

MG-RAST  
Technical report and manual  
for version 3.3.6 – rev. 1

Andreas Wilke<sup>1,2</sup>, Elizabeth M. Glass<sup>1,2</sup>, Jared Bischof<sup>2,1</sup>, Daniel Braithwaite<sup>1,2</sup>,  
Mark DSouza<sup>2,1</sup>, Wolfgang Gerlach<sup>2,1</sup>, Travis Harrison<sup>2,1</sup>, Kevin Keegan<sup>1,2</sup>,  
Hunter Matthews<sup>1,2</sup>, Tobias Paczian<sup>2,1</sup>, Wei Tang<sup>1,2</sup>, William L. Trimble<sup>2,1</sup>, Jared  
Wilkening<sup>1,2</sup>, Narayan Desai<sup>1,2</sup>, and Folker Meyer<sup>1,2</sup>

<sup>1</sup>Argonne National Laboratory

<sup>2</sup>University of Chicago

May 17, 2013

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Version history . . . . .	7
1.2	Comparison of versions 2 and 3 . . . . .	7
1.3	The MG-RAST team . . . . .	9
<b>2</b>	<b>Under the hood: The MG-RAST technology platform and pipeline</b>	<b>12</b>
2.1	Data model . . . . .	12
2.2	Details on the new MG-RAST pipeline . . . . .	14
2.2.1	Preprocessing . . . . .	14
2.2.2	Dereplication . . . . .	17
2.2.3	DRISEE . . . . .	17
2.2.4	Screening . . . . .	17
2.2.5	Gene calling . . . . .	17
2.2.6	AA clustering . . . . .	18
2.2.7	Protein identification . . . . .	18
2.2.8	Annotation mapping . . . . .	18
2.2.9	Abundance profiles . . . . .	18
2.3	The rRNA pipeline . . . . .	20
2.3.1	rRNA detection . . . . .	20
2.3.2	rRNA clustering . . . . .	20
2.3.3	rRNA identification . . . . .	20
2.4	Quality assessment . . . . .	21
2.4.1	DRISEE . . . . .	21
2.4.2	Kmer profiles . . . . .	21
2.4.3	Nucleotide histograms . . . . .	21
2.5	Best hit, representative hit, and lowest common ancestor interpretation . . . . .	23
2.5.1	Best hit . . . . .	23

2.5.2	Representative hit . . . . .	23
2.5.3	Lowest Common Ancestor (LCA) . . . . .	24
2.5.4	Comparison of methods . . . . .	24
2.6	Numbers of annotations vs. number of reads . . . . .	24
2.7	Metadata, Publishing, and Sharing . . . . .	25
2.7.1	Metadata . . . . .	25
2.7.2	Publishing . . . . .	26
2.7.3	Sharing . . . . .	26
2.7.4	Identifiers . . . . .	26
2.7.5	Linking to MG-RAST . . . . .	27
<b>3</b>	<b>MG-RAST v3 web interface</b>	<b>28</b>
3.1	Technical details . . . . .	30
3.1.1	Browser requirements . . . . .	30
3.1.2	Downloading figures . . . . .	30
3.2	Sitemap for MG-RAST . . . . .	30
3.3	Upload page . . . . .	31
3.4	Browse page – Metadata-enabled data discovery . . . . .	32
3.5	Project page . . . . .	32
3.6	Overview page . . . . .	35
3.6.1	Technical part of the Overview page – Details on sequencing and analysis .	35
3.6.2	Biological part of the Overview page . . . . .	40
3.7	Download page . . . . .	42
3.8	Search Page . . . . .	44
3.9	Analysis page . . . . .	46
3.9.1	Normalization . . . . .	48
3.9.2	Rarefaction . . . . .	49
3.9.3	KEGG mapper . . . . .	51
3.9.4	Recruitment plots . . . . .	51
3.9.5	Bar charts . . . . .	53
3.9.6	Tree diagram . . . . .	53
3.9.7	Heatmap/Dendrogram . . . . .	58
3.9.8	Ordination . . . . .	60
3.9.9	Table . . . . .	60
3.9.10	Workbench . . . . .	63

<b>4 User Manual</b>	<b>64</b>
4.1 Privacy, Identifiers, Sharing, and Publication . . . . .	64
4.2 Uploading to MG-RAST . . . . .	64
4.2.1 Assembled data with read abundance information . . . . .	65
4.2.2 Steps for submission via the web interface . . . . .	65
4.2.3 Cmd-line uploader . . . . .	69
4.2.4 Managing the Inbox . . . . .	69
4.2.5 Generating metadata for the submission . . . . .	72
4.3 Working with Projects and Collections . . . . .	76
4.4 Understanding Datasets . . . . .	77
4.5 Drilling Down with the Workbench . . . . .	81
4.6 Downloads from the Workbench . . . . .	84
4.7 Viewing Evidence . . . . .	86
4.8 MG-RAST Output . . . . .	86
4.8.1 Data products on the website . . . . .	89
4.8.2 FTP server . . . . .	89
4.8.3 Downloads . . . . .	90
<b>5 Putting It All in Perspective</b>	<b>91</b>
5.1 MG-RAST” A community resource . . . . .	91
5.2 Future Work . . . . .	92
5.2.1 Roadmap . . . . .	93
<b>A The downloadable files for each data set</b>	<b>95</b>
<b>B Terms of Service</b>	<b>99</b>
<b>C Tools and data used by MG-RAST</b>	<b>101</b>
C.1 Databases . . . . .	101
C.1.1 Protein databases . . . . .	101
C.1.2 Ribosomal RNA databases: . . . . .	102
C.2 Software . . . . .	102
C.2.1 Bioinformatics codes: . . . . .	102
C.2.2 Web/UI tools: . . . . .	102
C.2.3 Behind the scenes: . . . . .	103
Glossary and Acronyms . . . . .	108



# Chapter 1

## Introduction

The National Human Genome Research Institute (NHGRI), a division of the National Institutes of Health, publishes information (see Figure 1) describing the development of computing costs and DNA sequencing costs over time [25]. The dramatic gap between the shrinking costs of sequencing and the more or less stable costs of computing is a major challenge for biomedical researchers trying to use next-generation DNA sequencing platforms to obtain information on microbial communities. Wilkening *et al.* [43] provide a real currency cost for the analysis of 100 gigabasepairs of DNA sequence data using BLASTX on Amazon’s EC2 service: \$300,000.<sup>1</sup> A more recent study by University of Maryland researchers [1] estimates the computation for a terabase of DNA shotgun data using their CLOVR metagenome analysis pipeline at over \$5 million per terabase.

Nevertheless, the growth in data enabled by next-generation sequencing platforms also provides an exciting opportunity for studying microbial communities: 99% of the microbes in which have not yet been cultured [33]. Cultivation-free methods (often summarized as metagenomics) offer novel insights into the biology of the vast majority of life on Earth [37].

Three types of metagenomics experiments are commonly used:

1. Environmental clone libraries (“functional metagenomics): use of Sanger sequencing (frequently) instead of more cost-efficient next-generation sequencing
2. Amplicon studies (single gene studies, 16s rDNA): next-generation sequencing of PCR amplified ribosomal genes providing a single reference gene-based view of microbial community ecology
3. Shotgun metagenomics: use of next-generation technology applied directly to environmental samples

---

<sup>1</sup>This includes only the computation cost, no data transfer cost, and was computed by using 2009 prices.

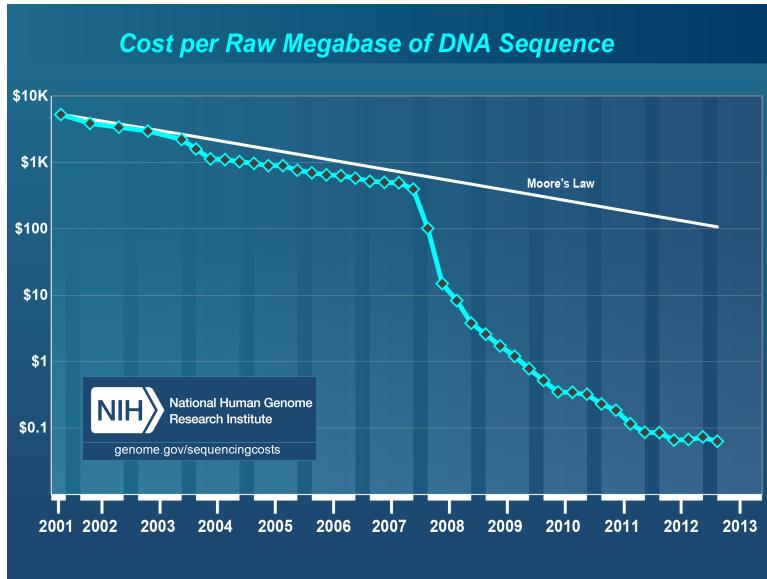


Figure 1.1: Chart showing shrinking cost for DNA sequencing. This comparison with Moore’s law roughly describing the development of computing costs highlights the growing gap between sequence data and the available analysis resources. Source: NHGRI

Each of these methods has strengths and weaknesses (see [37]), as do the various sequencing technologies (see [26]).

To support user-driven analysis of all types of metagenomic data, we have provided MG-RAST [24] (<http://metagenomics.anl.gov>). MG-RAST enables researchers to study the function and composition of microbial communities. The MG-RAST portal offers automated quality control, annotation, comparative analysis, and archiving services. At the time of writing (June 14, 2013) MG-RAST has completed the analysis of over 25 terabasepairs of DNA data in over 78,000 datasets contributed by thousands of researchers worldwide.

The MG-RAST system provides answers to the following scientific questions:

- Who is out there? Identifying the composition of microbial composition either by using amplicon data for single genes or by deriving community composition from shotgun metagenomic data using sequence similarities.
- What are they doing? Using shotgun data (or metatranscriptomic data) to derive the functional complement of a microbial community using similarity searches against a number of databases.
- Who is doing what? Based on sequence similarity searches, identifying the organisms en-

coding specific functions.

## 1.1 Version history

### Version 1

The original version of MG-RAST was developed in 2007 by Folker Meyer, Andreas Wilke, Daniel Paarman, Bob Olson, and Rob Edwards. It relied heavily on the SEED [28] environment and allowed upload of preprocessed 454 and Sanger data.

### Version 2

Version 2, released in 2008, had numerous improvements. It was optimized to handle full-sized 454 datasets and was the first version of MG-RAST that was not fully SEED based. Version 2.0 used BLASTX analysis for both gene prediction and functional classification [24].

### Version 3

While version 2 of MG-RAST was widely used, it was limited to datasets smaller than a few hundred megabases, and comparison of samples was limited to pairwise comparisons. Version 3 is not based on SEED technology; instead, it uses the SEED subsystems as a preferred data source. Starting with version 3, MG-RAST moved to github.

## 1.2 Comparison of versions 2 and 3

In the 3.0 version, datasets of tens of gigabases can be annotated, and comparison of taxa or functions that differ between samples is now limited only by the available screen real estate. Figure 1.2 shows a comparison of the analytical and computational approaches used in MG-RAST v2 and v3. The major changes are the inclusion of a dedicated gene-calling stage using FragGenescan [32], clustering of predicted proteins at 90% identified by using uclust [9], and the use of BLAT [18] for the computation of similarities. Together with changes in the underlying infrastructure, this version has allowed dramatic scaling of the analysis with the limited hardware available.

Similar to version 2.0, the new version of MG-RAST does not pretend to know the correct parameters for the transfer of annotations. Instead, users are empowered to choose the best parameters for their datasets.

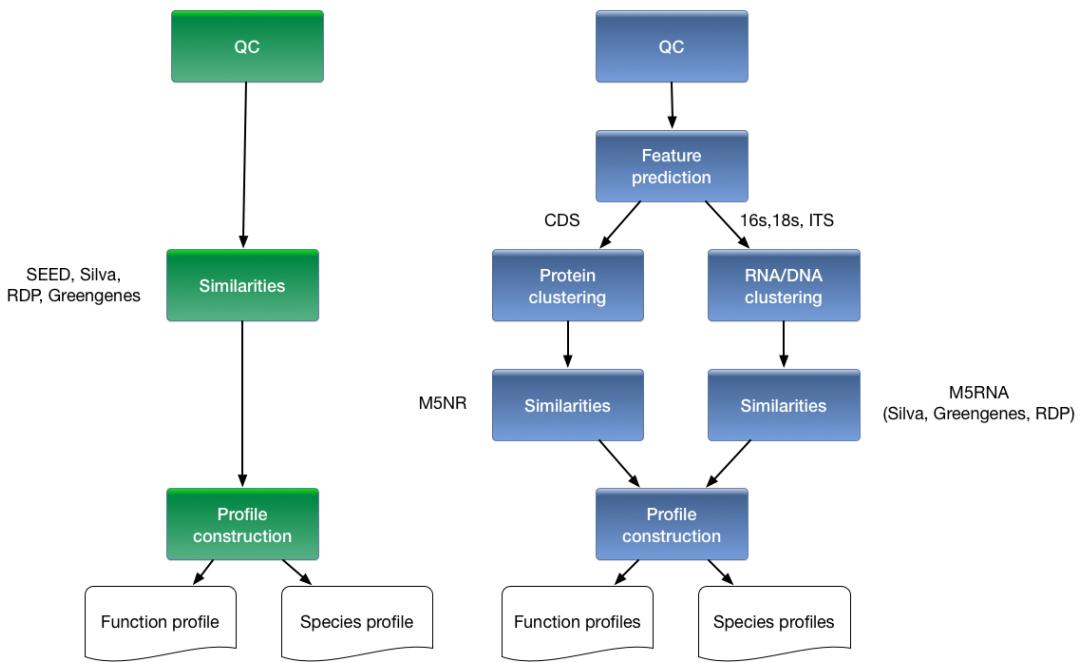


Figure 1.2: Overview of processing pipeline in (left) MG-RAST v2 and (right) MG-RAST v3. In the old pipeline, metadata was rudimentary, compute steps were performed on individual reads on a 40-node cluster that was tightly coupled to the system, and similarities were computed by BLAST to yield abundance profiles that could then be compared on a per sample or per pair basis. In the new pipeline, rich metadata can be uploaded, normalization and feature prediction are performed, faster methods such as BLAT are used to compute similarities, and the resulting abundance profiles are fed into downstream pipelines on the cloud to perform community and metabolic reconstruction and to allow queries according to rich sample and functional metadata.

The new version of MG-RAST represents a rethinking of core processes and data products, as well as new user interface metaphors and a redesigned computational infrastructure. MG-RAST supports a variety of user-driven analyses, including comparisons of many samples, previously too computationally intensive to support for an open user community.

Scaling to the new workload required changes in two areas: the underlying infrastructure needed to be rethought, and the analysis pipeline needed to be adapted to address the properties of the newest sequencing technologies.

## 1.3 The MG-RAST team

MG-RAST was started by Rob Edwards and Folker Meyer in 2007. The MG-RAST team has significantly expanded in the past few years. The current team is listed below.

- Andreas Wilke
- Elizabeth M. Glass
- Jared Bischof
- Daniel Braithwaite
- Mark DSouza
- Wolfgang Gerlach
- Travis Harrison
- Kevin Keegan
- Hunter Matthews
- Tobias Paczian
- Wei Tang
- William L. Trimble
- Jared Wilkening
- Narayan Desai
- Folker Meyer

The image shows a light gray rectangular box containing the email address "mg-rast@mcs.anl.gov". The text is in a standard sans-serif font, with "mg-rast" in blue, "@" in black, "mcs" in orange, and ".anl.gov" in blue.

Figure 1.3: The email address for the MG-RAST project. Note that this is inserted into the document as an image, you will have to type it.

## Contacting the MG-RAST team and help desk

The MG-RAST project uses a ticket system to manage interactions with users. Mark D'Souza is managing the help desk interaction with the users. To email him, please use the email address for the MG-RAST project:

We recommend including as much detail as possible into your emails to the help-desk, details like account names, MG-RAST identifiers will help us identify any issues and speed up resolving them.

Below is an example of the types of details we'd like to receive:

- your name
- your account name for MG-RAST
- a clear text description of your problem
- any MG-RAST identifiers (those are the 444xxxx.3 numbers)
- any project numbers
- the browser and which version you are using in it, if the problem relates to the web site
- what platform your data was created on
- if your data was a failure in the web site, what time the failure occurred

## Past members of the MG-RAST team

The following people were associated with MG-RAST in the past:

- Daniel Paarmann, 2007-2008

- Rob Edwards, 2007-2008
- Mike Kubal, 2007-2008
- Alex Rodriguez, 2007-2008
- Bob Olson, 2007-2009
- Daniela Bartels, 2007-2011
- Yekaterina Dribinsky, 2011

MG-RAST was started by Rob Edwards and Folker Meyer in 2007.

# Chapter 2

## Under the hood: The MG-RAST technology platform and pipeline

One key aspect of scaling MG-RAST to large numbers of modern NGS datasets is the use of cloud computing,<sup>1</sup> which decouples MG-RAST from its previous dedicated hardware resources. Using our task server AWE [42] and the SHOCK data management system developed alongside it, we have updated our underlying computational platform using a purpose-built software platform optimized for large-scale sequence analysis.

The new analytical pipeline for MG-RAST v3 is encapsulated and separated from the data store, enabling far greater scalability.

Combined, these changes in infrastructure and pipeline have made version 3 of MG-RAST 750 times faster than version 2.

### 2.1 Data model

The MG-RAST data model (see Figure 2.1) has changed dramatically in order to handle the size of modern next-generation sequencing datasets. In particular, we have made a number of choices that reduce the computational and storage burden.

We note that the size of the derived data products for a next-generation dataset in MG-RAST is typically about 10x the size of the actual dataset. Individual datasets now may be as large as a terabase,<sup>2</sup> with the on-disk footprint significantly larger than the basepair count because of the inefficient nature of FASTQ files, which basically double the on-disk size for FASTQ representations.

---

<sup>1</sup>We use the term *cloud* as a shortcut for Infrastructure as a Service (IaaS).

<sup>2</sup>This would be for several metagenomes that are part of the JGI Prairie pilot.

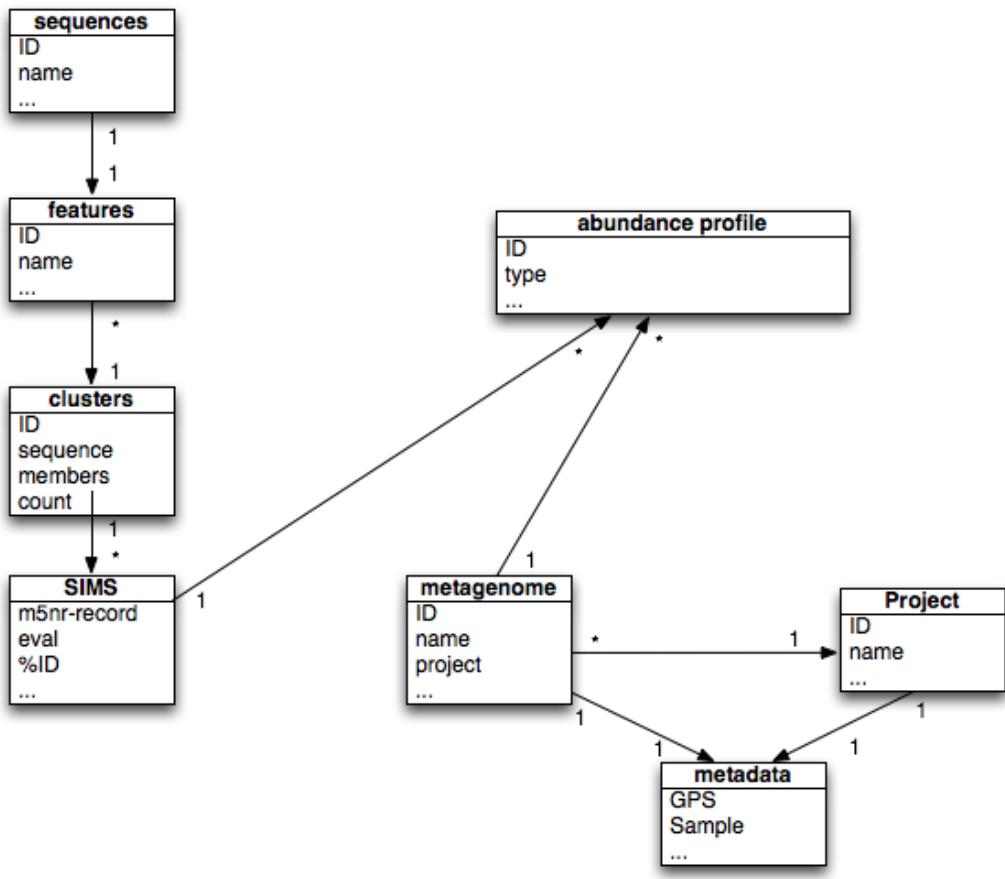


Figure 2.1: MG-RAST v3 data model.

- Abundance profiles. Using abundance profiles, where we count the number of occurrences of function or taxon per metagenomic dataset, is one important factor that keeps the datasets manageable. Instead of growing the dataset sizes (often with several hundred million individual sequences per dataset), the data products now are more or less static in size.
- Single similarity computing step per feature type. By running exactly one similarity computation for proteins and another one for rRNA features, we have limited the computational requirements.
- Clustering of features. By clustering features at 90% identity, we reduce the number of times we compute on similar proteins. Abundant features will be even more efficiently clustered, leading to more compression among for abundant species.

As shown in Figure 2.1, MG-RAST relies on abundance profiles to capture information for

each metagenome. The following abundance profiles are calculated for every metagenome.

- MD5s – number of sequences (clusters) per database entry in the M5nr.
- functions – summary of all the MD5s that match a given function.
- ontologies – summary of all the MD5s that match a given hierarchy entry.
- organisms – summary of all MD5s that match a given taxon entry.
- lowest common ancestors

The static helper tables (show in blue in Figure 2.2) help keep the main tables smaller, by normalizing and providing integer representations for the entities in the abundance profiles.

## 2.2 Details on the new MG-RAST pipeline

The pipeline shown in Figure 2.3 contains a significant number of improvements over version 3.0. Using the M5NR [41] (an MD5 nonredundant database), the new pipeline computes results against many reference databases instead of only SEED. Several key algorithmic improvements were needed to support the flood of user-generated data (see Figure 2.4). Using dedicated software to perform gene prediction instead of using a similarity-based approach reduces runtime requirements. The additional clustering of proteins at 90% identity reduces data while preserving biological signals. We also restrict the pipeline annotations only to protein coding genes and ribosomal RNA (rRNA) genes.

Below we describe each step of the pipeline in some detail. All datasets generated by the individual stages of the processing pipeline are made available as downloads. Appendix A lists the available files for each dataset.

### 2.2.1 Preprocessing

After upload, data is preprocessed by using SolexaQA [7] to trim low-quality regions from FASTQ data. Platform-specific approaches are used for 454 data submitted in FASTA format: reads more than two standard deviations away from the mean read length are discarded following [13]. All sequences submitted to the system are available, but discarded reads will not be analyzed further.

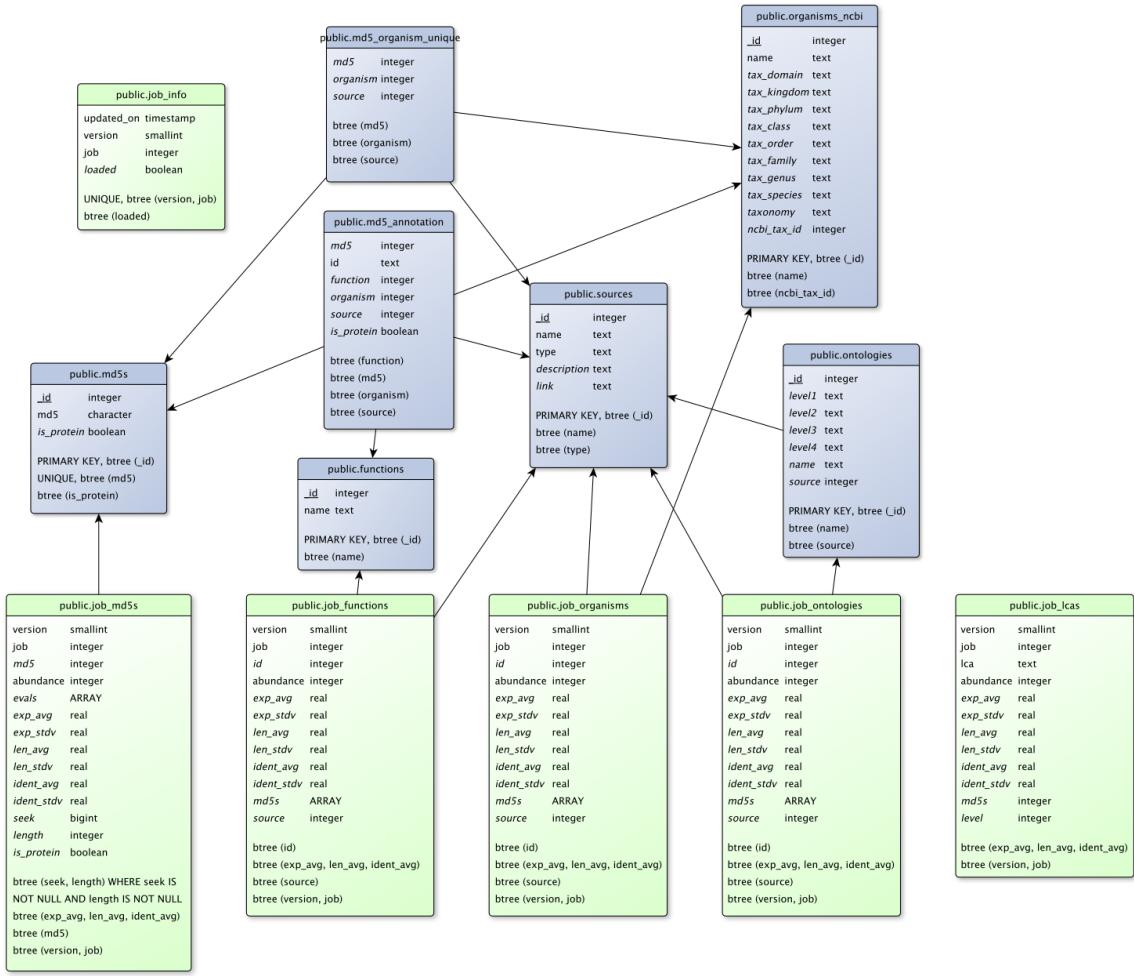


Figure 2.2: Analysis database schema: static objects (blue) and per metagenome (variable) objects (green).

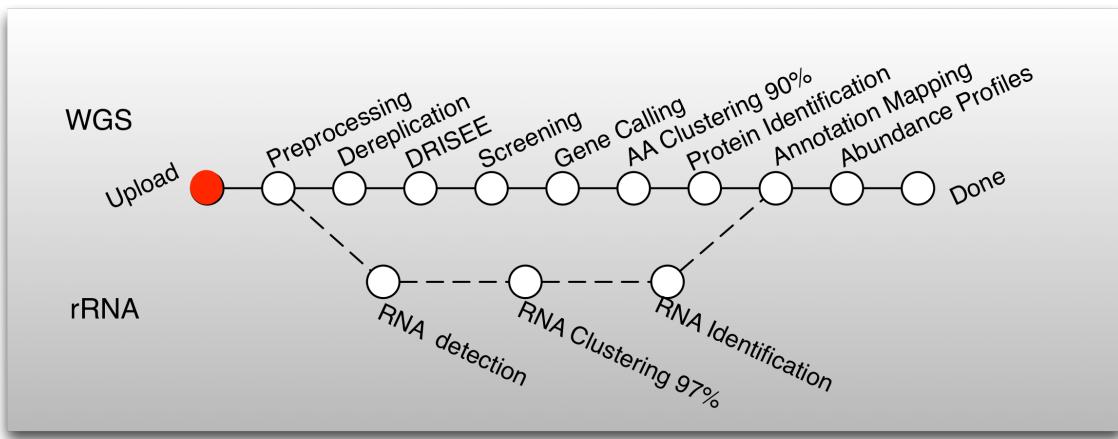


Figure 2.3: Details of the analysis pipeline for MG-RAST version 3

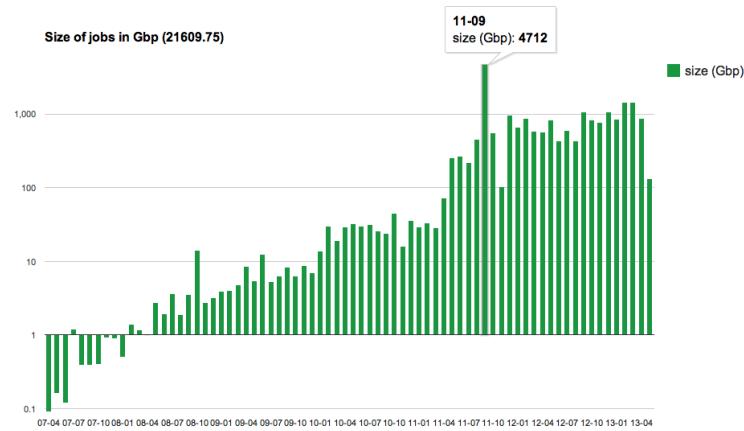


Figure 2.4: Sizes of MG-RAST jobs per month in gigabasepairs from 2007 to 2013.

