

# Metagenomics: Characterization of Microbial Communities using NGS

Part 2

Brandi Cantarel, PhD  
Bioinformatic Scientist  
Baylor Scott and White

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

# Taxonomic Assessment using 16S

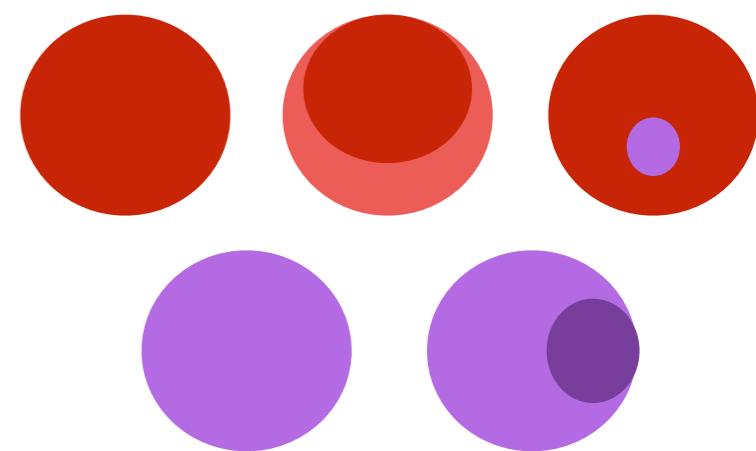
- 16S is targeted sequencing for a single gene which acts as a marker for organisms
- Pros
  - Well established
  - Relatively inexpensive ~ \$100/sample
  - Amplifies only bacteria not host or environmental fungi, plants, etc
- Cons
  - Amplifies only bacteria not viruses, microbial fungi, archaea, 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

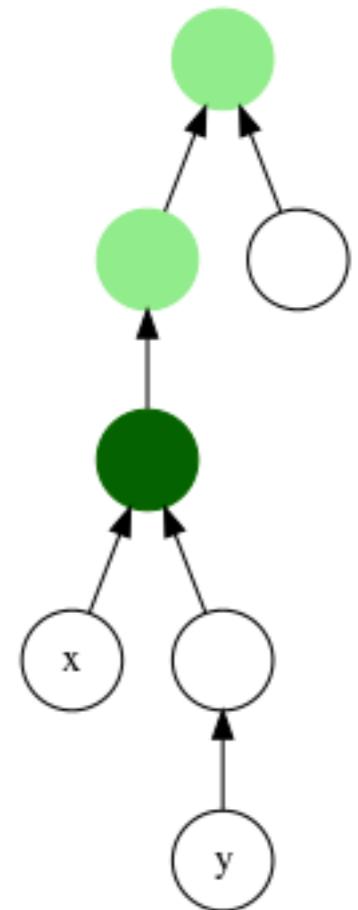
# Taxonomic Assignment: Complex Analysis

- 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

The screenshot shows the homepage of the JGI GOLD database. At the top, there's a navigation bar with links for JGI HOME, LOGIN, Home, Search, Distribution Graphs, Biogeographical Metadata, Statistics, References, Team, Help, and News. Below this is a sidebar with counts for Studies (22,229), Biosamples (69,168), Sequencing Projects (69,463), and Analysis Projects (57,213). A button to 'Download Excel Data file' is also present. The main content area features a welcome message for 'GOLD Release v.5' and three main steps: 1. Register (with a 'Register' button), 2. Annotate (with a 'Annotate' button), and 3. Publish (with a 'Publish' button). The 'Publish' step is associated with the 'SIGS Standards in Genomic Sciences' logo.

The screenshot shows the NCBI genome search interface. The top navigation bar includes links for Resources, How To, and a search bar set to 'Genome'. The main content area is titled 'Reference and representative genomes'.

The screenshot shows the HMP Reference Genomes page. The top navigation bar includes links for Overview, Reference Genomes, Microbiome Analysis, Health & Ethics, Resources, Outreach, and Data Browser. A 'Login' button and a search bar are also at the top. The main content area is titled 'Microbial Reference Genomes' and contains a paragraph about the HMP's plan to sequence 3000 reference genomes. It includes 'GET DATA' and 'GET TOOLS' buttons. A sidebar on the left lists 'Current News' items: 'January 2015 Metagenome Analysis Workshop March 3-6'.

The screenshot shows the EnsemblBacteria website. The top navigation bar includes links for Sequence Search, BLAST, Tools, and Download. The main content area is titled 'Access to over 20,000 Bacterial Genomes' and contains two search boxes: one for genes ('Search for a gene') and one for genomes ('Search for a genome'). Below these are instructions: 'e.g. ftsZ or uridine\*' and 'e.g. type esc to find Escherichia'. A decorative graphic of bacterial DNA is in the bottom right corner.

# 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

- MetaPhlAn — web based
- PhyloSift
- 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.

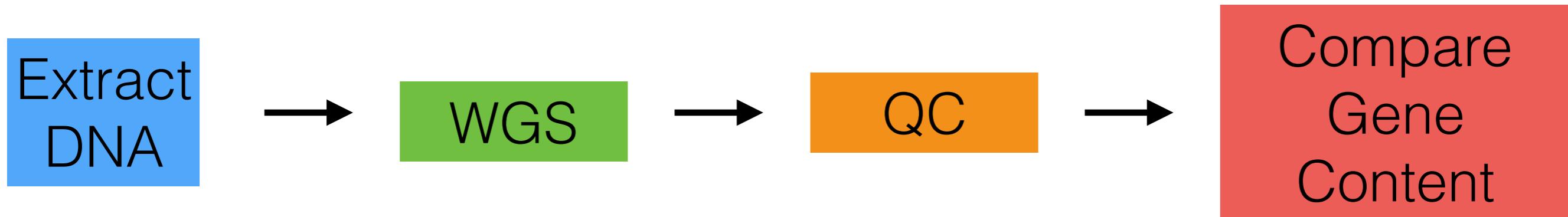
# WGS Taxonomy Assignment and Visualization

- 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
- STAMP
  - <http://kiwi.cs.dal.ca/Software/STAMP>
  - Tools for statistical analysis and visualization
- LefSe
  - <http://huttenhower.sph.harvard.edu/galaxy/>
  - A method for metagenomic biomarker discovery by way of class comparison, tests of biological consistency and effect size estimation.
- Plotting Tools including R, Excel, Matlab

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

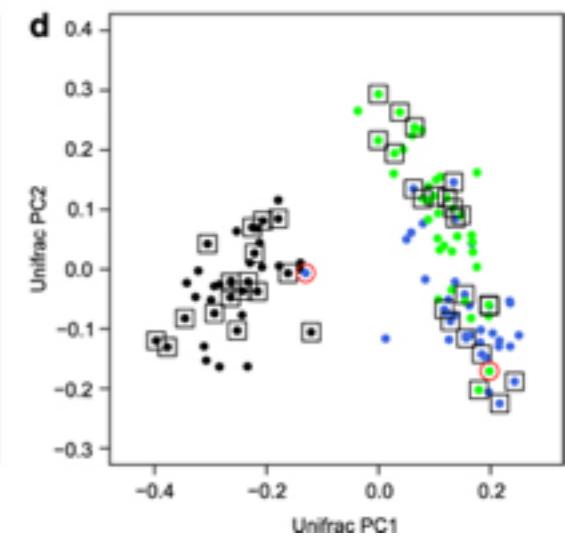
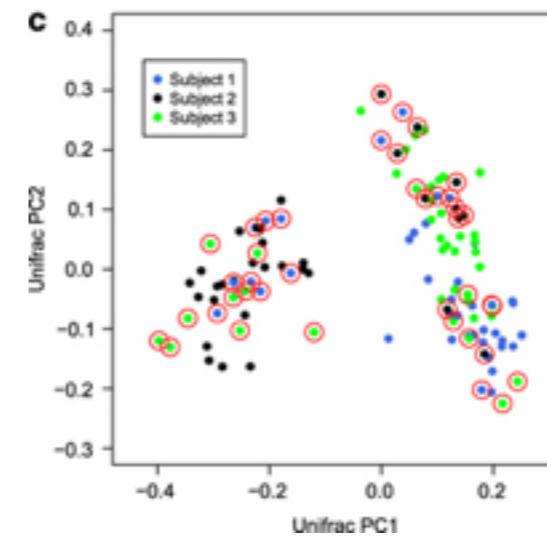
# What is the Functional Capability

