensure group trends and identify sample mislabeling and possible “compromised” samples



Knights D, Kuczynski J, Koren O, Ley RE, Field D, Knight R, DeSantis TZ, Kelley ST. Supervised classification of microbiota mitigates mislabeling errors. ISME J. 2011 Apr;5(4):570-3. doi: 10.1038/ismej.2010.148. Epub 2010 Oct 7. PubMed PMID: 20927137; PubMed Central PMCID: PMC3105748.

# Host/Environmental Contamination

- In the human body — in human stool composes < 5% of reads, but the skin can be > 80% human reads
- Fungal, plant and soil bacteria can contaminate environmental samples.
- When you are collecting samples from “inside” of a habitat, it can be easy to contaminate the site with another site ie a colon biopsy with rectal microbiome.
- The natural environment can also contaminate samples, even the lab.

## Metagenome Databases



NIH HUMAN  
MICROBIOME  
PROJECT



# Comprehensive Functional Databases

eggNOG  
version 3.0



- KEGG
- eggNOG/COG
- PFAM
- SEED used by MG-RAST
- MetaCyc
- Uniref

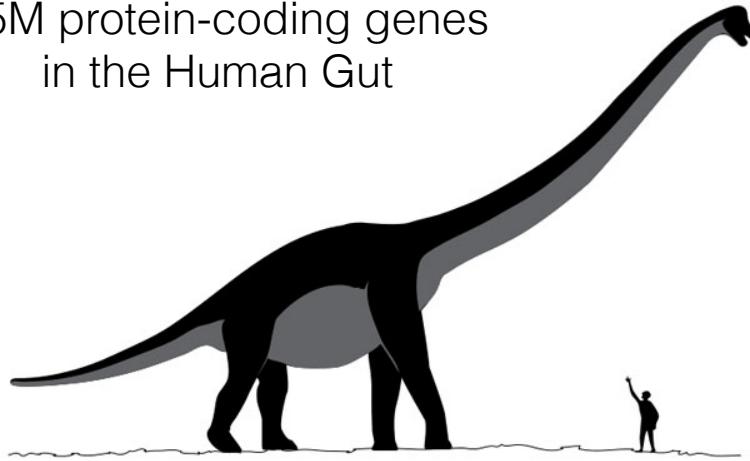


## Specialized Functional Databases

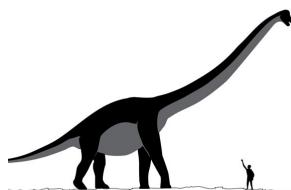
- Antibiotic resistant genes
  - <http://ardb.cbcn.umd.edu/>
  - <https://card.mcmaster.ca/>
- Virulent factors
  - <http://www.mgc.ac.cn/VFs/main.htm>
- Carbohydrate Active Enzymes
  - [www.cazy.org](http://www.cazy.org)
- Phage
- Proteases
  - <http://merops.sanger.ac.uk/>
- Transporters
  - <http://www.membranetransport.org/>

# Microbial Gene Content

3-5M protein-coding genes  
in the Human Gut



~25K Genes in the Human Genome



## Metagenomic Datasets Tend to Be Big

- Depending on taxonomic diversity, sequencing depth for each sample averages from 1M - 100M reads
- Analysis programs such as assembly and some alignment algorithms require >100 GB of RAM
- High performance computing platform is necessary
  - There are some publicly available resources for analysis

# Available Web-based Analysis Pipelines

- MG-RAST
  - Preference given to “public” datasets
  - Every easy to use
- EBI Metagenomics
  - Includes data visualization and customizable samples comparisons
  - DIAG
- JGI Integrated Microbial Genomes
  - Includes data visualization and customizable samples comparisons
- CloVR
  - Cloud-based workflow manager
  - Can run pipelines on your desktop
  - Available on the Academic Cloud



**MG-RAST**  
metagenomics analysis server

**Warning:** This application has been optimized for the Firefox browser. Since you are using Chrome, many features will not be available and / or behave incorrectly.  
Firefox is freely available [here](#).

[Browse Metagenomes](#)

[About](#) [Register](#) [Contact](#) [Help](#) [Upload\\*](#) [News](#)

MG-RAST (the Metagenomics RAST) server is an automated analysis platform for metagenomes providing quantitative insights into microbial populations based on sequence data.

The server primarily provides upload, quality control, automated annotation and analysis for prokaryotic metagenomic shotgun samples. MG-RAST was launched in 2007 and has over 12,000 registered users and 212,065 data sets. The current server version is 3.6. We suggest users take a look at MG-RAST for the impatient. Also available for download is the MG-RAST manual.

# of metagenomes 212,065  
# base pairs 85.9 Tbp  
# of sequences 683.67 billion  
# of public metagenomes 30,034

• MG-RAST newsletter, August 2015  
• Upcoming change to MG-RAST upload (early August 2015)  
• MG-RAST API available  
• MG-RAST newsletter, September 2014

\* login required

This project has been funded in part with Federal funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN272200900040C.

This work was supported in part by the Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, under Contract DE-AC02-06CH11357.

[cite MG-RAST](#) [cite MG-RAST API](#)

**MG-RAST**  
metagenomics analysis server

Warning: This application has been optimized for the Firefox browser. Since you are using Chrome, many features will not be available and/or behave incorrectly.  
Firefox is freely available [here](#).

Browse Metagenomes    Register    Contact    Help    Uploads    News

About    Log in    Log out    Search    Advanced search

MS-RAST (the Metagenomics RAST) server is an automated analysis platform for metagenomics samples. It provides quality control, automated assembly, and annotation of metagenomic data sets. The server primarily provides upload, quality control, automated assembly, and annotation of metagenomic data sets. MG-RAST was launched in 2007 and has over 12,000 registered users and 212,065 data sets. The current server version is 3.8. We suggest that you use the latest version of the MG-RAST software for your analysis. Also available for download is the MG-RAST manual.

# of metagenomes: 212,065    # of samples: 854,795    # of requests: 683,471    # of public metagenomes: 34,045

MG-RAST newsletter, August 2015  
Upcoming change in MG-RAST update (early August 2015)  
MG-RAST API available  
MG-RAST newsletter, September 2014

High res PDF    MG-RAST API

The project has been funded in part with Federal Funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN272200400404C. This work was supported by the Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, under Contract DE-AC02-05ER25435.