## **2.2.2 DerePLICATION**

For shotgun metagenome and shotgun metatranscriptome datasets we perform a derePLICATION step. We use a simple k-mer approach to rapidly identify all 20 character prefix identical sequences. This step is required in order to remove artificial duplicate reads (ADRs) [12]. Instead of simply discarding the ADRs, we set them aside and use them later.

We note that derePLICATION is not suitable for amplicon datasets that are likely to share common prefixes.

## **2.2.3 DRISEE**

MG-RAST v3 uses DRISEE (Duplicate Read Inferred Sequencing Error Estimation) [17] to analyze the sets of ADRs [12] and determine the degree of variation among prefix-identical sequences derived from the same template. See below for details.

## **2.2.4 Screening**

The pipeline provides the option of removing reads that are near-exact matches to the genomes of a handful of model organisms, including fly, mouse, cow, and human. The screening stage uses Bowtie [20] (a fast, memory-efficient, short read aligner), and only reads that do not match the model organisms pass into the next stage of the annotation pipeline.

Note that this option will remove all reads similar to the human genome and render them inaccessible. This decision was made in order to avoid storing any human DNA on MG-RAST.

## **2.2.5 Gene calling**

The previous version of MG-RAST used similarity-based gene predictions, an approach that is significantly more expensive computationally than de novo gene prediction. After an in-depth investigation of tool performance [38], we have moved to a machine learning approach: FragGeneScan [32]. Using this approach, we can now predict coding regions in DNA sequences of 75 bp and longer. Our novel approach also enables the analysis of user-provided assembled contigs.

We note that FragGeneScan is trained for prokaryotes only. While it will identify proteins for eukaryotic sequences, the results should be viewed as more or less random.

## 2.2.6 AA clustering

MG-RAST builds clusters of proteins at the 90% identity level using the uclust [9] implementation in QIIME [5] preserving the relative abundances. These clusters greatly reduce the computational burden of comparing all pairs of short reads, while clustering at 90% identity preserves sufficient biological signals.

## 2.2.7 Protein identification

Once created, a representative (the longest sequence) for each cluster is subjected to similarity analysis. Instead of BLAST we use sBLAT, an implementation of the BLAT algorithm [18], which we parallelized using OpenMPI [11] for this work.

Once the similarities are computed, we present reconstructions of the species content of the sample based on the similarity results. We reconstruct the putative species composition of the sample by looking at the phylogenetic origin of the database sequences hit by the similarity searches.

## 2.2.8 Annotation mapping

Sequence similarity searches are computed against a protein database derived from the M5NR [41], which provides nonredundant integration of many databases: GenBank, [3], SEED [28], IMG [22], UniProt [21], KEGG [16], and eggNOGs [15]. Unlike MG-RAST v2, which relied solely on SEED, MG-RAST now supports many complementary views into the data with one similarity search, including different functional hierarchies: SEED Subsystems, IMG terms, COG [36], eggNOGs [15], and ontologies such as GO (Gene Ontology Consortium, 2013). Users can easily change views without recomputation. For example, COG and KEGG views can be displayed, which both show the relative abundances of histidine biosynthesis in a dataset of four cow rumen metagenomes.

## 2.2.9 Abundance profiles

Abundance profiles are the primary data product that MG-RAST’s user interface uses to display information on the datasets.

Using the abundance profiles, the MG-RAST system defers making a decision on when to transfer annotations. Since there is no well-defined threshold that is acceptable for all use cases, the abundance profiles contain all similarities and require their users to set cut-off values.

The threshold for annotation transfer can be set by using the following parameters: e-value, percent identity, and minimal alignment length.

The taxonomic profiles use the NCBI taxonomy. All taxonomic information is projected against this data. The functional profiles are available for data sources that provide hierarchical information. These currently comprise the following.

- SEED Subsystems

The SEED subsystems [28] represent an independent reannotation effort that powers, for example, the RAST [2] effort. Manual curation of subsystems makes them an extremely valuable data source.

Subsystems represent a four-level hierarchy:

1. Subsystem level 1 – highest level
2. Subsystem level 2 –
3. Subsystem level 3 – similar to a KEGG pathway
4. Subsystem level 4 – actual functional assignment to the feature in question

The page at <http://pubseed.theseed.org//SubsysEditor.cgi> allows browsing the subsystems.

- KEGG Orthologs

We use the KEGG [16] enzyme number hierarchy to implement a four-level hierarchy.

1. KEGG level 1 – first digit of the EC number (EC:X.\*.\*.\*)
2. KEGG level 2 – first two digits of the EC number (EC:X.Y.\*.\*)
3. KEGG level 3 – first three digits of the EC number (EC:X:Y:Z:.\*)
4. KEGG level 4 – entire four digits EC number

We note that KEGG data is no longer available for free download. We thus have to rely on using the latest freely downloadable version of the data.

The high-level KEGG categories are as follows.

1. Cellular Processes
2. Environmental Information Processing
3. Genetic Information Processing
4. Human Diseases
5. Metabolism

## 6. Organisational Systems

- COG and EGGNOG Categories

The high-level COG and EGGNOG categories are as follows.

1. Cellular Processes
2. Information Storage and Processing
3. Metabolism
4. Poorly Characterized

We note that for most metagenomes the coverage of each of the four namespaces is quite different. The “source hits distribution” (see Section 3.6.1.2) provides information on how many sequences per dataset were found for each database.

## 2.3 The rRNA pipeline

rRNA reads are identified by using a simple rRNA detection pipeline and are searched in a separate flow in the pipeline.

### 2.3.1 rRNA detection

An initial BLAT [18] search against a reduced RNA database efficiently identifies RNA.

The reduced database is a 90% identity clustered version of the SILVA database and is used merely to rapidly identify sequences with similarities to ribosomal RNA.

### 2.3.2 rRNA clustering

The rRNA-similar reads are then clustered at 97% identity, and the longest sequence is picked as the cluster representative.

### 2.3.3 rRNA identification

A BLAT similarity search for the longest cluster representative is performed against the M5rna database, integrating SILVA [29], Greengenes [8], and RDP [6].

## 2.4 Quality assessment

The MG-RAST pipeline offers a variety of summaries of technical aspects of the sequence quality to enable sequence data triage. These tools include DRISEE for estimating sequence error, summaries of the spectra of long kmers, and visualizations of the base caller output.

### 2.4.1 DRISEE

DRISEE [17] is a method for measuring sequencing error in whole-genome shotgun metagenomic sequence data that is independent of sequencing technology and overcomes many of the shortcomings of Phred. It utilizes ADRs to generate internal sequence standards from which an overall assessment of sequencing error in a sample is derived. DRISEE data are presented on the Overview page for each MG-RAST sample for which a DRISEE profile can be determined. Total DRISEE Error presents the overall DRISEE-based assessment of the sample as a percent error:

$\text{TotalDRISEEError} = \text{base\_errors}/\text{total\_bases} * 100$  where ‘base\_errors’ refers to the sum of DRISEE-detected errors and ‘total\_bases’ refers to the sum of all bases considered by DRISEE.

The current implementation of DRISEE is not suitable for amplicon sequencing data or other samples that may contain natural duplicated sequences (e.g., eukaryotic DNA where gene duplication and other forms of highly repetitive sequences are common) in high abundance.

### 2.4.2 Kmer profiles

Kmer digests are an annotation-independent method for describing sequence datasets that can support inferences about genome size and coverage. Here the Overview page presents several visualizations, evaluated at  $k=15$ : the kmer spectrum, kmer rank abundance, and ranked kmer consumed. All three graphs represent the same spectrum, but in different ways. The kmer spectrum plots the number of distinct kmers against kmer coverage; the kmer coverage is equivalent to number of observations of each kmer. The kmer rank abundance plots the relationship between kmer coverage and the kmer rank answering the question “What is the coverage of the nth most-abundant kmer? Ranked kmer consumed plots the largest fraction of the data explained by the nth most-abundant kmers only.

### 2.4.3 Nucleotide histograms

Nucleotide histograms are graphs showing the fraction of base pairs of each type (A, C, G, T, or ambiguous base N) at each position starting from the beginning of each read. Amplicon datasets

(see Figure 2.5) should show biased distributions of bases at each position, reflecting both conservation and variability in the recovered sequences:

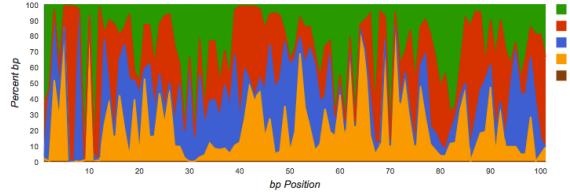


Figure 2.5: Nucleotide histogram with biased distributions typical for an amplicon dataset.

Shotgun datasets should have roughly equal proportions of A, T, G and C basecalls, independent of position in the read as shown in Figure 2.6.

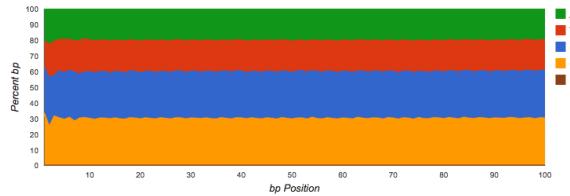


Figure 2.6: Nucleotide histogram showing ideal distributions typical for a shotgun metagenome.

Vertical bars at the beginning of the read indicate untrimmed (see Figure 2.7), contiguous barcodes. Gene calling via FragGeneScan [32] and RNA similarity searches are not impacted by the presence of barcodes. However, if a significant fraction of the reads is consumed by barcodes, it reduces the biological information contained in the reads.

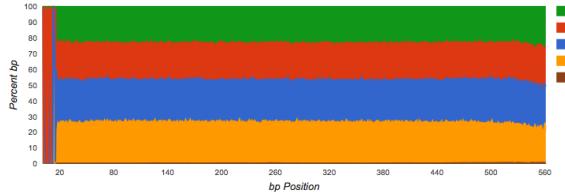


Figure 2.7: Nucleotide histogram with untrimmed barcodes.

If a shotgun dataset has clear patterns in the data, these indicate likely contamination with artificial sequences. The dataset shown in see Figure 2.8 had a large fraction of adapter dimers.

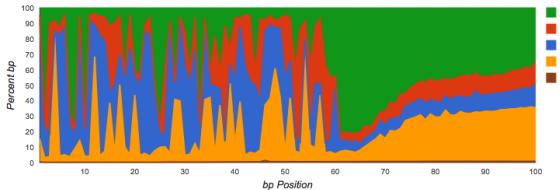


Figure 2.8: Nucleotide histogram with contamination.

## 2.5 Best hit, representative hit, and lowest common ancestor interpretation

MG-RAST searches the nonredundant M5NR and M5RNA databases in which each sequence is unique. These two databases are built from multiple sequence database sources, and the individual sequences may occur multiple times in different strains and species (and sometimes genera) with 100% identity. In these circumstances, choosing the “right” taxonomic information is not a straightforward process.

To optimally serve a number of different use cases, we have implemented three methods—best hit, representative hit, and lowest common ancestor—for end users to determine the number of hits (“occurrences of the input sequence in the database) reported for a given sequence in their dataset.

### 2.5.1 Best hit

The best hit classification reports the functional and taxonomic annotation of the best hit in the M5NR for each feature. In those cases where the similarity search yields multiple same-scoring hits for a feature, we do not choose any single “correct” label. For this reason we have decided to double count all annotations with identical match properties and leave determination of truth to our users. While this approach aims to inform about the functional and taxonomic potential of a microbial community by preserving all information, subsequent analysis can be biased because of a single feature having multiple annotations, leading to inflated hit counts. For users looking for a specific species or function in their results, the best hit classification is likely what is wanted.

### 2.5.2 Representative hit

The representative hit classification selects a single, unambiguous annotation for each feature. The annotation is based on the first hit in the homology search and the first annotation for that hit in our database. This approach makes counts additive across functional and taxonomic levels

and thus allows, for example, the comparison of functional and taxonomic profiles of different metagenomes.

### 2.5.3 Lowest Common Ancestor (LCA)

To avoid the problem of multiple taxonomic annotations for a single feature, we provide taxonomic annotations based on the widely used LCA method introduced by MEGAN [14]. In this method all hits are collected that have a bit score close to the bit score of the best hit. The taxonomic annotation of the feature is then determined by computing the LCA of all species in this set. This replaces all taxonomic annotations from ambiguous hits with a single higher-level annotation in the NCBI taxonomy tree.

### 2.5.4 Comparison of methods

Users should be aware that the number of hits might be inflated if the best hit filter is used or that a favorite species might be missing despite a similar sequence similarity result if the representative hit filter is used (in fact, even if a 100% identical match to a favorite species exists).

One way to consider both the best hit and representative hit is that they overinterpret the available evidence. With the LCA classifier function, on the other hand, any input sequence is classified only down to a trustworthy taxonomic level. While naively this seems to be the best function to choose in all cases because it classifies sequences to varying depths, the approach causes problems for downstream analysis tools that might rely on everything being classified to the same level.

## 2.6 Numbers of annotations vs. number of reads

The MG-RAST v3 annotation pipeline does not usually provide a single annotation for each submitted fragment of DNA. Steps in the pipeline map one read to multiple annotations and one annotation to multiple reads. These steps are a consequence of genome structure, pipeline engineering, and the character of the sequence databases that MG-RAST uses for annotation.

The first step that is not one-to-one is gene prediction. Long reads ( $>400$  bp) and contigs can contain pieces of two or more microbial genes; when the gene caller makes this prediction, the multiple predicted protein sequences (called fragments) are annotated separately.

An intermediate clustering step identifies sequences at 90% amino acid identity and performs one search for each cluster. Sequences that do not fall into clusters are searched separately. The “abundance column in the MG-RAST tables presents the estimate of the number of sequences that contain a given annotation, found by multiplying each selected database match (hit) by the number

of representatives in each cluster. The final step that is not one-to-one is the annotation process itself. Sequences can exist in the underlying data sources many times with different labels. When those sequences are the best hit similarity, we do not have a principled way to choose the “correct” label. For this reason we have decided to double count these annotations and leave determination of truth to our users. Note: Even when considering a single data source, double-counting can occur depending on the consistency of annotations. Also note: Hits refer to the number of unique database sequences that were found in the similarity search, **not** the number of reads. The hit count can be smaller than the number of reads because of clustering or larger due to double counting.

## 2.7 Metadata, Publishing, and Sharing

MG-RAST is both an analytical platform and a data integration system. To enable data reuse, for example for meta-analyses, we require that all data being made available to third parties contain at least minimal metadata. The MG-RAST team has decided to follow the minimal checklist approach used by the Genomics Standards Consortium (GSC) [10].

### 2.7.1 Metadata

While the GSC provides a GCDML [19] encoding, this XML-based format is more useful to programmers than to end users submitting data. We have therefore elected to use spreadsheets to transport metadata. Specifically we use MIxS (Minimum information about any (x) sequence (MIxS) and MIMARKS (Minimum Information about a MARKer gene Survey) to encode minimal metadata [44].

The metadata describe the origins of samples and provide details on the generation of the sequence data. While the GSC checklist aims at capturing a minimum of information, MG-RAST can handle additional metadata if supplied by the user. The metadata is stored in a simple key value format and is displayed on the Metagenome Overview page.

Once uploaded, the metadata spreadsheets are validated automatically, and users are informed of any problems.

The presence of metadata enables discovery by end users using contextual metadata. Users can perform searches such as “retrieve soil samples from the continental U.S. If the users have added additional metadata (domain specific extension), additional queries are enabled: for example, “restrict the results to soils with a specific pH.”

[» Back to the Metagenome Select](#)

 Job Information

Name - ID: 4447970.3 - CA\_05\_4.6  
Job: #1  
User: Pedro.Belda

share multiple metagenomes

To share the above job and its data with another user, please enter the email address of the user. Please note that you have to enter the email address which that person used to register at the MG-RAST service. The user will receive an email that notifies him how to access the data. Once you have granted the right to view one of your MG-RAST jobs to another user or group, the name will appear at the bottom of the page with the option to revoke it.

 Enter an email address

Enter an email address:

 This job is currently available to:

Figure 2.9: Dialogue showing the sharing mechanism. The mechanism requires a valid email address for the user with whom the data is to be shared. A list of users with access to the data is displayed at the bottom on the page.

## 2.7.2 Publishing

MG-RAST provides a mechanism to make data and analyses publicly accessible. Only the submitting user can make data public on MG-RAST. As stated above, metadata is mandatory for dataset publication.

## 2.7.3 Sharing

In addition to publishing, data and analysis can also be shared with specific users. To share data, users simply enter their email address via clicking sharing on the Overview page. The dialogue shown in Figure 2.9 will allow entering email addresses.

As shown in Figure 2.10, we tend to see dataset sharing between small groups of users.

## 2.7.4 Identifiers

MG-RAST automatically assigns a unique identifier to every dataset submitted. Upon completion of the automated pipeline, datasets can be viewed via the web interface by using the identifiers. The dataset identifiers are of the form integer\_prefix.revision. An example is 4440283.3.

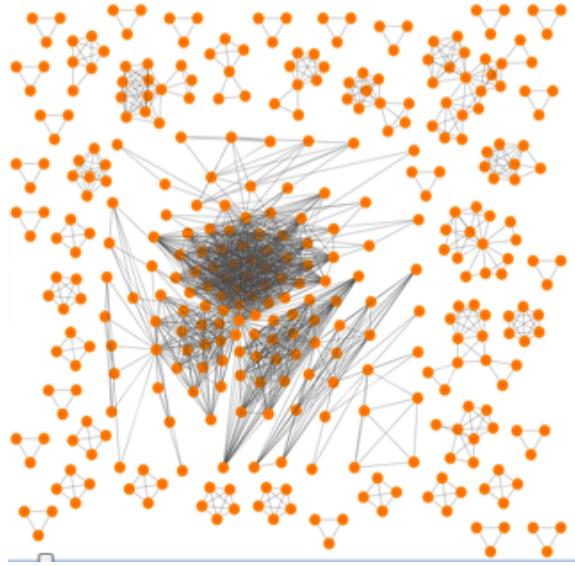


Figure 2.10: Data sets shared in MG-RAST by users (orange dots), shown as connecting edges.

In addition to datasets, MG-RAST supports projects (groups of datasets) that can be addressed with simple numerical project identifiers. An example is <http://metagenomics.anl.gov/linkin.cgi?project=128>.

### 2.7.5 Linking to MG-RAST

Because future versions of MG-RAST may change, we provide a link-in mechanism as a stable way of linking to MG-RAST. To link out to datasets in MG-RAST, users should use the `linkin.cgi`.

```
http://metagenomics.anl.gov/linkin.cgi?metagenome=
http://metagenomics.anl.gov/linkin.cgi?project=
```

Figure 2.11: Stable URLs provided by the `linkin.cgi` mechanism for linking to MG-RAST.

For example, for the metagenome ID 4440283.3 the URL is <http://metagenomics.anl.gov/linkin.cgi?metagenome=4440283.3>. This URL provides a stable method of linking to data that does not require the viewer to have an MG-RAST account. Note: One should not use the URL that is shown when browsing the site.

By default, a user's data is not visible to others; the user needs to explicitly grant permission for the data to be visible to anyone on the Internet, by making it public through the MG-RAST website.

# Chapter 3

## MG-RAST v3 web interface

The MG-RAST system provides a rich web user interface that covers all aspects of the metagenome analysis, from data upload to ordination analysis. The web interface can also be used for data discovery. Metagenomic datasets can be easily selected individually or on the basis of filters such as technology (including read length), quality, sample type, and keyword, with dynamic filtering of results based on similarity to known reference proteins or taxonomy. For example, a user might want to perform a search such as “phylum eq ‘actinobacteria’ and function in KEGG pathway Lysine Biosynthesis and sample in ‘Ocean’ ” to extract sets of reads matching the appropriate functions and taxa across metagenomes. The results can be displayed in familiar formats, including bar charts, trees that incorporate abundance information, heatmaps, or principal components analyses, or exported in tabular form. The raw or processed data can be recovered via download pages. Metabolic reconstructions based on mapping to KEGG pathways are also provided.

Sample selection is crucial for understanding large-scale patterns when multiple metagenomes are compared. Accordingly, MG-RAST supports MIxS and MIMARKS (Yilmaz, 2011) (as well as domain-specific plug-ins for specialized environments not extending the minimal GSC standards); several projects, including TerraGenome, HMP, TARA, and EMP, use these GSC standards, enabling standardized queries that integrate new samples into these massive datasets. An example query using the metadata browser, enabling the user to interrogate the existing pool of public datasets for a biome of interest (e.g., hot springs) and performing comparisons and a search for organisms encoding a specific gene function (e.g., beta-lactamase or aldo/keto reductase) is given in Figure 3.1.

One key aspect of the MG-RAST approach is the creation of smart data products enabling the user at the time of analysis to determine the best parameters for, for example, a comparison between samples. This is done without the need for recomputation of results.

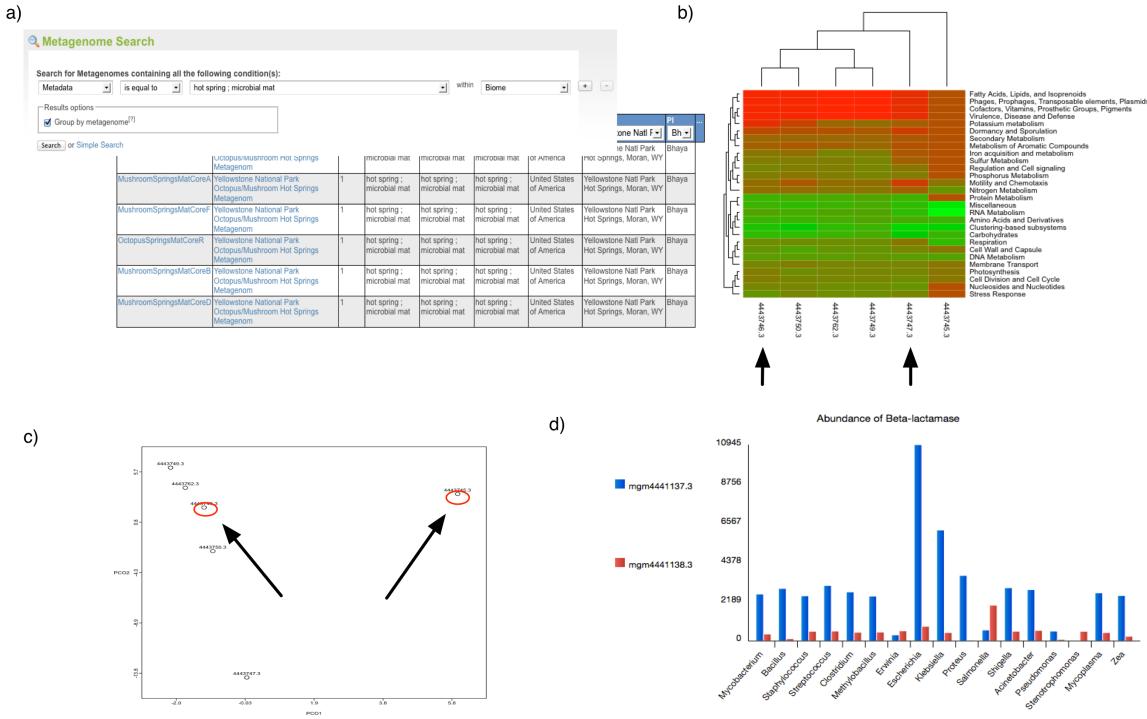


Figure 3.1: (a) Using the web interface for a search of metagenomes for microbial mats in hot-springs (GSC-MIMS-Keywords Biome=hotspring; microbial mat), we find 6 metagenomes (refs: 4443745.3, 4443746.3, 4443747.3, 4443749.3, 4443750.3, 4443762.3). (b) Initial comparison reveals some differences in protein functional class abundance (using SEED subsystems level 1). (c) From the PCoA plot using normalized counts of functional SEED Subsystem-based functional annotations (level 2) and Bray-Curtis as metric, we attempt to find differences between two similar datasets (MG-RAST-IDs: 4443749.3, 4443762.3). (d) Using exported tables with functional annotations and taxonomic mapping, we analyze the distribution of organisms observed to contain beta-lactamase and plot the abundance per species for two distinct samples.

## **3.1 Technical details**

This section briefly presents information about browsers and downloading.

### **3.1.1 Browser requirements**

The current web interface for MG-RAST is being developed for recent versions of Firefox. If you are using another browser, please understand that some or all of the web site will not function.

We realize that Firefox may not be your favorite (or institutionally prescribed) browser, but writing interactive web sites for many browsers is hard. While we are aiming to create a multi-browser version of the web interface, the current version is limited to Firefox.

### **3.1.2 Downloading figures**

Almost all figures and tables are downloadable into either graphics or spreadsheets. Please look for a download chart data link next to the graphic.

## **3.2 Sitemap for MG-RAST**

The MG-RAST web site (as shown in Figure 3.2) is complex and offers a lot of different options.

The site at <http://metagenomics.anl.gov> has five main pages and a home page, shown in blue in Figure 3.2.

- Download page – lists all publicly available data for download. The data is structured into projects.
- Browse page – allows interactive browsing of all datasets and is powered by metadata.
- Search page – allows identifier, taxonomy, and function-driven searches against all public data.
- Analysis page – enables in-depth analyses and comparisons between datasets.
- Upload page – allows users to provide their samples and metadata to MG-RAST. More details on uploading are below.
- Home (Metagene Overview) page – provides an overview for each individual dataset.

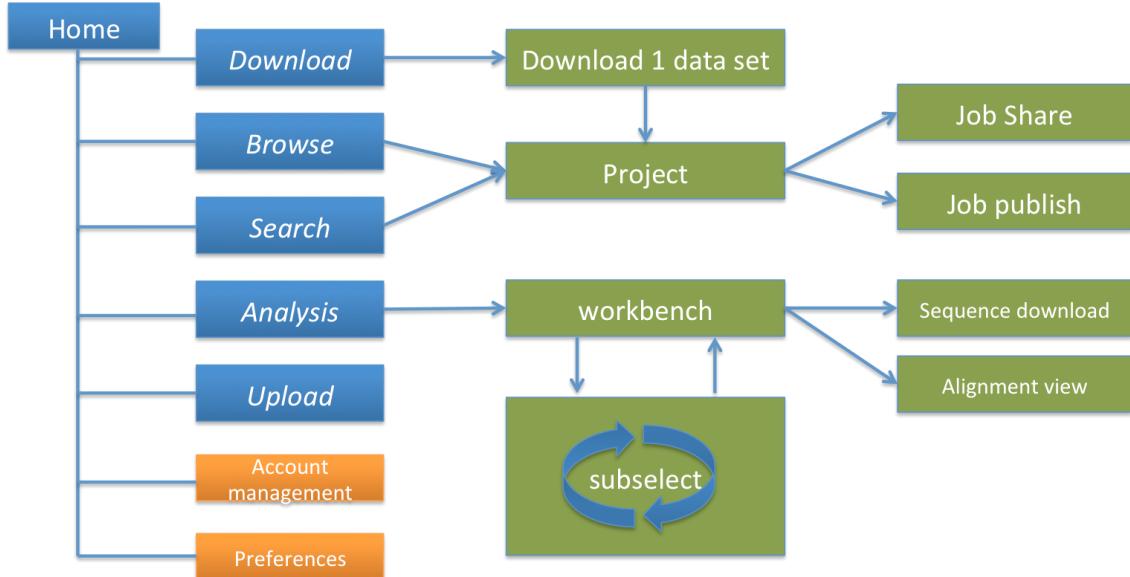


Figure 3.2: Sitemap for the MG-RAST version 3 web site. On the site map the main pages are shown in blue, management pages in orange. The green boxes represent pages that are not directly accessible from the home page.

