

The (IMG) Systems for Comparative Analysis of Microbial Genomes & Metagenomes

Victor M. Markowitz *

Nikos C. Kyrpides **

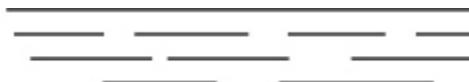
* Lawrence Berkeley National Lab

** Genome Biology Program, Joint Genome Institute

Genome sequence data processing & analysis



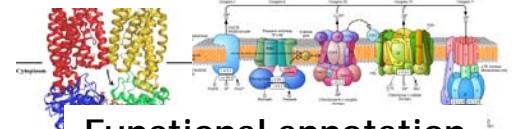
Genome data
processing



Assembly/Finishing



Gene prediction



Functional annotation



Single genome
analysis

Multi genome
annotation

Multi genome
analysis

Multi genome
annotation
review

Genome
annotation
revision

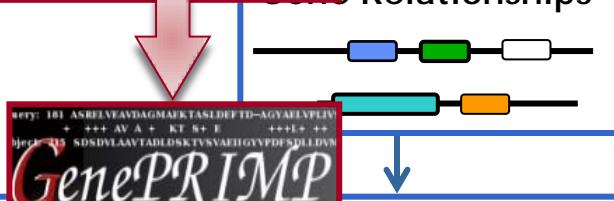


Examine Geno



Finishing QC
Gene Model
Validation

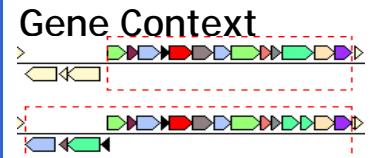
Examine Relationships



Compare Geno
• Scaffold alignment
• Genome clustering

Compare Genes
• Sequence similarity
• Phylogenetic profiler

Examine Functions



Compare Functions
• Function profiles
• Context analysis

Missing protein pro
☛ Candidate prod

Reports

- CRISPR elements
- Overlaps
- Short genes
- Long genes
- Unique genes
- Dubious genes
- Broken genes
- Transposases
- Missed genes