✓ qc_stats	10/7/2015, 9:19:23 AM
✓ preprocess	10/7/2015, 9:16:29 AM
✓ dereplication	10/7/2015, 9:16:46 AM
✓ screen	10/7/2015, 9:16:54 AM
✓ rna detection	10/7/2015, 9:17:11 AM
✓ rna clustering	10/7/2015, 9:17:29 AM
✓ rna sims blat	10/7/2015, 9:17:48 AM
✓ genecalling	10/7/2015, 9:17:48 AM
✓ aa filtering	10/7/2015, 9:17:53 AM
✓ aa clustering	10/7/2015, 9:19:26 AM
✓ aa sims blat	10/7/2015, 10:00:58 AM
✓ aa sims annotation	10/7/2015, 10:09:27 AM
✓ rna sims annotation	10/7/2015, 9:17:54 AM
✓ index sim seq	10/7/2015, 10:16:36 AM
✓ md5 annotation summary	10/7/2015, 10:19:23 AM
✓ function annotation summary	10/7/2015, 10:10:47 AM
✓ organism annotation summary	10/7/2015, 10:10:29 AM
✓ lca annotation summary	10/7/2015, 10:10:54 AM
✓ ontology annotation summary	10/7/2015, 10:10:51 AM
✓ source annotation summary	10/7/2015, 10:10:03 AM
✓ md5 summary load	10/7/2015, 10:32:53 AM
✓ function summary load	10/7/2015, 10:21:30 AM
✓ organism summary load	10/7/2015, 10:15:54 AM
✓ lca summary load	10/7/2015, 10:16:01 AM
✓ ontology summary load	10/7/2015, 10:17:31 AM
✓ done stage	10/7/2015, 10:35:28 AM
✓ notify job completion	10/7/2015, 10:35:31 AM

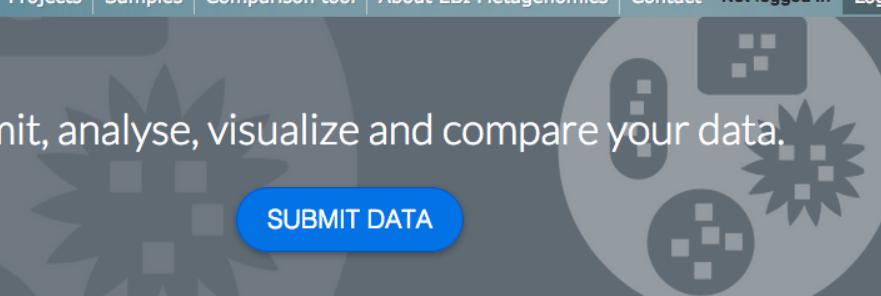
# EBI Metagenomics



Home    Submit data    Projects    Samples    Comparison tool    About EBI Metagenomics    Contact    Not logged in    Login

Submit, analyse, visualize and compare your data.

**SUBMIT DATA**



 **8192** data sets



**4167** metagenomics  
**780** metatranscript  
**3178** amplicons  
**67** assemblies



**5608** runs  
**4879** samples  
Public    **138** projects



**2584** runs  
**2522** samples  
Private    **93** projects



CloVR | Automated Sequence Analysis from Your Desktop

Welcome      Protocols      Getting Started      Download      Developers      Blog

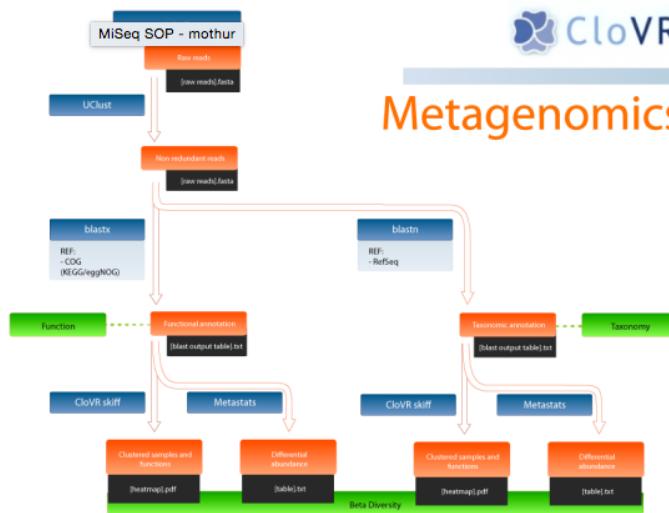
**Try CloVR**  
Read [CloVR tutorials](#) and run test applications on the [DIAG cloud](#).

**Get CloVR**  
Download and install CloVR to run supported microbial sequence analysis locally or on the cloud.

**About CloVR**  
The Cloud Virtual Resource supports user-friendly automated microbial sequence analysis applications.



## Metagenomics



We will sometimes refer to the protocol described above as *ClovR-Metagenomics (no-orfs)*, which is our default. For users who wish to first call open reading frames (ORFs) on their sequences, we provide an [alternative metagenomic analysis protocol](#) that utilizes *MetaGene* for ORF-calling prior to functional assignment.

**JGI** **IMG/EDU**

INTEGRATED MICROBIAL GENOMES / EDUCATION SITE

ALL Genomes Quick Genome Search:  Go

IMG Home Find Genomes Find Genes Find Functions Compare Genomes Analysis Cart My IMG Data Marts Using

**IMG Content**

Datasets

Bacteria	25871
Archaea	532
Eukarya	190
Plasmids	1186
Viruses	3888
Genome Fragments	1192
Total Datasets	32859

Genome by Metadata Project Map Metagenome Projects Map System Requirements

Hands on training available at the [Microbial Genomics & Metagenomics Workshop](#)

The Integrated Microbial Genomes (IMG) system ([Nucleic Acids Research, Volume 42 Issue D1](#)) serves as a community resource for comparative analysis and annotation of all publicly available genomes from three domains of life in a uniquely integrated context. Plasmids that are not part of a specific microbial genome sequencing project and phage genomes are also included into IMG in order to increase its genomic context for comparative analysis.