We note that if you want to create links to the MG-RAST web site, you should use the *linkin* mechanism instead of linking to any web page directly. All pages intended for users to create external links provide the linkin feature.

### 3.3 Upload page

Data and metadata can be uploaded in the form of spreadsheets along with the sequence data by using both the ftp and the http protocols. The web uploader will automatically split larger files and allow parallel uploads.

MG-RAST supports datasets that are augmented with rich metadata using the standards and technology developed by the GSC. Each user has a temporary storage location inside the MG-RAST system. This inbox provides temporary storage for data and metadata to be submitted to the system. Using the inbox, users can extract compressed files, convert a number of vendor specific formats to MG-RAST submission-compliant formats, and obtain an MD5 checksum for verifying that transmission to MG-RAST has not altered the data.

The web uploader has been optimized for large datasets of over 100 gigabasepairs, often re-

sulting in file sizes in excess of 150 GB.

## 3.4 Browse page – Metadata-enabled data discovery

The Browse page lists all datasets visible to the user.<sup>1</sup> This page also provides an overview of the nonpublic datasets submitted by the user or shared with users.

Figure 3.3 shows the interactive metagenome browse table, which provides an interactive graphical means to discover data based on technical data (e.g., sequence type or dataset size) or metadata (e.g., location or biome).

## 3.5 Project page

Shown in Figure 3.4, the project page provides a list of datasets and metadata for a project. The table at the bottom of the Project page provides access to the individual metagenomes by clicking on the identifiers in the first column. In addition, the final column provides downloads for metadata, submitted data, and the analysis results via the three labelled arrows.

For the dataset owners the Project page provides an editing capability using a number of menu entries at the top of the page. Figure 3.5 shows the available options.

- Share Project – make the data in this project available to third parties via sending them access tokens.
- Add Jobs – add additional datasets to this project.
- Edit Project Data – edit the contents of this page.
- Upload Info – upload information to be displayed on this page.
- Upload MetaData – upload a metadata spreadsheet for the project.
- Export MetaData<sup>2</sup> – export the metadata spreadsheet for this project.

---

<sup>1</sup>Datasets in MG-RAST are private by default, but the submitting user has the option of sharing datasets with specific users or making datasets public.

<sup>2</sup>This option is available to non-dataset owners.

## ALL METAGENOMES

group by project

### Current table counts

public (12535) private (0) shared (0)

metagenomes	projects	biomes	features	materials	altitudes	depths	locations	ph's	countries	temperatures	pi's
12535	373	102	101	102	115	296	477	116	71	1005	144

[clear table filters](#)

[add selected to a collection](#)

display  items per page

displaying 1 - 10 of 12535

[next»](#) [last»](#)

project ▲▼	name ▲▼	bps ▲▼	sequences ▲▼	biome	feature	material	sequencing type ▲▼	select	...
		< □	< □	all ▼	all ▼	all ▼	all ▼	□	all
The oral metagenome in health and disease	<a href="#">CA_05_4.6</a>	27669924	70503	human-associated habitat	human-associated habitat	human-associated habitat	WGS	public	<input type="checkbox"/>
The oral metagenome in health and disease	<a href="#">CA_06_1.6</a>	37519874	97722	human-associated habitat	human-associated habitat	human-associated habitat	WGS	public	<input type="checkbox"/>
cDNA - Plymouth Marine Lab Coastal Waters project	<a href="#">1-19-DNA-fix</a>	59316369	344216	marine habitat	marine habitat	marine habitat	WGS	public	<input type="checkbox"/>
cDNA - Plymouth Marine Lab Coastal Waters project	<a href="#">6-19-DNA-fix</a>	68187679	304020	marine habitat	marine habitat	marine habitat	WGS	public	<input type="checkbox"/>
Northern Line Islands	<a href="#">FannLIMic20050811</a>	30909241	290844	marine habitat	marine habitat	marine habitat	WGS	public	<input type="checkbox"/>
Northern Line Islands	<a href="#">FannLIVir20050811</a>	39607682	380355	marine habitat	marine habitat	marine habitat	WGS	public	<input type="checkbox"/>
Soudan Mine Metagenome	<a href="#">RedSoudMineMic20050331</a>	35439683	334386	mine drainage	mine drainage	mine drainage	WGS	public	<input type="checkbox"/>
Soudan Mine Metagenome	<a href="#">BlackSoudMineMic20050331</a>	38502057	388627	mine drainage	mine drainage	mine drainage	WGS	public	<input type="checkbox"/>
Chicken Cecum Microbiome	<a href="#">Chicken Cecum A</a>	32296796	310801	animal-associated habitat	animal-associated habitat	animal-associated habitat	WGS	public	<input type="checkbox"/>
Chicken Cecum Microbiome	<a href="#">Chicken_Cecum_B</a>	26378422	254712	animal-associated habitat	animal-associated habitat	animal-associated habitat	WGS	public	<input type="checkbox"/>

displaying 1 - 10 of 12535

[next»](#) [last»](#)

Figure 3.3: Browse page, enabling sorting and data search. Users can select the metadata they wish to view and search. Some of the metadata is hidden by default and can be viewed by clicking on the last column header on the right side of the table and selecting the desired columns; this can also be used to hide unwanted columns.

**THE ORAL METAGENOME IN HEALTH AND DISEASE (ID 128)** [metagenomes](#) [project metadata](#)

Visibility Public  
Static Link <http://metagenomics.anl.gov/linkin.cgi?project=128>

Share Project | Add Jobs | Edit Project Data | Upload Info | Upload MetaData | Export MetaData

#### DESCRIPTION

The oral cavity of humans is inhabited by hundreds of bacterial species and some of them have a key role in the development of oral diseases, mainly dental caries and periodontitis. We describe for the first time the metagenome of the human oral cavity under health and diseased conditions, with a focus on supragingival dental plaque and cavities. Direct pyrosequencing of eight samples with different oral-health status produced 1 Gbp of sequence without the biases imposed by PCR or cloning. These data show that cavities are not dominated by Streptococcus mutans (the species originally identified as the etiological agent of dental caries) but are in fact a complex community formed by tens of bacterial species, in agreement with the view that caries is a polymicrobial disease. The analysis of the reads indicated that the oral cavity is functionally a different environment from the gut, with many functional categories enriched in one of the two environments and depleted in the other. Individuals who had never suffered from dental caries showed an over-representation of several functional categories, like genes for antimicrobial peptides and quorum sensing. In addition, they did not have mutants streptococci but displayed high recruitment of other species. Several isolates belonging to these dominant bacteria in healthy individuals were cultured and shown to inhibit the growth of cariogenic bacteria, suggesting the use of these commensal bacterial strains as probiotics to promote oral health and prevent dental caries.

#### FUNDING SOURCE

Spanish MICINN: SAF2009-13032-C02-02 from the I+D program, BIO2008-03419-E from the EXPLORA program and MICROGEN CSD2009-00006 from the Consolider- Ingenio program.

#### CONTACT

Administrative  
**Alex Mira(CSISP)**  
Avda. Cataluña, 21. Valencia, Spain

Technical  
**Pedro Belda-Ferre(Center for Advanced Research in Public Health, Department of Genomics and Health)**  
Avda. Cataluña, 21 ; 46020 ; Valencia ; Comunidad Valenciana, Spain

#### ADDITIONAL DATA

administrative-contact_PI_lastname	Mira
project-description_Internal_project_ID	The oral metagenome in health and disease
administrative-contact_PI_email	mira_ale@gva.es
administrative-contact_PI_organization	Center for Advanced Research in Public Health, Department of Genomics and Health
administrative-contact_PI_organization_country	Spain
administrative-contact_PI_organization_address	Avda. Cataluña, 21 ; 46020 ; Valencia ; Comunidad Valenciana
administrative-contact_PI_organization_url	www.csisp.gva.es/web/csisp
administrative-contact_PI_firstname	Alex

#### METAGENOMES

There are 8 metagenomes in this project.

Export Jobs Table												
MG-RAST ID	Metagenome Name ▲▼	bp Count ▲▼	Sequence Count ▲▼	Biome ▲▼	Feature ▲▼	Material ▲▼	Location ▲▼	Country ▲▼	Coordinates ▲▼	Sequence Type ▲▼	Sequence Method ▲▼	Download
4447943.3	CA_04P	142,374,233	339,503	human-associated habitat	human-associated habitat	human-associated habitat	Valencia	Spain	39.481448, 0.353066	WGS	454	<a href="#">metadata</a> <a href="#">submitted</a> <a href="#">analysis</a>
4447192.3	NOCA_01P	77,538,485	204,218	human-associated habitat	human-associated habitat	human-associated habitat	Valencia	Spain	39.481448, 0.353066	WGS	454	<a href="#">metadata</a> <a href="#">submitted</a> <a href="#">analysis</a>
4447103.3	CA1_01P	203,711,161	464,594	human-associated habitat	human-associated habitat	human-associated habitat	Valencia	Spain	39.481448, 0.353066	WGS	454	<a href="#">metadata</a> <a href="#">submitted</a> <a href="#">analysis</a>
4447102.3	NOCA_03P	100,125,112	244,881	human-	human-	human-	Valencia	Spain	39.481448	WGS	454	<a href="#">metadata</a> <a href="#">submitted</a> <a href="#">analysis</a>

Figure 3.4: Project page, providing a summary of all data in the project and an interface for downloads.

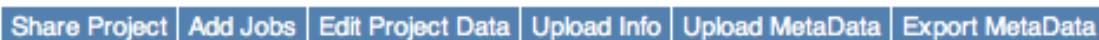


Figure 3.5: Buttons displayed by Project page to dataset owner.

The screenshot shows the top portion of the MG-RAST Metagenome Overview page. At the top left is the title "Metagenome Overview". Below it, the MG-RAST ID is listed as "4447970.3" with options to "Download", "Analyze", and "Search". To the right of the ID are links for "NCBI Project ID -", "GOLD ID -", and "PubMed ID 21716308". Below these links, there is a table with the following data:

Metagenome Name	CA_05_4.6
PI	Alex Mira
Organization	CSISP
Visibility	Public
Static Link	<a href="http://metagenomics.anl.gov/linkin.cgi?metagenome=4447970.3">http://metagenomics.anl.gov/linkin.cgi?metagenome=4447970.3</a>

Figure 3.6: Top of the metagenome Overview page.

## 3.6 Overview page

MG-RAST automatically creates an individual summary page for each dataset. This metagenome overview page provides a summary of the annotations for a single dataset. The page is made available by the automated pipeline once the computation is finished. This page is a good starting point for looking at a particular dataset. It provides a significant amount of information on technical details and biological content.

The page is intended as a single point of reference for metadata, quality, and data. It also provides an initial overview of the analysis results for individual datasets with default parameters. Further analyses are available on the Analysis page.

### 3.6.1 Technical part of the Overview page – Details on sequencing and analysis

The Overview page provides the MG-RAST ID for a data set, a unique identifier that is usable as accession number for publications. Additional information such as the name of the submitting PI and organization and a user-provided metagenome name are displayed at the top of the page as well. A static URL for linking to the system that will be stable across changes to the MG-RAST web interface is provided as additional information (Figure 3.6).

We provide an automatically generated paragraph of text describing the submitted data and the results computed by the pipeline. By means of the project information we display additional information provided by the data submitters at the time of submission or later.

One of the key diagrams in MG-RAST is the sequence breakdown pie chart (Figure 3.7) classifying the submitted sequences submitted into several categories according to their annotation status. As detailed in the description of the MG-RAST v3 pipeline above, the features annotated in MG-RAST are protein coding genes and ribosomal proteins.

We note that for performance reasons no other sequence features are annotated by the default pipeline. Other feature types such as small RNAs or regulatory motifs (e.g., CRISPRS [4]) not only will require significantly higher computational resources but also are frequently not supported by

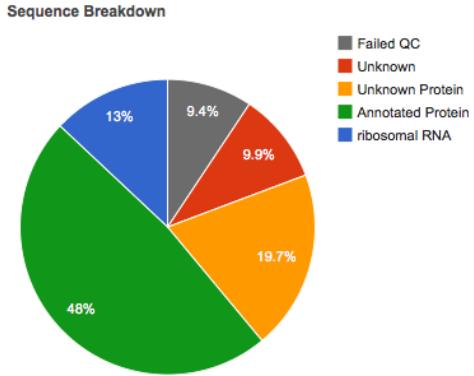


Figure 3.7: Sequences to the pipeline are classified into one of five categories: grey = failed the QC, red = unknown sequences, yellow = unknown function but protein coding, green = protein coding with known function, and blue = ribosomal RNA. For this example over 50% of sequences were either filtered by QC or failed to be recognized as either protein coding or ribosomal.

the unassembled short reads that constitute the vast majority of todays metagenomic data in MG-RAST. The quality of the sequence data coming from next-generation instruments requires careful design of experiments, lest the sensitivity of the methods is greater than the signal-to-noise ratio the data supports.

The overview page also provides metadata for each dataset to the extent that such information has been made available. Metadata enables other researchers to discover datasets and compare annotations. MG-RAST requires standard metadata for data sharing and data publication. This is implemented using the standards developed by the Genomics Standards Consortium. Figure 3.8 shows the metadata summary for a dataset.

All metadata stored for a specific dataset is available in MG-RAST; we merely display a standardized subset in this table. A link at the bottom of the table (“More Metadata”) provides access to a table with the complete metadata. This enables users to provide extended metadata going beyond the GSC minimal standards. A mechanism to provide community consensus extensions to the minimal checklists and the environmental packages are explicitly encouraged but not required when using MG-RAST.

### 3.6.1.1 Metagenome quality control

The analysis flowchart and analysis statistics provide an overview of the number of sequences at each stage in the pipeline (Figure ??). The text block next to the analysis flowchart presents the numbers next to their definitions.

GSC MIxS INFO	
<i>Investigation Type</i>	metagenome
<i>Project Name</i>	The oral metagenome in health and disease
<i>Latitude and Longitude</i>	39.481448, 0.353066
<i>Country and/or Sea, Location</i>	Spain Valencia
<i>Collection Date</i>	2010-03-01 10:00:00 UTC
<i>Environment (Biome)</i>	human-associated habitat
<i>Environment (Feature)</i>	human-associated habitat
<i>Environment (Material)</i>	human-associated habitat
<i>Environmental Package</i>	human-oral
<i>Sequencing Method</i>	454
<a href="#">More Metadata</a>	

Figure 3.8: Information from the GSC MIxS checklist providing minimal metadata on the sample.

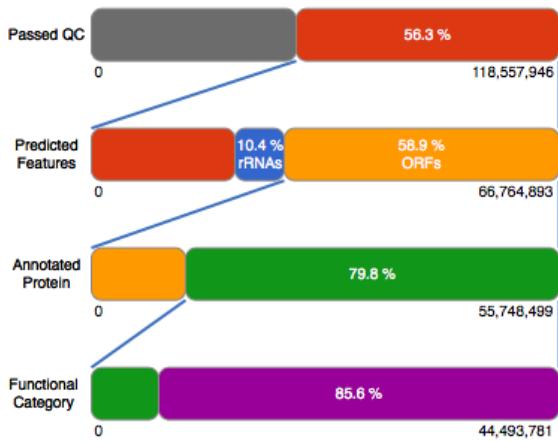


Figure 3.9: Analysis flowchart providing an overview of the fractions of sequences surviving the various steps of the automated analysis. In this case about 20% of sequences were filtered during quality control. From the remaining 37,122,128 sequences, 53.5% were predicted to be protein coding, 5.5% hit ribosomal RNA. From the predicted proteins, 76.8% could be annotated with a putative protein function. Of 32 million annotated proteins, 24 million have been assigned to a functional classification (SEED, COG, EggNOG, KEEG), representing 84% of the reads.

labelfig:analysis-flowchart

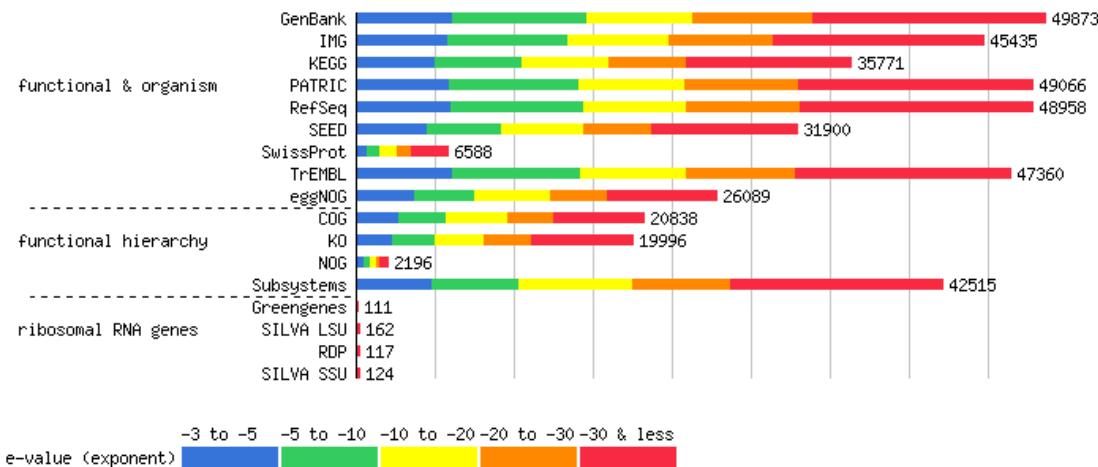


Figure 3.10: Graph showing the number of features in this dataset annotated by the different databases. The bars representing annotated reads are colored by e-value range. Different databases have different numbers of hits but can also have different types of annotation data.

### 3.6.1.2 Source hits distribution

The source hits distribution shows what percentage of the predicted protein features could be annotated with similarity to a protein of known function and which database those functions were from. In addition, ribosomal RNA genes are mapped to the rRNA databases.

Figure 3.10 shows the number of features in this dataset that were annotated by the different databases. These include protein databases, protein databases with functional hierarchy information, and ribosomal RNA databases.

In addition this display will print the number of records in the M5NR protein database and in the M5RNA ribosomal databases.

### 3.6.1.3 Other statistics

MG-RAST also provides a quick link to other statistics.

For example, the Analysis Statistics and Analysis Flowchart provide sequence statistics for the main steps in the pipeline from raw data to annotation, describing the transformation of the data between steps. Sequence length and GC histograms display the distribution before and after quality control steps. Metadata is presented in a searchable table that contains contextual metadata describing sample location, acquisition, library construction, and sequencing using GSC compliant metadata. All metadata can be downloaded from the table.

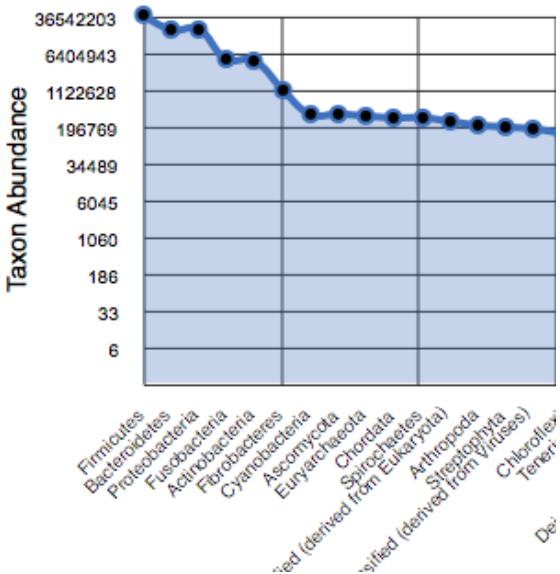


Figure 3.11: Sample rank abundance plot by phylum.

In addition, a breakdown of BLAT hits is provided per data source (e.g., hits to RefSeq [30], UniProt [21], or SEED [28]).

### 3.6.2 Biological part of the Overview page

The taxonomic hit distribution display divides taxonomic units into a series of pie charts of all the annotations grouped at various taxonomic ranks (domain, phylum, class, order, family, genus). The subsets are selectable for downstream analysis; this also enables downloads of subsets of reads, for example, those hitting a specific taxonomic unit.

#### 3.6.2.1 Rank abundance

The rank abundance plot (Figure 3.11) provides a rank-ordered list of taxonomic units at a user-defined taxonomic level, ordered by their abundance in the annotations.

#### 3.6.2.2 Rarefaction

The rarefaction curve of annotated species richness is a plot (see Figure 3.12 of the total number of distinct species annotations as a function of the number of sequences sampled. The slope of the right-hand part of the curve is related to the fraction of sampled species that are rare. On the left, a steep slope indicates that a large fraction of the species diversity remains to be discovered. If the

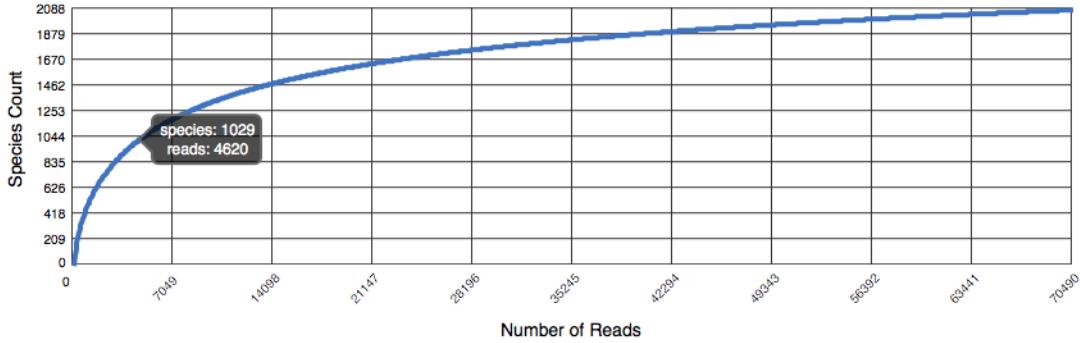


Figure 3.12: Rarefaction plot showing a curve of annotated species richness. This curve is a plot of the total number of distinct species annotations as a function of the number of sequences sampled.

curve becomes flatter to the right, a reasonable number of individuals is sampled: more intensive sampling is likely to yield only few additional species. Sampling curves generally rise quickly at first and then level off toward an asymptote as fewer new species are found per unit of individuals collected.

The rarefaction curve is derived from the protein taxonomic annotations and is subject to problems stemming from technical artifacts. These artifacts can be similar to the ones affecting amplicon sequencing [31], but the process of inferring species from protein similarities may introduce additional uncertainty.

### 3.6.2.3 Alpha diversity

In this section we display an estimate of the alpha diversity based on the taxonomic annotations for the predicted proteins. The alpha diversity is presented in context of other metagenomes in the same project (see Figure 3.13).

The alpha diversity estimate is a single number that summarizes the distribution of species-level annotations in a dataset. The Shannon diversity index is an abundance-weighted average of the logarithm of the relative abundances of annotated species.

We compute the species richness as the antilog of the Shannon diversity:

$$\text{Richness} = 10^{-\sum_i p_i \log(p_i)}$$

where  $p_i$  are the proportions of annotations in each of the species categories. Shannon species richness has units of the effective number of species. Each  $p$  is a ratio of the number of annotations

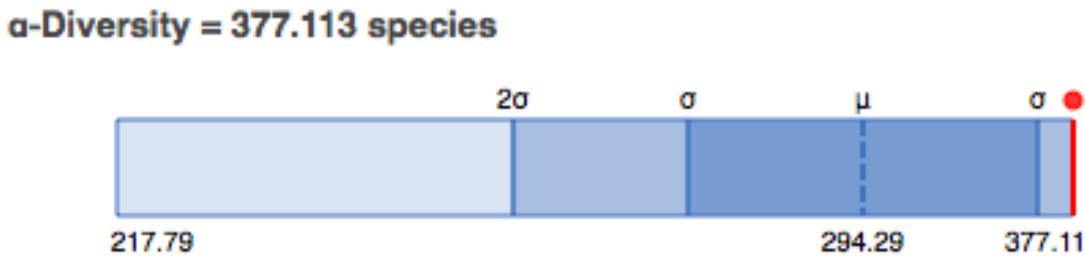


Figure 3.13: Alpha diversity plot showing the range of  $\alpha$ -diversity values in the project the data set belongs to. The min, max, and mean values are shown, with the standard deviation ranges ( $\sigma$  and  $2\sigma$ ) in different shades. The  $\alpha$ -diversity of this metagenome is shown in red.

for each species to the total number of annotations. The species-level annotations are from all the annotation source databases used by MG-RAST. The table of species and number of observations used to calculate this diversity estimate can be downloaded under “download source data” on the Overview page.

#### 3.6.2.4 Functional categories

This section contains four pie charts providing a breakdown of the functional categories for KEGG [16], COG [36], SEED Subsystems [28], and EggNOGs [15]. Clicking on the individual pie chart slices will save the respective sequences to the workbench. The relative abundance of sequences per functional category can be downloaded as a spreadsheet, and users can browse the functional breakdowns via the Krona tool [27] integrated in the page.

A more detailed functional analysis, allowing the user to manipulate parameters for sequence similarity matches, is available from the Analysis page.

## 3.7 Download page

The Download page provides all publicly available datasets for download. Three types of data are available for download.

- Metadata – data describing data in GSC-compliant format.
- Submitted data – the original user submission.

[Subsystems](#) [Download chart data](#)  
 has 42,515 predicted functions  
 79.8% of predicted proteins  
 104.4% of annotated proteins  
[View Subsystems interactive chart](#)

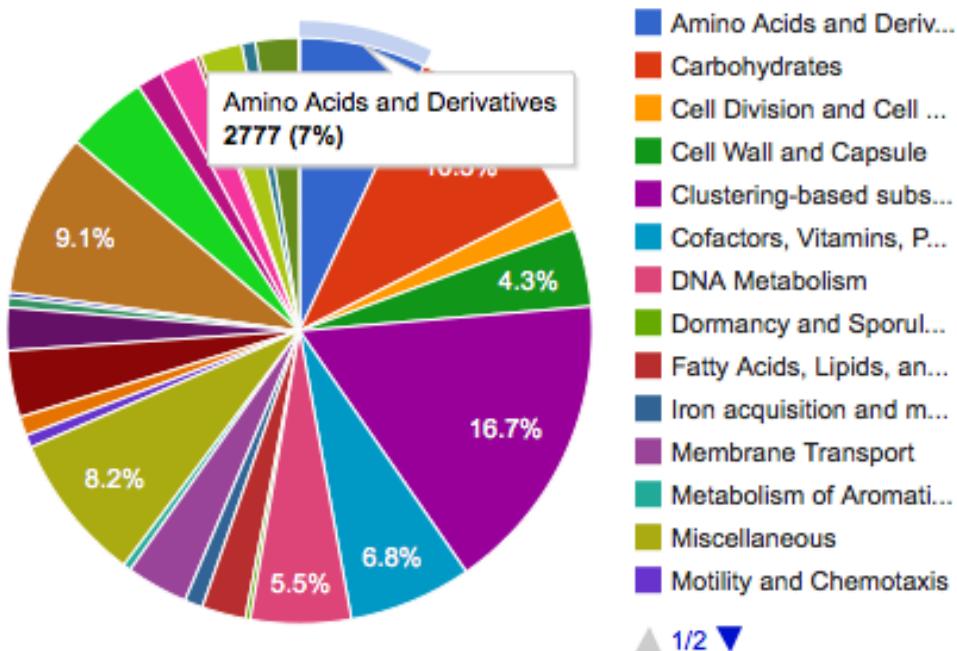


Figure 3.14: The Subsystems function piechart, showing reads classified into SEED subsystem level-one functions. In contrast to the COG, EGGNOG, and KEGG classification schemes, there are over 20 top-level subsystem categories, creating a more highly resolved “fingerprint” for the metagenome.