- 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



EBI Metagenomics



# Comprehensive Functional Databases

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

eggNOG  
version 3.0

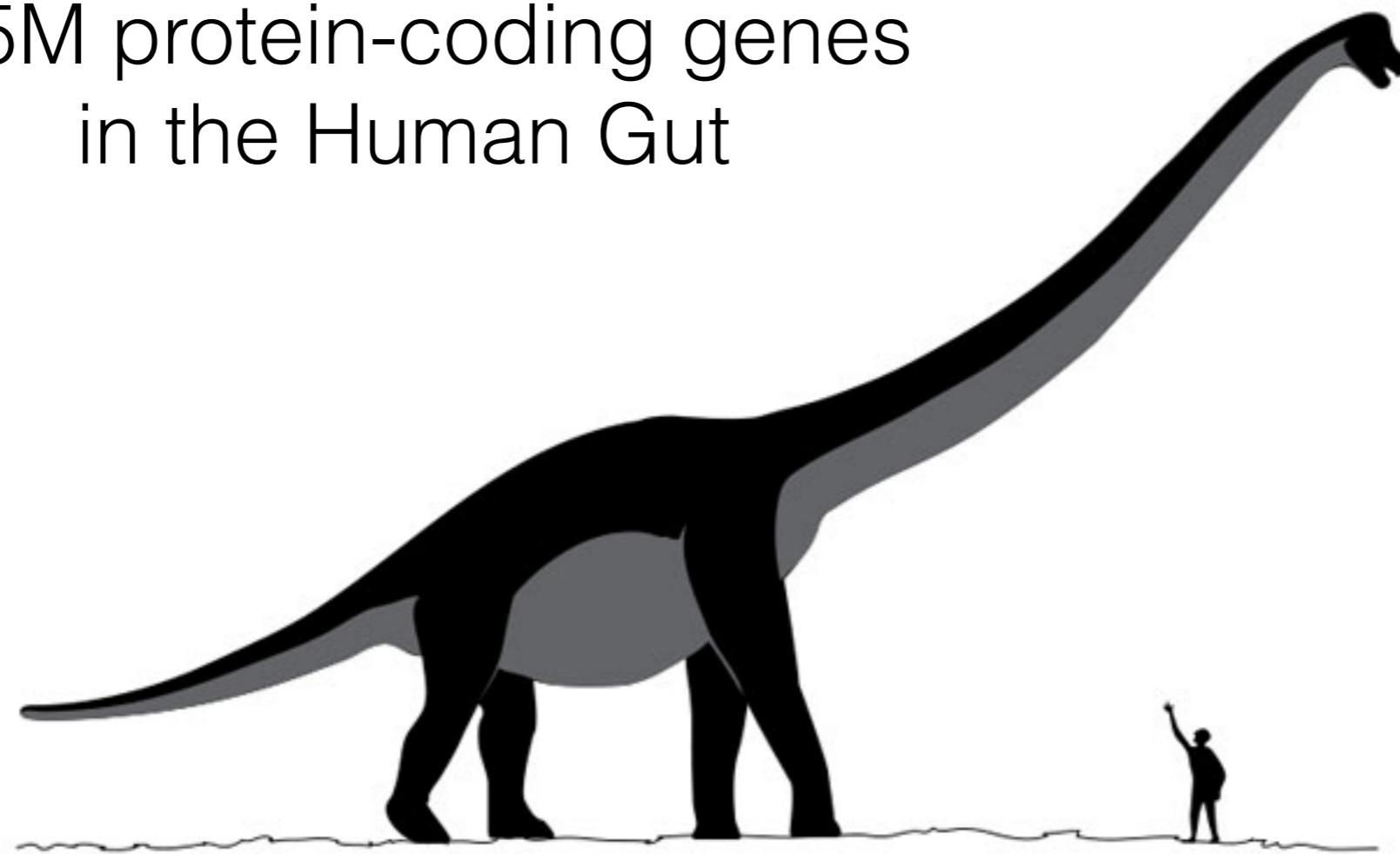


# Specialized Functional Databases

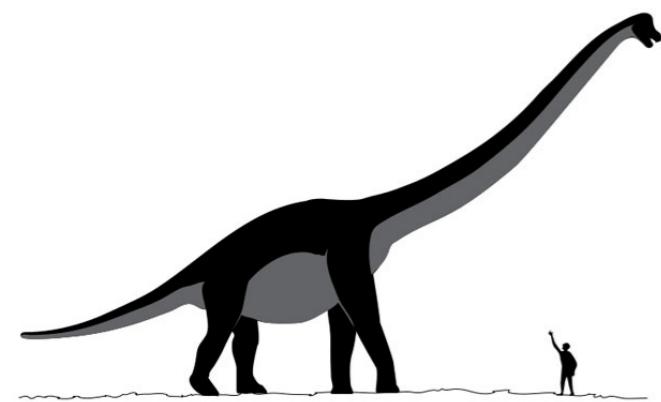
- Antibiotic resistant genes
  - <http://ardb.cbcb.umd.edu/>
- Virulent factors
  - <http://www.mgc.ac.cn/VFs/main.htm>
- Carbohydrate Active Enzymes
  - [www.cazy.org](http://www.cazy.org)
- Phage
- Proteases
  - <http://merops.sanger.ac.uk/>
- Transporters
  - <http://www.membranetransport.org/>

# 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

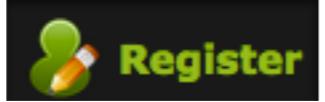
metagenomics analysis server

**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



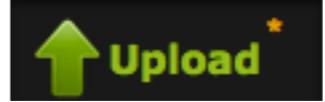
Register



Contact



Help



Upload\*



News

About

search for metagenomes



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

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.

- 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

About Register Help News

The 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. HHSN272200900404C.

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

The job pig59\_colon (234800) was submitted as part of the project [PigColon](#) at 10/7/2015, 9:16:17 AM.

The current status is completed, the computation is finished. It took 1 hours 19 minutes from job submission until completion.

The result data is available for download on the [download page](#). You can take a look at the overview analysis data on the [metagenome overview page](#).

✓ 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



Public

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



Private

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

[Back to query page](#)

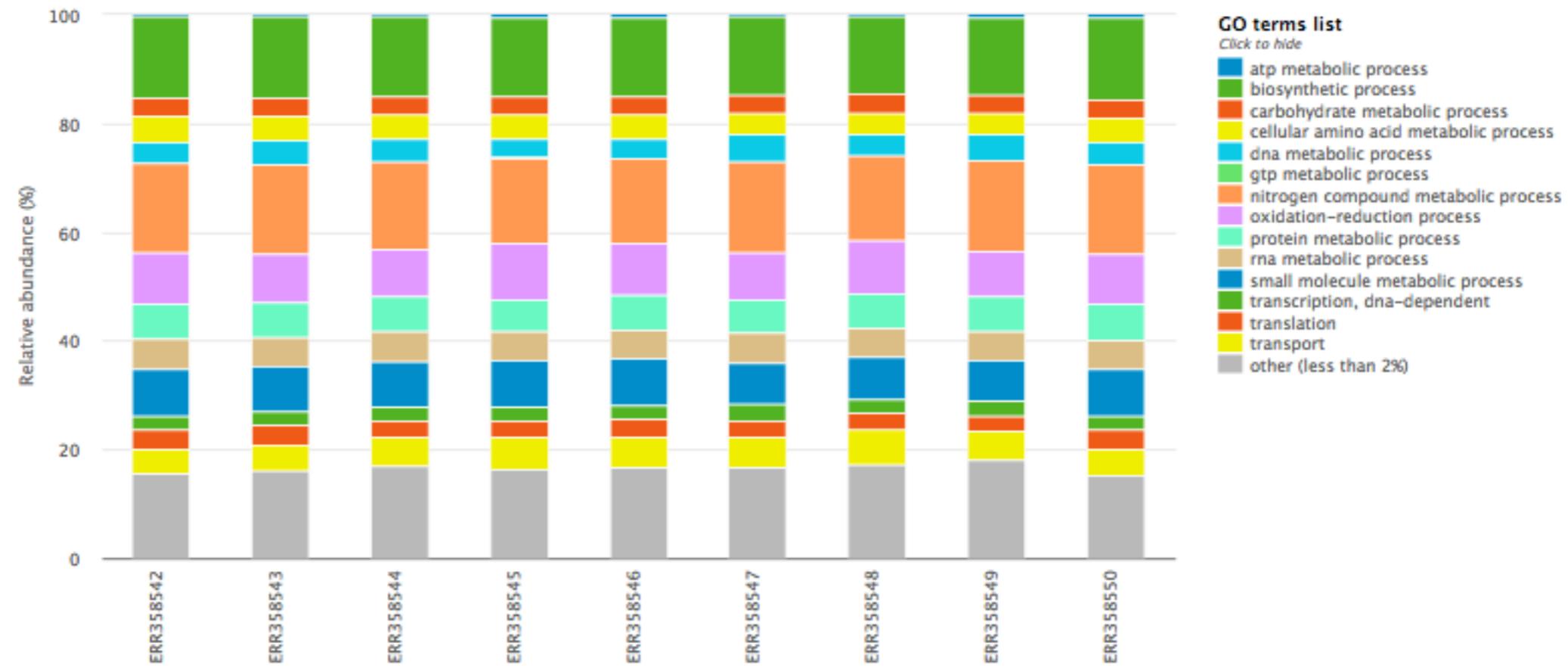
## Sample comparison tool: GO term annotation (functional analysis)

- Runs: [ERR358542](#), [ERR358543](#), [ERR358544](#), [ERR358545](#), [ERR358546](#), [ERR358547](#), [ERR358548](#), [ERR358549](#), [ERR358550](#)
- Project: Comparative freshwater metagenomics of Swedish and American lakes ([ERP004168](#))

[Barcharts](#)[Stacked columns](#)[Heatmap](#)[Principal Component Analysis](#)[Table](#)Jump to: [Biological process](#) | [Molecular function](#) | [Cellular component](#)[Export](#) ▾

### Biological process

Most frequent GO terms (biological process)





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.

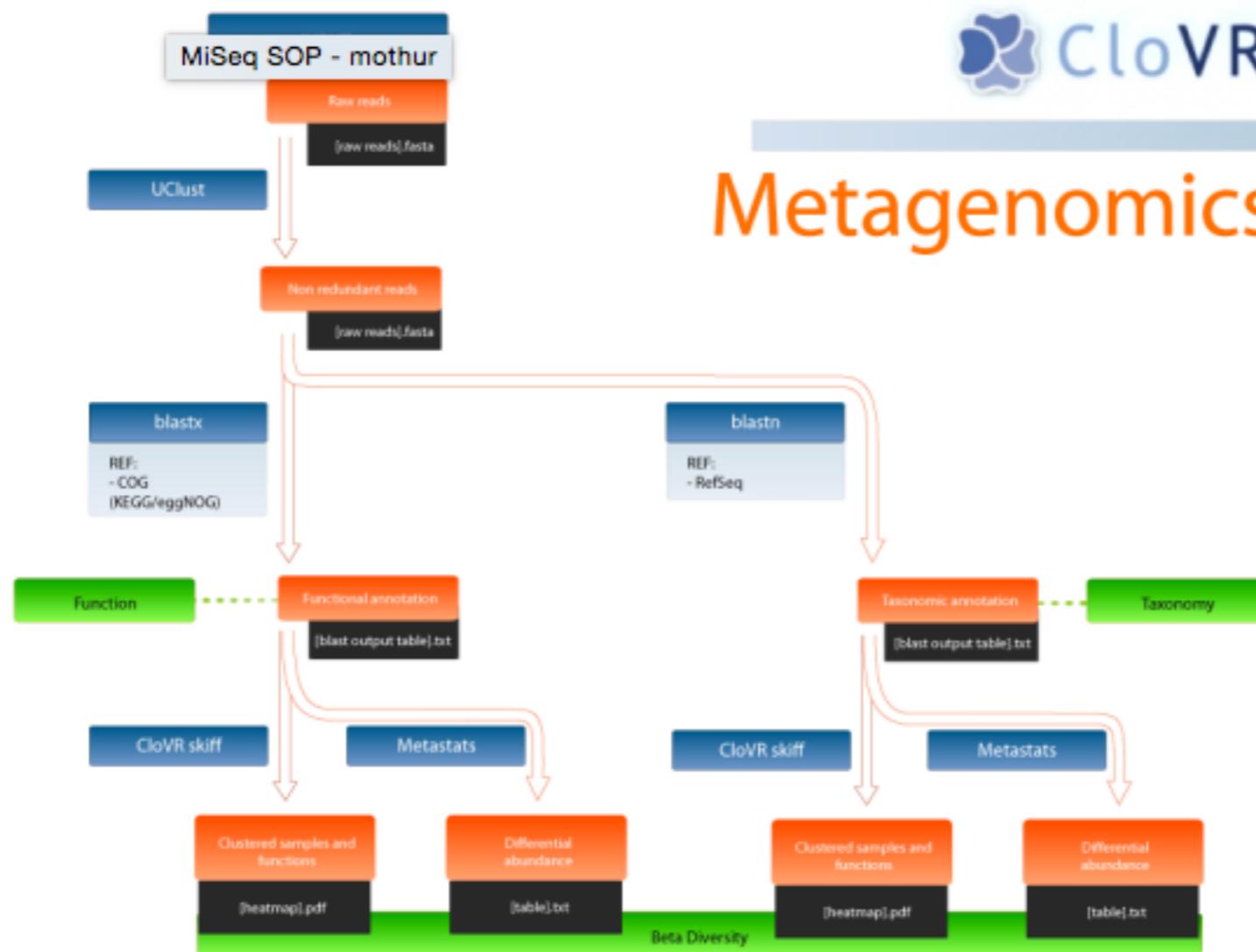


### About CloVR

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.