Count	Total
DNA, number of bases	<a href="#">135,697,930,103</a>
Total Genes	<a href="#">98,482,933</a>
Total Genomes	<a href="#">32,859</a>

**IMG Statistics**

All Genomes

Legend:

- Bacteria
- Archaea
- Eukarya
- Plasmids
- Viruses
- Genome Fragments

**News**

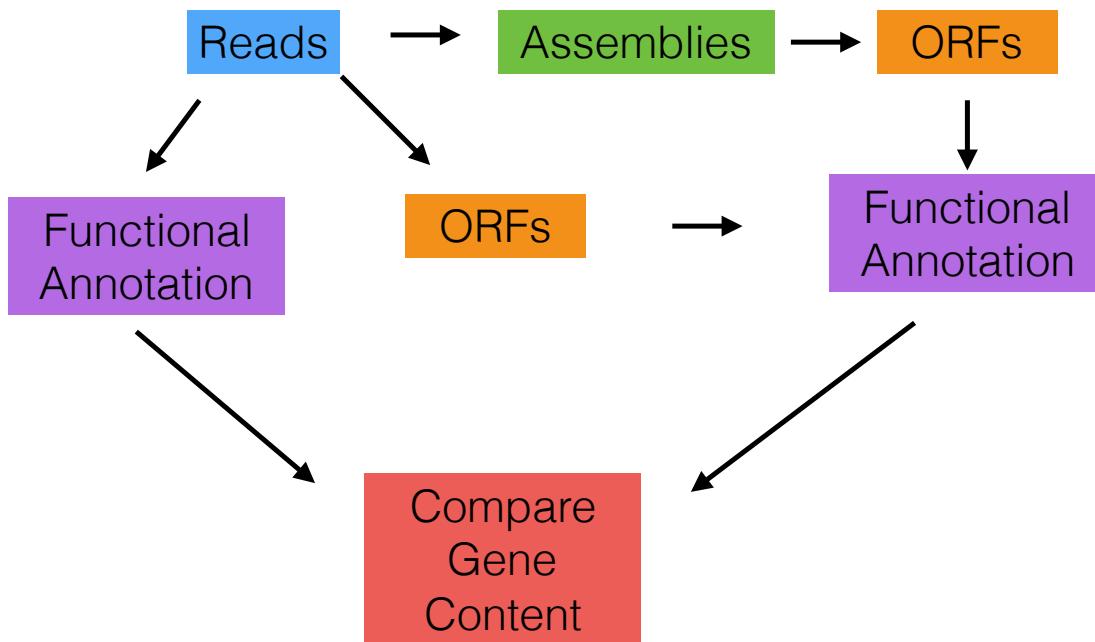
- Oct 5 2015 After 10 Years, IMG Still Revolutionizing Genomics
- Sep 2015 IMG ABC Data Mart
- Sep 2015 MGM Workshops
- Aug 11 2015 IMG Maintenance
- July 9 2015 ANI News Release
- July 8 2015 IMG Data Marts Changes
- June 15 2015 ProDeGe News Release
- June 11 2015 Plotting IMG's Next 10 Years
- May 2015 IMG accounts deprecated
- Apr 2015 BLAST in Workspace
- Mar 2015 IMG using GOLD's new metadata

[Read more...](#)

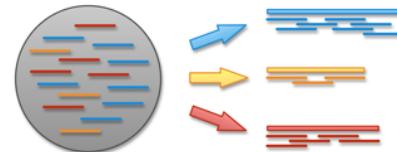
10 min break

Analysis Strategies

# Many Paths for Functional Annotations



## Assembly



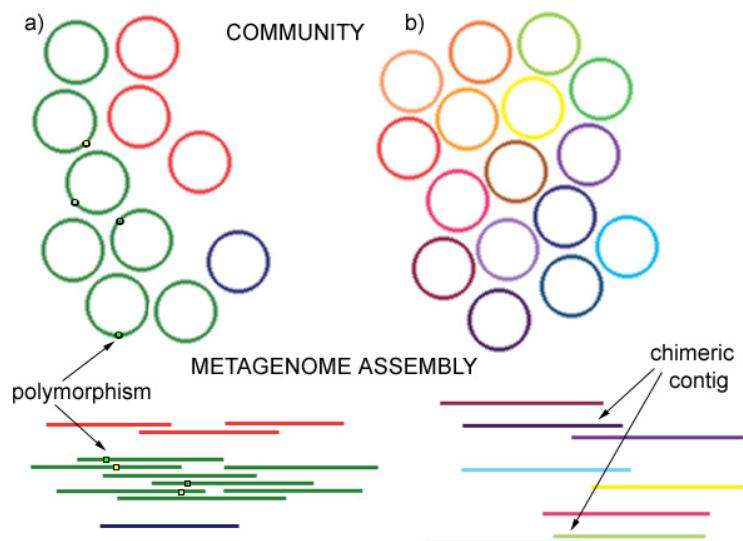
- Assembly can reduce the “amount of data” to optimize the annotation for function
- Assemblies in metagenomics can combine closely related strains or species
- Assemblies are high memory operations so there are some “pre-clustering” software to help reduce the data

# khmer: A Data Reduction Strategy

- khmer is a k-mer based dataset analysis and transformation toolkit
- It can be used to reduce the size of a dataset by:
  - abundance filtering and error trimming
  - graph-size filtering by removing disconnected reads
  - partitioning by splitting reads into disjoint sets.

## Assembly

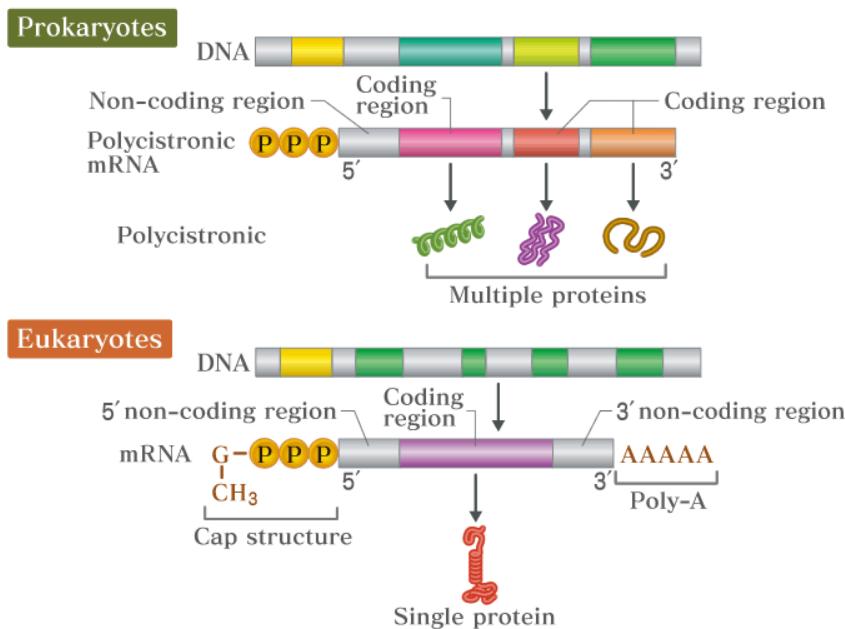
- Velvet/metaVelvet
- MetaAmos
- Mira
- Newbler (454 and hybrid assemblies)
- SOAPdenovo
- Meta-IDBA
- SPAdes