Found 11 metagenomes containing 3 metadata matches.

[map metagenomes](#) [create collection](#) [download results table](#) [clear all filters](#)

display  items per page  
displaying 1 - 11 of 11

Metagenome ▾	MG-RAST ID	Project ▾	Public	Match Counts ▾	Biome ▾	Feature ▾	Material ▾	Country ▾	Location ▾	PI ▾	...
			all	al	all	all	all	all	all	al	
NOCA_03P	4447102.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
NOCA_01P	4447192.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
February coral	4445755.3	Project for: February coral	1	1	animal-associated habitat	animal-associated habitat	animal-associated	Australia	Nelly Bay Magnetic Island		
CA_06P	4447903.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
CA_06_1.6	4447971.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
CA_05_4.6	4447970.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
CA_04P	4447943.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
CA1_02P	4447101.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
CA1_01P	4447103.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated	Spain	Valencia	Mira	
October coral	4445829.3	Project for: October coral	0	2	animal-associated habitat	animal-associated habitat	animal-associated	Australia	Nelly Bay Magnetic Island		
October coral	4445756.3	Project for: October coral	1	2	organism-associated habitat	organism-associated habitat	organism-associated	Australia	Nelly Bay Magnetic Island		

displaying 1 - 11 of 11

Figure 3.15: Searching for “oral health” returns 11 data sets for two projects.

- Analysis results – results of running the MG-RAST pipeline. The list includes all intermediate data products and is intended to serve as a basis for further analysis outside the MG-RAST pipeline.

Details on the individual files are in Appendix A.

## 3.8 Search Page

In addition to the Browse page, one can use the Search page to find datasets in MG-RAST. The basic function of the Search page is to find data sets that (1) contain a search string in the metadata (dataset name, project name, project description, GSC metadata), (2) contain specific functions (e.g., SEED functional roles, SEED subsystems, or GenBank annotations), or (3) contain specific organisms. The default search uses all three kinds of data.

In addition to a Google-like search that searches all data fields, we provide specialized searches in one of the three data types. Figure 3.15 shows the result of a metadata search for “oral health.” Figure 3.16 shows the results from Figure 3.15 after sorting by metagenome ID.

Found 11 metagenomes containing 3 metadata matches.

<a href="#">Metagenome ▾▼</a>	<a href="#">MG-RAST ID</a>	<a href="#">Project ▾▼</a>	<a href="#">Public</a>	<a href="#">Match Counts ▾▼</a>	<a href="#">Biome ▾▼</a>	<a href="#">Feature ▾▼</a>	<a href="#">Material ▾▼</a>	<a href="#">Country ▾▼</a>	<a href="#">Location ▾▼</a>	<a href="#">PI ▾▼</a>	<a href="#">...</a>
October coral	4445829.3	Project for: October coral	0	2	animal-associated habitat	animal-associated habitat	animal-associated habitat	Australia	Nelly Bay Magnetic Island		
October coral	4445756.3	Project for: October coral	1	2	organism-associated habitat	organism-associated habitat	organism-associated habitat	Australia	Nelly Bay Magnetic Island		
NOCA_03P	4447102.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
NOCA_01P	4447192.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
February coral	4445755.3	Project for: February coral	1	1	animal-associated habitat	animal-associated habitat	animal-associated Australia	Australia	Nelly Bay Magnetic Island		
CA_06P	4447903.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
CA_06_1.6	4447971.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
CA_05_4.6	4447970.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
CA_04P	4447943.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
CA1_02P	4447101.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	
CA1_01P	4447103.3	The oral metagenome in health and disease	1	1	human-associated habitat	human-associated habitat	human-associated Spain	Spain	Valencia	Mira	

displaying 1 - 11 of 11

Figure 3.16: Search results from the previous search sorted by projects.

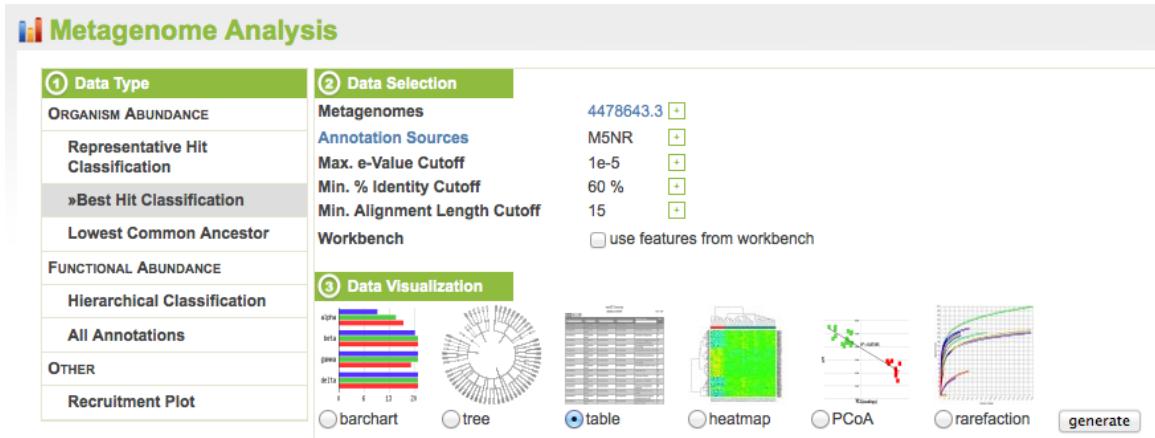


Figure 3.17: Three-step process in using the Analysis page: (1) select a profile and hit (see text) type; (2) select a list of metagenomes and set annotation source and similarity parameters; (3) choose a comparison.

## 3.9 Analysis page

The MG-RAST annotation pipeline produces a set of annotations for each sample; these annotations can be interpreted as functional or taxonomic abundance profiles. The analysis page can be used to view these profiles for a single metagenome or to compare profiles from multiple metagenomes using various visualizations (e.g., heatmap) and statistics (e.g., PCoA, normalization).

The page is divided into three parts following a typical workflow (Figure 3.17).

### 1. Data type

Selection of an MG-RAST analysis scheme, that is, selection of a particular taxonomic or functional abundance profile mapping. For taxonomic annotations, since there is not always a unique mapping from hit to annotation, we provide three interpretations: best hit, representative hit, and lowest common ancestor, as explained in Section 2.5.

We note that when choosing the LCA annotations, not all downstream tools are available. The reason is fact that for the LCA annotations not all sequences will be annotated to the same level: classifications are returned on different taxonomic levels.

Functional annotations can be grouped into mappings to functional hierarchies or can be displayed without a hierarchy. In addition, the recruitment plot displays the recruitment of protein sequences against a reference genome.

Each selected data type has data selections and data visualizations specific for it.

## 2. Data selection

Selection of sample and parameters. This dialog allows the selection of multiple metagenomes that can be compared individually or selected and compared as groups. Comparison is always relative to the annotation source, e-value, and percent identity cutoffs selectable in this section. In addition to the metagenomes available in MG-RAST, sets of sequences previously saved in the workbench can be selected for visualization.

## 3. Data visualization

Data visualization and comparison. Depending on the selected profile type, the profiles for the metagenomes can be visualized and compared by using barcharts, trees, spreadsheet-like tables, heatmaps, PCoA, rarefaction plots, circular recruitment plot, and KEGG maps.

The data selection dialogue provides access to data sets in four ways. The four categories can be selected from a pulldown menu.

- **private** data – list of private or shared data sets for browsing under available metagenomes.
- **collections** – defined sets of metagenomes grouped for easier analysis. This is the recommended way of working with the analysis page.
- **projects** – global groups of datasets grouped by the submitting user. The project name will be displayed.
- **public** data – display of all public datasets.

When using collections or projects, data can also be grouped into one set per collection or project and subsequently compared or added.

Once a category is selected, the data browser underneath “available metagenomes” will display data of the selected category. The text field under “available metagenomes” displays the available datasets or group identifiers.

The use of MG-RAST identifiers (e.g., 4447971.3) is possible in the text field underneath “available metagenomes.”

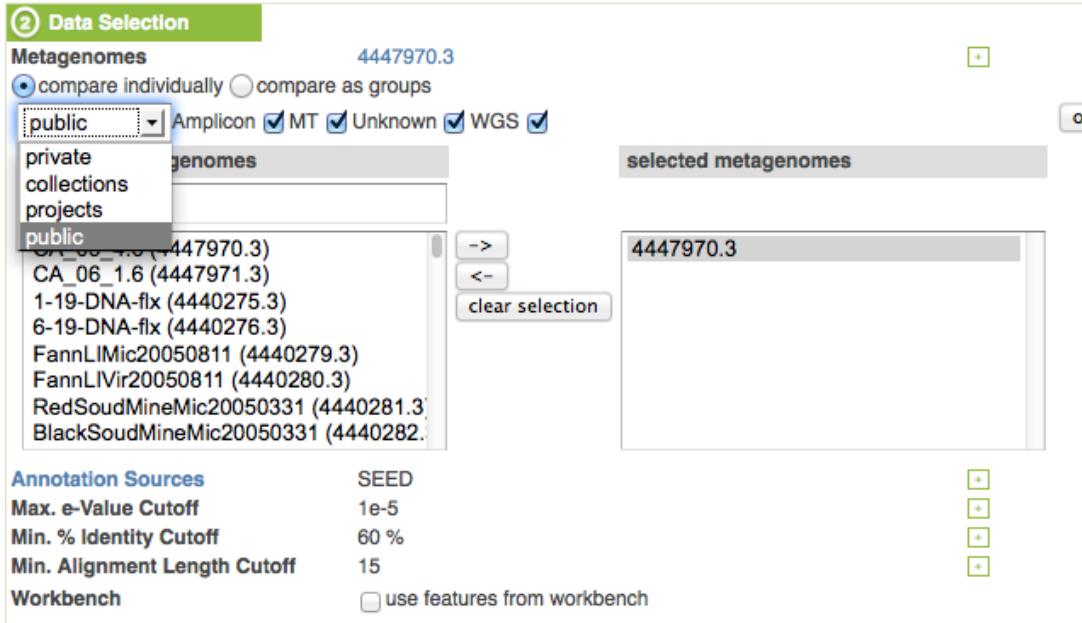


Figure 3.18: View of the data selection dialogue, with the list of four data categories expanded.

### 3.9.1 Normalization

Normalization refers to a transformation that attempts to reshape an underlying distribution. A large number of biological variables exhibit a log-normal distribution, meaning that when the data is transformed with a log transformation, the values exhibit a normal distribution. Log transformation of the counts data makes a normalized data product that is more likely to satisfy the assumptions of additional downstream tests such as ANOVA or t-tests.

Standardization is a transformation applied to each distribution in a group of distributions so that all distributions exhibit the same mean and the same standard deviation. This removes some aspects of intersample variability and can make data more comparable. This sort of procedure is analogous to commonly practiced scaling procedures but is more robust in that it controls for both scale and location.

The Analysis page calculates the ordination visualizations with either raw or normalized counts, at the users option. The normalization procedure is as follows.

$$\text{normalized\_value}_i = \log_2(\text{raw\_counts}_i + 1)$$

The standardized values then are calculated from the normalized values by subtracting the mean of each samples normalized values and dividing by the standard deviation of each samples normalized values.

$$\text{standardized}_i = (\text{normalized}_i - \text{mean}(\text{normalized}_i)) / \text{stddev}(\text{normalized}_i)$$

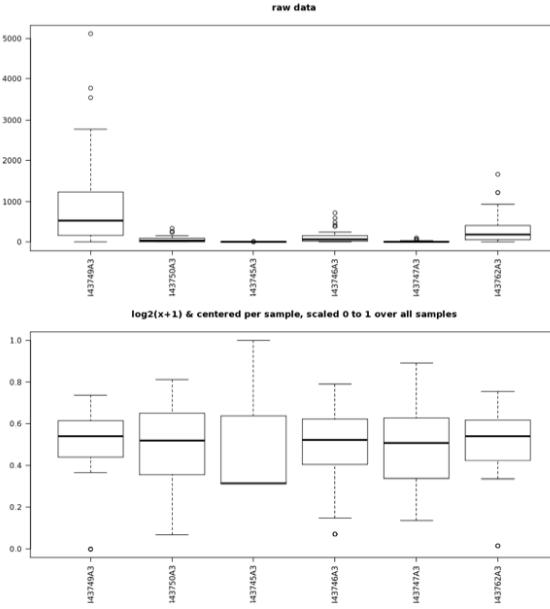


Figure 3.19: Boxplots of the abundance data for raw values (top) as well as values that have undergone the normalization and standardization procedures (bottom) described in the text. After normalization and standardization, samples exhibit value distributions that are much more comparable and that have a normal distribution; the normalized and standardized data are suitable for analysis with parametric tests; the raw data are not.

More about these procedures is available in a number of texts. We recommend Terry Speed's Statistical Analysis of Gene Expression in Microarray Data [35].

When data exhibit a nonnormal, normal, or unknown distribution, nonparametric tests (e.g., Man-Whitney or Kurskal-Wallis) should be used. Boxplots are easy to use, and the MG-RAST analysis page provides boxplots of the standardized abundance values for checking the comparability of samples (Figure 3.19).

### 3.9.2 Rarefaction

The rarefaction view is available only for taxonomic data. The rarefaction curve of annotated species richness is a plot (see Figure 3.20) of the total number of distinct species annotations as a function of the number of sequences sampled. As shown in the figure, ?? multiple data sets can be included.

The slope of the right-hand part of the curve is related to the fraction of sampled species that are rare. When the rarefaction curve is flat, more intensive sampling is likely to yield only a few

This data was calculated for metagenomes 4447970.3, 4447943.3, 4447192.3, 4447103.3, 4447102.3, 4447101.3, 4447971.3 and 4447903.3. The data was compared to M5NR using a maximum e-value of 1e-5, a minimum identity of 60 %, and a minimum alignment length of 15 measured in aa for protein and bp for RNA databases.

Metagenome 4447103.3 contains no organism data for the above selected sources and cutoffs. They are being excluded from the analysis.

The image is currently dynamic. To be able to right-click/save the image, please click the static button [static](#)

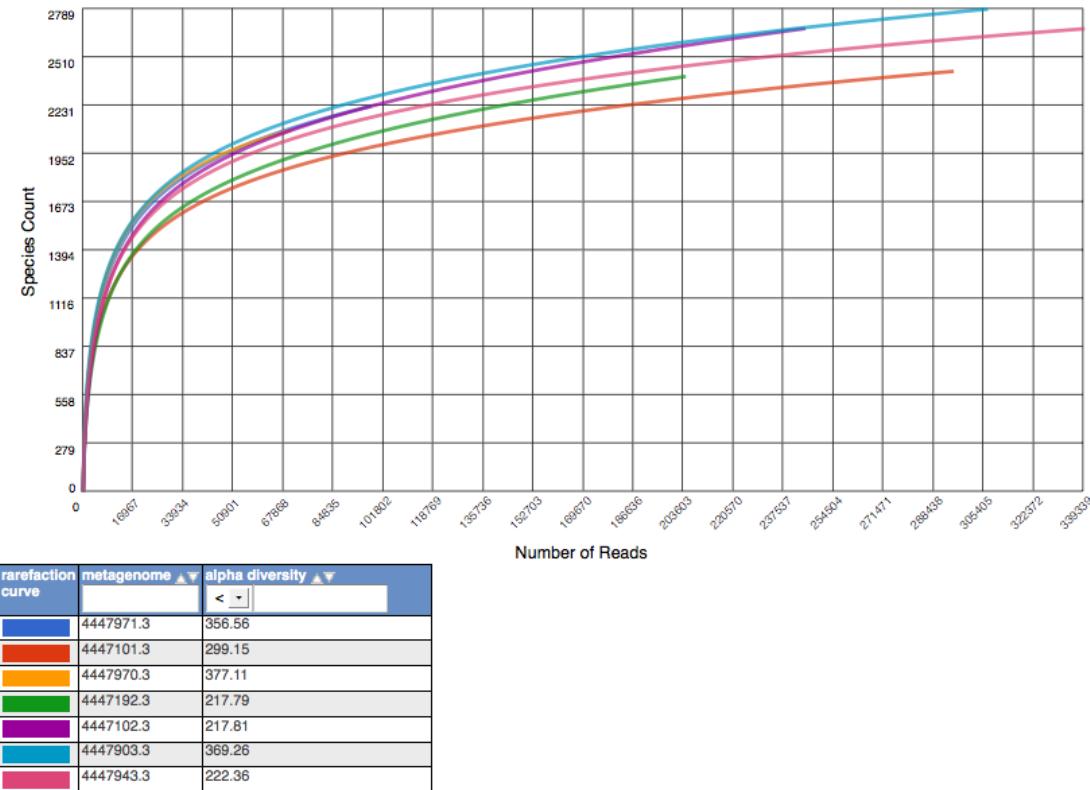


Figure 3.20: Rarefaction plot showing a curve of annotated species richness. This curve is a plot of the total number of distinct species annotations as a function of the number of sequences sampled.

additional species. The rarefaction curve is derived from the protein taxonomic annotations and is subject to problems stemming from technical artifacts. These artifacts can be similar to the ones affecting amplicon sequencing [31], but the process of inferring species from protein similarities may introduce additional uncertainty.

On the Analysis page the rarefaction plot serves as a means of comparing species richness between samples in a way independent of the sampling depth.

On the left, a steep slope indicates that a large fraction of the species diversity remains to be discovered. If the curve becomes flatter to the right, a reasonable number of individuals is sampled: more intensive sampling is likely to yield only a few additional species.

Sampling curves generally rise very quickly at first and then level off toward an asymptote as

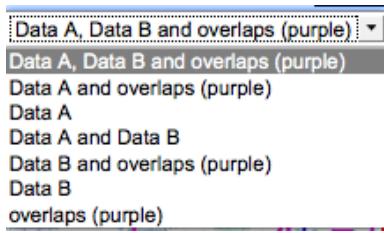


Figure 3.21: Options available for coloring the KEGG maps.

fewer new species are found per unit of individuals collected. These rarefaction curves are calculated from the table of species abundance. The curves represent the average number of different species annotations for subsamples of the the complete dataset.

### 3.9.3 KEGG mapper

The KEGG map tool allows the visual comparison of predicted metabolic pathways in metagenomic samples. It maps the abundance of identified enzymes onto a KEGG [16] map of functional pathways; note that the mapper is available only for functional data). Users can select from any available KEGG pathway map. Different colors indicate different metagenomic datasets.

The KEGG mapper works by providing two buffers that users can assign datasets to. After loading the buffers with the intended datasets, the KEGG mapper can highlight parts of the KEGG map that are present in the dataset. Several combinations of the two datasets can be displayed, as shown in Figure 3.21. Metagenomes can be assigned into one of two groups, and those groups can be visually compared (see Figure 3.22).

### 3.9.4 Recruitment plots

The recruitment page allows mapping of protein sequences in a single metagenome onto the complete genome sequences that are represented in the M5NR. Once the metagenome is selected, the page will provide a list of genomes, sorted by the number of hits per genome, from which the user can choose a genome to display (see Figure 3.23).

A circular genome plot or a table will be printed. See Figure 3.24 for an example. The following elements are contained in the figure:

- outmost circle: forward strand genes (red: protein, black: RNA)
- 2nd circle: contigs for the reference genome

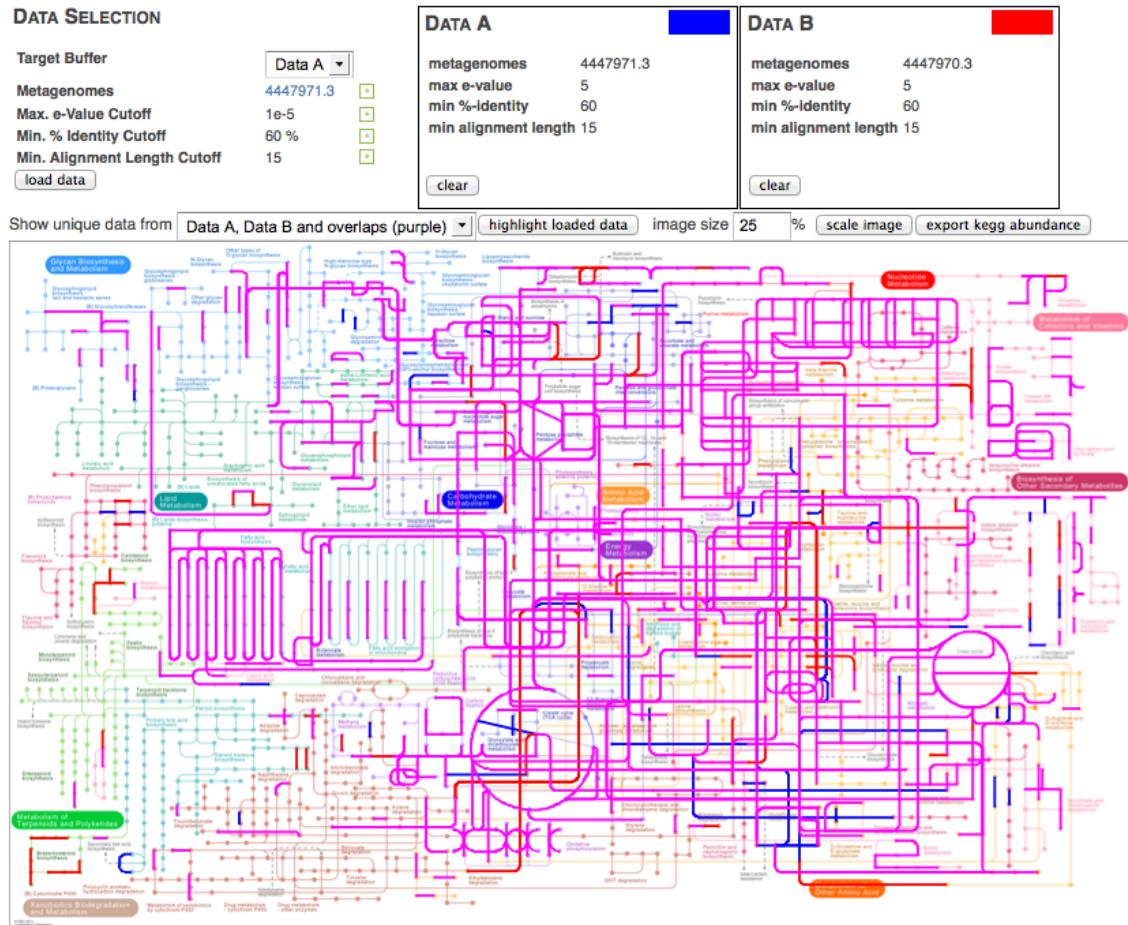


Figure 3.22: Comparison of two datasets using the KEGG mapper. Parts of metabolism common are shown in purple; unique to A are in blue; unique to B are in red.

- 3rd circle: reverse strand genes (red: protein, black: RNA)
- innermost circle: abundance information (color coded for E value)

The table view has the same information as the circular view and can easily be downloaded into a local spreadsheet. We use RefSeq [30] identifiers for the table as well as RefSEQ functions because the underlying contig information is present in the GenBank [3] downloads.

The recruitment plot uses the best hit approach.

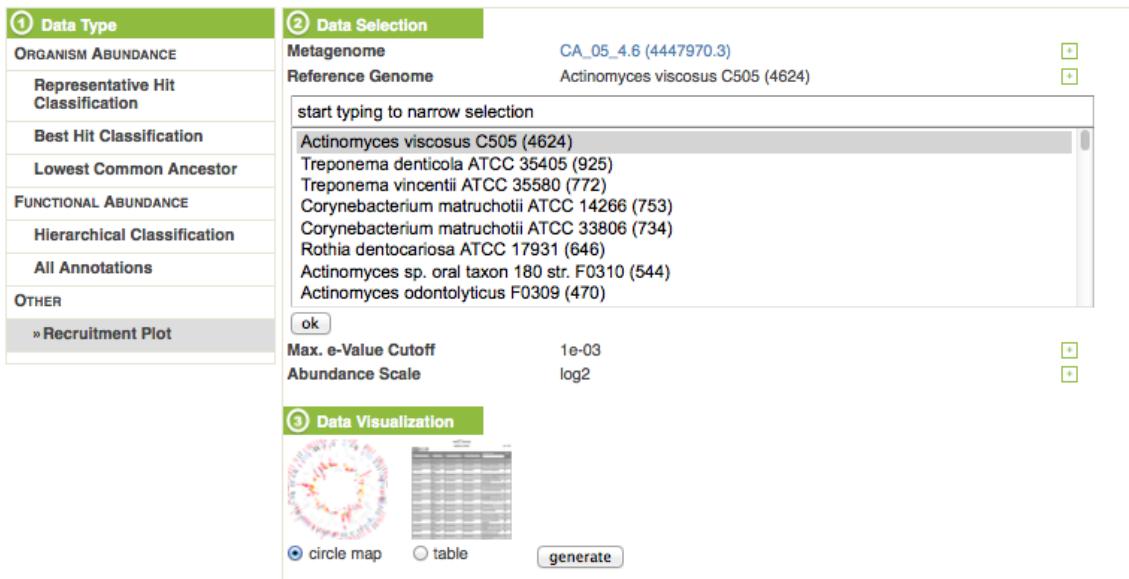


Figure 3.23: Selection of a genome for display, sorted by number of hits per genome.

### 3.9.5 Bar charts

Figure 3.25 shows the bar chart visualization option on the Analysis page. One important property of the page is the built-in ability to drill down by clicking on a specific category. In this example we have expanded the domain Bacteria to show the normalized abundance (adjusted for sample sizes) of bacterial phyla. The abundance information displayed can be downloaded into a local spreadsheet. Once a subselection has been made (e.g., the domain Bacteria selected), data can be sent to the workbench for detailed analysis. In addition, reads from a specific level can be added into the workbench.

### 3.9.6 Tree diagram

Figure 3.26 shows the tree diagram option on the Analysis page.

The tree diagram allows comparison of datasets against a hierarchy (e.g., Subsystems or the NCBI taxonomy). The hierarchy is displayed as a rooted tree, and the abundance (normalized for dataset size or raw) for each dataset in the various categories is displayed as a bar chart for each category. By clicking on a category (inside the circle), detailed information can be requested for that node; see Figure 3.27.

The tree offers several other capabilities, as shown in Figure 3.28.