## IMG Content

## Datasets

Bacteria	<a href="#">25871</a>
Archaea	<a href="#">532</a>
Eukarya	<a href="#">190</a>
Plasmids	<a href="#">1186</a>
Viruses	<a href="#">3888</a>
Genome Fragments	<a href="#">1192</a>
Total Datasets	<a href="#">32859</a>

[Genome by Metadata](#)[Project Map](#)[Metagenome Projects Map](#)[System Requirements](#)

Hands on  
training  
available at the

[Microbial Genomics &  
Metagenomics Workshop](#)

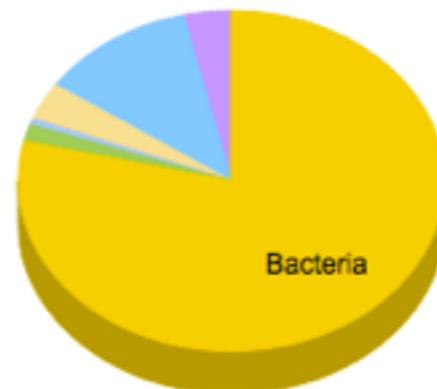
The Integrated Microbial Genomes (IMG) system ([Nucleic Acids Research, Volume 42 Issue D1](#)) serves as a community resource for comparative analysis and annotation of all publicly available genomes from three domains of life in a uniquely integrated context. Plasmids that are not part of a specific 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



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

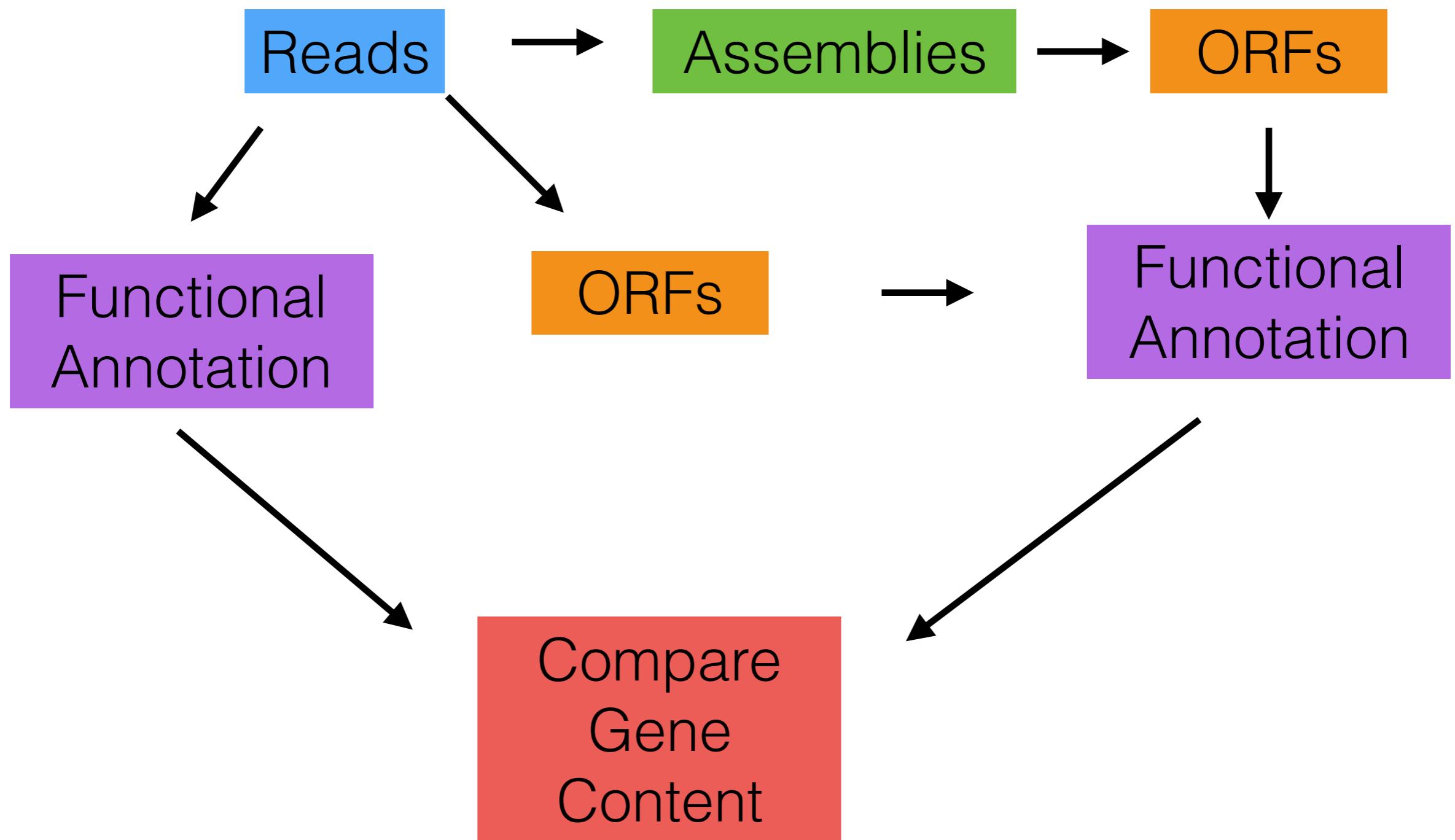
## News

- [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...](#)

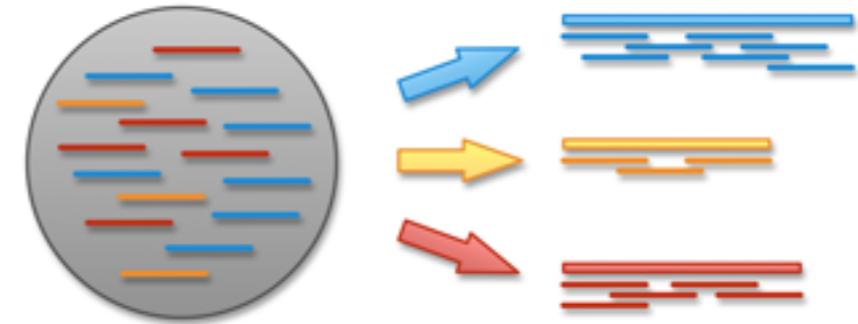
# Coffee Break

# Analysis Strategies

# Many Paths for Functional Annotations



# Assembly



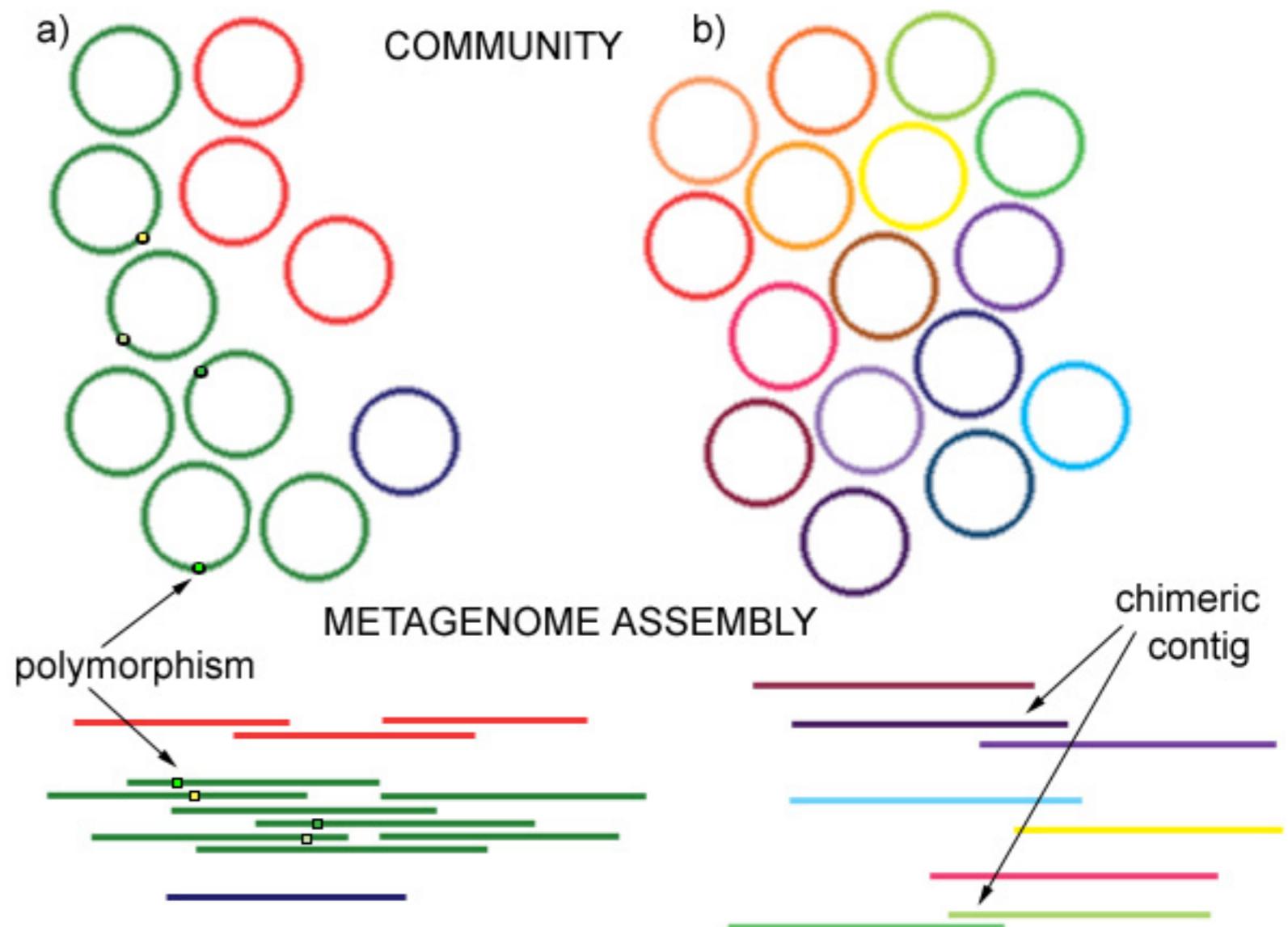
- Assembly can reduce the “amount of data” to optimize the annotation for function
- Assemblies in metagenomics can combine closely related strains or species
- Assemblies are high memory operations so there are some “pre-clustering” software to help reduce the data