# ORF Detection

- Most aligners can perform translated alignment which can be more sensitive and “overcome” sequencing errors
- These alignments can be slower than protein alignments (6-frame translations)
- ORF detection can:
  - Reduce computations for functional profiling
  - Provide “de-novo” genes
  - Allow for a complete sets of genes for gene clustering and sample comparison

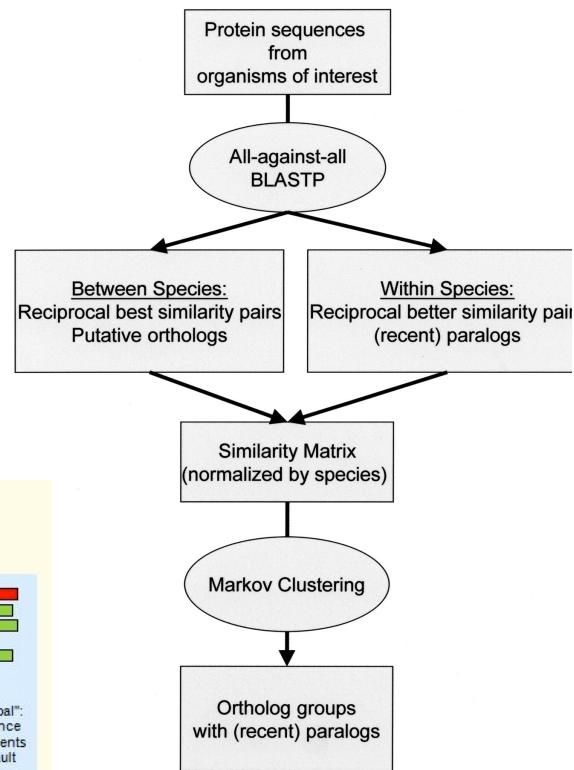
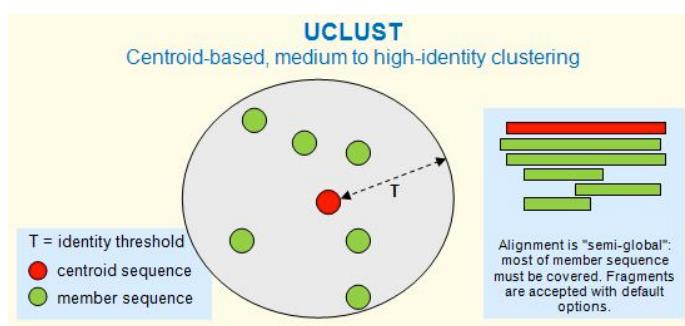
# ORF Detection



# Gene Finding Packages

- Most Metagenomic gene finders are modified prokaryotic gene finders
- MetaGeneMark
- FragGeneScan (on reads)
- Glimmer MG
- Orphelia

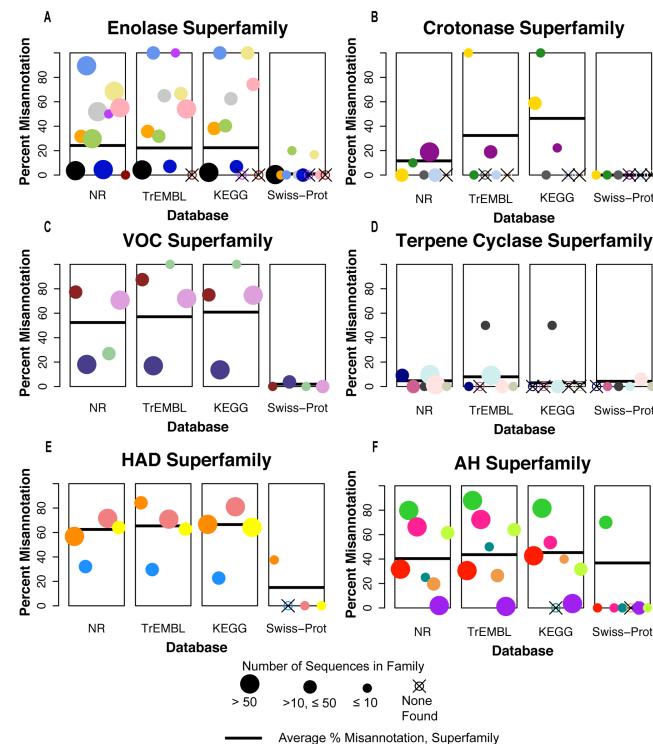
## Orthologous Clustering



# Functional Profiling

- High Throughput functional profiling comparison allows for gross comparisons of the functional capability of samples
  - Broad functional categories tend to be very similar in an ecological niche
- Profiling relies on alignments to functionally characterized proteins
- Homologous proteins tend to have similar broad “enzymatic function” i.e. kinase, hydrolase, transferase
  - However: Homology ≠ Same Biological Function