- Export of a high-resolution image – For publication purposes we provide an SVG version of

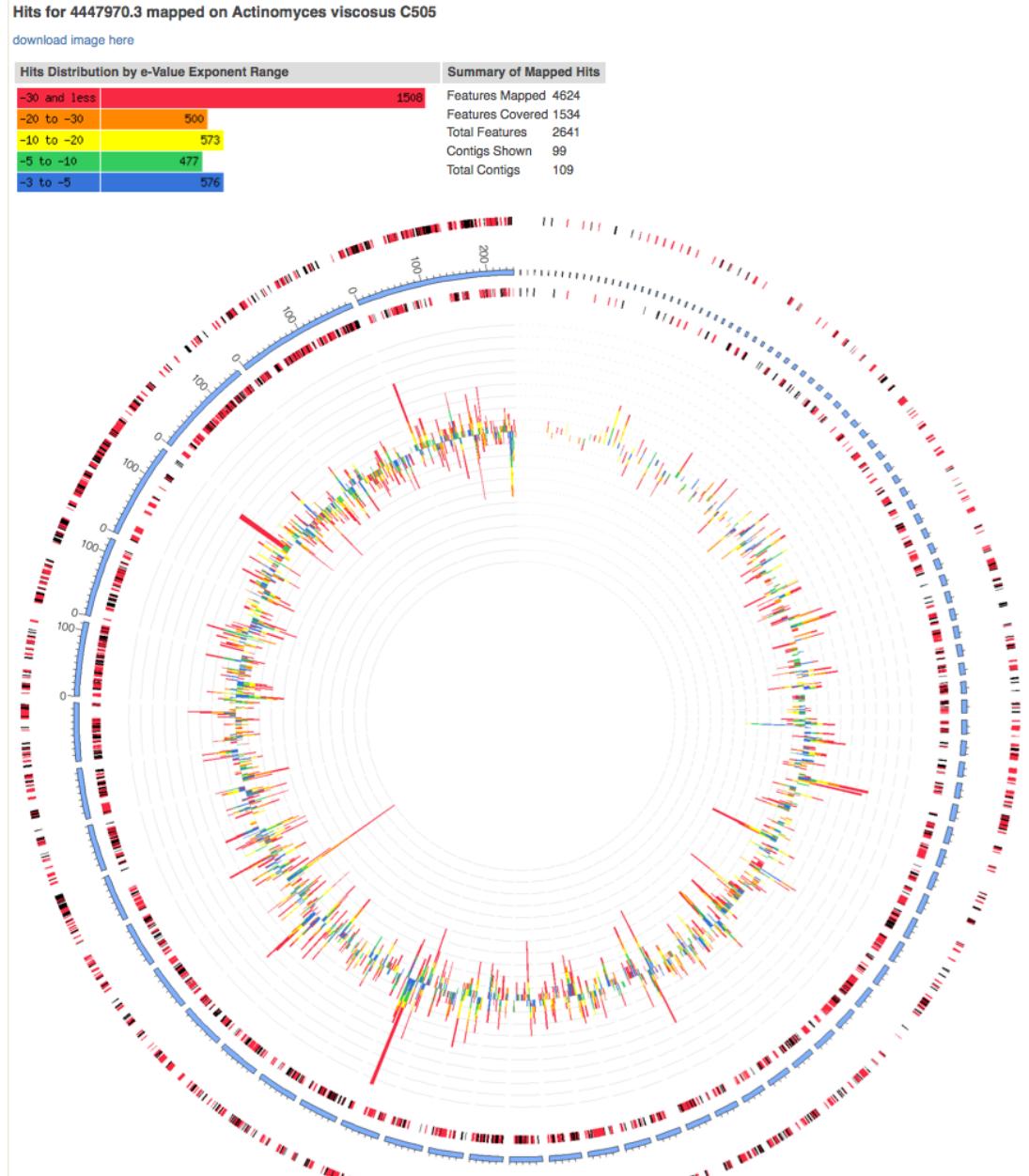


Figure 3.24: Example recruitment plot with the parameters from the previous figure for *Actinomyces viscosus* C505.



Figure 3.25: Bar chart view comparing normalized abundance of taxa. We have expanded the Bacteria domain to display the next level of the hierarchy.

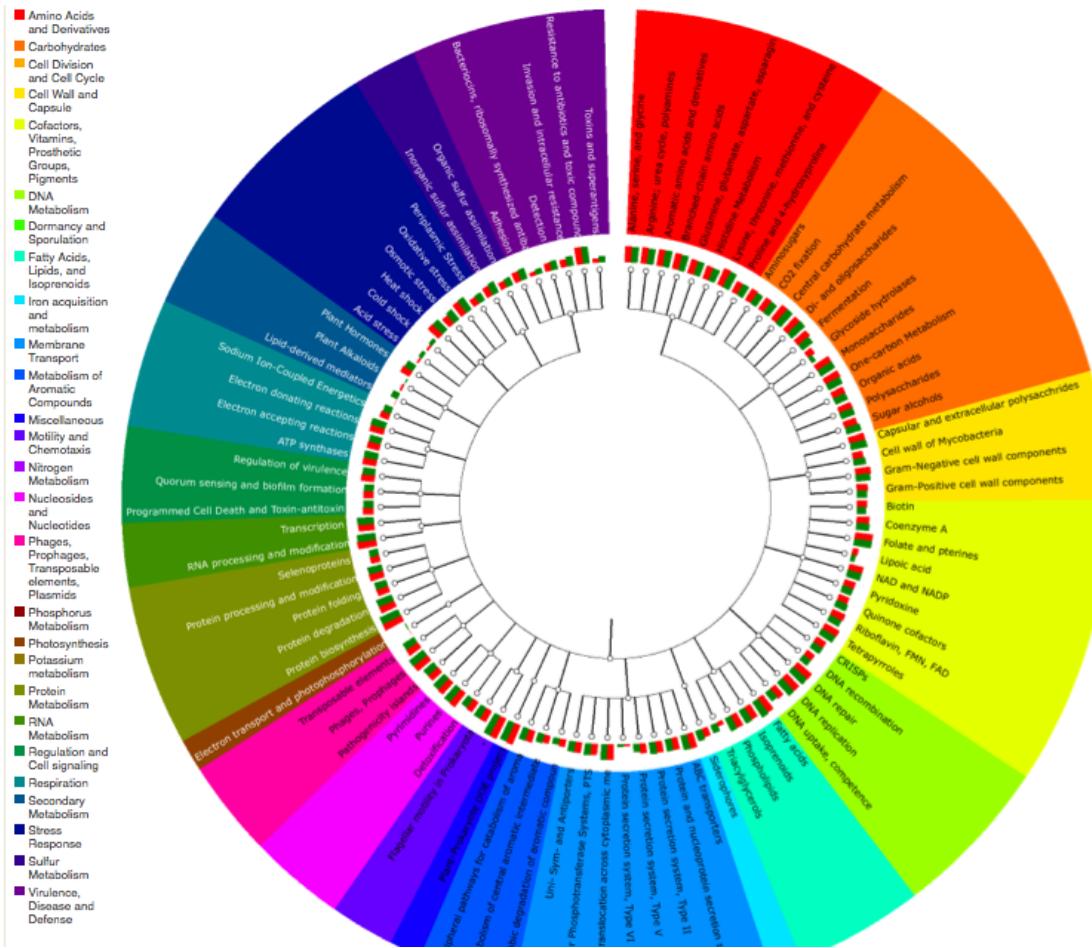


Figure 3.26: Tree diagram visualization option on the Analysis page.

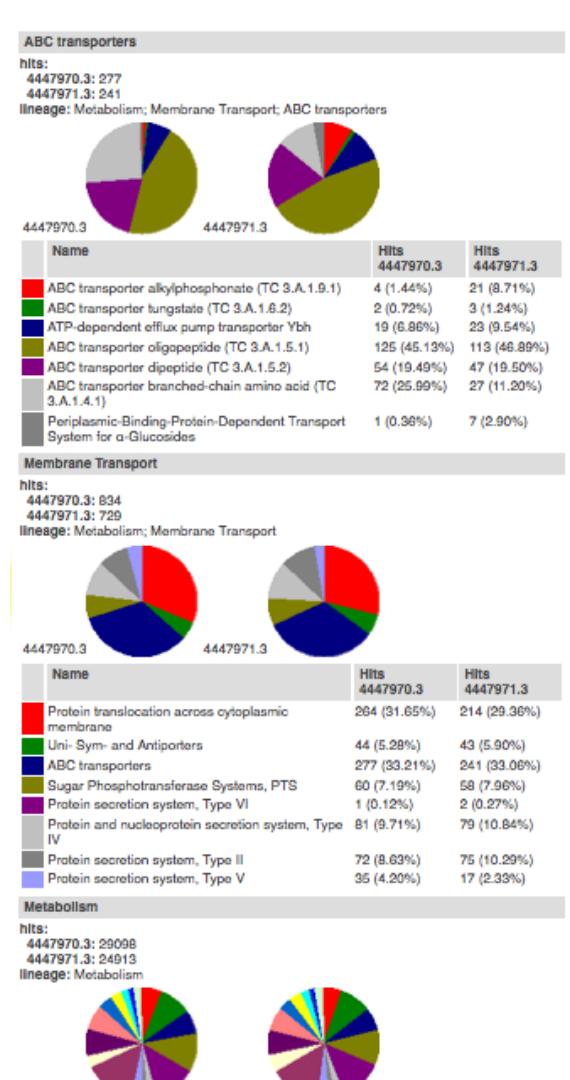


Figure 3.27: Tree diagram provision for detailed information: clicking on a node in the tree diagram will display addition information to the right of the tree display.

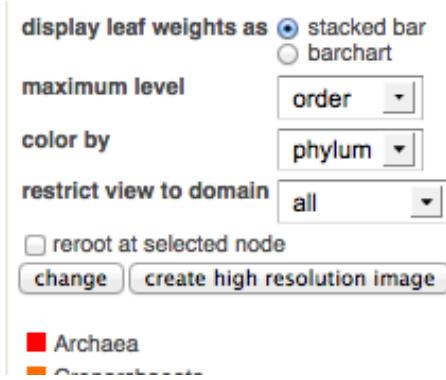


Figure 3.28: Options for the tree view.

the image.

- Rerooting – The tree display allows zooming in by rerooting the tree display. The user selects a node inside the tree (turning it red), and then selects “reroot” at the selected node. See Figure 3.28.
- Bar chart or stacked chart – The hierarchy entries can be displayed as bar charts per node or as a stacked graph.
- Restrict to domain – Is identical to rerooting the tree for a specific domain.
- Maximum level – This setting determines the depth of the tree being displayed.
- Color by – This setting determines the color (if any) used for the outer circle of the display.

Figure 3.29 shows the result of changing the display depth and coloring options. The color is used to group organisms visually into order level groups.

### 3.9.7 Heatmap/Dendrogram

The heatmap/dendrogram (Figure 3.30) allows an enormous amount of information to be presented in a visual form that is amenable to human interpretation. Dendograms are trees that indicate similarities between annotation vectors. The MG-RAST heatmap/dendrogram has two dendograms, one indicating the similarity/dissimilarity among metagenomic samples (x-axis dendrogram) and another indicating the similarity/dissimilarity among annotation categories (e.g., functional roles; the y-axis dendrogram). A distance metric is evaluated between every possible pair of sample

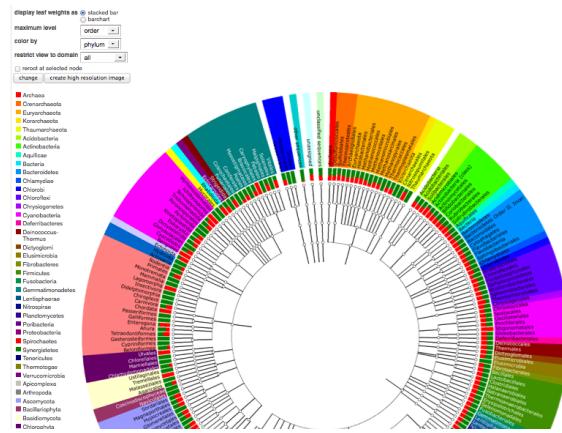


Figure 3.29: Tree view at order level with coloring set to phylum level.

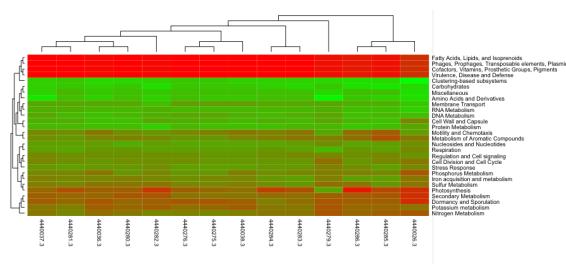


Figure 3.30: Heatmap/dendrogram example in MG-RAST. The MG-RAST heatmap/dendrogram has two dendrograms, one indicating the similarity/dissimilarity among metagenomic samples (x axis dendrogram) and another indicating the similarity/dissimilarity among annotation categories (e.g., functional roles; the y-axis dendrogram).

abundance profiles. A clustering algorithm (e.g., ward-based clustering) then produces the dendrogram trees. Each square in the heatmap dendrogram represents the abundance level of a single category in a single sample. The values used to generate the heatmap/dendrogram figure can be downloaded as a table by clicking on the download button.

### 3.9.8 Ordination

MG-RAST uses Principle Coordinate Analysis (PCoA) to reduce the dimensionality of comparisons of multiple samples that consider functional or taxonomic annotations. Dimensionality reduction is a process that allows the complex variation found in a large datasets (e.g., the abundance values of thousands of functional roles or annotated species across dozens of metagenomic samples) to be reduced to a much smaller number of variables that can be visualized as simple two- or three-dimensional scatter plots. The plots enable interpretation of the multidimensional data in a human-friendly presentation. Samples that exhibit similar abundance profiles (taxonomic or functional) group together, whereas those that differ are found farther apart.

A key feature of PCoA-based analyses is that users can compare components not just to each other but to metadata recorded variables (e.g., sample pH, biome, DNA extraction protocol) to reveal correlations between extracted variation and metadata-defined characteristics of the samples. It is also possible to couple PCoA with higher-resolution statistical methods in order to identify individual sample features (taxa or functions) that drive correlations observed in PCoA visualizations. This coupling can be accomplished with permutation-based statistics applied directly to the data before calculation of distance measures used to produce PCoAs; alternatively, one can apply conventional statistical approaches (e.g., ANOVA or Kruskal-Wallis test) to groups observed in PCoA-based visualizations.

### 3.9.9 Table

The table tool creates a spreadsheet-based abundance table that can be searched and restricted by the user. Tables can be generated at user-selected levels of phylogenetic or functional resolution. Table data can be visualized by using Krona [27] or can be exported in BIOM [23] format to be used in other tools (e.g., QIIME [5]). The tables also can be exported as tab-separated text.

Abundance tables serve as the basis for all comparative analysis tools in MG-RAST, from PCoA to heatmap/dendograms.

Consider the following example showing how to use the taxonomic information derived from an analysis of protein similarities found for the data set 4447970.3. We use the best hit classification, SEED, 10<sup>5</sup>, 60% identity, and a minimal alignment length of 15 amino acids. We select table

output. The results are shown in Figure 3.31.

The following control elements are connected to the table:

- group by – allows summarizing entries below the level chosen here to be subsumed.
- download table – downloads the entire table as a spreadsheet.
- Krona – invokes KRONA [27] with the table data.
- QIIME – creates a BIOM [23] format file with the data being displayed in the table.
- table size – changes the number of elements to display for the web page.

This data was calculated for metagenome 4447970.3. The data was compared to SEED using a maximum e-value of 1e-5, a minimum identity of 60 %, and a minimum alignment length of 15 measured in aa for protein and bp for RNA databases.

metagenome	source	domain	phylum	class	abundance	avg eValue	avg % ident	avg align len	# hits	to workbench
4447970.3	SEED	Archaea	Crenarchaeota	Thermoprotei	8	-12.00	68.08	51.60	7	
4447970.3	SEED	Archaea	Euryarchaeota	Archaeoglobi	2	-5.00	66.67	36.00	2	
4447970.3	SEED	Archaea	Euryarchaeota	Halobacteria	3	-8.50	68.16	43.25	3	
4447970.3	SEED	Archaea	Euryarchaeota	Methanobacteria	31	-14.75	67.90	55.97	16	
4447970.3	SEED	Archaea	Euryarchaeota	Methanococci	9	-10.39	68.21	49.39	8	
4447970.3	SEED	Archaea	Euryarchaeota	Methanomicrobia	59	-10.01	68.94	47.00	59	
4447970.3	SEED	Archaea	Euryarchaeota	Thermococci	15	-17.69	66.70	65.76	15	
4447970.3	SEED	Archaea	Euryarchaeota	Thermoplasmata	4	-5.00	68.10	33.40	4	
4447970.3	SEED	Archaea	Thaumarchaeota	unclassified (derived from Thaumarchaeota)	3	-10.00	60.04	50.67	3	
4447970.3	SEED	Bacteria	Acidobacteria	Solibacteres	19	-15.24	71.74	57.14	19	
4447970.3	SEED	Bacteria	Acidobacteria	unclassified (derived from Acidobacteria)	13	-21.64	62.83	78.64	13	
4447970.3	SEED	Bacteria	Actinobacteria	Actinobacteria (class)	13339	-20.36	70.30	69.27	8111	
4447970.3	SEED	Bacteria	Aquifcae	Aquifcae (class)	23	-19.74	67.38	69.09	23	
4447970.3	SEED	Bacteria	Bacteroidetes	Bacteroidia	2350	-27.79	77.40	75.62	1718	
4447970.3	SEED	Bacteria	Bacteroidetes	Cytophagia	102	-13.50	68.32	53.87	102	

Figure 3.31: View of the analysis page table.

Below we explain the columns of the table and the functions available for them. For each column we allow sorting the table by clicking on the upward- and downward-pointing triangles.

- metagenome

In the case of multiple datasets being displayed, this column allows sorting by metagenome ID or selecting a single metagenome.

- source

This displays the annotation source for the data being displayed.

- domain

The domain column allows subselecting from Archaea, Bacteria, Eukarya, and Viruses.

- phylum, class

Since we have selected to group results at the class level, only phylum and class are being displayed. The text fields in the column headers allow subsection (e.g., by entering Acidobacteria or Actinobacteria in the phylum field). The searches are performed inside the web browser and are efficient.

Any subselection will narrow down all datasets being displayed in the table.

Users can elect to have the results grouped by other taxonomy levels (e.g., genus), creating more columns in the table view.

- abundance

This indicates the number of sequences found with the parameters selected matching this taxonomic unit. (Note that the parameters chosen are displayed on top of the table.) Clicking on the abundance displays another page displaying the BLAT alignments underlying the assignments.

The abundance is calculated by multiplying the actual number of database hits found for the clusters by the number of cluster members.

- avg. evalue, avg percent identity, average alignment length

These indicate the average values for E value, percent identity, and alignment length.

- hits

This is the number of clusters found for this entity (function or taxon) in the metagenome.

- ...

This option allows extending the table to add additional columns.

### **3.9.10 Workbench**

The workbench was designed to allow users to select subsets of the data for comparison or export. Specifically, the workbench supports selecting sequence features and submitting them to further analysis or other analysis. A number of use cases are described below.

An important limitation with the current implementation is that data sent to the workbench exists only until the current session is closed.

# **Chapter 4**

## **User Manual**

### **4.1 Privacy, Identifiers, Sharing, and Publication**

Data in MG-RAST is private unless published to everyone or shared with specific users by the submitter.

Once data is submitted to the pipeline, a unique identifier is assigned (see 2.7.4 for details).

The web interface allows sharing and publication of data, requiring the presence of minimal metadata (see 2.7.1). Data can be shared only after the computation has finished.

### **4.2 Uploading to MG-RAST**

MG-RAST was designed to allow users to upload sequence data directly from next-generation sequencing machines. Data can be in FASTA, FASTQ, or SFF format. All uploaded sequence files must have one of the following extensions.

- .fasta
- .fna
- .fastq
- .fq
- .sff

Compressing large files will reduce the upload time and the chances of a failed upload. Users can use Zip (.zip) and gzip (.gz) as well as tarred gzipped files (.tgz).

We suggest uploading raw data (in FASTQ or SFF format) and letting MG-RAST perform the quality control step because this approach will allow us to identify any issues with the sequencing run. Frequently, local quality control will identify some issues but mask others.

It is not necessary to assemble data prior to upload to MG-RAST. The system has been optimized for short reads and can handle uploads of many hundreds of gigabytes.

### 4.2.1 Assembled data with read abundance information

For assembled data (in FASTA format) uploaded to MG-RAST, read abundance information for contigs can be imported as well. The “assembled” option for the pipeline will attempt to retrieve read abundance information from the sequence files using the following simple format:

```
>sequence_number_1_[cov=2]
CTAGCGCACATAGCATTCAAGCGTAGCAGTCACTAGTACGTAGTACGTACC
>sequence_number_2_[cov=4]
ACGTAGCTCACTCCAGTAGCAGGTACGTCGAGAAGACGTCTAGTCATCAT
. . .
```

The abundance information must be appended without spaces to the end of the sequence name (also without whitespace) in the format `_[cov=n]`, where n is the coverage or abundance of each contig.

### 4.2.2 Steps for submission via the web interface

To start uploading data to MG-RAST through the website, click on the up arrow. Doing so opens the Upload page. On this page you can upload files, modify the files where needed, add metadata, and submit files for analysis.

The page is split into two sections: “Prepare Data,” to upload, manipulate, and assemble all the files required for a submission, and “Submission,” to create the MG-RAST job(s), set analysis parameters, and start the analysis. Each section contains subsections that you can click to expand.

#### 4.2.2.1 Prepare data

**Download metadata spreadsheet template.** We provide a spreadsheet template that can be filled out with all the available metadata information for a dataset. The metadata can be modified later to add information or to correct errors. While the number of fields in the template is large, the number of required fields, labeled in red in the template, is small. The template file can be used to upload metadata for one or multiple samples and submit them to MG-RAST as a single project.

Table 4.1: Summary of upload times

Technology	Rate (bit/s)	Time for 1GB Upload
Modem 14.4 (2400 baud)	14.4 kbit/s	154 hours
ADSL Lite	1.5 Mbit/s	1.5 hours
Ethernet	10 Mbit/s	13.33 minutes
T3	44.736 Mbit/s	3 minutes
Fast Ethernet	100 Mbit/s	1.33 minutes

#### 4.2.2.2 Upload files

All files uploaded to MG-RAST should be named by using alphanumeric and .-\_ characters without spaces. Files larger than 50 MB should be compressed before upload, using gzip (preferable) or Zip. Compression will reduce the time taken for the upload of the file, which in turn reduces the chance that the upload will fail.

Files are divided into three types:

- Sequence files - FASTA, FASTQ, or SFF formats
- Metadata files - filled-out spreadsheet
- Barcode files - plain text ASCII containing lines with a barcode sequence followed by a unique filename separated by a tab, with as many lines as necessary for the barcodes in the sequence file you are submitting.

Click on the Browse button to select the file or files. The upload will begin automatically after the files are selected.

**Uploading.** For the actual uploading we use an HTML5 feature [40] that automatically breaks the files into chunks on the client side and sends them. Note: This is one of the reasons we request that you use a recent version of Firefox as older versions might be slower.

Table 1 summarizes observed upload times that might help users estimate how long the upload should take.

Based on observed values, upload times per 1 GB ( $10^{10}$  bytes) vary from 2 minutes to over an hour, with typical times being 10 to 15 minutes. Your experience will vary depending on the speed of your connection to the internet and the quality of service in your region.

#### Verifying the integrity of the uploaded files.

When the upload of your files has completed, you will be prompted with the MD5-sum of each file. You should generate an MD5 sum for each uploaded file on your machine, paste it into the

appropriate box in the prompt, and click the “check” button (see Figure 4.1). A popup will show you whether there is a match; additionally the check button will turn green upon success and red upon failure. Click the “Close” button if you have completed the checks or if you wish to skip this step.

Checking the integrity especially of large files is important because it will give you immediate feedback about whether your upload was successful. If not detected at upload time, a damaged file will lead to errors later in the pipeline, wasting both valuable compute cycles and, even more important, your time. To generate an MD5 sum of your file, you can use the “md5” shell command on a Mac, the “md5sum” shell command on Unix systems, or the freely available md5-sum tools on Windows (e.g., from <http://www.winmd5.com/>).

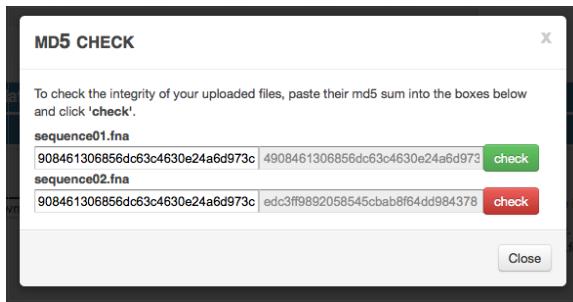


Figure 4.1: Dialogue requesting the user to put in a locally generated MD5 checksum for the files to identify any data corruption during the upload.

**File filters in place for uploading.** Since MG-RAST has been designed to work with metagenomic and metatranscriptomic datasets, there is a filter in place trying to identify datasets not suitable for MG-RAST. Those datasets will be colored red and cannot be submitted. Following are the criteria for rejection:

- Protein sequences – MG-RAST is optimized to perform translation from DNA to proteins.
- Reads shorter than 75 basepairs – The gene prediction stage performance deteriorates significantly with shorter reads.
- genomes – Submissions with complete genomes or a small number of contigs are rejected as well. Here our sister service RAST at <http://rast.nmpdr.org> should be used instead of MG-RAST.
- Files that are too small (sequence data less than 1 Mbp) – Files that are too small for MG-RAST to properly function are rejected at the submission stage. The minimal size requirement is 1 metagebase.

- Corrupted files – If the number of unique identifiers is not matching the number of sequence records in a file, the file is considered corrupt.
- Alignments – We cannot identify proteins from sequences containing alignment information.
- Colorspace – The tool chain does not function for ABIsolid sequences in colorspace. Please translate to standard FASTA.
- rar compressed files and Zip files over 4 GB – We cannot decompress these files.

In addition we will filter at the upload stage any Word documents, Rich Text Format files, and all files without the extension .fna, .fasta, .fq, .fastq, or .sff in their name.

**Note:** We recommend computing an MD5 checksum and verifying that the checksum computed by MG-RAST is identical to the locally computed checksum. This is the best way to ensure data integrity.

for Demultiplexing for 454 and similar datasets: MG-RAST performs demultiplexing based on the presence of the barcode sequence at the beginning of the reads. Assume you have a sequence file testseq.fasta and your barcode file has tab-separated lines like the following.

```
AAAAAAA   fileA
CCCCCCC   fileC
```

The demultiplexing step will split your sequence file into three files: fileA.fasta containing all the reads that begin with AAAAAAAA, fileC.fasta containing all the reads that begin with CCCCC-CCC, and testseq\_no\_MID\_tag.fasta containing reads that do not match either of the two.

We note that demultiplexing for Illumina needs to be done outside the MG-RAST system. Illumina barcodes work differently from 454 barcodes.

**Managing the inbox.** All files uploaded to MG-RAST will be displayed in your inbox and you can perform certain functions operations on them. Compressed and/or archived files can be unpacked, SFF files can be converted to FASTQ, and sequence data can be demultiplexed by using barcodes contained in uploaded files. Files can also be deleted. When a sequence file is selected, some information about the sequence data is displayed. It is a good idea to check that the uploaded files in your inbox match your local copies. The file MD5 checksum, file size, sequence count, and basepair count can be used for this purpose.

#### 4.2.2.3 Submission

The submission process allows you to create MG-RAST jobs using files in your inbox. It is designed to facilitate the creation of a large number of jobs easily.

**Select metadata file.** We recommend you supply metadata for all your samples. We assign a higher processing priority to samples with metadata. Metadata files that will have successfully passed our validation step will be displayed for selection.

**Select project.** All jobs created in MG-RAST will need to be placed in a project—either an existing project or a new one created by you during the submission. The project can be specified in one of three ways: in the metadata file if supplied or selected from the existing projects you have access to, or a new project name can be entered into the text box.

**Select sequence file(s).** All the sequence files in your inbox will be displayed for selection and job submission. Each sequence file can be used to create a single MG-RAST job.

**Choose pipeline options.** The MG-RAST analysis can be influenced by the options selected here, which affect dereplication, screening, and quality of the filtering of the reads. The options selected are applied to all the sequence files selected.

**Submit job.** This is the final step after which the analysis pipeline takes over and the processing begins. Once the job or jobs have been submitted, all the files required to create them will be removed from your inbox.

You can monitor the progress of your jobs in "My Data Summary" on the Browse Metagenomes page.

### 4.2.3 Cmd-line uploader

The following syntax will allow uploading to MG-RAST from the command line.

```
curl -H "auth: webkey" -X POST -F "upload=@/path_to_file/
metagenome.fasta" "http://api.metagenomics.anl.gov/
inbox/upload" > curl_output.txt
```

where you need to substitute “webkey” with the unique string of text generated by MG-RAST for your account. Your webkey is valid for a limited time period and ensures that the uploads you perform from the command line are recognized as belonging to your MG-RAST account and are placed in the correct inbox.