# 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

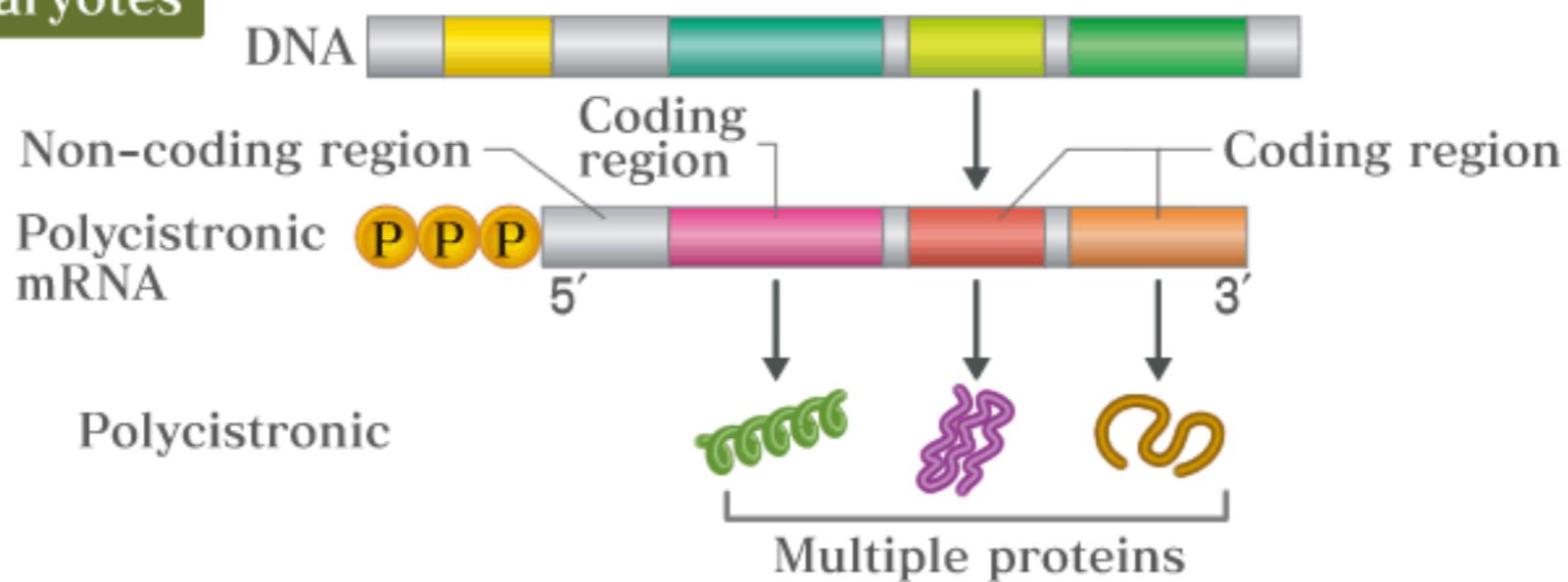


# ORF Detection

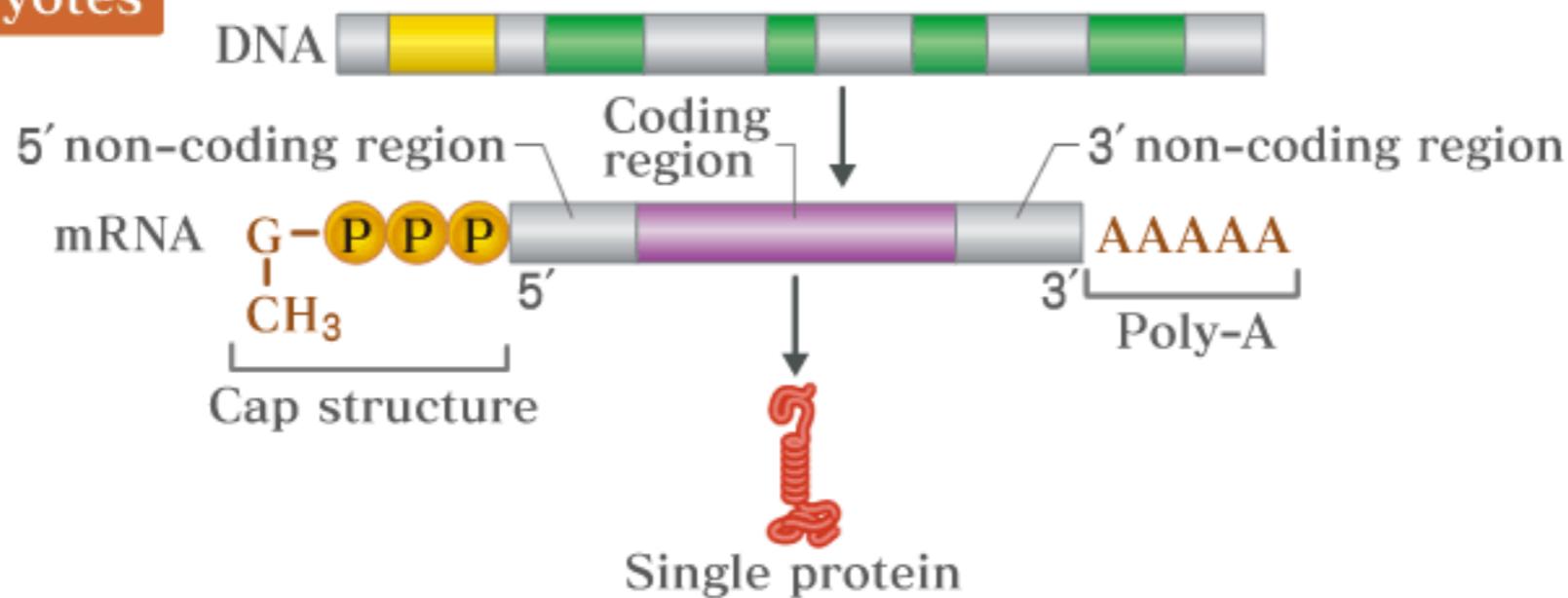
- Most aligners can perform translated alignment which can be more sensitive and “overcome” sequencing errors
- These alignments can be slower than protein alignments (6-frame translations)
- ORF detection can:
  - Reduce computations for functional profiling
  - Provide “de-novo” genes
  - Allow for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

# ORF Detection

## Prokaryotes



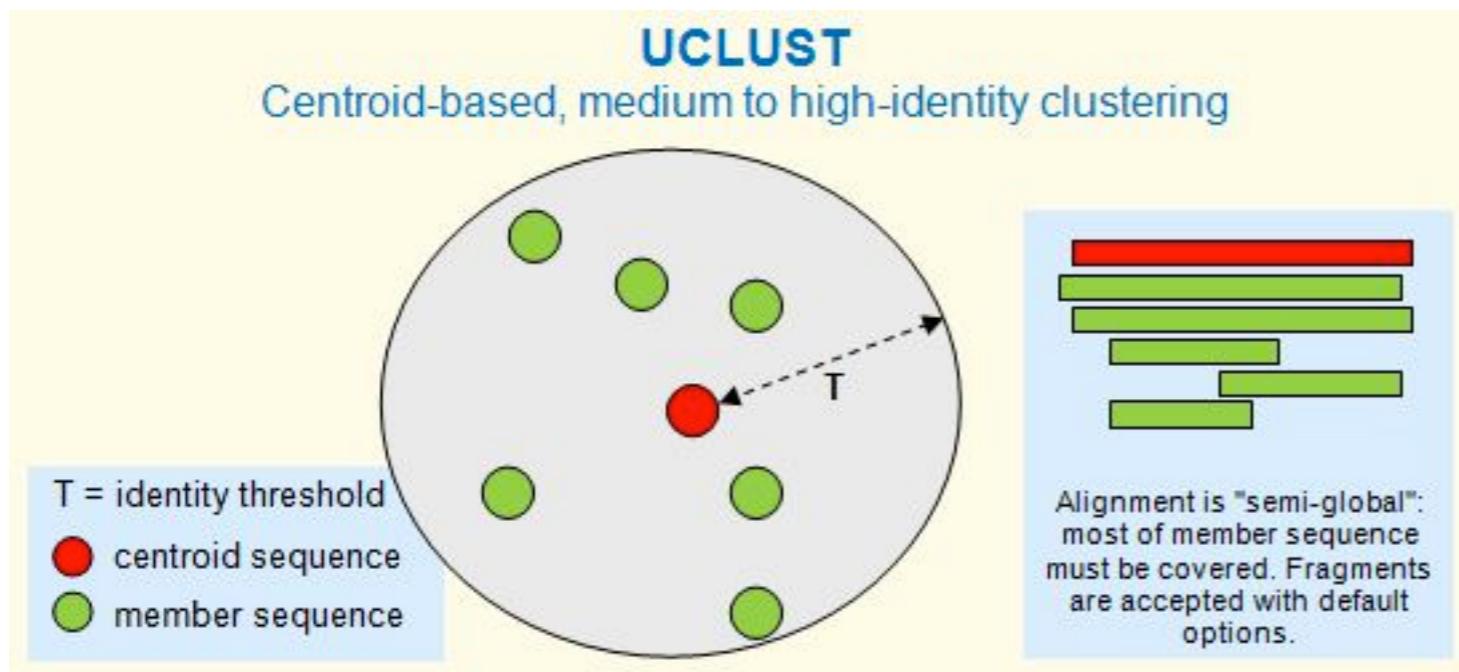
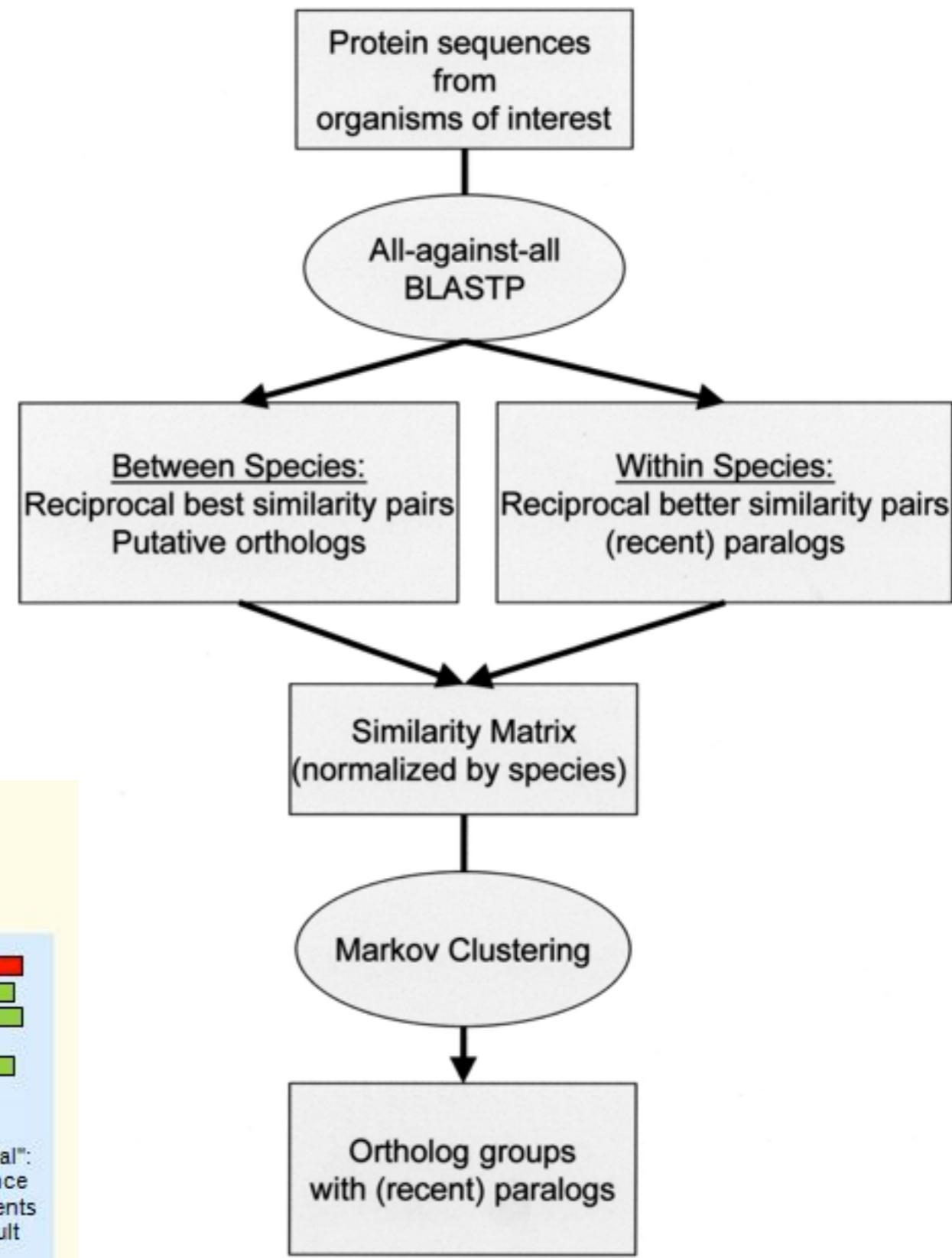
## Eukaryotes