Functional  
Annotation  
Error are  
Common



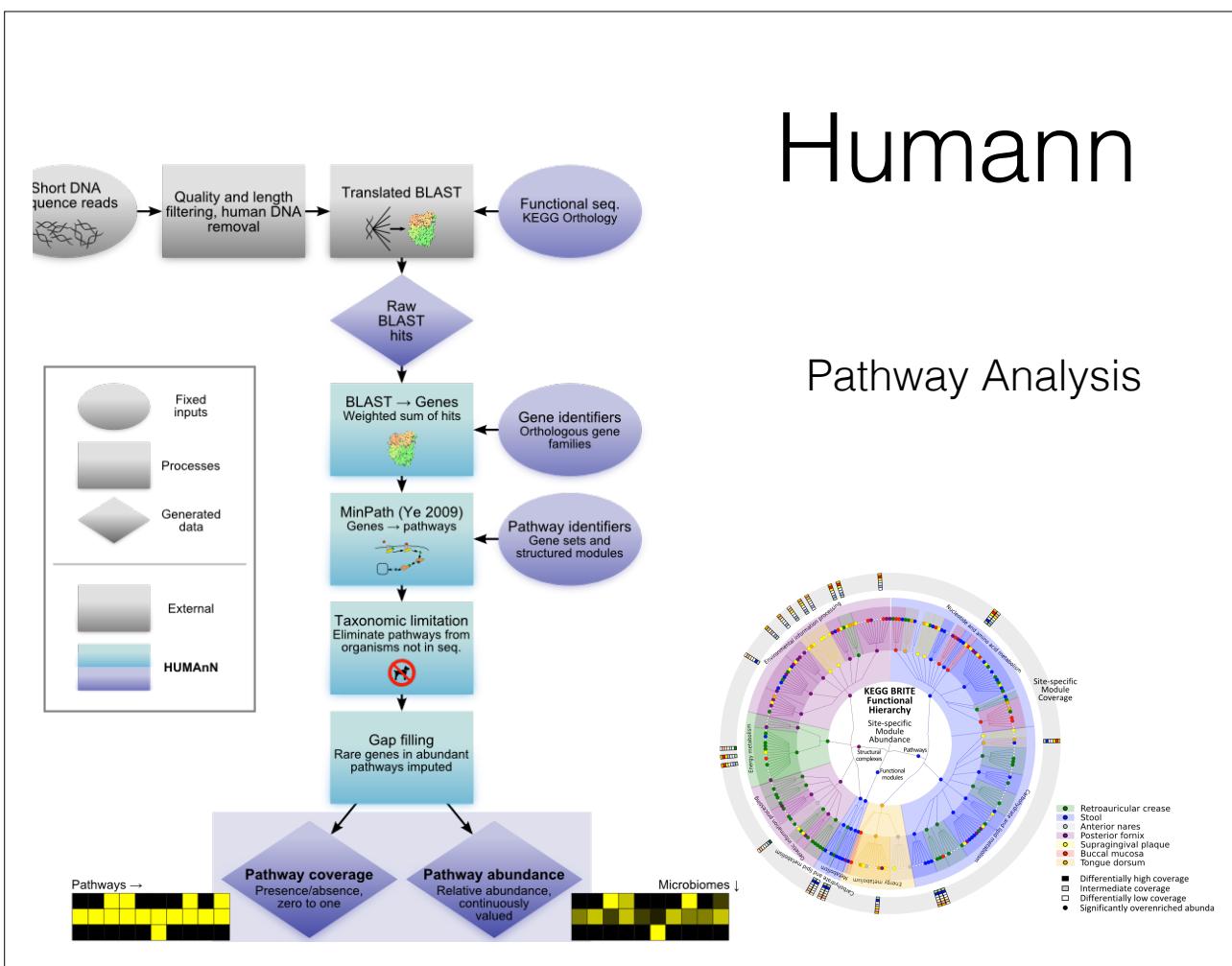
# Alignment Strategies

- BLASTP or BLASTX — very slow
- MALT — Requires > 100GB of memory
- USEARCH — Requires paid license for 64 bit version; memory requirement too high for 32 bit version
- VSEARCH —Free version of USEARCH, lacks sensitivity
- DIAMOND — Much more sensitive than VSEARCH, low memory requirement and fast

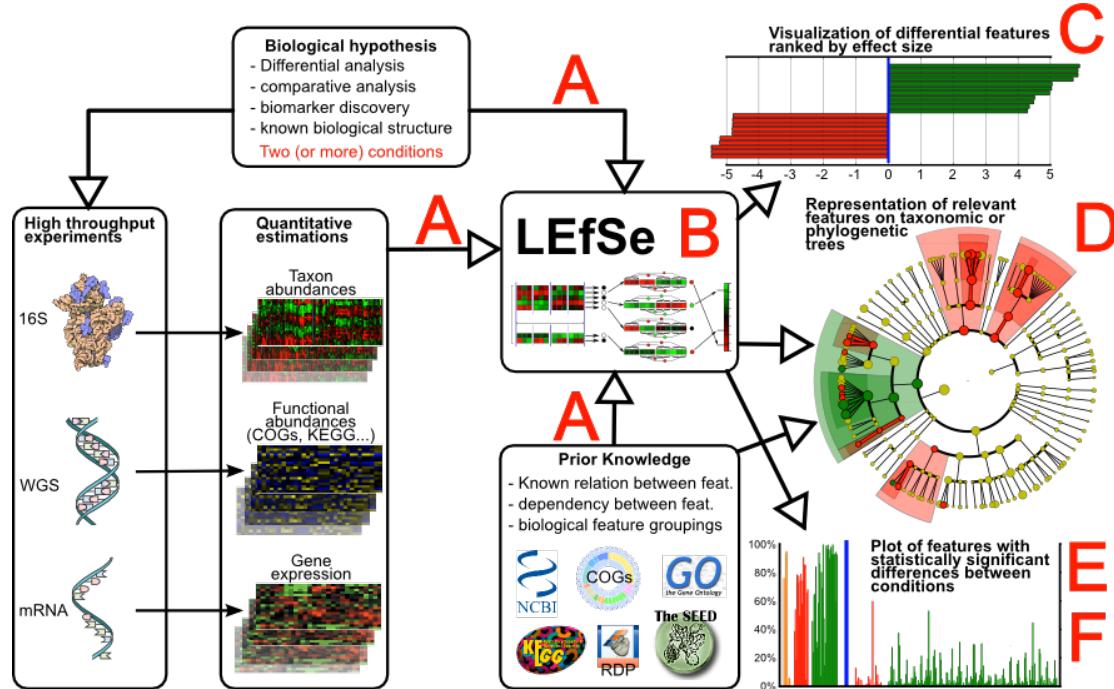
# Post-Alignment

- Using Alignments (Translated or Protein) — functional assignment is based on broad functional categories or pathways of annotated hits.
- Available Packages for functional assignment and pathway profiling:
  - Humann
  - Megan