### 4.2.4 Managing the Inbox

The Inbox is a temporary storage location for sequence and metadata files prior to submission to the pipeline. To protect us from any misuse of the facility, we have limited the Inbox to metadata spreadsheets and sequence files.

Files are visible only to the uploading user and will automatically be deleted after 72 hours.

You can unpack, delete, convert and demultiplex files from your inbox below. Metadata files will automatically appear in the 'select metadata file' section below. Sequence files will automatically appear in the 'select sequence file(s)' section below after sequence statistics are calculated (may take anywhere from seconds to hours depending on file size).

Filenames in gray in your Inbox are undergoing analysis and cannot be moved or submitted to a different process until analysis is complete. Filenames in red have encountered an error.

#### File Processing Operations

unpack selected

Unpacks selected zip, gzip, or tar files.

demultiplex

Demultiplexes selected files.

convert sff to fastq

Converts selected sff files to fastq format.

join paired-ends

Joins overlapping paired-end reads.

```
200_2M_09Nov09_lane_8_1.txt
200_2M_09Nov09_lane_8_1.txt.zip
E4GC80A02.fastq
E4GC80A02.ff
E4GC80A02.sff.fasta
E4GC80A02.sff.qual
E4GC80A02.xml
(uploading) E8N68KH02.sff
metadata_spreadsheet_biogas_reads.xls
MGRAST_MetaData_template_1.0.xlsx
```

#### Directory Management Operations

update inbox

Refreshes the contents of your inbox.

move selected

Moves the selected files into or out of a directory.

delete selected

Deletes the selected files.

create directory

Creates a new directory in your inbox.

delete directory

Allows you to select and delete an empty directory.

Figure 4.2: Temporary storage provided in Inbox before submitting data and limited editing features.

#### **4.2.4.1 File-processing options in the Inbox**

The following file-processing options are available.

- unpack selected – unpacks selected zip, gzip, or tar files.
- convert sff to fastq – converts selected sff files to fastq format (only FASTQ and FASTA files can be submitted to the system).
- demultiplex – demultiplexes selected files.

Note that this is suitable only for 454 type barcodes that are actual prefixes of the reads. This approach does not work for the Illumina barcode approach (basically a third read for each paired end read).

- join paired ends – joins overlapping paired-end reads.

**Please note:** After the actual upload is complete, the system will compute the statistics shown in Figure 4.3. Computing this information takes some time, so your data will not immediately be visible after you uploaded it.

#### **4.2.4.2 Directory management operations for the Inbox**

The following operations are available for managing the directory.

- update inbox – refreshes the contents of your inbox
- move selected – moves the selected files into or out of a directory
- delete selected – deletes the selected files
- create directory – creates a new directory in your inbox
- delete directory – allows you to select and delete an empty directory.

Users should always double check the MD5 checksum for files that are uploaded to the system in order to verify the integrity. Figure 4.3 shows the MD5 fingerprint that is computed upon upload for each file.

File Information	
sequence content	DNA
unique id count	118196
sequence type	WGS
standard deviation gc content	4.387
standard deviation length	23.360
standard deviation gc ratio	0.101
sequencing method guess	454
bp count	29996553
ambig sequence count	10432
length max	509
suffix	1
file size	59.5 MB
ambig char count	31973
sequence count	118196
length min	50
average gc content	69.145
average gc ratio	0.452
average ambig chars	0.271
file checksum	d5f9cd37554e6a858c84154aa0d2047
average length	253.787
type	ASCII text
creation date	2013 May 09 08:44:33
file type	fastq

Figure 4.3: Information displayed by the inbox for one file (once selected).

#### 4.2.5 Generating metadata for the submission

MG-RAST uses questionnaires to capture metadata for each project with one or more samples. Users download and fill out the questionnaire and then submit it. Questionnaires are validated automatically by MG-RAST for completeness and compliance with the controlled vocabularies for certain fields.

MG-RAST has implemented the use of Minimum Information about any (X) Sequence (MIXS) [44] developed by the Genomic Standards Consortium. In addition to the minimal checklists, more detailed data can be captured in optional environmental packages.

We use simple spreadsheets to capture metadata, with a minimal number of required fields (in red in the spreadsheets) and a number of optional fields. The spreadsheet is separated into multiple tabs representing the different metadata categories. The MG-RAST metadata spreadsheet template is available on the MG-RAST upload page or at [ftp://ftp.metagenomics.anl.gov/data/misc/metadata/MGRAST\\_MetaData\\_template\\_1.3.xlsx](ftp://ftp.metagenomics.anl.gov/data/misc/metadata/MGRAST_MetaData_template_1.3.xlsx).

A filled-out version of the spreadsheet is available at [ftp://ftp.metagenomics.anl.gov/data/misc/metadata/MGRAST\\_MetaData\\_template\\_example.xlsx](ftp://ftp.metagenomics.anl.gov/data/misc/metadata/MGRAST_MetaData_template_example.xlsx).

In Figure 4.4 we show the template tab for project and the required field labels (in red) (in essence, your contact information). Figure 4.5 shows the various tabs in the spreadsheet.

Note: Use the third line in the spreadsheet and as shown in Figure 4.6 to enter your data. Do not attempt to alter the first two lines or delete them; they are read only. The first line contains

	A	B	C	D	E	F	G
1	project_name	project_description	project_fund	project_id	PI_email	PI_firstname	PI_lastname
2	Name of the project	Description of the project	Funding source	Internal ID of the project	Administrative contact email	Administrative contact first name	Administrative contact last name
3							
4							
5							
6							
7							
8							
9							
10							
11							
12							
13							

Figure 4.4: Project spreadsheet. In red are required fields. Note that the 2nd row contains information on how to fill out the form.



Figure 4.5: The various tabs in the spreadsheet. Project, sample and one of library metagenome or library mimarks survey are required.

the field labels, and the second line contains descriptions that can help explain how to fill out the fields, along with what unit to use (e.g., temperature in Celsius and distance in meters).

**Required sheets** You need to fill out four sheets to describe your metadata:

1. Project – This sheet has only one row, and describes a set of samples uploaded together; the other sheets have one row per sample.
2. Sample – This sheet includes either the filename or metagenome name used for matching.
3. Library – This sheet includes either the metagenome (for WGS and WXS) or mimarks-survey (for 16s amplicon).
4. Environmental package – Several packages of suggested standard metadata are available. Choose the package that best describes your dataset (e.g., water, human-skin, soil).

	A	B	C	D	E	F	G	H
1	sample_name	sample_id	latitude	longitude	continent	country	location	depth
2	Unique name	Internal ID of	The geograph	Depth is				
3	sample1							
4	sample2							
5	sample3							
6								
7								

Figure 4.6: Sample tab with 3 new samples (sample1, sample2, and sample3) added. Again red text in the first row indicates required fields. Rows 1 and 2 cannot be altered.

The sample section (below) requires minimal information (including the sample name) about where and when the sample was taken. Note that some fields in the spreadsheet must be filled out with terms from a controlled vocabulary or in a certain way. Country and environment (biome, feature, material) fields require entries from curated ontologies, gazetteer and environmental ontology, respectively.

Figure 4.6 shows the sample tab with three new samples (sample1, sample2, and sample3) added. Again red text in the first row indicates required fields.

**Mandatory fields** Five fields must be completed.

- Country – United States of America, Netherlands, Australia, Uruguay
- Latitude and longitude – 106.84517, -104.60667, 28 42.306N, 88 24.099W, 45.30 N, 73.35 W
- Biome – small lake biome, tropical humid forests, mangrove biome. This term must be one of the terms from the bioportal ontology. Terms that are not on this list are not valid.
- Feature – city, fish farm, livestock-associated habitat, marine habitat, ocean basin, microbial mat. This term must be one of the terms from the bioportal ontology. Terms that are not on this list are not valid.
- Material – air, dust, volcanic soil, saliva, blood, dairy product, surface water, piece of gravel. This term must be one of the terms from the bioportal ontology. Terms that are not on this list are not valid.

**Library section** The library section captures technical data on the preparation and sequencing done. You should choose the library tab to fill out (“metagenome” for shotgun sequencing or “mimarks-survey” for amplicon) based on the type of sequencing done. These are separated as different sequencing techniques involving different metadata fields. Each row describes one library for one sample. The samples need to have the identical sample name you used in the sample tab before.

The library\_metagenome tab shows the required fields in red. The **file\_name** field holds the filename of the sequence file uploaded, or the filename to use for creating the demultiplexed file if you uploaded a multiplexed sequence file and have barcode sequences in the spreadsheet. This is used for mapping sequence files to metadata.

The **metagenome\_name** field holds the name of the metagenome you are submitting. If the **file\_name** field is empty, it will be used for mapping metadata to sequence files, in this case it would need to match the uploaded sequence filename (not including file extension).

The **investigation\_type** field is required to be “metagenome” for shotgun metagenome samples and “mimarks-survey” for amplicon studies (reflecting what tab was filled out).

The type of sequencing instrument used is another required field. Values are, for example, Illumina, 454, Ion Torrent, Sanger, or assembled.

Again, only a limited number of fields are required. However, the more info you provide, the easier it is for you and others to understand any potential uses of your data and to understand why results appear in a particular way. It might, for example, allow understanding of specific biases caused by technology choices or sampled environments.

You can fill out one or more environmental metadata packages. Currently we provide support for the following GSC environmental packages:

- Air
- Built Environment
- Host-associated
- Human-associated
- Human-oral
- Human-skin
- Human-vaginal
- Microbial mat/biofilm

- Miscellaneous natural or artificial environment
- Plant-associated
- Sediment
- Soil
- Wastewater sludge
- Water

We strongly encourage users to submit rich metadata, but we understand the effort required in providing it. Using the environmental packages (which were designed and are used by practitioners in the respective field) should make it reasonably simple to report the essential metadata required to analyze the data. If there is no environmental package to report metadata for your specific sample, please contact MG-RAST staff: we will work with the GSC [10] to create the required questionnaire.

## 4.3 Working with Projects and Collections

Collections provide an efficient way to create multiple sets of metagenomes for analysis. For example, if you want to compare human gut with cow rumen samples, you probably want to see a dialogue like that in Figure 4.7.

MG-RAST v3 provides a mechanism to make this happen. Users can create collections that are persistent across multiple sessions. Below, we show how to define a collection that allows comparison of multiple datasets. Please note that collections are just shortcuts to the actual samples; they cannot be shared at this time.

Step 1: you start with the metadata browser (either on the front page or in the little menu block in the top right-hand corner), and click on the globe symbol. See Figure 4.8 for the symbol.

Step 2: This will take you to the Browser dialogue, showing a large number of metagenomes (Figure 4.9).

Step 3: Clicking on biome will allow selecting a specific Biome (here we pick Animal associated). This results in a list of metagenomes shown in Figure 4.10.

The list of samples still shows too many samples when restricted to just animal associated metagenomes.

Step 4: To downselect, search for Twin to further restrict to samples from Peter Turnbaugh and Jeffrey Gordons Human Twin study (see Figure 4.12).

collection	job name ▲▼	select ...
all		<input type="checkbox"/> all
all	ObeseMouseCecumMic2005	<input type="checkbox"/>
CF	LeanMouseCecumMic2005	<input type="checkbox"/>
CF2	CFLungPat001Rep1SDVir20060505	<input type="checkbox"/>
CFLung	CFLungPat001Rep2SDVir20060505	<input type="checkbox"/>
human	CFLungPat001Rep3SDVir20060505	<input type="checkbox"/>
marine		
Northern Line	FXPY	<input type="checkbox"/>
null	FYGT	<input type="checkbox"/>
plant virus	HealSputRep2SDVir20060707	<input type="checkbox"/>
plant2	HealSputRep3SDVir20060707	<input type="checkbox"/>
Seawater	BGlgutGeneSet	<input type="checkbox"/>
St Louis - human samples	human In-R	<input type="checkbox"/>
Unspecified Biome 4/7/2011	human In-M	<input type="checkbox"/>
Unspecified Biome :-(	human In-E	<input type="checkbox"/>
human	human In-D	<input type="checkbox"/>
human	human In-B	<input type="checkbox"/>
human	human In-A	<input type="checkbox"/>
human	human F2-Y	<input type="checkbox"/>
human	human F2-X	<input type="checkbox"/>
human	human F2-W	<input type="checkbox"/>
human	human F2-V	<input type="checkbox"/>

Figure 4.7: View of the browse table with the collection column enabled. Clicking on the “...” at the right end of the table allows expanding the table columns.

Step 5: Clicking on the black shopping-cart symbol in the top right-hand corner will allow the creation of a new collection entry. The next step is naming the collection (see Figure 4.13. Here we name the collection Twin Study and hit OK.

Step 6: Once the collection is added, the new collection will appear in the list of collections (in “Your Data Summary”).

Step 7: Use collection in the Metagenome selection on the Analysis page. It is possible to analyze individual metagenomes or compare whole groups of metagenomes (see Figure 4.15).

## 4.4 Understanding Datasets

Unfortunately not every sequencing run works equally well. Users of MG-RAST have provided data with many different sources of error, allowing us to provide a number of tools to identify the

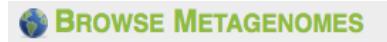


Figure 4.8: Symbol for the MG-RAST metagenome browser

All Metagenomes  
group by project

Current table counts  
metagenomes 22081 projects 101 biomes 51 altitudes 82 depths 110 locations 306 pH's 0 countries 41 temperatures 145 sequencing methods 5 pi's 58

public (4279) private (3) shared (17799)

clear table filters display 25 items per page add selected to a collection

project	name	biome	type	select
unknown	CA_005_4_6_repl_filt_bact	human-associated habitat	WGS	<input checked="" type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
unknown	ConsensReass	unknown	WGS	<input checked="" type="checkbox"/> shared <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
unknown	ConsensReass2	unknown	WGS	<input checked="" type="checkbox"/> shared <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
unknown	MRFinal	unknown	Amplicon	<input checked="" type="checkbox"/> shared <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
unknown	CA_006_1_6_repl_filt_bact_fa	human-associated habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
cDNA - Plymouth Marine Lab Coastal Waters project	1-19-DNA-fbx	marine habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
cDNA - Plymouth Marine Lab Coastal Waters project	6-19-DNA-fbx	marine habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
unknown	db_mg_case1	unknown	WGS	<input checked="" type="checkbox"/> shared <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Northern Line Islands	FannLMic20050811	marine habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Northern Line Islands	FannLVir20050811	marine habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Soudan Mine Metagenome	RedSoudanMineMic20050331	mine drainage	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Soudan Mine Metagenome	BlackSoudanMineMic20050331	mine drainage	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Chicken Cecum Microbiome	Chicken Cecum A	animal-associated habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Chicken Cecum Microbiome	Chicken_Cecum_B	animal-associated habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
Chicken Cecum	Chicken Cecum A Contigs	animal-associated habitat	WGS	<input checked="" type="checkbox"/> public <input type="checkbox"/> all <input type="checkbox"/> job # <input type="checkbox"/> id <input type="checkbox"/> project <input type="checkbox"/> name <input type="checkbox"/> biome <input type="checkbox"/> type <input type="checkbox"/> altitude <input type="checkbox"/> depth <input type="checkbox"/> location <input type="checkbox"/> pH <input type="checkbox"/> country <input type="checkbox"/> temperature <input type="checkbox"/> sequencing method <input type="checkbox"/> pi
...				

Figure 4.9: MG-RAST metagenome browser

most common errors.

The quality assessment tools described in ?? provide a good tool set for an initial data quality analysis. While many potential sources of error exist, some common problems can be easily identified with just the nucleotide histograms. If your data exhibits patterns like the ones described in ??, there were likely to be problems with sequencing.

We frequently find datasets with high numbers of reads filtered out by the quality control. Below we list the major reasons for filtered reads:

## 1. Artificial duplicate reads

The presence of significant amounts of ADRs hints at problems with a PCR step. Frequently a problem with 454 runs occurs, but problems are also seen with other platforms. For example, a set of DNA templates may be copied many times during a PCR step, often consuming up to 80% of the entire dataset.

## 2. Filtered reads

If you have selected screening against a host organism (e.g., the human genome), reads

All Metagenomes  
group by project

Current table counts      public (4280) private (3) shared (17799)

metagenomes	projects	biomes	altitudes	depths	locations	ph's	countries	temperatures	sequencing methods	pi's
22082	102	51	82	110	306	0	41	145	5	58

clear table filters      add selected to a collection

display 25 items per page      displaying 1 - 25 of 22082      next » last »

project	name	biome	type	select ...	all
The oral metagenome in health and disease	CA_05_4.6	all			
unknown	ConsensReass	air			
unknown	ConsensReass2	animal manure			
unknown	MRfinal	animal manure, animal manure			
The oral metagenome in health and disease	CA_06_1.6	animal-associated habitat			
cDNA - Plymouth Marine Lab Coastal Waters project	1-19-DNA-fbx	animal-associated habitat, feces			
cDNA - Plymouth Marine Lab Coastal Waters project	6-19-DNA-fbx	animal-associated habitat, feces, feces			
aquatic habitat		aquatic habitat			
aquatic habitat, freshwater habitat		aquatic habitat, freshwater habitat			
aquatic habitat, marine habitat		aquatic habitat, marine habitat			
aquatic habitat, marine habitat, Aphotic zone		aquatic habitat, marine habitat, Aphotic zone			
unknown	db_mg_case1	aquatic habitat, sediment			
Northern Line Islands	FannLIMic20050811	biofilm, saline marsh			
Northern Line Islands	FannLIVir20050811	biofilm, sludge, waste water			
Soudan Mine Metagenome	RedSoudMineMic20050331	clouds			
Soudan Mine Metagenome	BlackSoudMineMic20050331	extreme habitat ; hypersaline			
Chicken Cecum Microbiome	Chicken Cecum A	extreme habitat, hydrothermal vent, hot spring			
Chicken Cecum Microbiome	Chicken_Cecum_B	extreme habitat, hydrothermal vent, microbial mat			
feces		feces			
Chicken Cecum Microbiome	Chicken Cecum A Contigs	feces, feces			
Chicken Cecum Microbiome	Chicken Cecum B Contigs	animal-associated habitat	WGS	public	<input type="checkbox"/>
unknown	Bray Reclaimed Water	animal-associated habitat	WGS	shared	<input checked="" type="checkbox"/>
unknown	nloke	unknown	WGS	shared	<input type="checkbox"/>
unknown	Vaginal Microbiome#1	unknown	WGS	shared	<input type="checkbox"/>
unknown	VE	unknown	WGS	shared	<input checked="" type="checkbox"/>

Figure 4.10: Filtering by BIOME information.

All Metagenomes  
group by project

Current table counts      public (114) private (0) shared (19)

metagenomes	projects	biomes	altitudes	depths	locations	p's	countries	temperatures	sequencing methods	p's
133	18	1	9	7	35	0	10	10	5	15

clear table filters      add selected to a collection

display  items per page      displaying 1 - 25 of 133      next > last >

project ▾	name ▾	biome ▾	type ▾	select ...
Chicken Cecum Microbiome	Chicken Cecum A	animal-associated habitat	WGS	public <input type="checkbox"/>
Chicken Cecum Microbiome	Chicken_Cecum_B	animal-associated habitat	WGS	public <input checked="" type="checkbox"/>
Chicken Cecum Microbiome	Chicken Cecum A Contigs	animal-associated habitat	WGS	public <input type="checkbox"/>
Chicken Cecum Microbiome	Chicken Cecum B Contigs	animal-associated habitat	WGS	public <input type="checkbox"/>
Human Lung Healthy vs Cystic Fibrosis Metagenome	CLungPat001Rep1SDVir20060505	animal-associated habitat	WGS	public <input type="checkbox"/>
Mosquito Metagenome	Mosq1SDVir20060125	animal-associated habitat	WGS	public <input type="checkbox"/>
Human Lung Healthy vs Cystic Fibrosis Metagenome	CLungPat001Rep2SDVir20060505	animal-associated habitat	WGS	public <input type="checkbox"/>
Mosquito Metagenome	MosqDigSDVir20060606	animal-associated habitat	WGS	public <input type="checkbox"/>
Mosquito Metagenome	Mosq2SDVir20060609	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealGutKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealGutKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealSlimKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealSlimKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealSlimKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealSlimKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealSlimKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Aquacultured Fish (Kent State)	FishHealSlimKentSTMic20060504	animal-associated habitat	WGS	public <input type="checkbox"/>
Stressed Coral Holobionts	BocasPAMic20050921	animal-associated habitat	WGS	public <input type="checkbox"/>
Stressed Coral Holobionts	DOCPorCompHawMic200602	animal-associated habitat	WGS	public <input type="checkbox"/>
Stressed Coral Holobionts	pHPorCompHawVir200602	animal-associated habitat	WGS	public <input type="checkbox"/>
Stressed Coral Holobionts	ConPorCompHawVir200602	animal-associated habitat	WGS	public <input type="checkbox"/>
Stressed Coral Holobionts	DOCPorCompVirHaw200602	animal-associated habitat	WGS	public <input type="checkbox"/>

Figure 4.11: A reduced list of metagenomes for one BIOME.

**BROWSE METAGENOMES**

**Your Data Summary**

Available for analysis <sup>[?]</sup>	3
In Progress <sup>[?]</sup>	1649
Shared with you <sup>[?]</sup>	17799
Collections <sup>[?]</sup>	13

Click on the blue links above to browse just your data. For more information visit the [Support page](#).

**All Metagenomes**  
group by project

**Current table counts**

metagenomes	projects	biomes	altitudes	depths	locations	ph's	countries	temperatures	sequencing methods	pi's
16	1	1	0	0	2	0	1	0	1	1

public (16) private (0) shared (0)

[clear table filters](#) [add selected to a collection](#)

display  items per page

displaying 1 - 16 of 16

project ▾▼	name ▾▼	biome ▾▼	type ▾▼	select ...
Twin		animal-associated habitat	WGS	<input checked="" type="checkbox"/> p <input type="checkbox"/> all
Twin Gut Microflora Study	TS1	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS2	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS4	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS19	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS20	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS28	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS49	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS50	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS3	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS21	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS51	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS6	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS7	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS8	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS9	animal-associated habitat	WGS	<input type="checkbox"/>
Twin Gut Microflora Study	TS30	animal-associated habitat	WGS	<input type="checkbox"/>

displaying 1 - 16 of 16

Figure 4.12: Selecting a specific project to further reduce the number of datasets being displayed.

matching that particular genome will be removed from the dataset. This action frequently consumes a large fraction of the sequence run. There are in vitro techniques to minimize the amount of host DNA in your sample. see [37].

### 3. Filter rRNA reads in metatranscriptomes

If you are trying to analyze a metatranscriptome, reducing the amount of ribosomal is essential. Unless an in vitro rRNA knockdown method is applied, up to 97% of all reads will be ribosomal.

## 4.5 Drilling Down with the Workbench

One of the new features of MGRAST v3 is the workbench. It is the main mechanism for exchanging subsets of data between analysis views. It also allows you to download the FASTA files of a selection of proteins.

When you initially go to the Analysis page (see 3.9), your workbench will be empty. It is displayed as the leftmost tab in the data tabular view. So how do you get data into the workbench?

The screenshot shows the MG-RAST 'BROWSE METAGENOMES' interface. On the left, there are two summary boxes: 'Your Data Summary' and 'Public Data Summary'. The 'Your Data Summary' box shows metrics like 3 datasets available for analysis, 1649 in progress, and 17799 shared with you. The 'Public Data Summary' box shows metrics like 4280 metagenomes, 88 projects, and 164 million sequences. In the center, there's a table titled 'All Metagenomes' with columns for metagenomes, projects, biomes, altitudes, depths, locations, ph's, countries, temperatures, sequencing methods, and p's. A red arrow points from the 'add selected to a collection' button in the top right to a modal dialog box. This dialog box has a question mark icon and asks 'Enter a name for this collection'. The input field contains 'Twin Study'. Below the input field are 'OK' and 'Cancel' buttons. To the right of the dialog is a table of 16 rows, each representing a twin study with columns for project, ID, habitat, and type (WGS, public). The first row is highlighted.

Figure 4.13: Saving a collection

There are two simple ways to select data subsets: from any generated table or from the drilldown of a barchart.

Try this example: Start by selecting the lean and obese mouse cecum samples (MG-RAST IDs 4440463.3 and 4440464.3) [39] in the data selection and creating a table. To do so, go to the Analysis page, and select the analysis view Organism Classification. Expand the metagenome selection by clicking the plus symbol next to metagenomes. Select public from the dropdown-box (to view only public data sets), and type “mouse” into the filter box. Select the two samples, and click the button with the right arrow, then the ok button. The default data visualization is “table,” so you can click the “generate” button (Figure 4.16).

After a short wait, a new tab will appear in the tabview, showing the data table with organism classifications for the two samples. The last column of this table will have a button labeled “to workbench” as the column header. Each cell in that column will have a checkbox. Checking a checkbox and clicking the “to workbench” button will send the proteins identified by that row to the workbench. Note that you have only one workbench and that putting a new set of proteins into it will replace the current content. What if you want to select all Bacteria: do you really need to click through all those checkboxes? No: you can use the grouping feature of the table, so you have

The screenshot shows the MG-RAST 'Browse Metagenomes' interface. On the left, there are two summary boxes: 'Your Data Summary' and 'Public Data Summary'. The 'Your Data Summary' box shows metrics like 'Available for analysis' (3), 'In Progress' (1654), 'Shared with you' (17799), and 'Collections' (14). The 'Public Data Summary' box provides an overview of the public dataset count (4280 metagenomes, 88 projects, etc.). Below these is a note about making datasets public. The main area is titled 'Your Collections' and includes links to 'back to all metagenomes' and 'delete selected entries'. It features a dropdown for 'display' (set to 20) and a page number 'displaying 1 - 20 of 2432'. A 'select all' checkbox is available. The collection table lists items grouped by domain, such as 'all', 'CF', 'CF2', 'CFLung', 'human', 'marine', 'Northern Line', 'null', 'plant virus', 'plant2', 'Seawater', 'St Louis - human samples', 'Twin Study', 'Unspecified Biome 4/7/2011', and 'Unspecified Biome :-('. Each row has a checkbox in the 'select' column.

collection	job name ▾	select ...
all	ObeseMouseCecumMic2005	<input type="checkbox"/>
CF	LeanMouseCecumMic2005	<input type="checkbox"/>
CF2	CFLungPat001Rep1SDVir20060505	<input type="checkbox"/>
CFLung	CFLungPat001Rep2SDVir20060505	<input type="checkbox"/>
human	CFLungPat001Rep3SDVir20060505	<input type="checkbox"/>
marine		
Northern Line		
null		
plant virus		
plant2		
Seawater		
St Louis - human samples		
Twin Study		
Unspecified Biome 4/7/2011		
Unspecified Biome :-()		
human	BGIgutGeneSet	<input type="checkbox"/>
human	human In-R	<input type="checkbox"/>
human	human In-M	<input type="checkbox"/>
human	human In-E	<input type="checkbox"/>
human	human In-D	<input type="checkbox"/>
human	human In-B	<input type="checkbox"/>
human	human In-A	<input type="checkbox"/>
human	human F2-Y	<input type="checkbox"/>
human	human F2-X	<input type="checkbox"/>
human	human F2-W	<input type="checkbox"/>
human	human F2-V	<input type="checkbox"/>

Figure 4.14: List of collections

to click only one checkbox per metagenome.

Above the table you will find a dropdown-box labeled “group table by” (Figure 4.17). Select “domain,” and the table will be grouped, so there is only one row per metagenome and domain.

Now check the two boxes in the “Bacteria” rows, and click the “to workbench” button (see Figure 4.18).

A pop-up message will appear, telling you how many proteins have been sent to the workbench. If you look at the tabular view now, you will notice that the workbench tab shows the number of proteins it currently contains (see Figure 4.19). If you click on that tab, you will get information about what the workbench contains. On this tab you will also find a “download as FASTA” button.

