# Next Generation Sequencing (NGS) Data Analysis for the Oral Microbiome
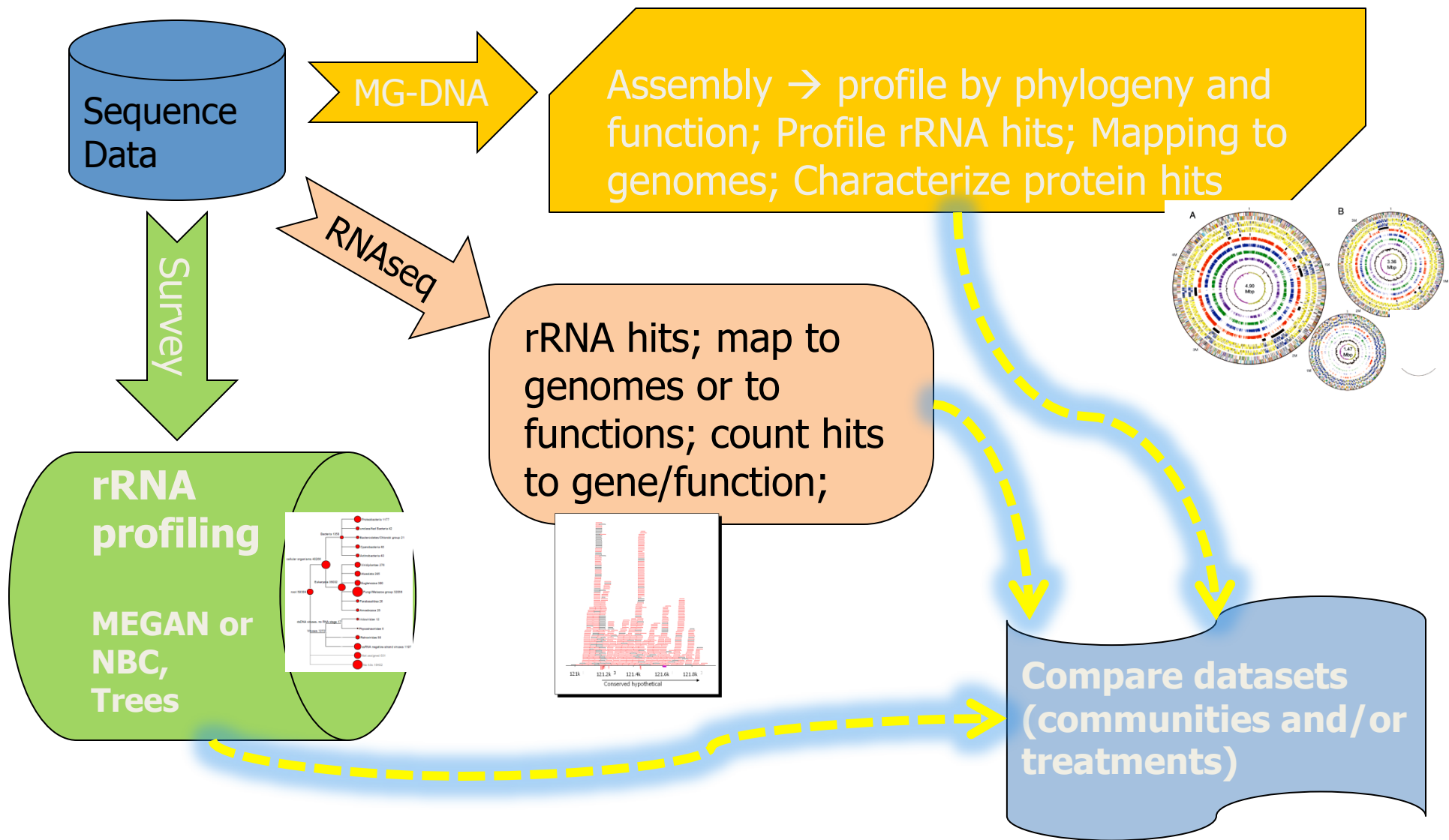


## Part 2: oral metagenome annotation and analysis

Gary Xie

xie@lanl.gov

National Institute of
Dental and Craniofacial Research
National Institutes of Health