# Gene Finding Packages

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

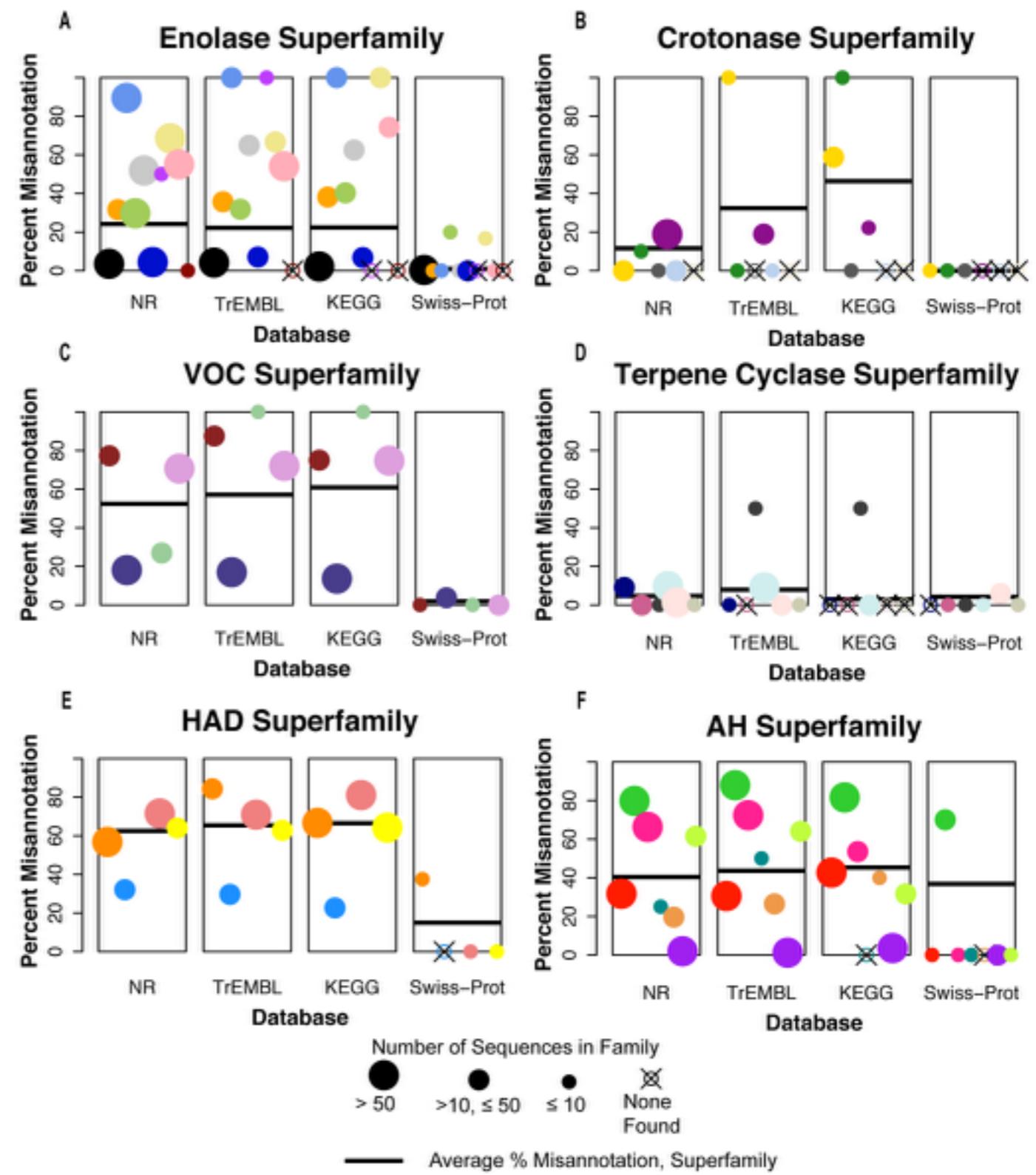
# Orthologous Clustering



# 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



Schnoes AM, Brown SD, Dodevski I, Babbitt PC. Annotation error in public databases: misannotation of molecular function in enzyme superfamilies. PLoS Comput Biol. 2009 Dec;5(12):e1000605. doi: 10.1371/journal.pcbi.1000605. Epub 2009 Dec 11. PubMed PMID: 20011109; PubMed Central PMCID: PMC2781113.

# Alignment Strategies

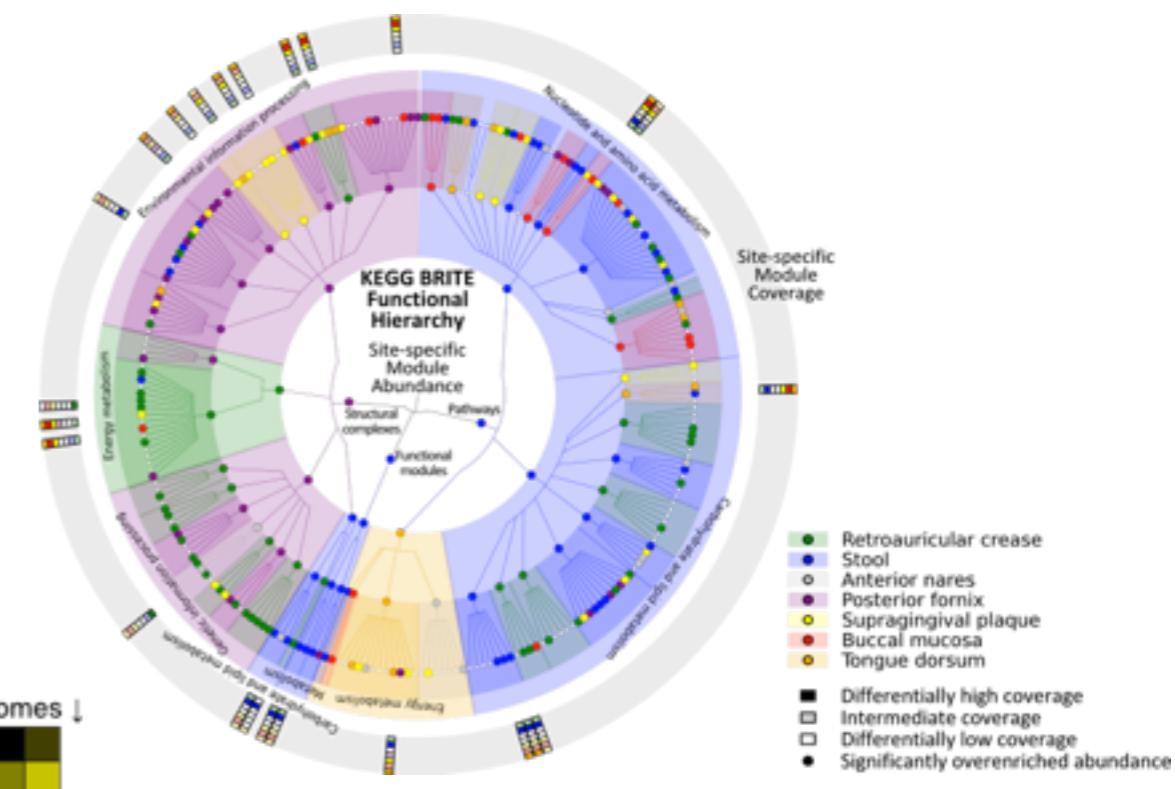
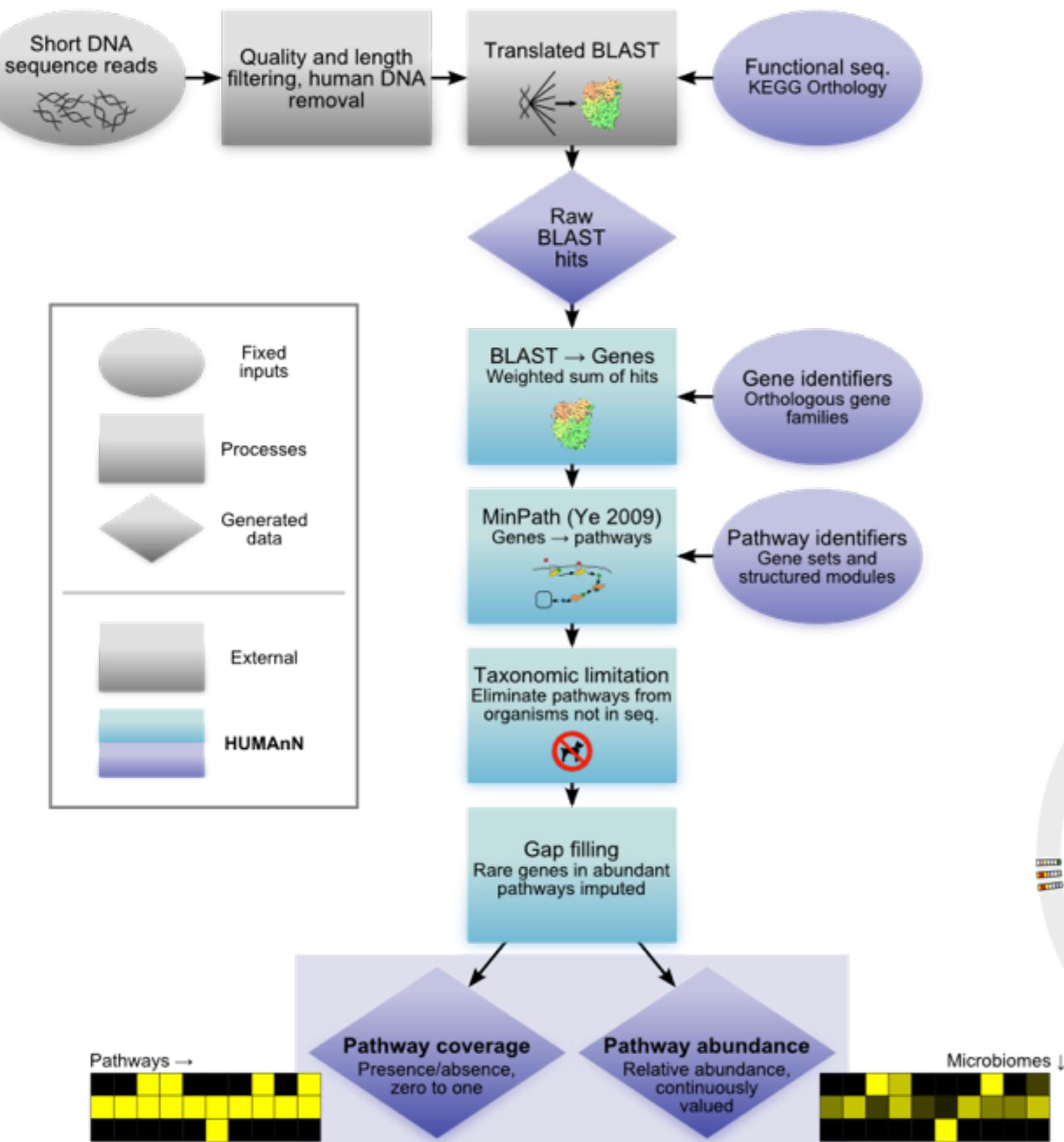
- BLASTP or BLASTX — very slow
- MALT — Requires > 100GB of memory
- USEARCH — Requires paid license for 64 bit version; memory requirement too high for 32 bit version
- VSEARCH —Free version of USEARCH, lacks sensitivity
- DIAMOND — Much more sensitive than VSEARCH, low memory requirement and fast

