# Metagenomic Analysis

Analyzing multibacterial (multispecies) samples from the environment

November 22, 2016

# Metagenomics

- Metagenomics deals with samples taken directly from the environment
  - Soil, water, hot spring, oil sands, human gut, stool
  - Also called environmental genomics

- Necessarily more complex than genomics
  - Mixture of multiple organisms
  - Many have never been looked at on the molecular level at all

# Metagenomic analysis

- Who's there?
  - Taxonomy analysis

- How many/much of them are there?
  - Relative abundance of organisms

- What do they do?
  - Functional analysis

# Basic principles

- Simple analysis using a reference gene 16SrRNA
  - Amplification of reference gene with general primers
  - Overall microbial community composition (presence-absence)
- Whole genome sequencing (WGS)
  - Sequencing + mapping to known genomes (or to specific marker database)
  - Overall microbial community composition (quantitative)
  - Dominant functions
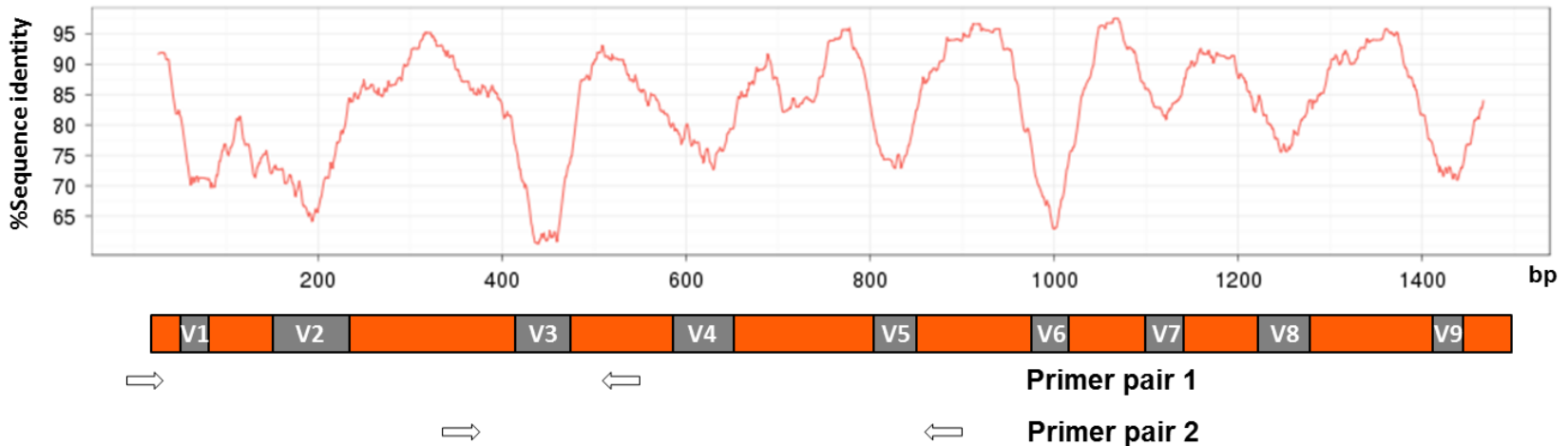  - It can also be done with assembly

# Comparative metagenomic analysis

- Across different environments
  - Individual taxon abundance
  - Overall microbial community composition
  - Dominant functions
  - Environment-specific functions of same taxon

- Within same environment
  - Changes over time

# 16S rRNA sequencing

- Early and still common method
- Highly conserved yet unique to individual (mostly bacterial) species
  - Consists of variable and conserved regions
- Targeted sequencing using primer pairs



http://**www.gatc-biotech.com/index.php?id=1025&L=1**

# 16s rRNA based community analysis procedures

**DNA Extraction**

**PCR Amplification**

**From Environment Samples**
**From Enrichment Cultures**

**Gene: 16s rRNA**
**Primer sets: 1392R454A, 926F454B**

**Data Analysis**

**Pyrosequencing**

**Microbial Community Analysis**

**454 Titanium Pyrosequencing**

# SILVA Database

# GreenGenes Database

# A 16S rRNA analysis pipeline

# Metagenomic Tools (Local)

- Qiime (Quantitative Insights Into Microbial Ecology)
  - Consists of Python scripts
  - Taxonomy and diversity statistics/visualization
- Mothur
  - Commands written in C++ programs
  - Taxonomy and diversity statics/visualization
- MEGAN (Metagenome Analyzer)
  - Provides functional as well as taxonomy analysis
  - GUI with tree-based visualization

# QUIIME (Local Installation)

# MOTHUR (Local Installation)

# Megan 5 (Local Installation)

## MEGAN5 - MEtaGenome ANalyzer

(Download here)

**MEGAN5**



MEGAN5 was written by D. H. Huson, with ideas or supporting code contributed by S.C. Schuster, S. Mitra, D.C. Richter, P. Rupek, H.-J. Ruscheweyh, R. Tappu and N. Weber.