Los Alamos
NATIONAL LABORATORY
— EST. 1943 —

# Toward a comprehensive bioinformatic workbench for the genomics community

# Community Composition Profiling

- 16s based Community Composition Profiling
  - Pyrotags
  - RDP pyro-pipeline
  - Mothur
  - Naïve Bayesian classifier
- Random shotgun metagenome based Community Composition Profiling
  - Read based
    - MEGAN
    - Sequence Recruitment
  - Contig/protein based
    - AMPHORA

# 16s Based Community Composition Profiling

Limitations:
 PCR introduced bias
 Access to the DNA (preparation work)

- Phylochip
  - Limited by probes on the microarray
  - Relative abundance estimation

- Pyrotags (Pyrosequencing PCR amplified 16s rRNA)
  - Homopolymer quality issues
  - Shorter reads: targeted variable regions

- i-tags (Illumina sequencing PCR amplified 16s rRNA)
  - hard to find universally conserved flanking region immediately adjacent to phylogenetically informative variable region
  - region of high-quality sequence is very short (<76bp)
    - reduced phylogenetic resolution
    - more dependent on reference 16s database

# Resources in 16S rRNA sequence analysis

- Various 16S rRNA databases and resources available:
  - RDPII Project
  - Greengenes
  - SILVA rRNA database project

Available 16S rRNA databases (January 2008):

| Database | B[1] | A[2] | E[3] | Alignment positions | Number of sequences |
|---|---|---|---|---|---|
| RDPII (http://rdp.cme.msu.edu/) | + | - | - | ~53,000 | ~472,000[4] |
| ARB SILVA (http://www.arb-silva.de/) | + | + | + | ~46,000 | ~504,000[5] |
| greengenes (http://greengenes.lbl.gov/) | + | + | (+) | 7,682 | ~180,000[6] |

[1] *Bacteria*, [2] *Archaea*, [3] *Eukarya*; [4] including partial sequences; [5] >300 nt; [6] >1250 nt



http://rdp.cme.msu.edu/



http://www.arb-silva.de/



http://greengenes.lbl.gov/cgi-bin/nph-index.cgi

# Overview of barcoded pyrosequencing workflow for rRNA-based community composition profiling



Amplify each sample, introducing barcode into each sequence using tagged PCR primers

Mix samples and sequence on pyrosequencer (e.g. GS FLX)

Use barcodes to assign each sequence to the sample it came from, dropping low-quality reads

K-mer based

Trimming

Use community clustering techniques (either OTU-based or tree-based) to relate samples to one another

Build phylogenetic tree using one representative of each OTU: track which parts of the tree came from which sample

Group related sequences into OTUs for downstream analyses

Trim barcodes and build multiple sequence alignment based on reference sequences

OTU based          Tree based

Homology based

Clustering

Alignment

# Analysis Toolbox

## K-mer based approach

- NBC
- Markov

## OTU-based approach

- DOTUR:  define & count OUT (distance matrix as input)
- SONS:  measuring overlap between communities (OTU designation as input )
- Estimate S - beta diversity index (Need abundance file for each OTU)

## Tree-based approach

- S-LIBSHUFF                   (distance matrix as input)
- Treeclimber                   (phylogenetic tree as input)
- Unifrac                          (phylogenetic tree as input)