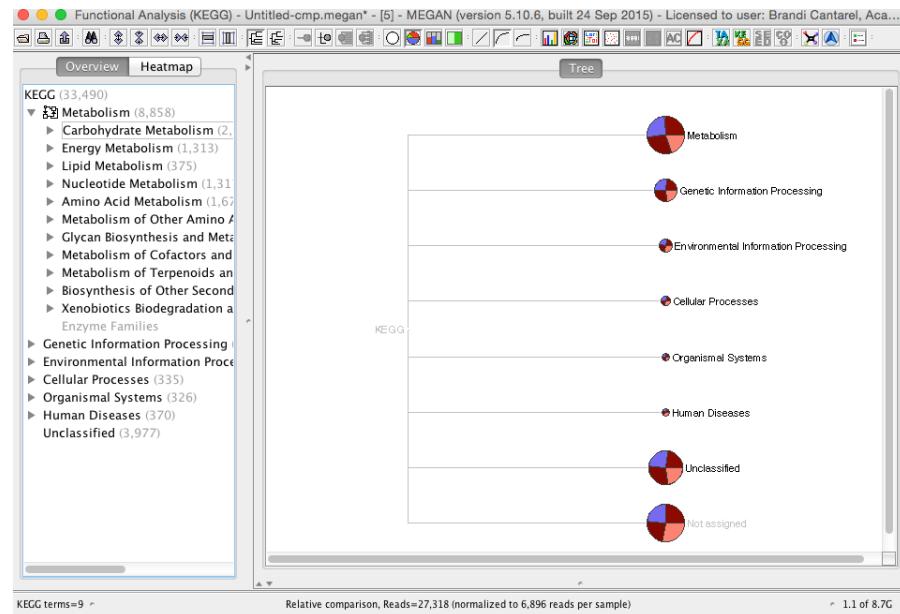
# Humann



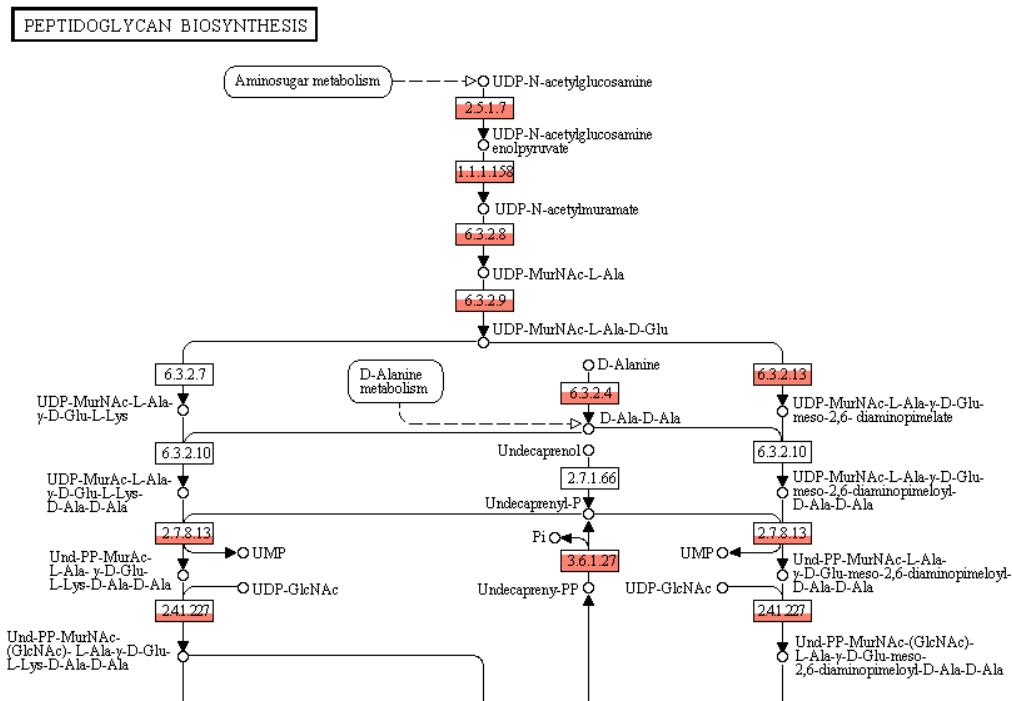
# LefSe



# MEGAN Broad Functional Comparisons



## Pathway Exploration



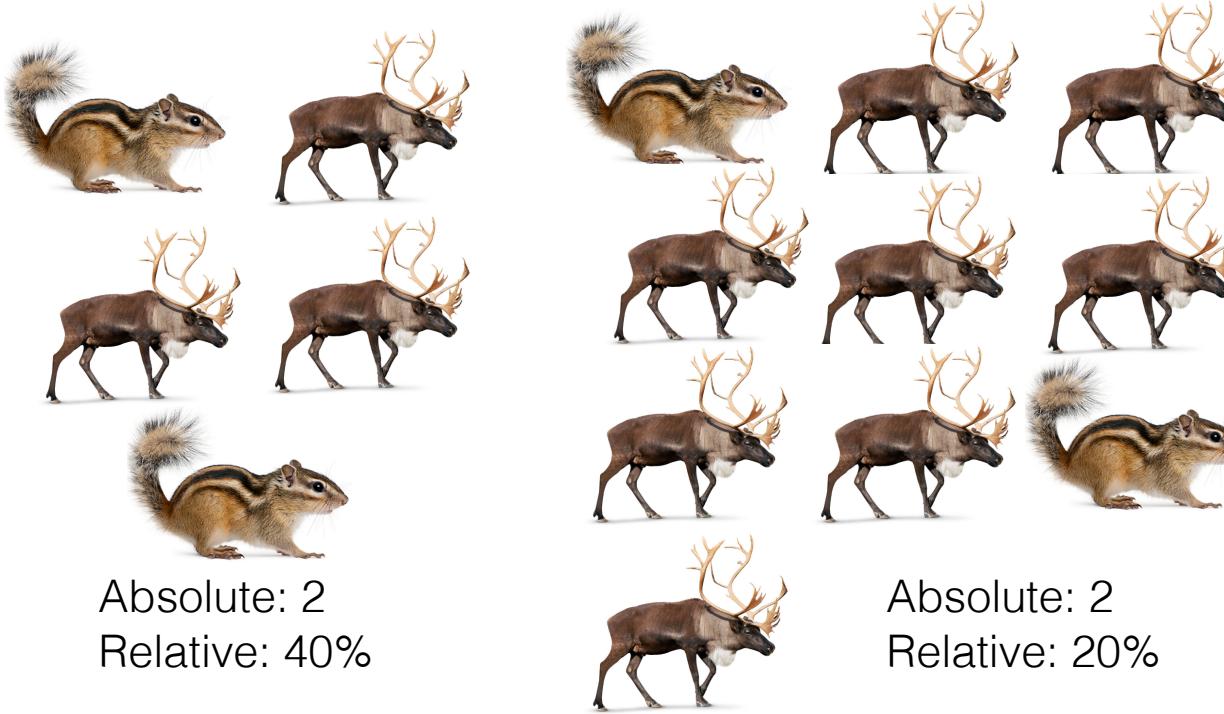
# Statistical Analysis

- MetagenomeSeq
  - Normalization of counts based on relative abundance
  - 2-sample t-statistic
- STAMP
  - ANOVA for multiple groups differential testing
  - Kruskal-Wallis H-test — nonparametric method test to determine differences in medians

## Relative Abundance vs Absolute Abundance

- Absolute abundance is a quantitate measure of the feature in the sample (qPCR)
- Relative abundance is the measure of a feature relative to all other features (sequencing)
- We can sequence every molecule in a sample, therefore abundances in microbiome studies are based on the number of measures (reads) and the proportion of the feature in the sample.