## Introduction

In metagenomics, the aim is to understand the composition and operation of complex microbial consortia in environmental samples through sequencing and analysis of their DNA. Similarly, metatranscriptomics and metaproteomics target the RNA and proteins obtained from such samples. Technological advances in next-generation sequencing methods are fueling a rapid increase in the number and scope of environmental sequencing projects. In consequence, there is a dramatic increase in the volume of sequence data to be analyzed.

# Operational Taxonomic Units

- OTUs are basic units of taxonomy analysis
- 16S rRNA reads are clustered into OTUs
  - At 97% or 95% sequence identity (usually)
- OTUs are mapped to a taxon

# What do the analysis pipelines do?

**QIIME**
- Quality control
- Chimera detection
- OTU clustering
- Pick representative sequences
- Assign taxonomy
- Taxonomy table

**mothur**
- Quality control
- Align sequences
- Clean alignment
- Pre-cluster sequences
- Chimera detection
- Classify sequences
- Remove non bacterial sequences
- Generate distance matrix
- OTU clustering
- Classify OTU
- Taxonomy table

**MG-RAST**
- Upload sequences and metadata
- Quality control
- RNA identification
- RNA clustering
- Assign taxonomy
- Taxonomy table

# Phoenix 2: Pyrotag Analysis Pipeline

551,995 raw reads in 58 samples

369,434 good reads

40,929 good unique reads

Sample 1, 2,...,n

QC

Deduplicating & Defragmenting

Multiple Sequence Alignment

70 partitions

Distance Based OTU Merging

OTU Assignment

Distance Calculation

Partitioning Based on Sequence Identities

Taxonomic Annotation

Alpha Diversity

Beta Diversity

Hypothesis Testing

Graphs, Tables, Trees etc.

- ❖ **Capacity: > 60 samples at once, which Mothur cannot handle**
- ❖ **OTU options: Single, average, complete linkage**
- ❖ **Parallelized: 44x speed up VS. Mothur based on a 5 sample test**
- ❖ **The results are almost the same as the results from Mothur**

# Phoenix 2 Web Interface

**Phoenix 2: SSU rRNA Analysis Pipeline**
offered by the Visual Genomics Centre at the University of Calgary

Please use this form to submit a Phoenix 2 analysis job. If your request is successful and the analysis is finished, you will receive an email with a link to a Web page, where you can view or download your analysis results. This can take a while depending on your dataset size and the server load.

**Submit Phoenix 2 analysis job**
\* indicates required input                                  About Phoenix 2

| | |
|---|---|
| Upload style* | ◉ Single archive file (recommended) <br> ○ Multiple pairs of fasta and quality files |
| Archive file* | [_____] Browse… Help |
| Analysis name* | [_____] Most result files will bear this name. |
| Email address* | [_____] Notifications will be sent to this email. |
| User name | [_____] Email will be used if not given. |
| Strand direction | ◉ Reverse        ○ Forward |

Default primers

Forward
```
aaacttaaaggaattgacgg
aaacttaaatgaattgacgg
aaactcaaatgaattgacgg
aaactcaaaggaattgacgg
```

Reverse
```
acgggcggtgtgtac
acgggcggtgtgtgc
```
Help

Sample-specific primers file

Forward
[_____] Browse…

Reverse
[_____] Browse… Help

Quality control    Quality cutoff 27    Min length 200    Max length 451    Help

Clustering
Method
Average neighbor ▾

Distance cutoff
☑ 0.03 ☑ 0.05 ☐ 0.07 ☐ 0.09    Help

Representative sequence    ◉ Consensus    ○ MOF    Help

Design file    [_____] Browse…    Help

Rare OTU filtering    ◉ Do not filter    ○ Filter by frequency    Help

**Submit Job**                    Clear Form

Hide options

Please be patient after submitting - uploading files can take several minutes.