**Patrick Schloss' Mothur has most of these packages**

SA Eichorst

# Comparing community memberships and structures



**Patrick Schloss' Mothur**

http://www.mothur.org/

RDP has better GUI and plug-in

# RDP Pyrosequencing Pipeline



http://pyro.cme.msu.edu/

RDP  Pyrosequencing Pipeline

Pro:
- Including the initial process: sorting, trimming tag and primers, removing low quality
- Aligner based on the 2$^{nd}$ structure (Nawrocki & Eddy, 2007)
- Including complete-linkage clustering, formatted for downstream analysis
- Including package for study community richness & diversity:
    - Generating Shannon/Chao1
    - Jaccard/sorensen index
    - Rarefaction curve
- Including the K-mer analysis: NBC and library comparison

Con:
- Only includes "Bacteria" and "Archaea" sequences, no Eukaryote
- No oral specific training set
- Library comparison only handle two samples at a time

# K-mer based NBC and LANL-Markov classifier

* **No multiple sequence alignment** required

* Provide confidence estimates for each assignment.

* Tested on the leave-one-out cross validation (LOOCV)

* Different from RDP classifier, our Markov classifier assumes a dependency on sequence reads

This page uploads a "16S rRNA" file (in fasta format) from your machine and displays the results of the 16S Classifier.

Enter the file name to upload:

[                                                    ] ( Browse... )

Training set:
- ⊙ HOMD
- ○ RDP

http://oralgen.lanl.gov/oralgen-tng/missimp.html  ( Upload File )

# Toward a comprehensive bioinformatic workbench for the genomics community

# Community Composition Profiling

- 16s based Community Composition Profiling
  - Pyrotags
  - RDP pyro-pipeline
  - Mothur
  - Naïve Bayesian classifier
- Random shotgun based Community Composition Profiling
  - Read based
    - Sequence Recruitment
    - MEGAN
  - Contig/protein based
    - AMPHORA

# Read based approach: Sequence Recruitment

Sequence recruitment---align metagenomic reads against reference genomes or genome fragments

# Fragment Recruitment Service

http://oralgen.lanl.gov/oralgen-tng/FragRecruit.html

Reference Genome:  Selenomonas flueggei ATCC 43531
Query Metagenome:   HumanGutCommunitySubject7

tabular data



37273.NZ_GG694006-NZ_GG694014.nuc

# Mapping the metagenome reads to reference genomes

HMP reference genomes: a starting point for reconstruction of microbial genomes from metagenomic sequences

# Read based: lowest common ancestor vs top blast approach

**Top Blast hit:**
—The poor representation of microbial diversity by sequenced isolates
—remote matching to phylogenetically distant organisms or the absence of any hits.

Such as MG-RAST, HOMD

**Lowest common ancestor:**
—All reads count
—Trace back to the lowest common ancestor of the set of taxa.

Such as MEGAN, AMPHORA

# Contig/protein based: Automated pipeline for phylogenomic analysis (AMPHORA)



Two applications:
- Build a genome tree from 578 complete bacterial genomes

- Identify bacterial phylotypes from metagenomic data

Wu et al 2008

# AMPHORA Features

- Fully Automated: in a pipeline
- Can be used for phylogenetic analyses of single gene or whole microbiomes.
- 31 pre-build phylogenetic marker genes
  - Most are single copy genes within each genome
  - Housekeeping genes
- Profile HMM-based multiple sequence alignment
  - High quality alignment according to seed alignment
  - Reproducible
  - Speed

A) N1 (MG-RAST)

B) N1(IMG-M ER)

Legend: dnaG, frr, infC, nusA, pgk, pyrG, rplA, rplB, rplC, rplD, rplE, rplF, rplK, rplL, rplM, rplN, rplP, rplS, rplT, rpmA, rpoB, rpsB, rpsC, rpsE, rpsI, rpsJ, rpsK, rpsM, rpsS, smpB, tsf

# Summary: Community Composition Profiling

Approaches for studying the community composition
- Homology based
    - Lowest common ancestor (MEGAN)
    - Top blast hits (shotgun reads, 16s)
    - Phylogenetic tree (31 house keeping gene)

- K-mer character based
    - NBC
    - Markov

Input data used:
- Targeted:
    - 16s rRNA reads
    - Other housekeeping genes
- Random shotgun reads