Missing genes
Candidate genes

Missing enzymes
☛ Candidate enzymes

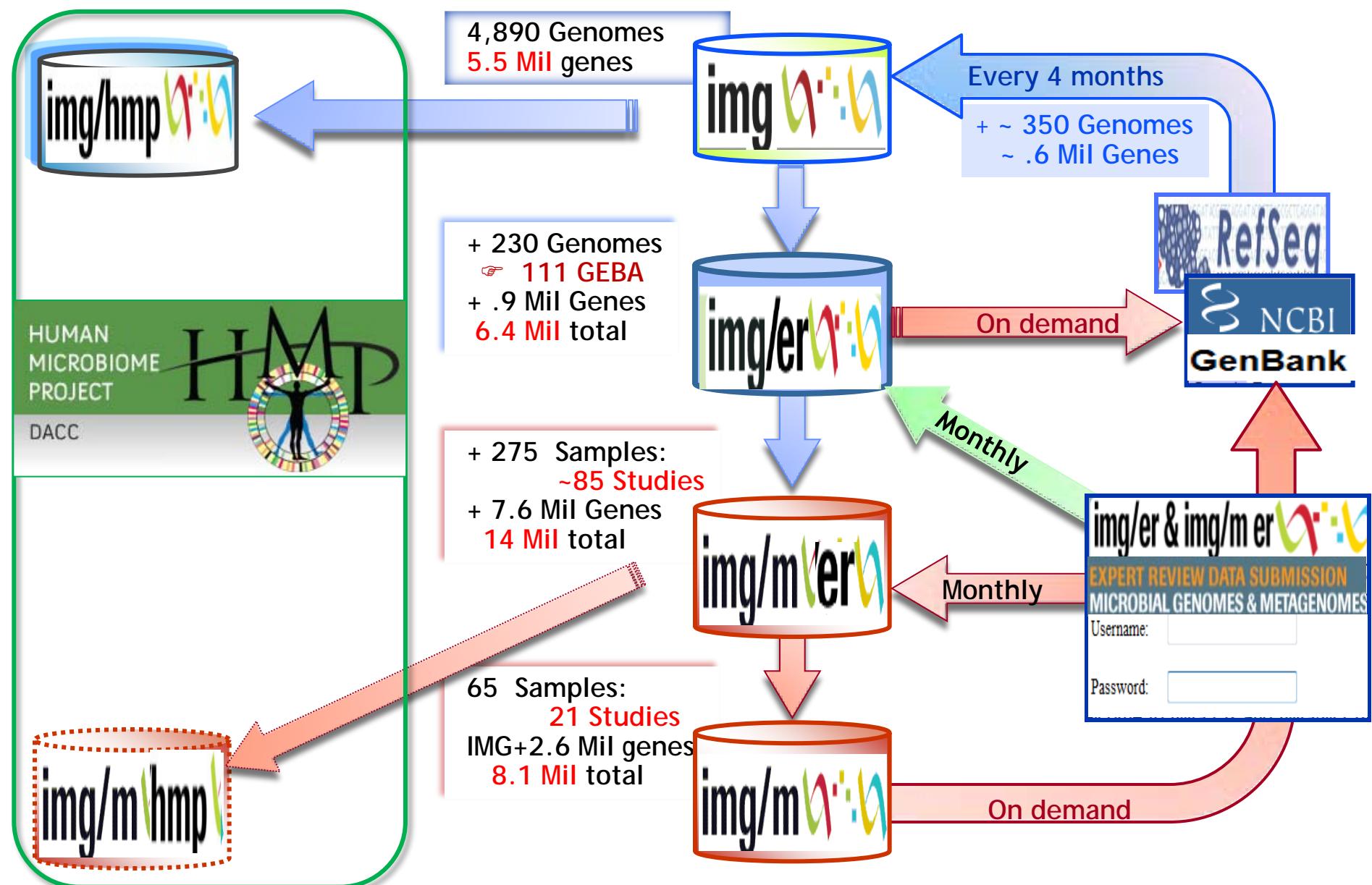
Protein products
☛ Change protein

Missing genes
Include new gene

Missing enzymes
☛ Identify enzyme



IMG systems: genome/metagenome data flow



IMG Synopsys

**img w:
INTEGRATED MICROBIAL GENOMES**

IMG Home Find Genomes Find Genes Find Functions Compare G...

IMG Genomes

	finished	draft	Total
Bacteria	781	503	1284
Archaea	56	3	59
Eukarya	19	30	49
Plasmids	974	0	974
Viruses	2524	0	2524
All Genomes	4354	536	4890

[Genome by Metadata](#)

IMG Statistics Project Map Content History

Hands on training available at the Microbial Genomics & Metagenomics Workshop

01 Archaea All None
02 Crenarchaeota All None
03 Thermoprotei All None
04 Desulfurococcales All None
05 Desulfurococcaceae All None
06 Aeropyrum All None

Metadata Categories

- [Oxygen Requirement](#) →
- [Body Sample Site](#)
- [Motility](#)
- [Sporulation](#)
- [Salinity](#)
- [Temperature Range](#)

Oxygen Requirement

Categories	Count
Aerobe	30
Anaerobe	1213
Facultative	22
Microaerophilic	10
Microaerobic	2
Obligate aerobic	61
Unknown	3

Temperature Range

Categories	Count
Hyperthermophile	30
Mesophile	1213
Psychrophile	22
Psychrotolerant	10
Psychrotrophic	2
Thermophile	61
Thermotolerant	2
Unknown	3

ornia

Summary of capabilities

- Examine individual genome annotations
 - Review, identify anomalies
 - Correct /curate using **IMG ER**
- Compare multiple genomes in terms of functional/metabolic capabilities
 - COG, Pfam, KEGG, EC#s, ...
- Analyze gene in terms of conservation, gene transfer
- Download genomes, set of contigs of interest with all their annotations

IMG: reasoning about annotations



IMG: Protein product review

Organism Information

Organism Name	Thermoplasma acidophilum
Taxon Object ID	638154521
NCBI Taxon ID	273075
NCBI Project ID	110
GOLD ID in IMG Database	Gc0003