18

# Pyrotag Download Table

File  Edit  View  History  Bookmarks  Tools  Help

16S Pyrotag Sequences Downl...

hmp.coe.ucalgary.ca/HMP/pyrotags/

Google

## 16S Pyrotag Sequence Downloads

### New sequences from Run 798 (uploaded October 7, 2013)

- Voordouw lab 23 samples (Batch 32) and Gieg lab 16 samples
- 18 other samples (1C1, 1C2, 1C3, 1I1, 1I2, 1I3, 1O1, 1O2, 1O3, 2C1, 2C2, 2C3, 2I1, 2I2, 2I3, 2O1, 2O2, 2O3)

### New sequences from Run 795 (uploaded September 25, 2013)

- Voordouw lab 119 samples (59 from Batch 30, 60 from Batch 31)

### Raw sequences

| Lab | Download | Number of samples | Available since | Download size |
|-----|----------|-------------------|-----------------|---------------|
| | Raw sequences from Run 798 (Batch 32, Gieg lab 16 samples, reads from regions 1 and 2 pooled per sample) | 39 | Oct 7, 2013 | 141 MB |
| | Raw sequences from Run 795 (Batches 30 and 31) | 119 | Sep 25, 2013 | 193 MB |
| | Raw sequences from Run 788 (Gieg lab) | 20 | Aug 9, 2013 | 52 MB |
| | Raw sequences from Run 777 (Batch 29, reads from regions 3 and 4 pooled per sample) | 59 | June 17, 2013 | 97 MB |
| | Raw sequences from Run 768 (Batches 27 and 28) | 116 | April 22, 2013 | 230 MB |
| | Raw sequences pooled from Runs 741 and 762 (Batch 26) | 35 | April 2, 2013 | 74 MB |
| | Raw sequences pooled from Runs 744 and 762 (Batch 25) | 60 | April 2, 2013 | 88 MB |
| | Raw sequences pooled from Runs 741 and 744 (Batch 24) | 60 | Mar 30, 2013 | 146 MB |
| | Raw sequences from Run 718 | 60 | Nov 9, 2012 | 125 MB |
| | Raw sequences from Run 695 | 60 | Aug 20, 2012 | 95 MB |
| | Raw sequences from Run 681 | 58 | June 7, 2012 | 122 MB |
| Gerrit Voordouw & Lisa Gieg | Raw sequences from Run 666 | 55 | Apr 19, 2012 | 126 MB |
| | Raw sequences from Run 657 | 24 | Mar 29, 2012 | 63 MB |
| | Raw sequences from Run 647 | 40 | Mar 26, 2012 | 94 MB |
| | Raw sequences from Run 630 | 80 | Feb 24, 2012 | 198 MB |

19

# Data Access: Phylogeny

# Metagenomic sequencing

- Whole metagenome with „binning"
  - Map genes to (annotated) genomic sequences
  - Count hits by taxa („bins") – gives taxonomic composition incl. Quantitation of taxa.
  - You can use marker database instead of full genomes. Faster but less sensitive.

- Whole metagenome with assembly
  - Assemble reads just like in genome assembly
  - Complicated because of multiple unknown sources of metagenomic reads, lower coverage on individual genomes.
  - Requires large computers.

1) The input reads are aligned against a dbase of genomes
2) The alignments are processed using the lowest common ancestor algorithm
3) Each read is assigned to a taxonomic level
4) The result is a list of reads with the assigned taxonomy



NGS read

strain1 ... strain2 ... strain3 ...

Database of all genomes

genome of strain1

genome of strain2

genome of strain3

Taxoner output:
   List of Taxonomies and the number of reads assigned to each one

# Taxoner: Principle

1) 
2) thm
3) 
4) 



Summary

| Taxonomy | Rank | No. of Reads |
|---|---|---|
| Staphylococcus aureus (1280) | species | 90185 |
| Staphylococcus aureus subsp. aureus (46170) | subspecies | 1040 |
| Staphylococcus aureus subsp. aureus USA300 (367830) | no rank | 565 |
| Staphylococcus (1279) | genus | 439 |
| Staphylococcus aureus subsp. aureus USA300_FPR3757 (451515) | no rank | 377 |
| Bacteria (2) | superkingdom | 189 |
| root (1) | no rank | 129 |
| Staphylococcus aureus subsp. aureus M013 (1118959) | no rank | 106 |
| Bacilli (91061) | class | 69 |
| Staphylococcus aureus subsp. aureus USA300_TCH959 (450394) | no rank | 62 |
| Staphylococcus aureus Bmb9393 (1321369) | no rank | 50 |
| Staphylococcus aureus subsp. aureus LGA251 (985006) | no rank | 49 |
| cellular organisms (131567) | no rank | 46 |
| Staphylococcus aureus CA-347 (1323661) | no rank | 43 |
| Staphylococcaceae (90964) | family | 27 |