# Combination of different strategies

# Toward a comprehensive bioinformatic workbench for the genomics community

***Bins*** are sets of metagenomic sequence fragments originating from one phylogenetic group, preferably from the same species

# binning methods

- Binning methods developed can fall into two categories (biology perspective):
  - similarity-based
  - composition-based

- From machine learning perspective, binning methods can also be divided into
  - supervised learning
  - unsupervised learning

# Biology perspective

- **similarity-based**
  - Assign metagenomic fragments to their closest phylogenetic neighbor based on coding-sequence identity.
  - ex. BLAST, dotter

- **composition-based**
  - Distinguish genomes from one another by intrinsic features of the sequence.
  - ex. olionucleotide frequencies, GC content

# machine learning perspective

- ## Supervised learning
  - o Methods that build a classifier using the knowledge of completed genomes
  - o (Chen et al., 2005 )Current amount of known genomes is insufficient to represent the almost limitless microbial.
- ## unsupervised learning
  - o Do not have dependence on training data
  - o Directly clustering metagenomic samples
  - o Focuses on the long fosmid-sized fragments

# Binning methods matrix

|  | Composition-based | Similarity-based |
|---|---|---|
| **supervised** | PhyloPythia (IBM), 2007<br>ClaMS<br><br>MEGAN, 2007<br><br>Naïve Basin Classifier, 2008 | BLAST |
| **unsupervised** | S-GSOM, 2008<br><br>TETRA, 2005<br>Fuzzy k-means classifier, 2008<br>CompostBin, 2008<br>GSOM, 2008 |  |

KL Liu

# Metagenome Annotation

IMG/M-ER – the combined approach of BLASTX similarity search and de novo gene prediction.

- 80-299 bp:
    - MultiBLASTx against IMG-NR
    - All frameshift fragments are joined afterwards.
- 300-699 bp:
    - MultiBLASTx > GeneMark>  Metagene.
- >700 bp:
    - GeneMark>Metagene