# Post-Alignment

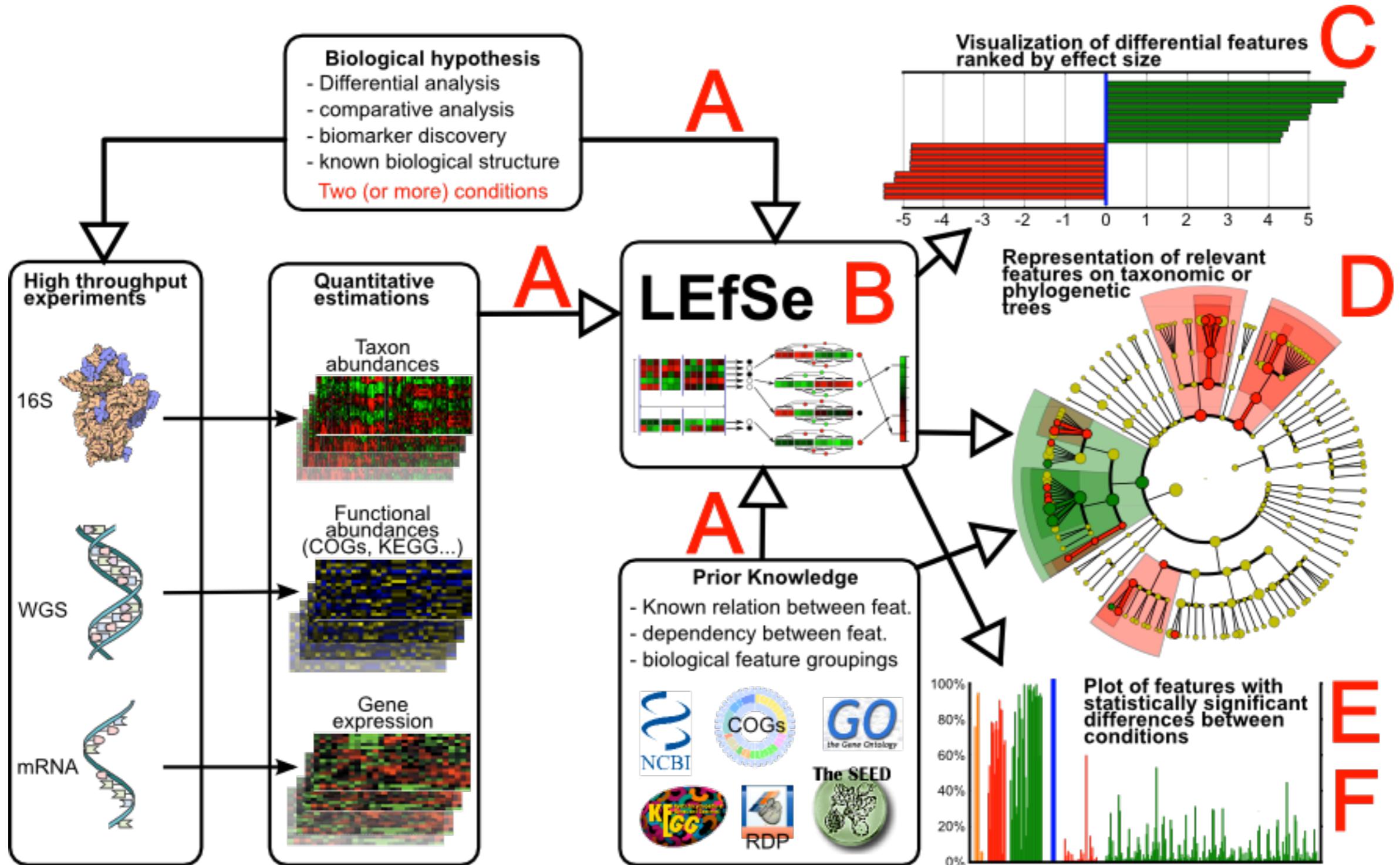
- Using Alignments (Translated or Protein) — functional assignment is based on broad functional categories or pathways of annotated hits.
- Available Packages for functional assignment and pathway profiling:
  - Humann
  - Megan