Taxoner output:
    List of Taxonomies and the number of reads assigned to each one

# Taxoner: Gene Assignment

1) Genes and functions are assigned to the lower taxonomic levels (species, subspecies and strain)
2) The algorithm is based on an integrated database created with JBioWH using Gene, PTT, COG and eggNOG databases
3) The result is a list of genes per taxonomy and COG-eggNOG functional classification



Gene and function assignment output:
List of COG groups and number of genes per taxonomies

# Taxoner: Gene Assignment

COG/eggNOG Top Classes

| Top Classes | Total |
|---|---|
| INFORMATION STORAGE AND PROCESSING | 675 |
| CELLULAR PROCESSES AND SIGNALING | 528 |
| METABOLISM | 1335 |
| POORLY CHARACTERIZED | 564 |

COG/eggNOG Functional Classification

| Functional Classes | One Letter | Total |
|---|---|---|
| RNA processing and modification | A | 0 |
| Translation, ribosomal structure and biogenesis | J | 243 |
| Chromatin structure and dynamics | B | 1 |
| Replication, recombination and repair | L | 202 |
| Transcription | K | 229 |
| Signal transduction mechanisms | T | 100 |
| Cell wall/membrane/envelope biogenesis | M | 179 |
| Extracellular structures | W | 0 |
| Cell motility | N | 2 |
| Nuclear structure | Y | 0 |

COG functional classification



| Taxonomy | Rank | No. of Genes | No. of Reads |
|---|---|---|---|
| Staphylococcus aureus subsp. aureus VC40 (1028799) | no rank | 2403 | 17677 |
| Staphylococcus aureus subsp. aureus T0131 (1006543) | no rank | 2256 | 11569 |
| Staphylococcus aureus subsp. aureus TW20 (663951) | no rank | 2259 | 11532 |
| Staphylococcus aureus subsp. aureus str. JKD6008 (546342) | no rank | 2174 | 11063 |
| Staphylococcus aureus subsp. aureus 11819-97 (1123523) | no rank | 2115 | 8125 |

Gene and function assignment output:
List of COG groups and number of genes per taxonomies

# Taxoner: Run times and examples

| | Dbase | Running time[1] | | |
|---|---|---|---|---|
| | | 1 thread | 4 threads | 12 threads |
| **MetaPhlAn** | own bacterial marker dbase[2] | 14 sec | 7 sec | 6 sec |
| **Taxoner** | NCBI nt Bacteria[3] | **165 sec** | **105 sec** | **90 sec** |
| **Taxoner** | NCBI nt full dbase[4] | 2446 sec | 2031 sec | 1866 sec |
| **MEGABLAST** | NCBI nt bacteria[3] | 8.3 h | n/a | 3.9 h |
| **MEGABLAST** | NCBI nt full dbase[4] | 37.6 h | n/a | 9.4h |

- MetaPhlAn was selected because of its speed and accuracy in estimating taxon composition
- Megablast was selected because of its reputation in alignments

Processor: Intel Xeon CPU: E5-2640
[1] Dataset: SRR292150, Reads 183 203 (27 MB)
[2] Own database with 366 988 039 nucleotides (367 MB)
[3] 15 400 949 699 nucleotides (15 GB)
[4] 52 380 339 934 nucleotides (54 GB)

# Marker databases make analysis fast but much less sensitive..



Number of reads necessary for positive identification
(for a "novel" anthrax species not included in the database)

Full dbase
165 sec
(Metaphlan)

Marker dbase
14 sec
(Metaphlan)

# HMP Metagenome Assembly Pipeline

# HMP Metagenome Analysis Flowchart

# Example: Tailings Pond Composition

Sand, clay, fines (< 44 $\mu$m)
Bitumen
Alkaline water
Hydrocarbon diluent
Naphthenic acids
Metals (V, Ni, Sr)

Microorganisms!