# Relative Abundance vs Absolute Abundance

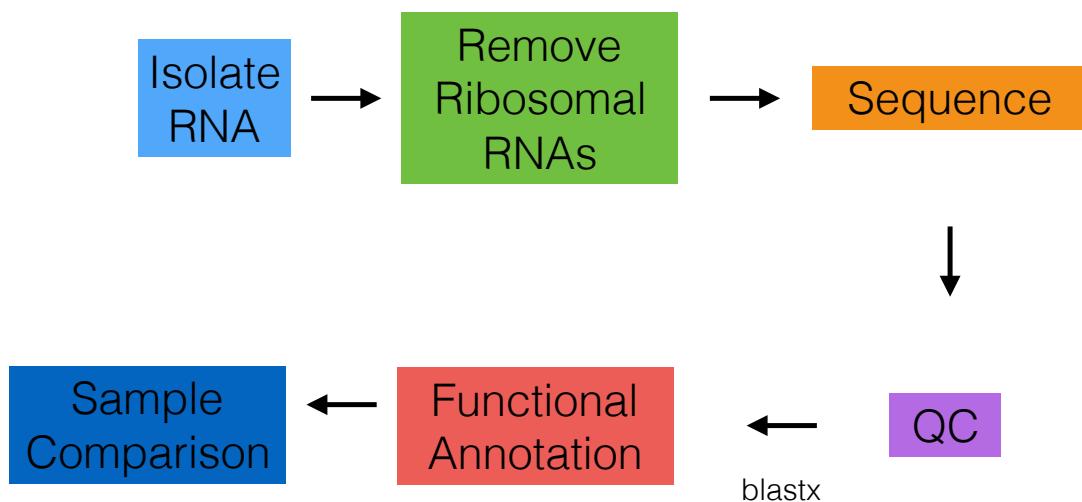


- Community Structure Based on WGS
  - 16S VS WGS
  - Least Common Ancestor
  - Reference Genome databases
  - Assignment Strategies
- Data Processing Workflow
  - QC
  - Functional Databases
  - Analysis Platforms
  - DIY Analysis
- Non-DNA based technologies: Metatranscriptomics, Metaproteomics and Metabolomics
  - Expression (RNASeq)
  - Translation (Proteomics)
  - Metabolites (Metabolomics)

# Metagenomics vs Metatranscriptomics

- Metagenomics can give insight into gene content.
- Metatranscriptomics can measure how expression (functional potential) changes in response to the environment
- Metatranscriptomics can also show which organisms are the most functionally active.