Genome Statistics

Compare Gene Annotation

[Compare Gene Annotations](#)

[Gene Object ID](#) **638180166**

Compare
View annotation
Select filter *

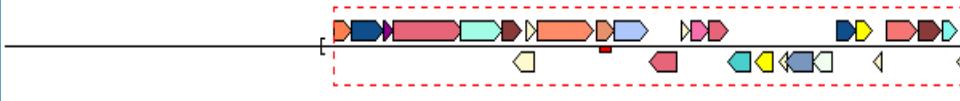
Gene Detail

Gene Information

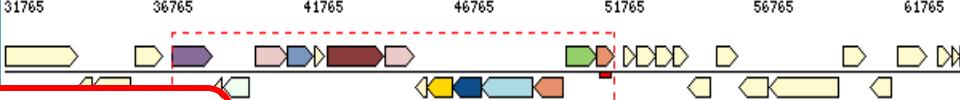
Gene Object ID
Gene Symbol
Locus Tag
Product Name
SwissProt Protein Product

Chromosomal Cassette By COG

Thermoplasma acidophilum DSM 1728: AL139299 MATCHES 24
 -10736 -5736 -736 4264 9264 14264 19264



Ferroplasma acidarmanus Fer1, unfinished sequence: NZ_AABC04000005 MATCHES 2
 31765 36765 41765 46765 51765 56765 61765



Evidence For Function Prediction

Conserved Neighborhood
Ortholog Neighborhood Viewer
Chromosomal Cassette Viewer By **COG**

FaciDRAFT_1021 : SNO glutamine amidotransferase(EC:2.6.-.) 51477 .. 52049 (190 aa) (COG0311)

Candidate Product Names for Query Gene (OID: 638180166)

Genome Name	Thermoplasma acidophilum DSM 1728
Gene Product Name	hypothetical protein
IMG Term(s)	pyridoxal phosphate synthase yaaE subunit
COG	Predicted glutamine amidotransferase involved in pyridoxine biosynthesis
Pfam	SNO glutamine amidotransferase family
KEGG Orthology (KO)	glutamine amidotransferase [EC:2.6.--]

Select	Homolog Gene	Original Product Name	IMG Term OID	IMG Term	D	C	Genome	Percent Identity	Alignment On Candidate	Alignment On Homolog	E-value
<input checked="" type="checkbox"/>	638394069	SNO glutamine amidotransferase	05446	pyridoxal phosphate synthase yaaE subunit	A	D	Ferroplasma acidarmanus Fer1	54.45			6.00e-

[Add to MyIMG Annotation](#)

MyIMG Annotation for Selected Genes

Select	Gene Object ID	Original Product Name	Annotated Product Name	Prot Description	EC Number
<input checked="" type="checkbox"/>	638180166	hypothetical protein	SNO glutamine amidotransferase		

MyIMG Annotation

Product Name	SNO glutamine amidotransferase
Prot Description	
EC Number	

[Update Annotation](#) [Delete Annotation](#)

IMG: Metabolic capability review- missing enzymes



KEGG Pathway Details

Details for *Lysine degradation*.

Enzymes in Pathway

Add Selected to Function Cart

Select	EC Number	
<input type="checkbox"/>	EC:1.1.1.35	3-hydroxya
<input type="checkbox"/>	EC:1.13.12.2	Lysine 2-m

View Pathway Map

Thermoplasma acidophilum DSM 17
Thermoplasma volcanium GSS1 (A)

Find missing enzymes

KEGG Map (for Finding Missing Enzymes)

Find Candidate Genes for Missing Function
Genome: Thermoplasma volcanium GSS1
Function: (EC:2.6.1.39) 2-aminoadipate transaminase.

LYSINE DEGRADATION

Using Homologs or Orthologs
Using KO
 Using Both

Go Reset

Using Homologs or Orthologs

Database Search Options:

- Currently selected genomes (fast)
- Whole database (slow)

1 distinct hits loaded. (6 total)

Candidate Genes for Missing Function

Genome: *Thermoplasma volcanium* GSS1

Function: (EC:2.6.1.39) 2-aminoadipate transaminase.