Besides being able to download the sequences of your selected proteins, you can also use

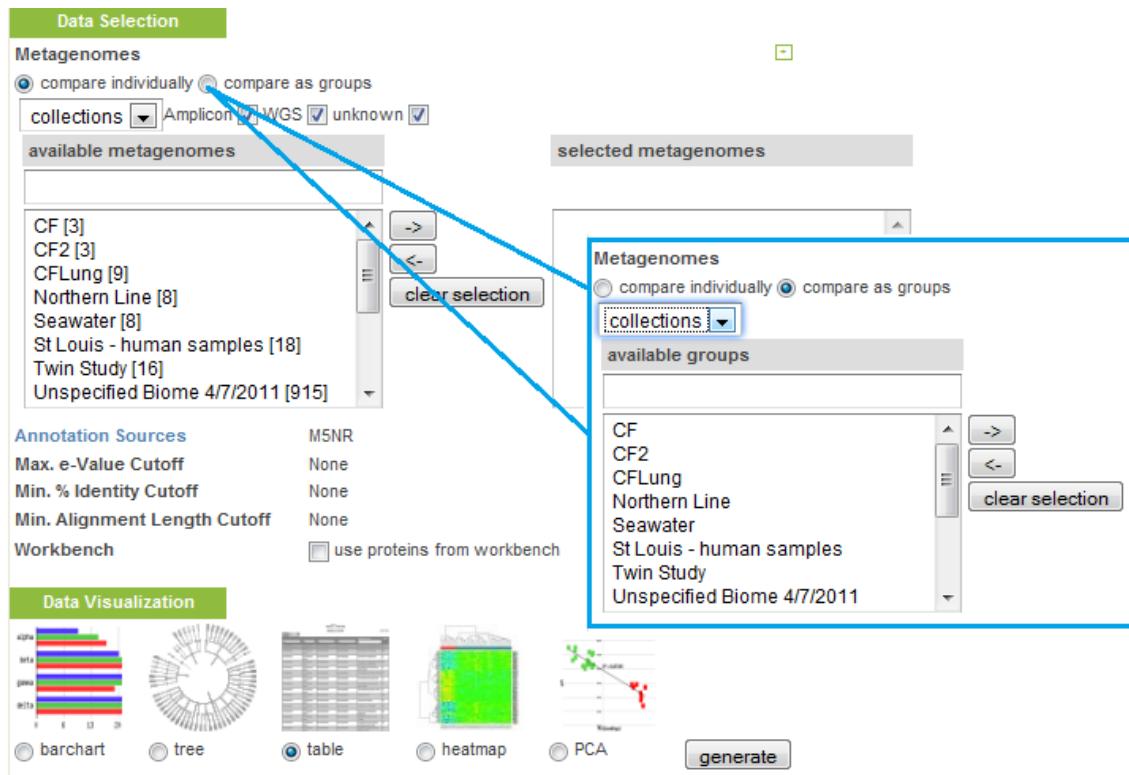


Figure 4.15: Comparing collections as groups.

them to generate other visualizations. For example, you can switch from organism to functional classification. To do so, simply check the “use proteins from workbench” box in the data selection when generating a new visualization (e.g., a circular tree using the proteins we just buffered).

The table is not the only visualization that allows you to put a subselection into the workbench. You can also use the barchart for this purpose (Figure 4.20). Simply click on the “to workbench” button next to the headline of a drilldown. Note that you cannot put the topmost barchart into the workbench because it is not yet a subselection of proteins.

## 4.6 Downloads from the Workbench

The workbench feature stores subselections of data and allows those to be used as input for further selection or displays: for example, select all *E. coli* reads and then display the functional categories present just in *E. coli* reads across multiple data sets. In addition, the workbench allows downloading the annotated reads for the subselection stored in the workbench as fasta (Figure 4.21).

Once processing datasets in MG-RAST is finished, a download page is created for the project.

**② Data Selection**

Metagenomes 4440463.3, 4440464.3

compare individually  compare as groups

public  Amplicon  MT  Unknown  WGS

available metagenomes	selected metagenomes
1-19-DNA-flx (4440275.3) 6-19-DNA-flx (4440276.3) FannLIMic20050811 (4440279.3) FannLIVir20050811 (4440280.3) RedSoudMineMic20050331 (4440281.3) BlackSoudMineMic20050331 (4440282.3) Chicken Cecum A (4440283.3) Chicken_Cecum_B (4440284.3)	LeanMouseCecumMic2005 (4440463.3) ObeseMouseCecumMic2005 (4440464.3)

**Annotation Sources** M5NR

**Max. e-Value Cutoff** 1e-5

**Min. % Identity Cutoff** 60 %

**Min. Alignment Length Cutoff** 15

**Workbench**  use features from workbench

**③ Data Visualization**

barchart  tree  table  heatmap  PCoA  rarefaction

Figure 4.16: Screenshot of the Analysis page and Workbench tab. Users can search and select metagenomes to analyze, the annotation sources and parameters to set, and the analysis and visualization they want to perform.

group table by  change  
download this

metagenome	source	domain	phy
all	M5NR	all	
4440464.3	M5NR	Archaea	Crei
4440463.3	M5NR	Archaea	Crei
4440463.3	M5NR	Archaea	Euri

Figure 4.17: Using the tables to group results.

metagenome	source	domain	abundance	avg eValue	avg % ident	avg align len	# hits	to workbench
all	M5NR	all	< 1	< -14.51	69.47	59.12	62	<input type="checkbox"/>
4440463.3	M5NR	Archaea	52	-14.51	69.47	59.12	62	<input type="checkbox"/>
4440463.3	M5NR	Archaea	108	-15.49	69.72	61.37	108	<input type="checkbox"/>
4440463.3	M5NR	Bacteria	16048	-25.93	69.66	79.16	16048	<input checked="" type="checkbox"/>
4440463.3	M5NR	Bacteria	16571	-28.57	70.35	84.17	16571	<input checked="" type="checkbox"/>
4440463.3	M5NR	Eukaryota	138	-25.53	74.19	73.10	69	<input type="checkbox"/>
4440463.3	M5NR	Eukaryota	160	-18.54	74.21	59.04	116	<input type="checkbox"/>

Figure 4.18: Use the table to select results you want to add to your workbench for further analyses.

On this page all data products created during the computation are made available as files. In addition, datasets that have been published in MG-RAST have links to an ftp site at the top of this page where you can download additional information.

## 4.7 Viewing Evidence

For individual proteins, the MG-RAST page allows users to retrieve the sequence alignments underlying the annotation transfers (see Figure 4.22). Using the M5NR [41] technology, users can retrieve alignments against the database of interest with no additional overhead.

## 4.8 MG-RAST Output

Users can access the data in MG-RAST in three ways.

- Through the website (which is authenticated)
- Through the MG-RAST API (which is also authenticated)

Workbench (31058 Features) | Getting Started | Organism table 1

workbench buffers  
31058 features

The workbench contains 31058 unique features.  
They were selected from the following table lines:

metagenome	source	domain	abundance	avg eValue	avg % ident	avg align len	# hits
4440464.3	M5NR	Bacteria	16048	-25.93	69.66	79.16	16048
4440463.3	M5NR	Bacteria	16571	-28.57	70.35	84.17	16571

download metagenome dna FASTA annotated by GenBank

display annotated hits below

Figure 4.19: View of the workbench with a summary of the proteins that have been added.

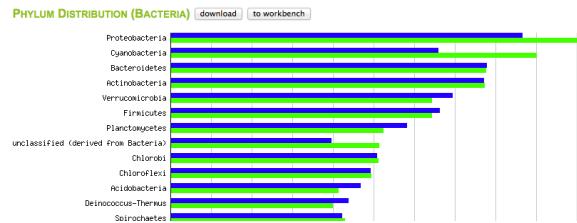


Figure 4.20: Using barcharts to download or add results to workbench.

Workbench (302 Features) | Getting Started | Organism bchart 1

workbench buffers  
302 features

The workbench contains 302 unique features.  
They were selected from a RefSeq phylogenetic bchart for the category Gammaproteobacteria

download metagenome dna FASTA annotated by RefSeq

display annotated hits below

Figure 4.21: Download of selected reads, via the workbench, using the name space of the selection.

## BLAT alignments

The sequence alignments underlying functional and organismal classification are stored in MG-RAST in an abbreviated format. This page allows re-creation of these alignments using the original parameters and tools.

**NOTE:** Since different annotation providers have different interpretation of the sequences, you can switch between **name spaces** when performing this query.

select annotation name space

fragment	organism	name space	name space ID	function	e-value <small>[?]</small>	score <small>[?]</small>	identity
4443749.3 CYOBK37TF	Calditerrivibrio nitroreducens DSM 19672	RefSeq	YP_004051066.1	alpha-glucan phosphorylase	3e-40	163 bits (420)	77/128 (60%)
				<input type="button" value="download database sequences"/>			
				<input type="button" value="download all sequences"/>			<input type="button" value="download predicted coding sequences"/>

>RefSeq: YP\_004051066.1 alpha-glucan phosphorylase [Calditerrivibrio nitroreducens DSM 19672]

Length = 850

Score = 163 bits (420), Expect = 3e-40  
Identities = 77/128 (60%), Positives = 88/128 (69%), Gaps = 0/128 (0%)

```
Query: 7 VAGVDVWLNNPLPRAESTGGMKAAANGQNLISILDGWWDEADYYQTGWPIGRGEYED 66
V GVDVWLNNP RP EASCTGMKAA NG SILDGWW E GW IG GEY D
Sbjct: 585 VRGVGVWLNNNRRPMEASCTGMKAAINGALNFSILDGWWVEGYKNNNGWSIGAGEYSD 644

Query: 67 RAYQDEVESNALYDLLEQEVAFLFYQRGSGLPHQWIQRMKQAIRLNCPQFSTQRMVLEY 126
YQD VE LYD LE E+ PLFY + GLP +W++ MX +I + C +FST RMV+EY
Sbjct: 645 PKYQDFVEGGQELYDKLENEIVPLFYAKDRSGLPREWLKMMKNSIFICCSFSTSRRMVMEY 704

Query: 127 VQRAYIPL 134
++ Y PL
Sbjct: 705 HEKYYTPL 712
```

Figure 4.22: BLAT hit details with alignment.

- Through the ftp site (which is not authenticated, and for which we put only data from public projects)

Access to private or shared data requires either password access via the web interface or a webkey generated for access via the API. All of the data files available on the website should be available via the urls returned by a query to <http://api.metagenomics.anl.gov/analysisset/mgm4450212.3>. This returns a JSON data structure with urls for all the data files. For instance, 900.abundance.organism.gz can be retrieved from <http://api.metagenomics.anl.gov/analysisset/mgm4450212.3-900-5>. Note: This is a public job; the private jobs can be accessed by providing the webkey in the auth field in

Table 4.2: Differences between the various access modalities

Service	Web Interface	API	FTP Server	Comment
public data access	Y	Y	Y	
private data access	Y	Y	N	
upload	Y	T	N	Unless specifically arranged by help desk

the GET request.

In general, we preserve the inputs and outputs for every stage of the pipeline for download and to ensure reproducibility. An example of why this is useful is the use of dereplicated reads for error estimation by DRISEE (see 2.2.3) or the LCA (see 2.5.3) algorithm being used to reinterpret the similarities for a given cluster.

### 4.8.1 Data products on the website

Users can obtain various data products from the website.

- Spreadsheets. Many of the web pages provide spreadsheets for download with the information rendered into tables or graphical displays. While most of the graphics can be downloaded directly, some require creating a static version for download (which can be achieved through a button next to the graphic). Note that the option to use a screenshot will provide images at screen resolution.
- BIOM file format exports. From the table on the Analysis page, users can download streams in BIOM [23] format, reflecting the parameter choices made. Using this approach, users can download abundance profiles in BIOM format, which can then be processed downstream with BIOM-compliant tools (e.g., QIIME [5]).
- Sequence files via the workbench. The workbench allows download of small subsets of sequences with annotations.

### 4.8.2 FTP server

All public data is made available on the FTP server. The FTP server provides a number of services:

- project – area where public data is made available for download, sorted into projects
- data – data created to enable MG-RAST (e.g., the M5NR is made available here)
- tools – tools developed by the MG-RAST team are made available here, in addition to github
- private – a private upload area, provided by MG-RAST help desk staff to users in certain situations

A project directory for `ftp://ftp.metagenomics.anl.gov/projects/128/` is shown in Figure 4.23.

## [Index of ftp://ftp.metagenomics.anl.gov/projects/128/](http://ftp.metagenomics.anl.gov/projects/128/)

[Up to higher level directory](#)

Name	Size	Last Modified
<a href="#">4447101.3</a>		5/29/12 12:00:00 AM
<a href="#">4447102.3</a>		5/29/12 12:00:00 AM
<a href="#">4447103.3</a>		5/29/12 12:00:00 AM
<a href="#">4447192.3</a>		5/29/12 12:00:00 AM
<a href="#">4447903.3</a>		5/29/12 12:00:00 AM
<a href="#">4447943.3</a>		5/29/12 12:00:00 AM
<a href="#">4447970.3</a>		5/29/12 12:00:00 AM
<a href="#">4447971.3</a>		5/29/12 12:00:00 AM
<a href="#">metadata.project-128.json</a>	100 KB	6/11/12 12:00:00 AM
<a href="#">metadata.project-128.xls</a>	15 KB	5/29/12 12:00:00 AM
<a href="#">metadata.project-128.xlsx</a>	15 KB	6/11/12 12:00:00 AM

Figure 4.23: Listing of the project directory for `ftp://ftp.metagenomics.anl.gov/projects/128/`

### 4.8.3 Downloads

One of the critical insights when developing MG-RAST version 3 was the need to make a maximum number of data products available for download for downstream analysis. For this purpose we have created the download page that contains all automatically created data products in a single location for each metagenome. In addition, a global download page provides access to all public datasets grouped by projects.

In the Appendix (see Appendix A), using a specific example (MG-RAST ID: 4465825.3), we list the data products available on the download page for each metagenome.

The general paradigm is to make all files available that are generated during the automated analysis. In addition, the user-submitted data and metadata are made available.

# Chapter 5

## Putting It All in Perspective

We have described MG-RAST, a community resource for the analysis of metagenomic sequence data. We have developed a new pipeline and environment for automated analysis of shotgun metagenomic data, as well as a series of interactive tools for comparative analysis. The pipeline is also being used for analyzing metatranscriptome data as well as amplicon data of various kinds. This service is being used by thousands of users worldwide, many contributing their data and analysis results to the community. We believe that community resources such as MG-RAST will fill a vital role in the bioinformatics ecosystem in the years to come.

### 5.1 MG-RAST” A community resource

MG-RAST has become a community clearinghouse for metagenomic data and analysis, with over 12,000 public datasets that can be freely used. Because analysis was performed in a uniform way, these datasets can serve as building blocks for new comparative analysis; so long as new datasets are analyzed similarly, results are robustly comparable between new and old dataset analysis. These datasets (and the resulting analysis data products) are made available for download and reuse as well.

Community resources like MG-RAST provide a clear value proposition to the metagenomics community. First, it enables low-cost meta-analysis. Users utilize the data products in MG-RAST as a basis for comparison without the need to reanalyze every dataset used in their studies. The high computational cost of analysis [43] makes precomputation a prerequisite for large-scale meta-analyses. In 2001, Angiuli et al. [1] determined the real currency cost of reanalysis for the over 12,000 datasets openly available on MG-RAST to be in excess of \$30 million if Amazons EC2 platform is used. This figure does not consider the 66,000 private datasets that have been analyzed with MG-RAST.

Second, it provides incentives to the community to adopt standards, in terms of both metadata and analysis approaches. Without this standardization, data products are not readily reusable, and computational costs quickly become unsustainable. We are not arguing that a single analysis is necessarily suitable for all users; rather, we are pointing out that if one particular type of analysis is run for all datasets, the results can be efficiently reused, amortizing costs. Open access to data and analyses foster community interactions that make it easier for researchers efforts to achieve consensus with respect to establishing best practices as well as identifying methods and analyses that could provide misleading results.

Third, community resources drive increased efficiency and computational performance. Community resources consolidate the demand for analysis resources sufficiently to drive innovation in algorithms and approaches. Because of this demand, the MG-RAST team has needed to scale the efficiency of their pipeline by a factor of nearly 1,000 over the past four years. This drive has caused improvements in gene calling, clustering, and sequence quality analysis, as well as many other areas. In less specialized groups with less extreme computational needs, this sort of efficiency gain would be difficult to achieve. Moreover, the large quantities of datasets that flow through the system have forced the hardening of the pipeline against a large variety of sequence pathology types that would not be readily observed in smaller systems.

We believe that our experiences in the design and operation of MG-RAST are representative of bioinformatics as a whole. The community resource model is critical if we are to benefit from the exponential growth in sequence data. This data has the potential to enable new insights into the world around us, but only if we can analyze it effectively. Only because of this approach have we been able to scale to the demands of our users effectively, analyzing over 200 billion sequences thus far.

We note that scaling to the required throughput by adding hardware to the system or simply renting time using an unoptimized pipeline on. For example, Amazons EC2 machine would not be economically feasible. The real currency cost on EC2 for the data currently analyzed in MG-RAST (26 terabasepairs) would be in excess of \$100 million using an unoptimized workflow such as CLOVR [1].

All of MG-RAST is open source and available on <https://github.com/MG-RAST>.

## 5.2 Future Work

While MG-RAST v3 is a substantial improvement over prior systems, much work remains to be done. Dataset sizes continue to increase at an exponential pace. Keeping up with this change remains a top priority, as metagenomics users continue to benefit from increased resolution of mi-

crobial communities. Upcoming versions of MG-RAST will include (1) mechanisms for speeding pipeline up using data reduction strategies that are biologically motivated; (2) opening up the data ecosystem via an API that will enable third-party development and enhancements; (3) providing distributed compute capabilities using user-provided resources; and (4) providing virtual integration of local datasets to allow comparison between local data and shared data without requiring full integration.

### **5.2.1 Roadmap**

We maintain a rough roadmap for future version of MG-RAST.

#### **version 3.5**

- provide a web services API
- develop an R client
- provide alpha version of MG-RAST remote compute client (using VMs)

#### **4.0**

- provide reviewer access tokens
- consolidate all SQL onto PostGRES
- provide beta version of MG-RAST remote compute client (using VMs)
- include IPython-based notebooks for analysis
- use AWE for all computations and SHOCK for all pipeline storage
- provide multi-metagenome recruitment plot
- convert all file access to SHOCK

#### **version 4.x**

- rewrite web interface to support many browsers
- provide BAM upload support
- provide BAM download support

- provide variation study supportr

## **version 5.0**

- provide federated SHOCK system
- provide an assembly based pipeline

## **Acknowledgments**

This work used the Magellan machine (U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research, under contract DE-AC02-06CH11357) at Argonne National Laboratory, and the PADS resource (National Science Foundation grant OCI-0821678) at the Argonne National Laboratory/University of Chicago Computation Institute. This work was supported in part by the U.S. Dept. of Energy under Contract DE-AC02-06CH11357, the Sloan Foundation (SLOAN #2010-12), NIH NIAID (HHSN272200900040C), and the NIH Roadmap HMP program (1UH2DK083993-01).

# Appendix A

## The downloadable files for each data set

### **Uploaded File(s) DNA (4465825.3.25422.fna)**

Uploaded nucleotide sequence data in FASTA format. Preprocessing

Depending on the options chosen, the preprocessing step filters sequences based on length, number of ambiguous bases and quality values if available.

#### **passed, DNA (4465825.3.100.preprocess.passed.fna)**

A FASTA formatted file containing the sequences which were accepted and will be passed on to the next stage of the analysis pipeline.

#### **removed, DNA (4465825.3.100.preprocess.removed.fna)**

A FASTA formatted file containing the sequences which were rejected and will not be passed on to the next stage of the analysis pipeline. Dereplication

The optional dereplication step removes redundant technical replicate sequences from the metagenomic sample. Technical replicates are identified by binning reads with identical first 50 base-pairs. One copy of each 50-base-pair identical bin is retained.

#### **passed, DNA (4465825.3.150.dereplication.passed.fna)**

A FASTA formatted file containing one sequence from each bin which will be passed on to the next stage of the analysis pipeline.

#### **removed, DNA (4465825.3.150.dereplication.removed.fna)**

A FASTA formatted file containing the sequences which were identified as technical replicates and will not be passed on to the next stage of the analysis pipeline. Screening

The optional screening step screens reads against model organisms using bowtie to remove reads which are similar to the genome of the selected species.

#### **passed, DNA (4465825.3.299.screen.passed.fna)**

A FASTA formatted file containing the reads which had no similarity to the selected genome and will be passed on to the next stage of the analysis pipeline. Prediction of protein

## coding sequences

Coding regions within the sequences are predicted using FragGeneScan, an ab-initio prokaryotic gene calling algorithm. Using a hidden Markov model for coding regions and non-coding regions, this step identifies the most likely reading frame and translates nucleotide sequences into amino acids sequences. The predicted coding regions, possibly more than one per fragment, are called features.

### **coding, Protein** (4465825.3.350.genecalling.coding.faa)

A amino-acid sequence FASTA formatted file containing the translations of the predicted coding regions.

### **coding, DNA** (4465825.3.350.genecalling.coding.fna)

A nucleotide sequence FASTA formatted file containing the predicted coding regions.  
RNA Clustering

Sequences from step 2 (before dereplication) are pre-screened for at least 60% identity to ribosomal sequences and then clustered at 97% identity using UCLUST. These clusters are checked for similarity against the ribosomal RNA databases (Greengenes [8], LSU and SSU from [29], and RDP [6]).

### **rna97, DNA** (4465825.3.440.cluster.rna97.fna)

A FASTA formatted file containing sequences that have at least 60% identity to ribosomal sequences and are checked for RNA similarity.

### **rna97, Cluster** (4465825.3.440.cluster.rna97.mapping)

A tab-delimited file that identifies the sequence clusters and the sequences that comprise them.

The columns making up each line in this file are:

Cluster ID, e.g. rna97\_998

Representative read ID, e.g. 11909294

List of IDs for other reads in the cluster, e.g. 11898451,11944918

List of percentage identities to the representative read sequence, e.g. 97.5%,100.0%

## **RNA similarities**

The two files labelled expand are comma- and semicolon- delimited files that provide the mappings from md5s to function and md5s to taxonomy:

### **annotated, Sims** (4465825.3.450.rna.expand.lca)

### **annotated, Sims** (4465825.3.450.rna.expand.rna)

Packaged results of the blat search against all the DNA databases with MD5 value of the database sequence hit followed by sequence or cluster ID, similarity information, annotation, organism, database name.

### **raw, Sims** (4465825.3.450.rna.sims)

This is the similarity output from BLAT. This includes the identifier for the query which is either the FASTA id or the cluster ID, and the internal identifier for the sequence that it hits.

The fields are in BLAST m8 format:

Query id (either fasta ID or cluster ID), e.g. 11847922

Hit id, e.g. lcl—501336051b4d5d412fb84afe8b7fdd87

percentage identity, e.g. 100.00

alignment length, e.g. 107

number of mismatches, e.g. 0

number of gap openings, e.g. 0

q.start, e.g. 1

q.end, e.g. 107

s.start, e.g. 1262

s.end, e.g. 1156

e-value, e.g. 1.7e-54

score in bits, e.g. 210.0

**filtered, Sims** (15:04 4465825.3.450.rna.sims.filter)

This is a filtered version of the raw Sims file above that removes all but the best hit for each data source. Gene Clustering

Protein coding sequences are clustered at 80% identity with UCLUST. This process does not remove any sequences but instead makes the similarity search step easier. Following the search, the original reads are loaded into MG-RAST for retrieval on-demand.

**aa90, Protein** (4465825.3.550.cluster\_aa90.faa)

An amino acid sequence FASTA formatted file containing the translations of one sequence from each cluster (by cluster ids starting with aa90\_) and all the unclustered (singleton) sequences with the original sequence ID.

**aa90, Cluster** (4465825.3.550.cluster\_aa90.mapping)

A tab-separated file in which each line describes a single cluster.

The fields are:

Cluster ID, e.g. aa90\_3270

protein coding sequence ID including hit location and strand, e.g. 11954908\_1\_121\_+

additional sequence ids including hit location and strand, e.g.  
11898451\_1\_119\_+,11944918\_19\_121\_+

sequence % identities, e.g. 94.9%,97.0%

Protein similarities

**annotated, Sims** (4465825.3.650.superblast.expand.lca)

The expand.lca file decodes the MD5 to the taxonomic classification it is annotated with.

The format is:

md5(s), e.g. cf036dfa9cdde3a8a4c09d7fabfd9ba5;1e538305b8319dab322b8f28da82e0a1

feature id (for singletons) or cluster id of hit including hit location and strand, e.g.

11857921\_1\_101\_-

alignment %, e.g. 70.97;70.97

alignment length, e.g. 31;31

E-value, e.g. 7.5e-05;7.5e-05

Taxonomic string, e.g. Bacteria;Actinobacteria;Actinobacteria  
(class);Coriobacterales;Coriobacteriaceae;Slackia;Slackia exigua;-

**annotated, Sims** (4465825.3.650.superblat.expand.protein)

Packaged results of the blat search against all the protein databases with MD5 value of the database sequence hit followed by sequence or cluster ID, similarity information, functional annotation, organism, database name.

Format is:

md5 (identifier for the database hit), e.g. 88848aa7224ca2f3ac117e7953edd2d9

feature id (for singletons) or cluster ID for the query, e.g. aa90\_22837

alignment % identity, e.g. 76.47

alignment length, e.g. 34

E-value, e.g. 1.3e-06

protein functional label, e.g. SsrA-binding protein

Species name associated with best protein hit, e.g. Prevotella bergensis DSM 17361 Ref-  
Seq 585502

**raw, Sims** (4465825.3.650.superblat.sims)

Blat output with sequence or cluster ID, md5 value for the sequence in the database and similarity information.

**filtered, Sims** (4465825.3.650.superblat.sims.filter)

Blat output filtered to take only the best hit from each data source.

# Appendix B

## Terms of Service

- MG-RAST is a web-based computational metagenome analysis service provided on a best-effort basis. We strive to provide correct analysis, privacy, but can not guarantee correctness of results, integrity of data or privacy. That being said, we are not responsible for any HIPPA regulations regarding human samples uploaded by users. We will try to provide as much speed as possible and will try to inform users about wait times. We will inform users about changes to the system and the underlying data.
- We reserve the right to delete non public data sets after 120 days.
- We reserve the right to reject data set that are not complying with the purpose of MG-RAST.
- We reserve the right to perform additional data analysis (e.g. search for novel sequence errors to improve our sequence quality detection, clustering to improve sequence similarity searches etc.) AND in certain cases utilize the results. We will NOT release user provided data without consent and or publish on user data before the user.
- User acknowledges the restrictions stated above and will cite MG-RAST when reporting on their work.
- User acknowledges the fact that data sharing on MG-RAST is meant as a pre-publication mechanism and we strongly encourage users to make data publicly accessible in MG-RAST once published in a journal (or after 120 days).
- User acknowledges that data (including metadata) provided is a) correct and b) user either owns the data or has the permission of the owner to upload data and or publish data on MG-RAST.

- We reserve the right to curate and update public meta data.
- We reserve the right at any time to modify this agreement. Such modifications and additional terms and conditions will be effective immediately and incorporated into this agreement. MG-RAST will make a reasonable effort to contact users via email of any changes and your continued use of MG-RAST will be deemed acceptance thereof.

# **Appendix C**

## **Tools and data used by MG-RAST**

The MG-RAST team is happy to acknowledge the use of the following great software and data products: Databases

MG-RAST uses a number of protein and ribosomal RNA databases integrated into the M5NR [41] (Wilke et al, BMC Bioinformatics 2012. Vol 13, No. 151) non-redundant database using the M5NR tools.