Aerobes at surface
Anaerobes at depth

# Depth Profile of an Active Tailings Pond

- The pond in operation since 2004
- Sampled in October 2008
- Samples collected from surface → 60 ft. deep
- "Soft" or pre-consolidated (pre-CT) tailings
- Pond is routinely treated with gypsum ($CaSO_4$)



Meddőhányó

Overview of Sequences Found at Phylum Level as a Function of Depth

# *Proteobacteria*



Increasing depth

**Surface**

α: *Hypomicrobium*
β: *Methyloversatilis*
δ: *Pelobacter*
γ: *Hydrocarboniphaga*

**5 ft.**

α: *Porphyrobacter*
β: *Rhodoferax, Hydrogenophaga*
δ: *Syntrophus, Desulfocapsa*
γ: *Coxiella*

**30 ft.**

β: *Rhodoferax, Azoarcus*
δ: *Syntrophus, Geobacter, Desulfatibacillum*
γ: *Coxiella*

**60 ft.**

β: *Brachymonas*
δ: *Syntrophus, Desulfuromonas*
γ: *Coxiella*

Nonmetric multidimensional scaling plot
by Phoenix analysis Highlevel_167samples (distance: 0.03, dissimilarity: Bray–Curtis)

**Phoenix 2 NMDS**

Oil_sands
Sample: V6_226, Axis 1: 0.381116, Axis 2: −0.410497

Legend:
- CBM_cores
- CBM_cuts
- CBM_dcuts
- CBM_water
- Oil_field
- Oil_sands
- Outgroup
- TP5
- TP6_1011
- TP6_2008
- TP_MLSB
- TP_surface

35

# Metagenomic Analysis Related Sites

- MG-RAST (Metagenomic Analysis Server)
  - https://metagenomics.anl.gov
  - 27,003 public metagenomes
- IMG (Integrated Microbial Genomes)
  - https://img.jgi.doe.gov
- SILVA and GreenGenes (Ribosomal RNA Collection)
  - http://www.arb-silva.de
  - http://greengenes.lbl.gov/cgi-bin/nph-index.cgi
  - Ribosomal RNA databases and tools

# MG-RAST Server

# Integrated Microbial Genomes

IMG MISSION | IMG ACCESS & USERS | DATA DISTRIBUTION | IMG DATA WAREHOUSE | IMG DATA MARTS | IMG OTHER | IMG REFERENCES

JGI HOME    CONTACT US

## IMG Data Management

### IMG Mission top

The **mission** of the **Integrated Microbial Genomes (IMG)** system is to support the annotation, analysis and distribution of microbial genome and metagenome datasets sequenced at **DOE's Joint Genome Institute** (JGI).

IMG is also open to **scientists worldwide** for the annotation, analysis, and distribution of their own genome and metagenome datasets, as long as they agree with the IMG **data release policy** and follow **the metadata requirements** for integrating data into IMG (see IMG submission site).

### IMG Access & Users top

IMG is committed to provide scientists worldwide **free** support for genome & metagenome data annotation & integration and open access comparative analysis of integrated genome and metagenomes.

IMG users need to register at: JGI Single Sign On (JGI SSO) in order to obtain a **login** and **password** for gaining access to IMG's data content and analysis tools. Logins/passwords allow users to (i) **submit** their own genomes/metagenomes and keep them "private" for up to two years while they review and revise annotations; (ii) **employ** IMG's **curation** tools for identifying and correcting annotation anomalies, such as protein products, for both private or public genomes-annotation revisions are recorded/saved in user specific "MyIMG" files on IMG's file system; (iii) **employ** IMG's **Workspace** which supports a persistent version of IMG's "Carts" and performing long running analysis computations; (iv) **download** IMG genome and metagenome **datasets** via JGI's Portals.

**As of Dec. 31, 2014**, IMG has **10,310** users from **88** countries across **6** continents. User Map

### Data Distribution & Distribution Policy top

Genome and metagenome datasets submitted for annotation and/or integration in **IMG** will be kept "private" for up to **two years** from the date they become available for analysis, then they will become **public**: isolate genome datasets will be kept private for 18 months, while single cell and metagenome datasets will be kept private for 24 months. A genome or metagenome dataset submitted to IMG can be replaced by newer versions of the same genome/metagenome dataset, but **cannot be removed** in order to avoid making them public.

For genome and metagenome datasets with **multiple submissions**, only the latest version will be kept in IMG, with older versions **automatically removed**.

Genome and metagenome datasets submitted for annotation and/or integration in **IMG** are distributed solely through individual genome and metagenome **data portals** and are limited to **assembled** and **annotated** datasets; no other type of data distribution (data downloads) is provided.

# What have we learnt?

- Metagenomics is the analysis of samples - usually environmental or gut microflora samples - with many thousand species

- Traditional approach uses one reference gene, 16S rRNA, amplified by PCR and NGS sequenced. Bacterial composition is obtained..

- In whole genome sequencing (WGS) the reads can be mapped to annotated genome sequences, bacterial composition and biological functions are obtained.

- Alternatively, WGS reads can be assembled at large computational overheads which makes analysis more accurate.

- Main programs: MG RAST, Megan, Mothur