Select	<u>Candidate Gene</u>	<u>Candidate Gene Product</u>	<u>Enzyme for Candidate Gene</u>	<u>Ortholog Gene</u>	Add Enzyme to Candidate Gene(s) in MyIMG Annotation				Loaded: 
<input type="checkbox"/>	638190918	aspartate aminotransferase		6412764					
	Add to MyIMG Annotation		Select						

IMG: Gene model review- missing genes

Phylogenetic Profiler for Single Genes

Find Genes In*	With Homologs In	Without Homologs In	Ignoring	Taxon Name
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Archaea
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Euryarchaeota
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Thermoplasma
<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<i>Thermoplasma acidophilum</i> DSM 1728
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<i>Thermoplasma volcanium</i> GS1

Phylogenetic Profiler for

Add Selected to Gene Cart

Missing Gene? TBLASTN of t
In Genomes

Select	Result Row	Gene Object ID	Loc Ta
<input type="checkbox"/>	1	638190495	TVG000
<input checked="" type="checkbox"/>	41	638190774	TVG280

BLAST against *Thermoplasma acidophilum* DSM 1728

TBLASTN 2.2.15 [Oct-15-2006]

Sequences producing significant alignments:

Score = 103 bits (257), Expect = 9e-25
 Identities = 48/49 (97%), Positives = 49/49 (100%)
 Frame = +2

Query: 1 MAFPEAVERRLNKKICMRCYARN SIRATRCRKCGY
 MAFPEAVERRLNKKICMRCYARN SIRATRCRKCGY

Sbjct: 1365728 MAFPEAVERRLNKKICMRCYARN SIRATRCRKCGY

new gene true coords: 1365728, 1366168*

List of Potential Missing Genes

Select	Query Gene OID	Query Start Coord	Query End Coord	Subject Subject Taxon Subject Start Subject End Frame Scaffold Bit Score	Coord	Score
<input type="checkbox"/>	638190774	1	49	638154521 Thermoplasma acidophilum DSM 1728 1365728 1365874 +2 AL139299 103 9e-25 1365728 1366168*		

Add Missing Gene

Add My Missing Gene Annotation

Update Missing Gene Annotation

Missing Gene OID: 34

Genome 638154521: *Thermoplasma acidophilum* DSM 1728

Product Name	50S ribosomal protein L40E
Scaffold	AL139299
Start Coord	1365928
End Coord	1366168
Strand	+

Update My Missing Gene Annotation Refresh

Sequence Viewer

indicates potential start codon region.
 indicates possible Shine-Dalgarno region.

F1	M	K	P	R	Y	F	F	Q	S	G	B	A	V
F2	*	N	R	G	I	F	F	S	S	Q	E	K	R
F3	E	T	E	V	F	F	S	V	R	R	S	G	
GC	45	45	46	47	46	47	48	47	46	47	47	47	47
1365928	A	T	G	A	A	C	G	A	G	T	A	G	C
1365928	T	A	C	T	T	G	G	C	T	C	A	G	T
GC	55	55	54	53	54	53	52	53	54	53	53	53	53
F6	H	F	R	P	I	K	K	L	*	S	F	R	
F5	F	G	L	Y	K	K	*	D	P	S	A		
F4	S	V	S	T	N	K	E	T	L	L	L	P	