### **C.1 Databases**

#### **C.1.1 Protein databases**

- The SEED [28] (Overbeek et al., NAR, 2005, Vol. 33, Issue 17)
- GenBank [3] (Benson et al., NAR, 2011, Vol. 39, Database issue)
- RefSeq [30] (Pruitt et al., NAR, 2009, Vol. 37, Database issue)
- IMG/M (Markowitz et al., NAR, 2008, Vol. 36, Database issue)
- UniProt [21] (Apweiler et al., NAR, 2011, Vol. 39, Database issue)
- eggNOGG [15] (Muller et al., NAR, 2010, Vol. 38, Database issue)
- KEGG [16] (Kanehisa et al., NAR, 2008, Vol. 36, Database issue)
- PATRIC [34] (Gillespie et al., Infect. Immun., 2011, Vol. 79, no. 11)

### **C.1.2 Ribosomal RNA databases:**

- greengenes [8] (DeSantis et al., Appl Environ Microbiol., 2006, Vol. 72, no. 7)
- SILVA [29] (Pruesse et al., NAR, 2007, Vol. 35, issue 21)
- RDP [6] (Cole et al., NAR, 2009, Vol. 37, Database issue)

## **C.2 Software**

### **C.2.1 Bioinformatics codes:**

- FragGeneScan [32] (Rho et al, NAR, 2010, Vol. 38, issue 20)
- BLAT [18] (J. Kent, Genome Res, 2002, Vol. 12, No. 4)
- QIIME [5] (Caporaso et al, Nature Methods, 2010, Vol. 7, No. 5) (we also use uclust that is part of QIIME)
- Biopython
- Bowtie [20] (Langmead et al., Genome Biol. 2009, Vol 10, issue 3)
- sff\_extract, Jose Blanca and Joaquin Caizares
- Dynamic Trim, part of SolexaQA, [7] (Cox et al., BMC Bioinformatics, 2011, Vol. 11, 485)
- FastqJoin

### **C.2.2 Web/UI tools:**

- Krona [27] (Ondov et. al. BMC Bioinformatics, 2011, Vol. 12, 385)
- Raphael JavaScript Library (Dmitry Baranovskiy)
- jQuery
- Circos (Krzywinski et al., Genome Res. 2009, Vol. 19)
- cURL

### C.2.3 Behind the scenes:

- Perl
- Python
- R
- Googles V8 JavaScript engine
- Node.js
- nginx
- OpenStack

# List of Figures

1.1	Chart showing shrinking cost for DNA sequencing. This comparison with Moore's law roughly describing the development of computing costs highlights the growing gap between sequence data and the available analysis resources. Source: NHGRI .	6
1.2	Overview of processing pipeline in (left) MG-RAST v2 and (right) MG-RAST v3. In the old pipeline, metadata was rudimentary, compute steps were performed on individual reads on a 40-node cluster that was tightly coupled to the system, and similarities were computed by BLAST to yield abundance profiles that could then be compared on a per sample or per pair basis. In the new pipeline, rich metadata can be uploaded, normalization and feature prediction are performed, faster methods such as BLAT are used to compute similarities, and the resulting abundance profiles are fed into downstream pipelines on the cloud to perform community and metabolic reconstruction and to allow queries according to rich sample and functional metadata. . . . .	8
1.3	The email address for the MG-RAST project. Note that this is inserted into the document as an image, you will have to type it. . . . .	10
2.1	MG-RAST v3 data model. . . . .	13
2.2	Analysis database schema: static objects (blue) and per metagenome (variable) objects (green). . . . .	15
2.3	Details of the analysis pipeline for MG-RAST version 3 . . . . .	16
2.4	Sizes of MG-RAST jobs per month in gigabasepairs from 2007 to 2013. . . . .	16
2.5	Nucleotide histogram with biased distributions typical for an amplicon dataset. . . . .	22
2.6	Nucleotide histogram showing ideal distributions typical for a shotgun metagenome. . . . .	22
2.7	Nucleotide histogram with untrimmed barcodes. . . . .	22
2.8	Nucleotide histogram with contamination. . . . .	23

2.9	Dialogue showing the sharing mechanism. The mechanism requires a valid email address for the user with whom the data is to be shared. A list of users with access to the data is displayed at the bottom on the page. . . . .	26
2.10	Data sets shared in MG-RAST by users (orange dots), shown as connecting edges. . . . .	27
2.11	Stable URLs provided by the <code>linkin.cgi</code> mechanism for linking to MG-RAST. . . . .	27
3.1	(a) Using the web interface for a search of metagenomes for microbial mats in hotsprings (GSC-MIMS-Keywords Biome=hotspring; microbial mat), we find 6 metagenomes (refs: 4443745.3, 4443746.3, 4443747.3, 4443749.3, 4443750.3, 4443762.3). (b) Initial comparison reveals some differences in protein functional class abundance (using SEED subsystems level 1). (c) From the PCoA plot using normalized counts of functional SEED Subsystem-based functional annotations (level 2) and Bray-Curtis as metric, we attempt to find differences between two similar datasets (MG-RAST-IDs: 444749.3, 4443762.3). (d) Using exported tables with functional annotations and taxonomic mapping, we analyze the distribution of organisms observed to contain beta-lactamase and plot the abundance per species for two distinct samples. . . . .	29
3.2	Sitemap for the MG-RAST version 3 web site. On the site map the main pages are shown in blue, management pages in orange. The green boxes represent pages that are not directly accessible from the home page. . . . .	31
3.3	Browse page, enabling sorting and data search. Users can select the metadata they wish to view and search. Some of the metadata is hidden by default and can be viewed by clicking on the last column header on the right side of the table and selecting the desired columns; this can also be used to hide unwanted columns. . . . .	33
3.4	Project page, providing a summary of all data in the project and an interface for downloads. . . . .	34
3.5	Buttons displayed by Project page to dataset owner. . . . .	34
3.6	Top of the metagenome Overview page. . . . .	35
3.7	Sequences to the pipeline are classified into one of five categories: grey = failed the QC, red = unknown sequences, yellow = unknown function but protein coding, green = protein coding with known function, and blue = ribosomal RNA. For this example over 50% of sequences were either filtered by QC or failed to be recognized as either protein coding or ribosomal. . . . .	36
3.8	Information from the GSC MIxS checklist providing minimal metadata on the sample. . . . .	37

3.9	Analysis flowchart providing an overview of the fractions of sequences surviving the various steps of the automated analysis. In this case about 20% of sequences were filtered during quality control. From the remaining 37,122,128 sequences, 53.5% were predicted to be protein coding, 5.5% hit ribosomal RNA. From the predicted proteins, 76.8% could be annotated with a putative protein function. Of 32 million annotated proteins, 24 million have been assigned to a functional classification (SEED, COG, EggNOG, KEEG), representing 84% of the reads. . . . .	38
3.10	Graph showing the number of features in this dataset annotated by the different databases. The bars representing annotated reads are colored by e-value range. Different databases have different numbers of hits but can also have different types of annotation data. . . . .	39
3.11	Sample rank abundance plot by phylum. . . . .	40
3.12	Rarefaction plot showing a curve of annotated species richness. This curve is a plot of the total number of distinct species annotations as a function of the number of sequences sampled. . . . .	41
3.13	Alpha diversity plot showing the range of -diversity values in the project the data set belongs to. The min, max, and mean values are shown, with the standard deviation ranges ( and 2) in different shades. The -diversity of this metagenome is shown in red. . . . .	42
3.14	The Subsystems function piechart, showing reads classified into SEED subsystem level-one functions. In contrast to the COG, EGGNOG, and KEGG classification schemes, there are over 20 top-level subsystem categories, creating a more highly resolved “fingerprint” for the metagenome. . . . .	43
3.15	Searching for “oral health” returns 11 data sets for two projects. . . . .	44
3.16	Search results from the previous search sorted by projects. . . . .	45
3.17	Three-step process in using the Analysis page: (1) select a profile and hit (see text) type; (2) select a list of metagenomes and set annotation source and similarity parameters; (3) choose a comparison. . . . .	46
3.18	View of the data selection dialogue, with the list of four data categories expanded. . . . .	48
3.19	Boxplots of the abundance data for raw values (top) as well as values that have undergone the normalization and standardization procedures (bottom) described in the text. After normalization and standardization, samples exhibit value distributions that are much more comparable and that have a normal distribution; the normalized and standardized data are suitable for analysis with parametric tests; the raw data are not. . . . .	49

3.20	Rarefaction plot showing a curve of annotated species richness. This curve is a plot of the total number of distinct species annotations as a function of the number of sequences sampled. . . . .	50
3.21	Options available for coloring the KEGG maps. . . . .	51
3.22	Comparison of two datasets using the KEGG mapper. Parts of metabolism common are shown in purple; unique to A are in blue; unique to B are in red. . . . .	52
3.23	Selection of a genome for display, sorted by number of hits per genome. . . . .	53
3.24	Example recruitment plot with the parameters from the previous figure for <i>Actinomyces viscosus</i> C505. . . . .	54
3.25	Bar chart view comparing normalized abundance of taxa. We have expanded the Bacteria domain to display the next level of the hierarchy. . . . .	55
3.26	Tree diagram visualization option on the Analysis page. . . . .	56
3.27	Tree diagram provision for detailed information: clicking on a node in the tree diagram will display addition information to the right of the tree display. . . . .	57
3.28	Options for the tree view. . . . .	58
3.29	Tree view at order level with coloring set to phylum level. . . . .	59
3.30	Heatmap/dendrogram example in MG-RAST. The MG-RAST heatmap/dendrogram has two dendograms, one indicating the similarity/dissimilarity among metagenomic samples (x axis dendrogram) and another indicating the similarity/dissimilarity among annotation categories (e.g., functional roles; the y-axis dendrogram). . . . .	59
3.31	View of the analysis page table. . . . .	61
4.1	Dialogue requesting the user to put in a locally generated MD5 checksum for the files to identify any data corruption during the upload. . . . .	67
4.2	Temporary storage provided in Inbox before submitting data and limited editing features. . . . .	70
4.3	Information displayed by the inbox for one file (once selected). . . . .	72
4.4	Project spreadsheet. In red are required fields. Note that the 2nd row contains information on how to fill out the form. . . . .	73
4.5	The various tabs in the spreadsheet. Project, sample and one of library metagenome or library mimarks survey are required. . . . .	73
4.6	Sample tab with 3 new samples (sample1, sample2, and sample3) added. Again red text in the first row indicates required fields. Rows 1 and 2 cannot be altered. . . . .	74
4.7	View of the browse table with the collection column enabled. Clicking on the “...” at the right end of the table allows expanding the table columns. . . . .	77
4.8	Symbol for the MG-RAST metagenome browser . . . . .	78

4.9	MG-RAST metagenome browser . . . . .	78
4.10	Filtering by BIOME information. . . . .	79
4.11	A reduced list of metagenomes for one BIOME. . . . .	80
4.12	Selecting a specific project to further reduce the number of datasets being displayed. . . . .	81
4.13	Saving a collection . . . . .	82
4.14	List of collections . . . . .	83
4.15	Comparing collections as groups. . . . .	84
4.16	Screenshot of the Analysis page and Workbench tab. Users can search and select metagenomes to analyze, the annotation sources and parameters to set, and the analysis and visualization they want to perform. . . . .	85
4.17	Using the tables to group results. . . . .	86
4.18	Use the table to select results you want to add to your workbench for further analyses. . . . .	86
4.19	View of the workbench with a summary of the proteins that have been added. . . . .	87
4.20	Using barcharts to download or add results to workbench. . . . .	87
4.21	Download of selected reads, via the workbench, using the name space of the selection. . . . .	87
4.22	BLAT hit details with alignment. . . . .	88
4.23	Listing of the project directory for <code>ftp://ftp.metagenomics.anl.gov/projects/128/</code> . . . . .	90

# Glossary

**16s** 16S ribosomal RNA (or 16S rRNA) is a component of the 30S small subunit of prokaryotic ribosomes. 4, 5, 74, 114

**ADR** Artificial duplicate read. 4, 16, 20, 114

**DNA** Deoxyribonucleic acid. 4, 114

**EC2** Amazon Elastic Compute Cloud. 4, 5, 114

**MD5** The MD5 message-digest algorithm is a widely used cryptographic hash function that produces a 128-bit (16-byte) hash value. Specified in RFC 1321, MD5 has been utilized in a wide variety of security applications, and is also commonly used to check data integrity.. 4, 32, 66–68, 71, 101–103, 112, 114

**RNA** Ribonucleic acid. 4, 21, 114

**rRNA** ribosomal ribonucleic acid. 4, 12, 15, 19, 114

**SEED** The SEED effort led by Ross Overbeek is a systematic annotation effort for prokaryotic genomes using Subsystems.. 4, 7, 15, 17, 18, 29, 38, 39, 42, 60, 106, 110, 111, 114

**Subsystem** A subsystem is a set of functional roles that an annotator has decided should be thought of as related. Frequently, subsystems represent the collection of functional roles that make up a metabolic pathway, a complex (e.g., the ribosome), or a class of proteins (e.g., two-component signal-transduction proteins within *Staphylococcus aureus*). Construction of a large set of curated populated subsystems is at the center of the SEED annotation efforts.. 4, 17, 18, 29, 42, 52, 110, 114

# Bibliography

- [1] S. V. Angiuoli, M. Matalka, A. Gussman, K. Galens, M. Vangala, D. R. Riley, C. Arze, J. R. White, O. White, and W. F. Fricke. Clovr: a virtual machine for automated and portable sequence analysis from the desktop using cloud computing. *BMC Bioinformatics*, 12:356, 2011.
- [2] Ramy Aziz, Daniela Bartels, Aaron Best, Matthew DeJongh, Terrence Disz, Robert Edwards, Kevin Formsma, Svetlana Gerdes, Elizabeth Glass, Michael Kubal, Folker Meyer, Gary Olsen, Robert Olson, Andrei Osterman, Ross Overbeek, Leslie McNeil, Daniel Paarmann, Tobias Paczian, Bruce Parrello, Gordon Pusch, Claudia Reich, Rick Stevens, Olga Vassieva, Veronika Vonstein, Andreas Wilke, and Olga Zagnitko. The rast server: Rapid annotations using subsystems technology. *BMC Genomics*, 9(1):75, 2008.
- [3] D. A. Benson, M. Cavanaugh, K. Clark, I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, and E. W. Sayers. Genbank. *Nucleic Acids Res*, 41(Database issue):D36–42, 2013.
- [4] A. Bolotin, B. Quinquis, A. Sorokin, and S. D. Ehrlich. Clustered regularly interspaced short palindrome repeats (crisprs) have spacers of extrachromosomal origin. *Microbiology*, 151(Pt 8):2551–61, 2005.
- [5] J. G. Caporaso, J. Kuczynski, J. Stombaugh, K. Bittinger, F. D. Bushman, E. K. Costello, N. Fierer, A. G. Pena, J. K. Goodrich, J. I. Gordon, G. A. Huttley, S. T. Kelley, D. Knights, J. E. Koenig, R. E. Ley, C. A. Lozupone, D. McDonald, B. D. Muegge, M. Pirrung, J. Reeder, J. R. Sevinsky, P. J. Turnbaugh, W. A. Walters, J. Widmann, T. Yatsunenko, J. Zaneveld, and R. Knight. Qiime allows analysis of high-throughput community sequencing data. *Nat Methods*, 7(5):335–6, 2010.
- [6] J. R. Cole, B. Chai, T. L. Marsh, R. J. Farris, Q. Wang, S. A. Kulam, S. Chandra, D. M. McGarrell, T. M. Schmidt, G. M. Garrity, J. M. Tiedje, and Ribosomal Database Project. The ribosomal database project (RDP-II): previewing a new autoaligner that allows regular

- updates and the new prokaryotic taxonomy. *Nucleic acids research*, 31(1):442–443, January 2003.
- [7] M. P. Cox, D. A. Peterson, and P. J. Biggs. Solexaqa: At-a-glance quality assessment of illumina second-generation sequencing data. *BMC Bioinformatics*, 11:485, 2010.
  - [8] T. Z. DeSantis, P. Hugenholtz, N. Larsen, M. Rojas, E. L. Brodie, K. Keller, T. Huber, D. Dalevi, P. Hu, and G. L. Andersen. Greengenes, a Chimera-Checked 16S rRNA gene database and workbench compatible with ARB. *Appl. Environ. Microbiol.*, 72(7):5069–5072, July 2006.
  - [9] R. C. Edgar. Search and clustering orders of magnitude faster than blast. *Bioinformatics*, 26(19):2460–1, 2010.
  - [10] D. Field, L. Amaral-Zettler, G. Cochrane, J.R. Cole, P. Dawyndt, G.M. Garrity, J. Gilbert, F.O. Glöckner, L. Hirschman, and I. Karsch-Mizrachi. The genomic standards consortium. *Plos Biology*, 9(6):e1001088, 2011.
  - [11] Edgar Gabriel, Graham E. Fagg, George Bosilca, Thara Angskun, Jack J. Dongarra, Jeffrey M. Squyres, Vishal Sahay, Prabhanjan Kambadur, Brian Barrett, Andrew Lumsdaine, Ralph H. Castain, David J. Daniel, Richard L. Graham, and Timothy S. Woodall. Open MPI: Goals, concept, and design of a next generation MPI implementation. In *Proceedings, 11th European PVM/MPI Users' Group Meeting*, pages 97–104, Budapest, Hungary, September 2004.
  - [12] V. Gomez-Alvarez, T. K. Teal, and T. M. Schmidt. Systematic artifacts in metagenomes from complex microbial communities. *ISME J*, 3(11):1314–7, 2009.
  - [13] S. M. Huse, J. A. Huber, H. G. Morrison, M. L. Sogin, and D. M. Welch. Accuracy and quality of massively parallel dna pyrosequencing. *Genome Biol*, 8(7):R143, 2007.
  - [14] D. H. Huson, A. F. Auch, J. Qi, and S. C. Schuster. Megan analysis of metagenomic data. *Genome Res*, 17(3):377–86, 2007.
  - [15] L. J. Jensen, P. Julien, M. Kuhn, C. von Mering, J. Muller, T. Doerks, and P. Bork. eggNOG: automated construction and annotation of orthologous groups of genes. *Nucleic Acids Res*, 36(Database issue):D250–4, 2008.
  - [16] M. Kanehisa. The kegg database. *Novartis Found Symp*, 247:91–101; discussion 101–3, 119–28, 244–52, 2002.

- [17] K. P. Keegan, W. L. Trimble, J. Wilkening, A. Wilke, T. Harrison, M. D’Souza, and F. Meyer. A platform-independent method for detecting errors in metagenomic sequencing data: Drisee. *PLoS Comput Biol*, 8(6):e1002541, 2012.
- [18] W. J. Kent. Blat—the blast-like alignment tool. *Genome Res*, 12(4):656–64, 2002.
- [19] Murphy S Kagan L Kravitz S Lombardot T Field D Glckner FO; Genomic Standards Consortium Kottmann R, Gray T. A standard migs/mims compliant xml schema: toward the development of the genomic contextual data markup language (gcdml). *OMICS*, 12(2):115–21, 2008.
- [20] B. Langmead, C. Trapnell, M. Pop, and S. L. Salzberg. Ultrafast and memory-efficient alignment of short dna sequences to the human genome. *Genome Biol*, 10(3):R25, 2009.
- [21] Michele Magrane and Uniprot Consortium. UniProt knowledgebase: a hub of integrated protein data. *Database : the journal of biological databases and curation*, 2011, January 2011.
- [22] V. M. Markowitz, N. N. Ivanova, E. Szeto, K. Palaniappan, K. Chu, D. Dalevi, I. M. Chen, Y. Grechkin, I. Dubchak, I. Anderson, A. Lykidis, K. Mavromatis, P. Hugenholtz, and N. C. Kyrpides. Img/m: a data management and analysis system for metagenomes. *Nucleic Acids Res*, 36(Database issue):D534–8, 2008.
- [23] D. McDonald, J.C. Clemente, J. Kuczynski, J. Rideout, J. Stombaugh, D. Wendel, A. Wilke, S. Huse, J. Hufnagle, and F. Meyer. The biological observation matrix (biom) format or: how i learned to stop worrying and love the ome-ome. *Gigascience*, 2012.
- [24] F Meyer, D Paarmann, M D’Souza, R Olson, EM Glass, M Kubal, T Paczian, A Rodriguez, R Stevens, A Wilke, J Wilkening, and RA Edwards. The metagenomics rast server - a public resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics*, 9(1):386, 2008.
- [25] NHGRI. Cost per raw megabase of dna sequence, 2012.
- [26] Timothy J. Dallman Chrystala Constantinidou Saheer E Gharbia John Wain Mark J. Pallen Nicholas J. Loman, Raju V Misra. Performance comparison of benchtop high-throughput sequencing platforms. *Nature Biotechnology*, (5):434439, 2012.
- [27] B. D. Ondov, N. H. Bergman, and A. M. Phillippy. Interactive metagenomic visualization in a web browser. *BMC Bioinformatics*, 12:385, 2011.

- [28] R. Overbeek, T. Begley, R.M. Butler, J.V. Choudhuri, N. Diaz, H.-Y. Chuang, M. Cohoon, V. de Crécy-Lagard, T. Disz, R Edwards, M Fonstein, E.D. Frank, S. Gerdes, E.M. Glass, A. Goesmann, L. Krause, B. Linke, A.C. McHardy, F. Meyer, A. Hanson, D. Iwata-Reuyl, R. Jensen, N. Jamshidi, M. Kubal, N. Larsen, H. Neuweiler, C. Rückert, G. J. Olsen, R. Olson, A. Osterman, V. Portnoy, G.D. Pusch, D.A. Rodionov, J. Steiner, R. Stevens, I. Thiele, O. Vassieva, Y. Ye, O. Zagnitko, and V. Vonstein. The subsystems approach to genome annotation and its use in the project to annotate 1000 genomes. *Nucleic Acids Res*, 33(17), 2005.
- [29] Elmar Pruesse, Christian Quast, Katrin Knittel, Bernhard M. Fuchs, Wolfgang Ludwig, Jörg Peplies, and Frank Oliver O. Glöckner. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic acids research*, 35(21):7188–7196, December 2007.
- [30] K. D. Pruitt, T. Tatusova, and D. R. Maglott. NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res*, 35(Database issue), January 2007.
- [31] J. Reeder and R. Knight. The 'rare biosphere': a reality check. *Nat Methods*, 6(9):636–7, 2009.
- [32] Mina Rho, Haixu Tang, and Yuzhen Ye. Fraggenescan: Predicting genes in short and error-prone reads,. *NAR*, (in print), 2009.
- [33] C. S. Riesenfeld, P. D. Schloss, and J. Handelsman. Metagenomics: genomic analysis of microbial communities. *Annu Rev Genet*, 38:525–52, 2004.
- [34] E. E. Snyder, N. Kampanya, J. Lu, E. K. Nordberg, H. R. Karur, M. Shukla, J. Soneja, Y. Tian, T. Xue, H. Yoo, F. Zhang, C. Dharmanolla, N. V. Dongre, J. J. Gillespie, J. Hamelius, M. Hance, K. I. Huntington, D. Jukneliene, J. Koziski, L. Mackasmie, S. P. Mane, V. Nguyen, A. Purkayastha, J. Shallom, G. Yu, Y. Guo, J. Gabbard, D. Hix, A. F. Azad, S. C. Baker, S. M. Boyle, Y. Khudyakov, X. J. Meng, C. Rupprecht, J. Vinje, O. R. Crasta, M. J. Czar, A. Dickerman, J. D. Eckart, R. Kenyon, R. Will, J. C. Setubal, and B. W. Sobral. PATRIC: the VBI PathoSystems resource integration center. *Nucleic Acids Res*, 35(Database issue), January 2007.
- [35] Terry Speed. *Statistical Analysis of Gene Expression Microarray Data*. Chapman and Hall/CRC, 2003.

- [36] R. L. Tatusov, N. D. Fedorova, J. D. Jackson, A. R. Jacobs, B. Kiryutin, E. V. Koonin, D. M. Krylov, R. Mazumder, S. L. Mekhedov, A. N. Nikolskaya, B. S. Rao, S. Smirnov, A. V. Sverdlov, S. Vasudevan, Y. I. Wolf, J. J. Yin, and D. A. Natale. The cog database: an updated version includes eukaryotes. *BMC Bioinformatics*, 4:41, 2003.
- [37] Torsten Thomas, Jack Gilbert, and Folker Meyer. Metagenomics - a guide from sampling to data analysis. *Microbial Informatics and Experimentation*, 2(1):3, 2012.
- [38] W. L. Trimble, K. P. Keegan, M. D'Souza, A. Wilke, J. Wilkening, J. Gilbert, and F. Meyer. Short-read reading-frame predictors are not created equal: sequence error causes loss of signal. *BMC Bioinformatics*, 13(1):183, 2012.
- [39] P. J. Turnbaugh, R. E. Ley, M. A. Mahowald, V. Magrini, E. R. Mardis, and J. I. Gordon. An obesity-associated gut microbiome with increased capacity for energy harvest. *Nature*, 444(7122):1027–31, 2006.
- [40] W3C. File api; w3c working draft 25 october 2012, 2012.
- [41] A. Wilke, T. Harrison, J. Wilkening, D. Field, E. M. Glass, N. Kyrpides, K. Mavrommatis, and F. Meyer. The m5nr: a novel non-redundant database containing protein sequences and annotations from multiple sources and associated tools. *BMC Bioinformatics*, 13:141, 2012.
- [42] A. Wilke, J. Wilkening, E.M. Glass, N. Desai, and F. Meyer. An experience report: porting the mg-rast rapid metagenomics analysis pipeline to the cloud. *Concurrency and Computation: Practice and Experience*, 23(17):22502257, 2011.
- [43] J. Wilkening, A. Wilke, Narayan Desai, and Folker Meyer. Using clouds for metagenomics: A case study. In *IEEE Cluster 2009*, 2009.
- [44] Pelin Yilmaz, Renzo Kottmann, Dawn Field, Rob Knight, James R. Cole, Linda Amaral-Zettler, Jack A. Gilbert, Ilene Karsch-Mizrachi, Anjanette Johnston, Guy Cochrane, Robert Vaughan, Christopher Hunter, Joonhong Park, Norman Morrison, Phillip Rocca-Serra, Peter Sterk, Mani Arumugam, Laura Baumgartner, Bruce W. Birren, Martin J. Blaser, Vivien Bonazzi, Tim Booth, Peer Bork, Frederic D. Bushman, Pier Luigi Buttigieg, Patrick Chain , Elizabeth K. Costello, Heather Huot-Creasy, Peter Dawyndt, Todd DeSantis , Noah Fierer, Jed Fuhrman, Rachel E. Gallery , Dirk Gevers , Richard A. Gibbs , Michelle Gwinn Giglio , Inigo San Gil , Antonio Gonzalez3 , Jeffrey I. Gordon, Robert Guralnick , Wolfgang Haneln , Sarah Highlander , Philip Hugenholtz, Janet Jansson , Scott T. Kelley , Jerry Kennedy

, Dan Knights , Omry Koren , Justin Kuczynski , Nikos Kyriakis , Robert Larsen , Christian L. Lauber , Teresa Legg , Ruth E. Ley , Catherine A. Lozupone , Wolfgang Ludwig , Donna Lyons , Eamonn Maguire , Barbara A. Methé , Folker Meyer , Sara Nakielski , Karen E. Nelson , Diana Nemergut , Lindsay K. Neubold , Josh D. Neufeld , Anna E. Oliver , Norman R. Pace , Giriprakash Palanisamy , Jörg Peplies , Jane Peterson , Joseph Petrosino , Lita Proctor , Elmar Pruesse , Christian Quast , Jeroen Raes , Sujeevan Ratnasingham , Jacques Ravel , David A. Relman , Susanna Assunta-Sansone , Patrick D. Schloss , Lynn Schriml , Erica Sodergren , Aymé Spor , Jesse Stombaugh , James M. Tiedje , Doyle V. Ward , George M. Weinstock , Doug Wendel , Owen White , Andrew Whately , Andreas Wilke , Jennifer Wortmann , and Frank Oliver Glöckner . The “minimum information about an environmental sequence” (miens) specification. *Nature Biotechnology*, 2010.

The submitted manuscript has been created by UChicago Argonne, LLC, Operator of Argonne National Laboratory (“Argonne”). Argonne, a U.S. Department of Energy Office of Science laboratory, is operated under Contract No. DE-AC02-06CH11357. The U.S. Government retains for itself, and others acting on its behalf, a paid-up nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government.