# Humann

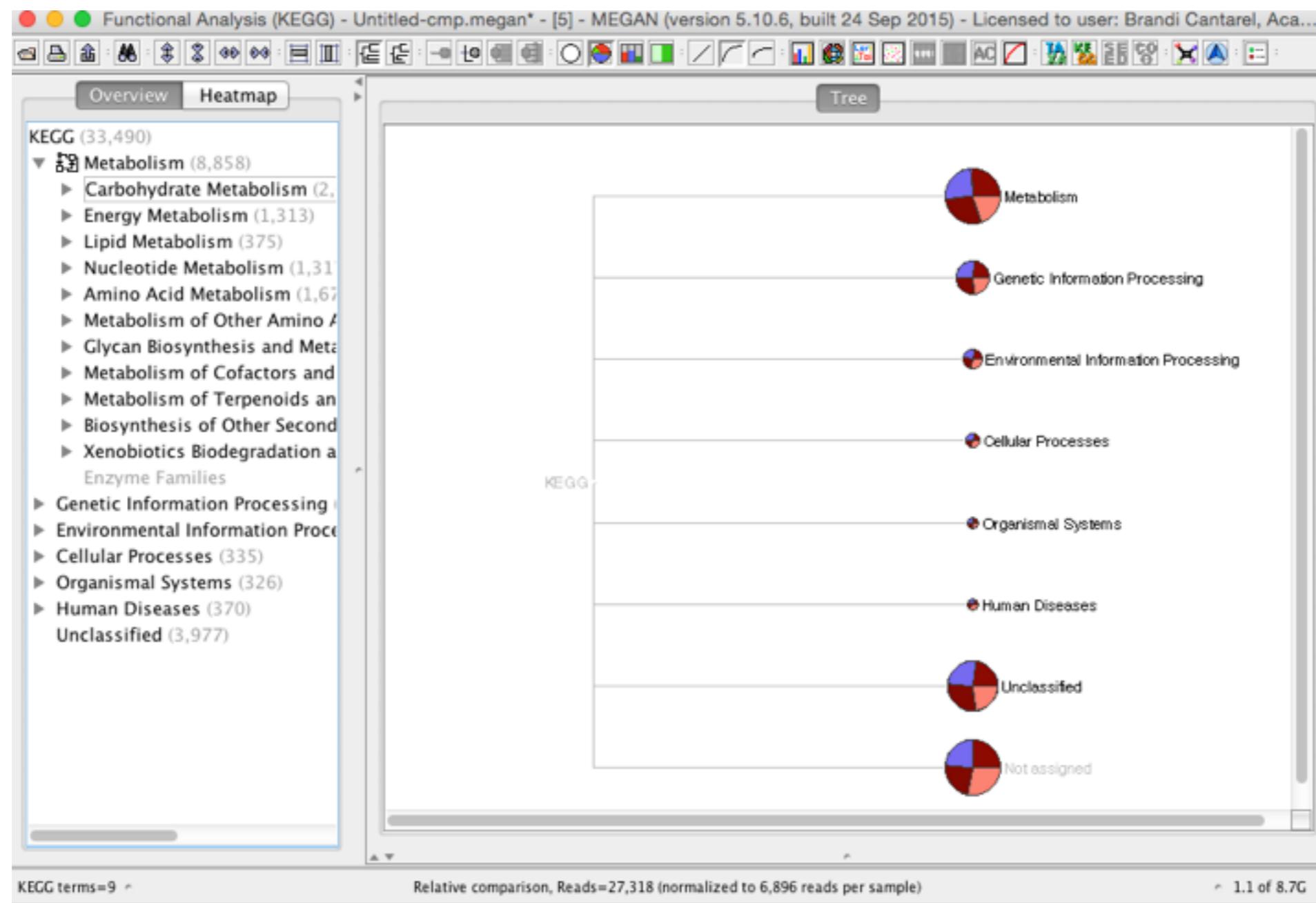
## Pathway Analysis



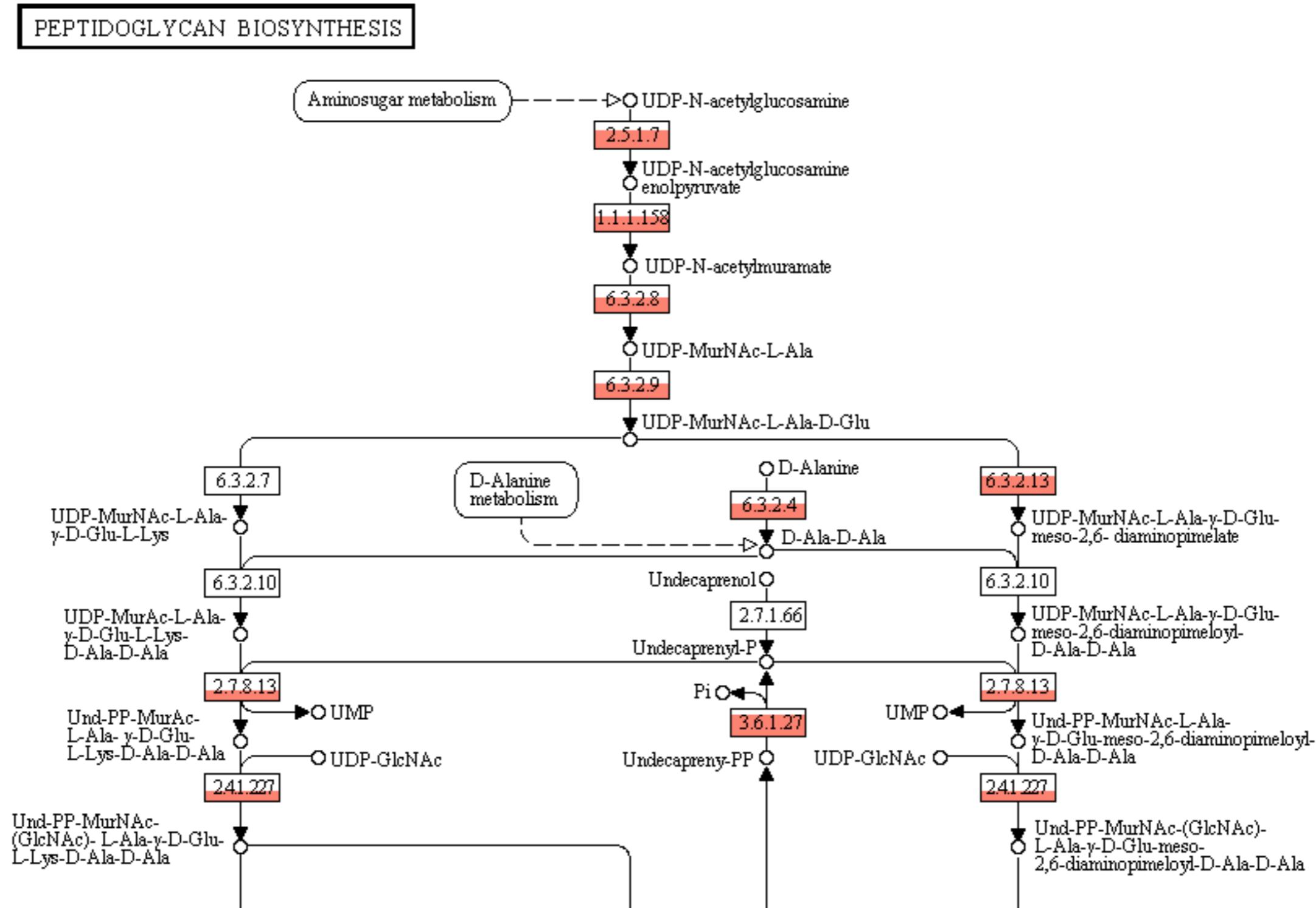
# LefSe



# MEGAN Broad Functional Comparisons



# Pathway Exploration



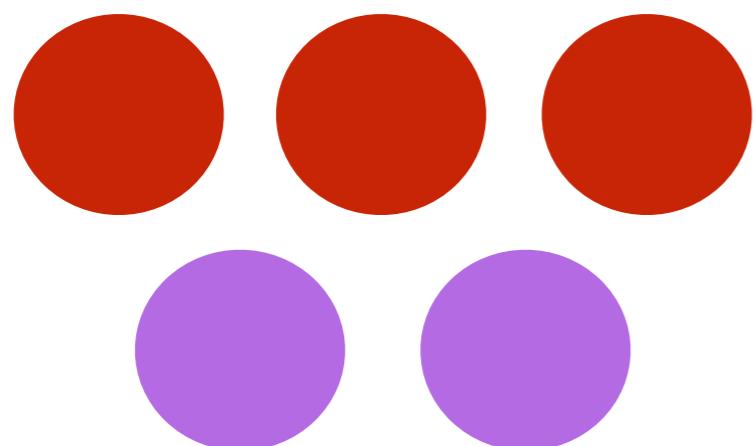
# Statistical Analysis

- Metastats
  - 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

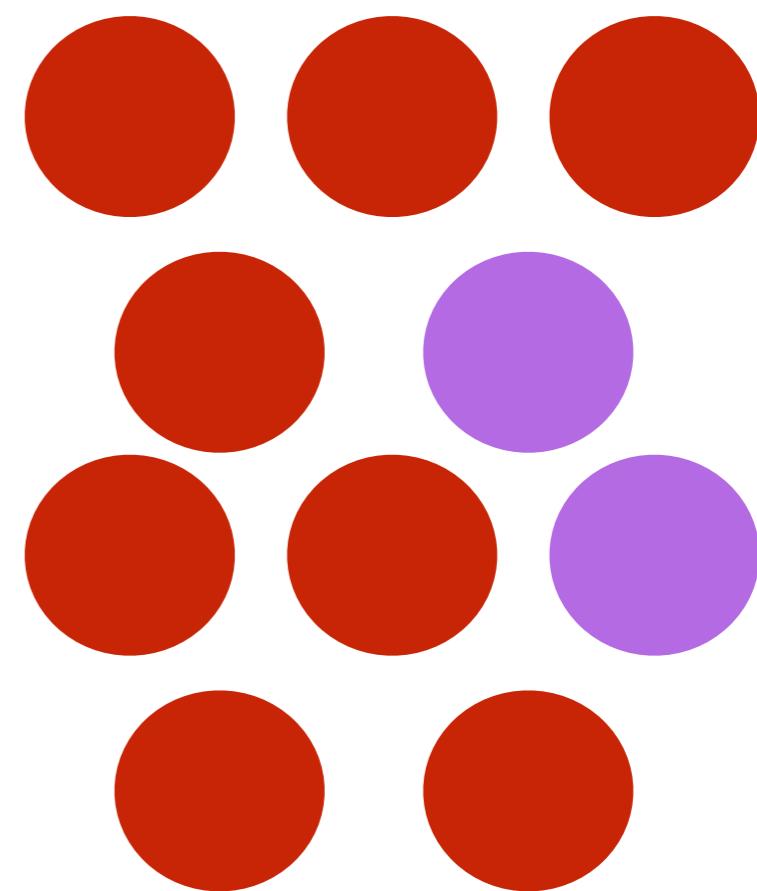
# Relative Abundance vs Absolute Abundance

- 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



Absolute: 2  
Relative: 40%



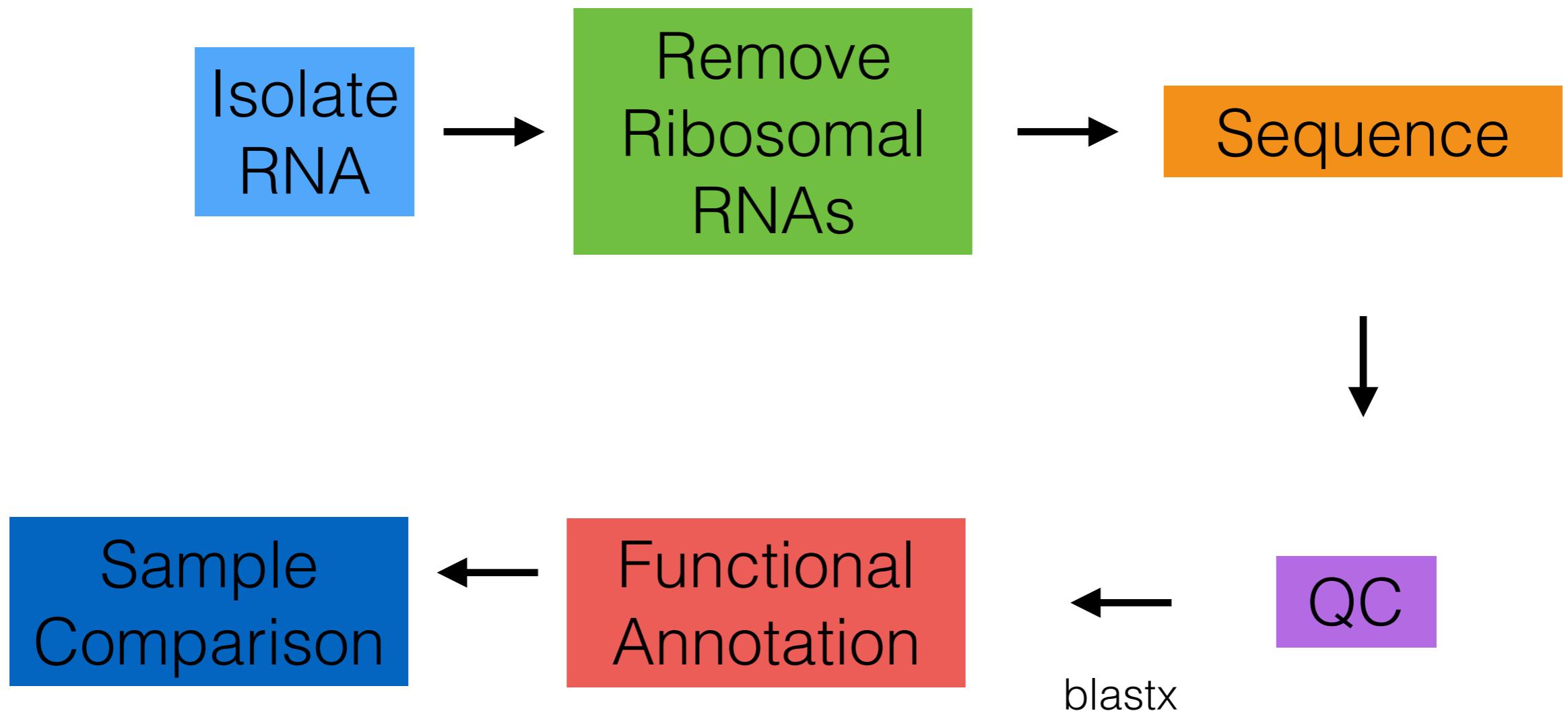
Absolute: 2  
Relative: 20%

- Community Structure Based on WGS
  - 16S VS WGS
  - Least Common Ancestor
  - Reference Genome databases
  - Assignment Stratgies
- Data Processing Workflow
  - QC
  - Functional Databases
  - Analysis Platforms
  - DYI 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 organism are the most functionally active.