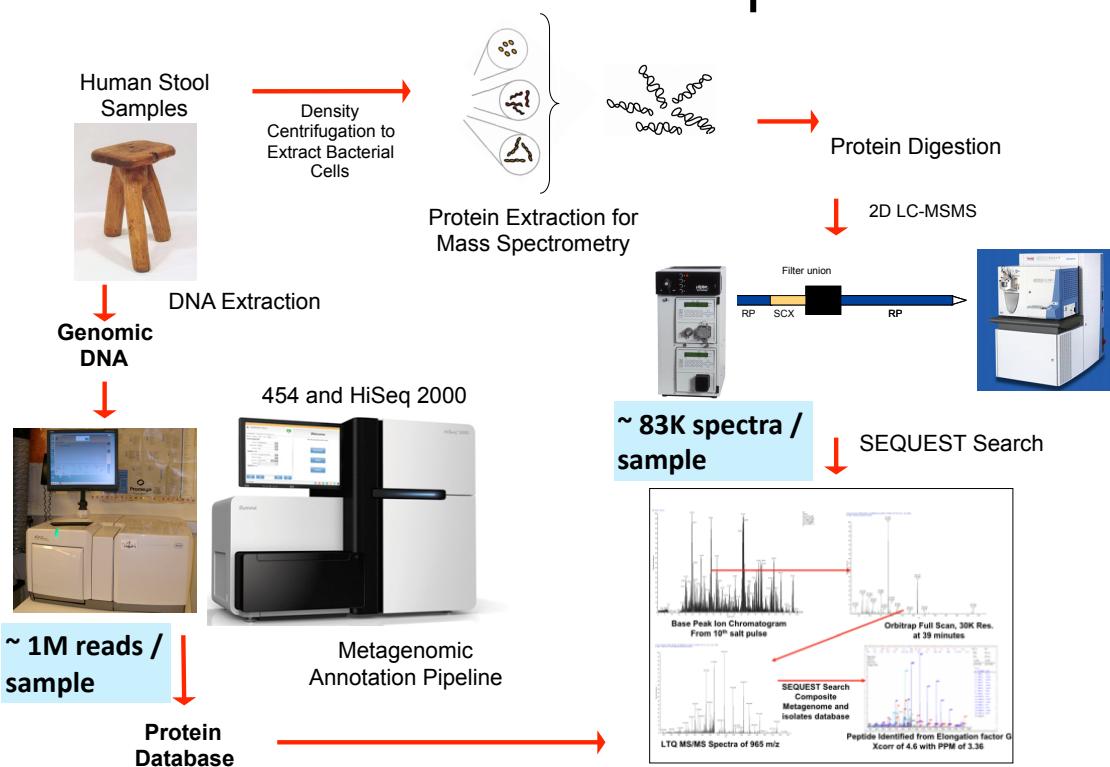
## Metatranscriptomics



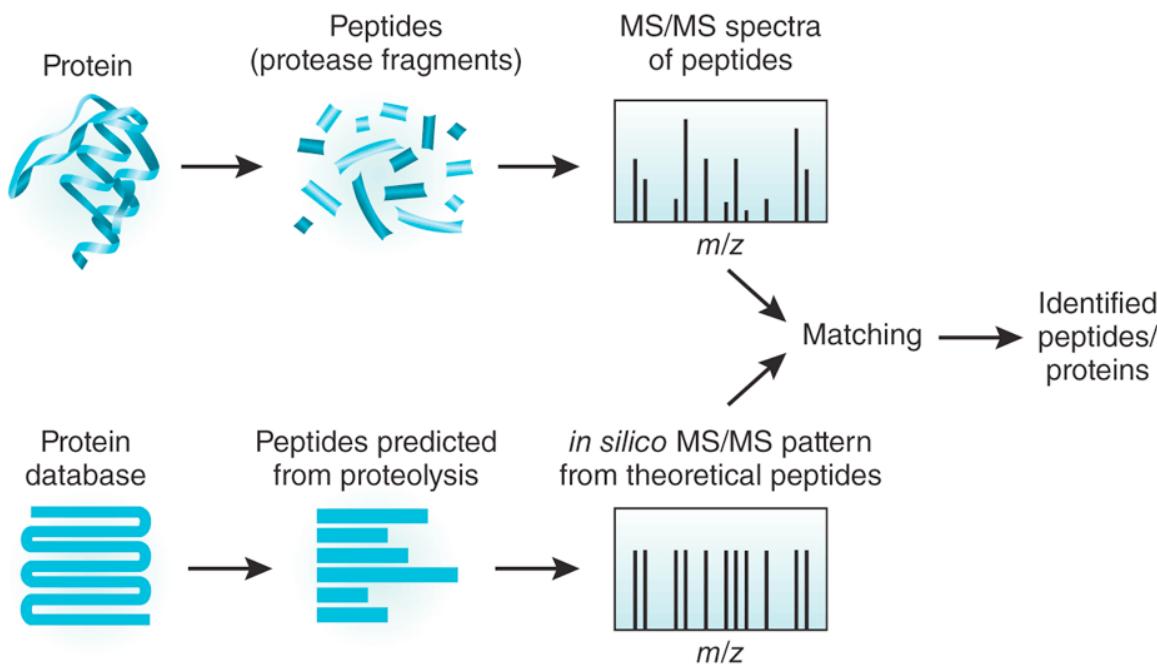
# Metaproteomics

- Like metagenomics and metatranscriptomics, metaproteomics is complicated by the lack of a complete reference set
- In order to determine the protein sequence of peptide fragments, a metagenomic or reference genome database is necessary.
- Unlike sequencing, denovo protein prediction from MS/MS is not trivial.
- Contains a mixture of environmental and microbiome proteins

## “Omics” Pipeline



# Peptide Spectral Matching

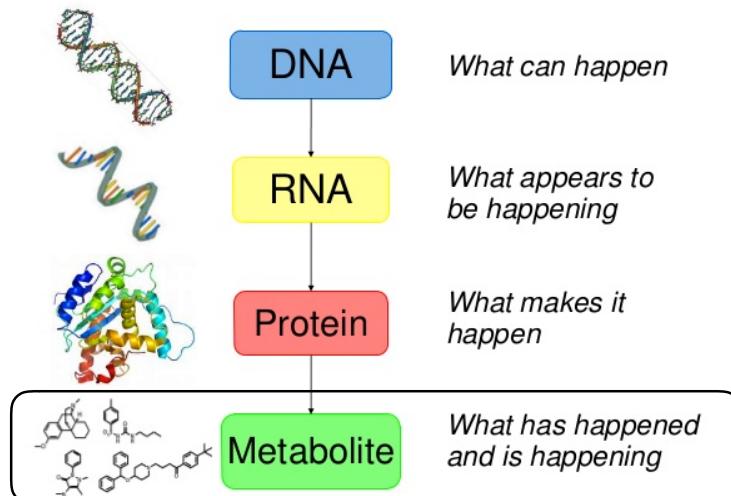


Duncan MW, Aebersold R, Caprioli RM. The pros and cons of peptide-centric proteomics. Nat Biotechnol. 2010 Jul;28(7):659-64.

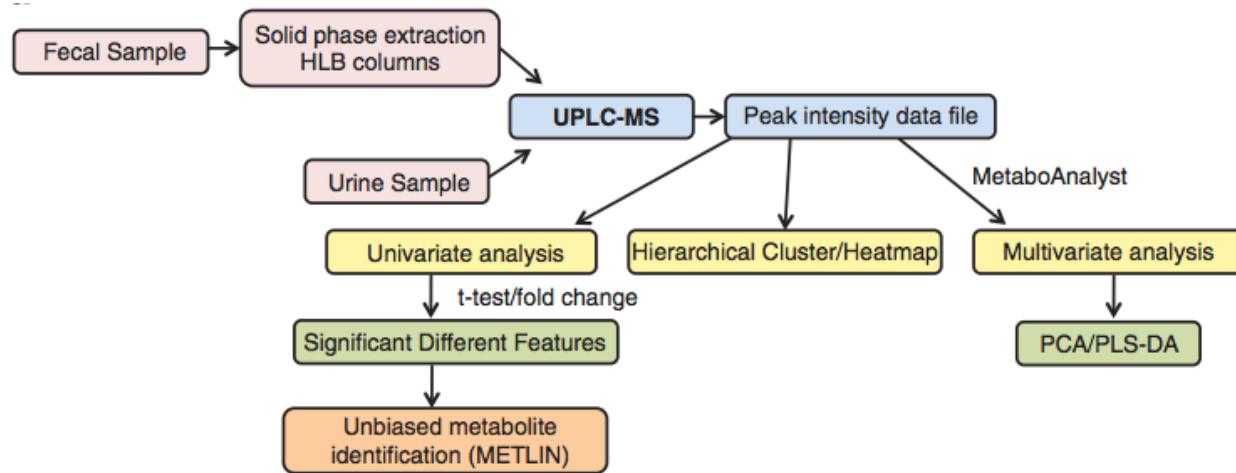
## meta-Metabolomics

Animal and environmental metabolomic studies are (meta)metabolomics — it is difficult to know “who” produced a particular metabolite.

### The central dogma of biology



# meta-Metabolomics



Marcabal A, Kashyap PC, Nelson TA, Aronov PA, Donia MS, Spormann A, Fischbach MA, Sonnenburg JL. A metabolomic view of how the human gut microbiota impacts the host metabolome using humanized and gnotobiotic mice. ISME J. 2013 Oct;7(10):1933-43. doi: 10.1038/ismej.2013.89. Epub 2013 Jun 6. PubMed PMID: 23739052; PubMed Central PMCID: PMC3965317.

## Workshop 2