MG-RAST –read based only
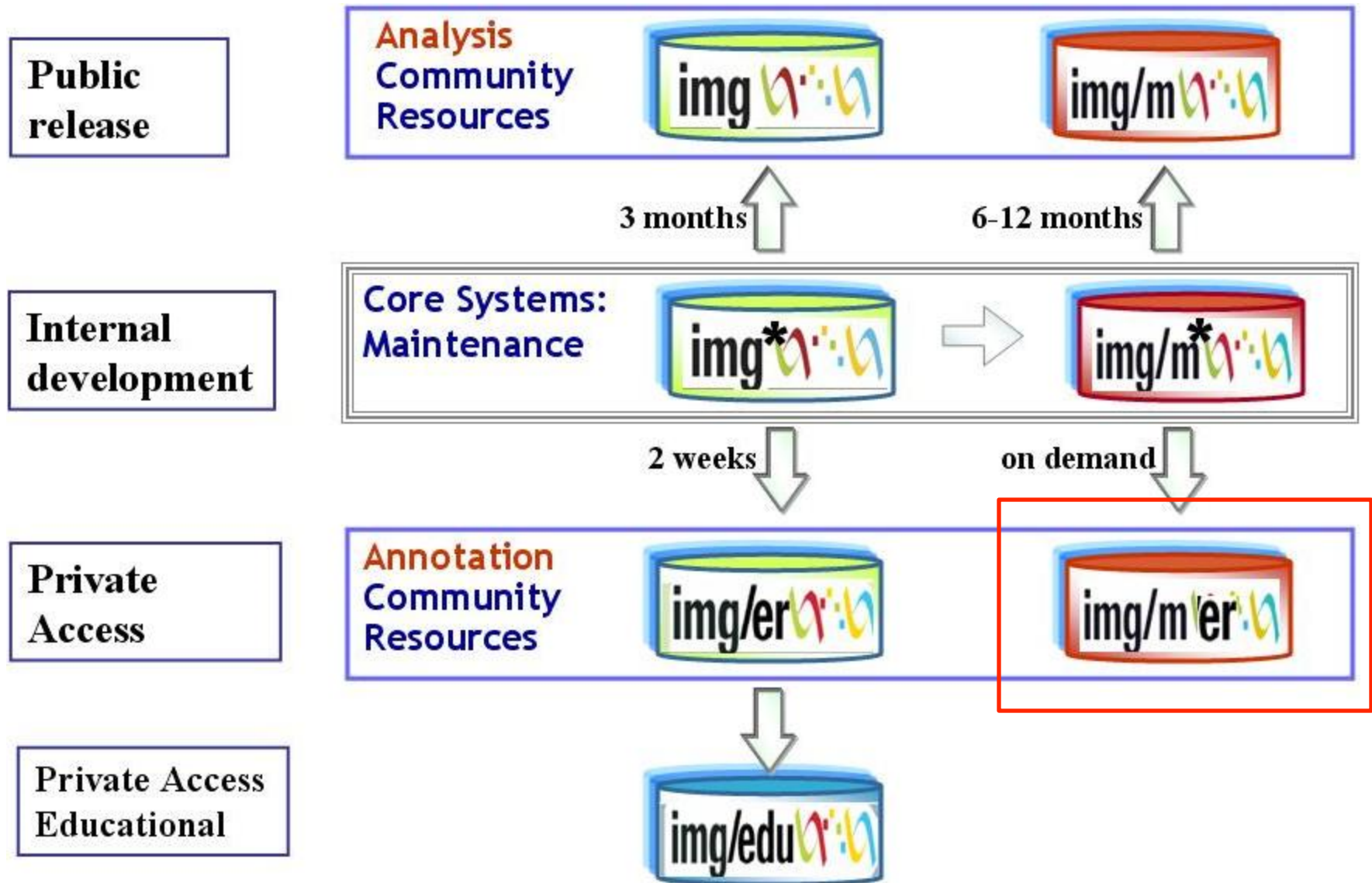Top blastx hit to SEED (reference db) only

# The Integrated Microbial Genomes (IMG) system
## http://img.jgi.doe.gov

http://metagenomics.nmpdr.org/

# MG-RAST
## Meta Genome Rapid Annotation using Subsystem Technology

version 1.2

The NMPDR, SEED-based, prokaryotic genome annotation service.
For more information about the SEED please visit theSEED.org.

The metagenomics RAST server (http://metagenomics.nmpdr.org) is a SEED-based environment that allows users to upload metagenomes for automated analyses. The server is built as a modified version of the RAST server. The RAST (Rapid Annotation using Subsystem Technology) technology was originally implemented to allow automated high-quality annotation of complete or draft microbial genomes using SEED data, and has been adapted for metagenome analysis.

Our freely available server provides the annotation of sequence fragments, their phylogenetic classification, functional classification of samples, and comparison between multiple metagenomes. The server also computes an initial metabolic reconstruction for the metagenome and allows comparison of metabolic reconstructions of metagenomes and genomes.

User submission and analysis are confidential. Although we do not guarantee a maximum turnover time, the current average processing time is about 24 hours. Currently the server handles 454 and Sanger sequence data. Data sets supplied by 454 can be uploaded directly.

The server relies on the technology and data established by FIG and the NMPDR team at Argonne National Laboratory and the University of Chicago.

In addition to SEED data we use the following ribosomal RNA databases for our analyses: GREENGENES, RDP-II and European ribosomal RNA database.

To be able to contact you once the computation is finished and in case user intervention is required, we request that you register with email address.

Login [                    ]

Password [                    ]  [ Login ]

Forgot your password?
Register a new account

# Credits

- NIH-NIDCR
- LANL-ORALGEN
  - Chienchi Lo
  - Pavel Senin
  - Andrew Liu
  - Patrick Chain