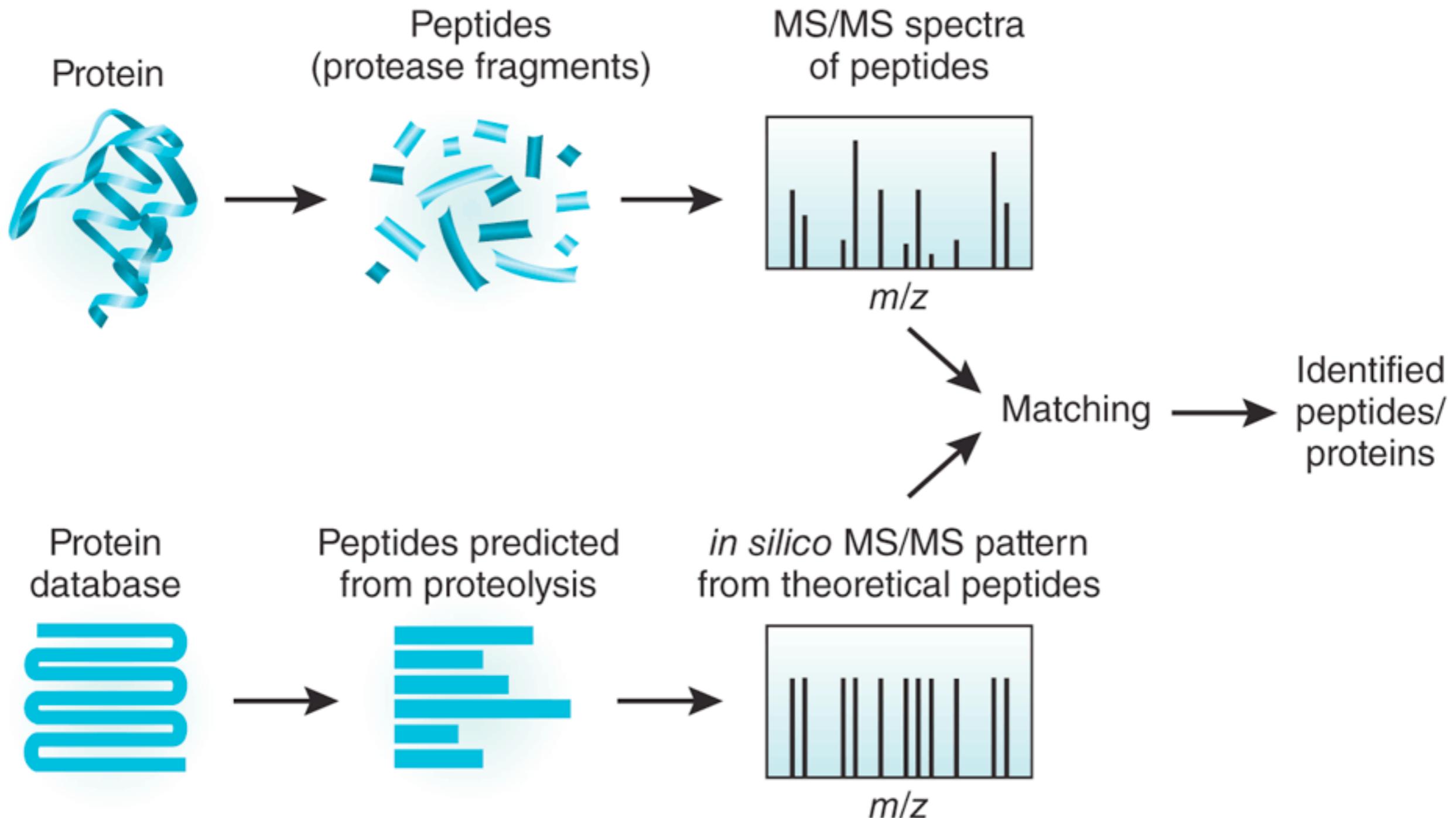
# Metatranscriptomics



# Metaproteomics

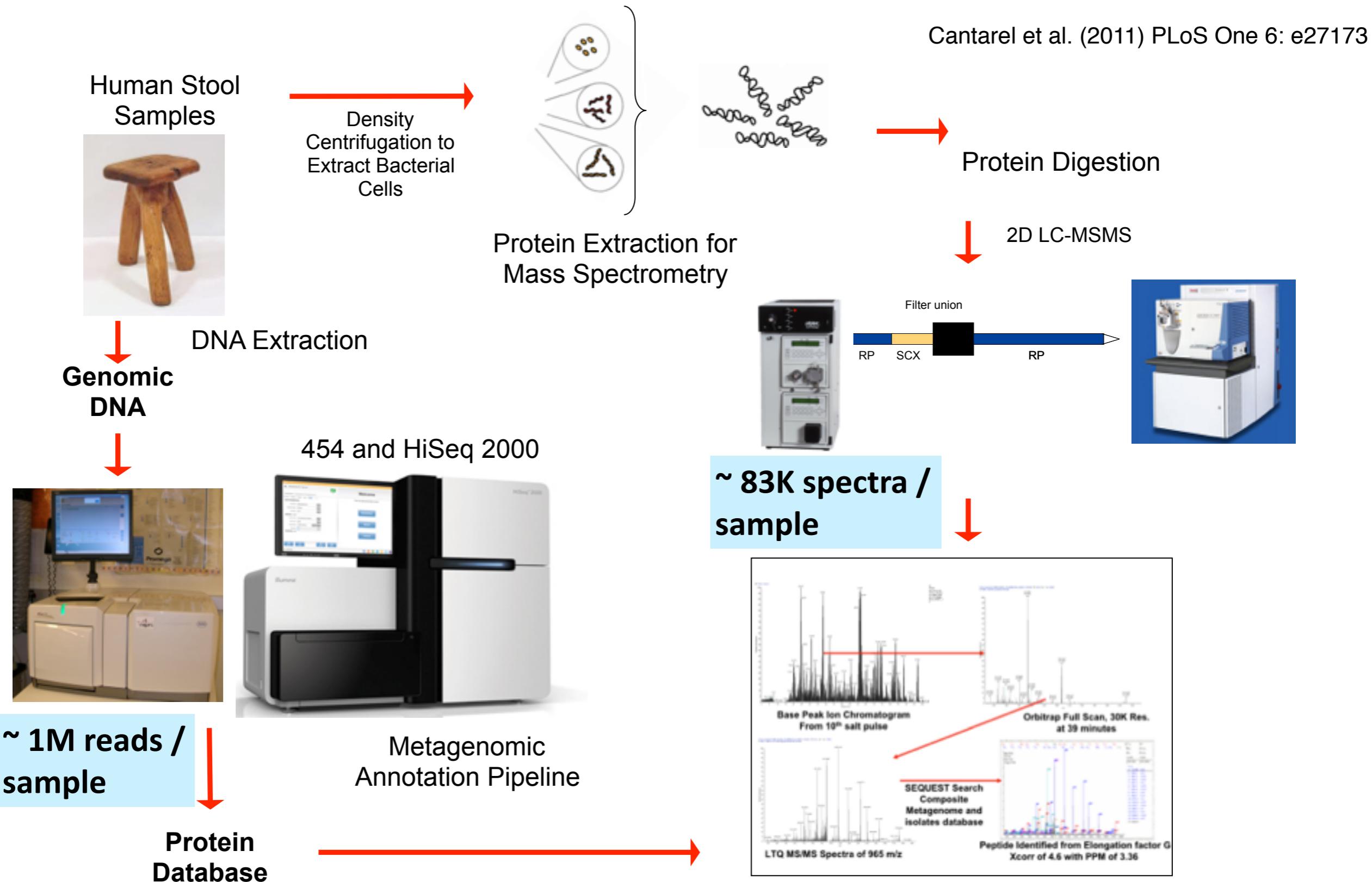
- Like metagenomics and metatranscriptomics, metaproteomics is complicated by the lack of a complete reference set
- Unlike sequencing, denovo protein prediction from MS/MS is not trivial.
- In order to determine the protein sequence of peptide fragments, a metagenomic or reference genome database is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

# Peptide Spectral Matching



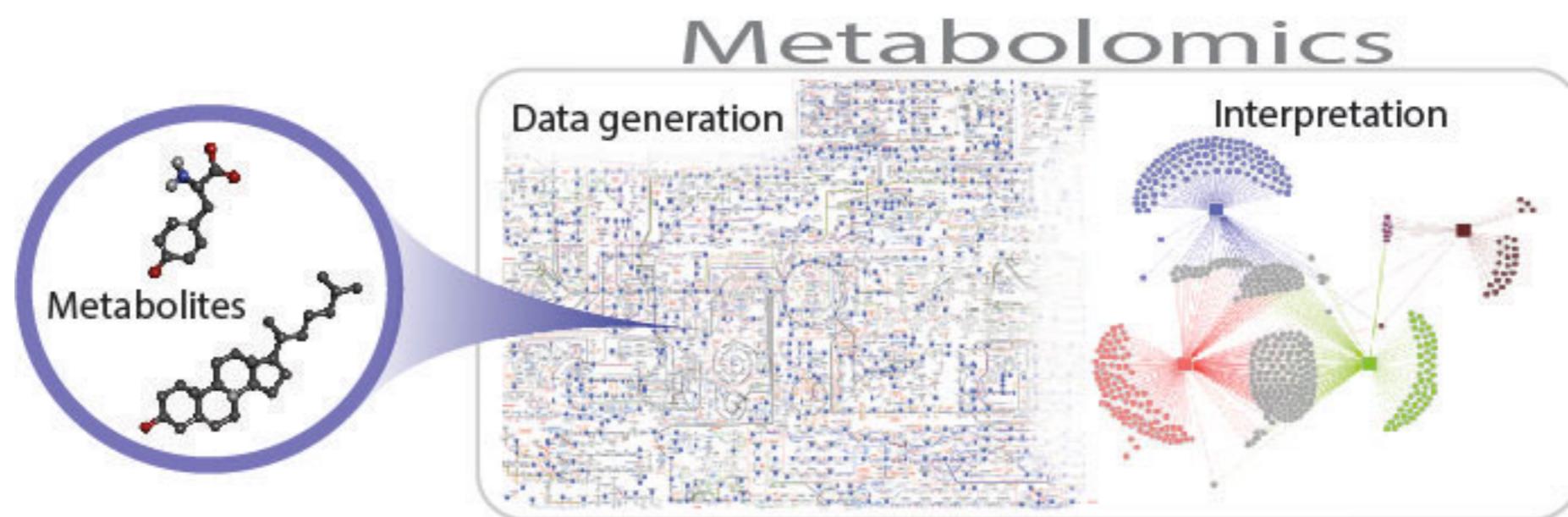
Duncan MW, Aebersold R, Caprioli RM. The pros and cons of peptide-centric proteomics. Nat Biotechnol. 2010 Jul;28(7):659–64.

# Metaproteomics



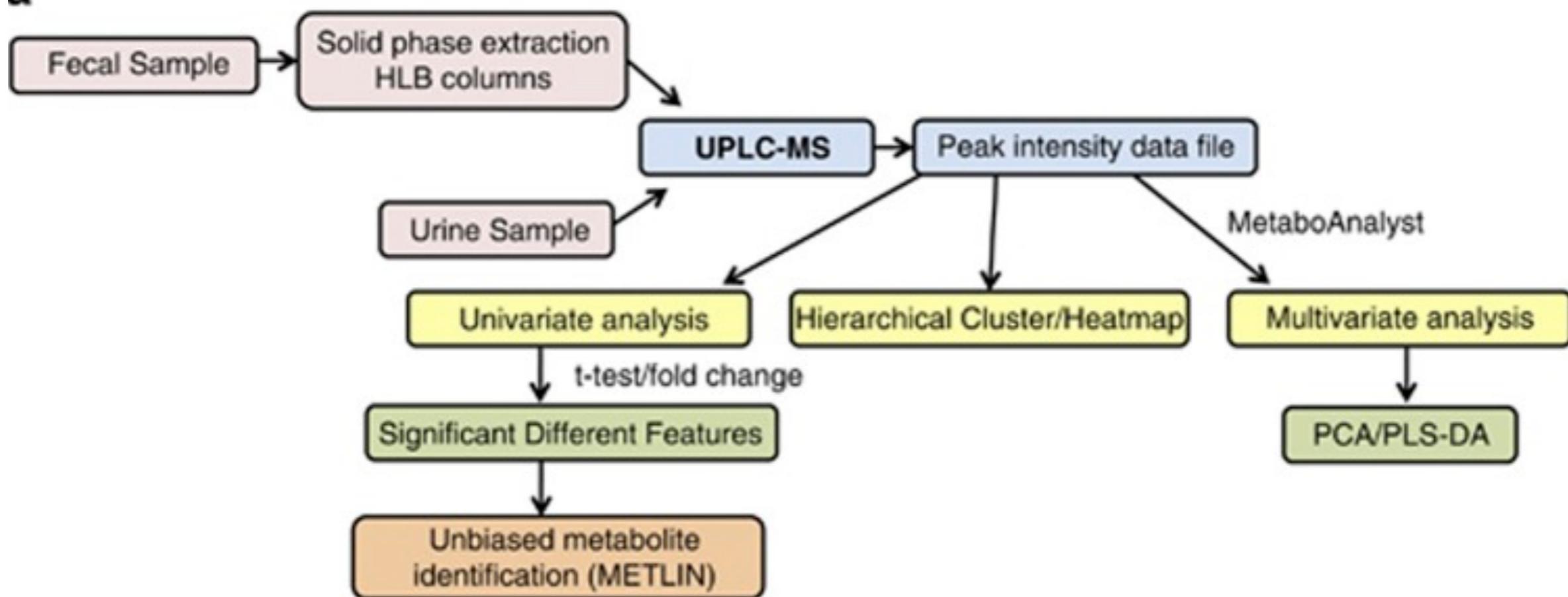
# (meta)Metabolomics

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



# Metabolomics

a



# Workshop 2