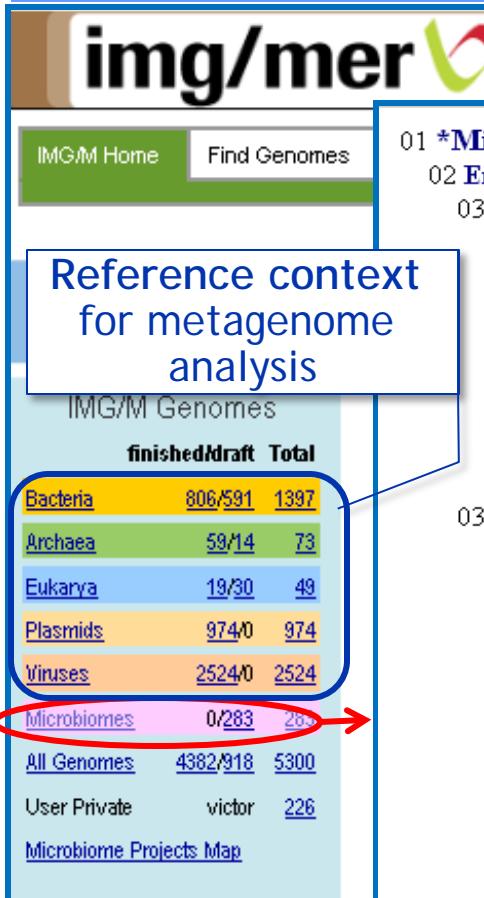
IMG/M Synopsis

INTEGRATED MICROBIAL GENOMES EXPERT REVIEW with MICROBIOME SAMPLES

Reference context for metagenome analysis

	finished	draft	Total
Bacteria	806	591	1397
Archaea	59	14	73
Eukarya	19	30	49
Plasmids	974	0	974
Viruses	2524	0	2524
Micromytes	0	283	283
All Genomes	4382	918	5300
User Private	victor		226
Microbiome Projects Map			

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



Summary of capabilities

- Examine individual metagenomes in the context of reference genomes
 - Review, identify anomalies
 - Correct /curate using **IMG/M ER**
- Compare multiple metagenomes in terms functional/metabolic capabilities
- Download metagenomes, set of contigs of interest with all their annotations

Key challenges

- Precision of metagenome annotations
 - Significance of functional comparisons
- Number/size /comparison of metagenomes



DOE Joint Genome Institute

Microbial Genomics & Metagenomics (MGM) Workshop



IMG Tutorial (Genome annotation and analysis)

[IMG - Genes and Genomes](#)

Microbial genome data analysis in IMG is set in the comparative context of multiple microbial genomes. IMG allows navigating the microbial genome data space along three key dimensions: genomes (organisms), functions (terms and pathways) and genes. In this section, IMG-based comparative analysis of gene families and genomes will be presented.

[Hands-on IMG \(exercises\)](#)

[Athanasios Lykidis](#)



Exercise solutions

[IMG - Functions and Pathways](#)

IMG has several ways for users to interact with protein functions and pathways, including Clusters of Orthologous Groups (COGs) and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathways. In addition, JGI is developing a controlled vocabulary for the representation of functions

[MyIMG](#)

The functional annotation for individual genes can be modified using the MyIMG Annotations features of MyIMG. In addition to curation of functional annotations, MyIMG provides support for uploading user genome selections that have been saved earlier from the Genome Browser or Genome Statistics and for setting systemwide user

[Gene Context Analysis in IMG](#)

[Hands-on IMG \(exercises\)](#)

[Iain Anderson](#)



[Iain Anderson](#)



[Kostas Mavrommatis](#)



[Users](#)

Exercise solutions

[Iain Anderson](#)



IMG/M Tutorial (metagenome analysis)

[IMG/M](#)

[Introduction to Metagenome analysis](#)

A snapshot of microbial community structure can be derived from analysis of metagenomic data. IMG/M methods and tools for establishing the taxonomic identity of community members will be presented along with tools for determining the fine population structure, genetic variation and genome dynamics of the dominant populations. Methods for assessing the diversity and abundance of microbial communities will be discussed.

[Natalia Ivanova](#)



[Statistical analysis of metagenomic datasets](#)

The systematic evaluation of the relative abundances of individual as well as sets of protein functions across various metagenomic datasets, can yield statistically significant deductions about over- and under-representation of protein function(s) and biological pathways in these communities. We can derive statistical methods for comparing the relative abundances of both individual as well as sets of protein families in 2 given metagenomic datasets. Statistical models for modeling individual abundances and methods for identifying protein families whose difference in abundances are statistically significant, will be presented.

[Amrita Pati](#)



[Hands-on IMG \(exercises\)](#)

[Users](#)

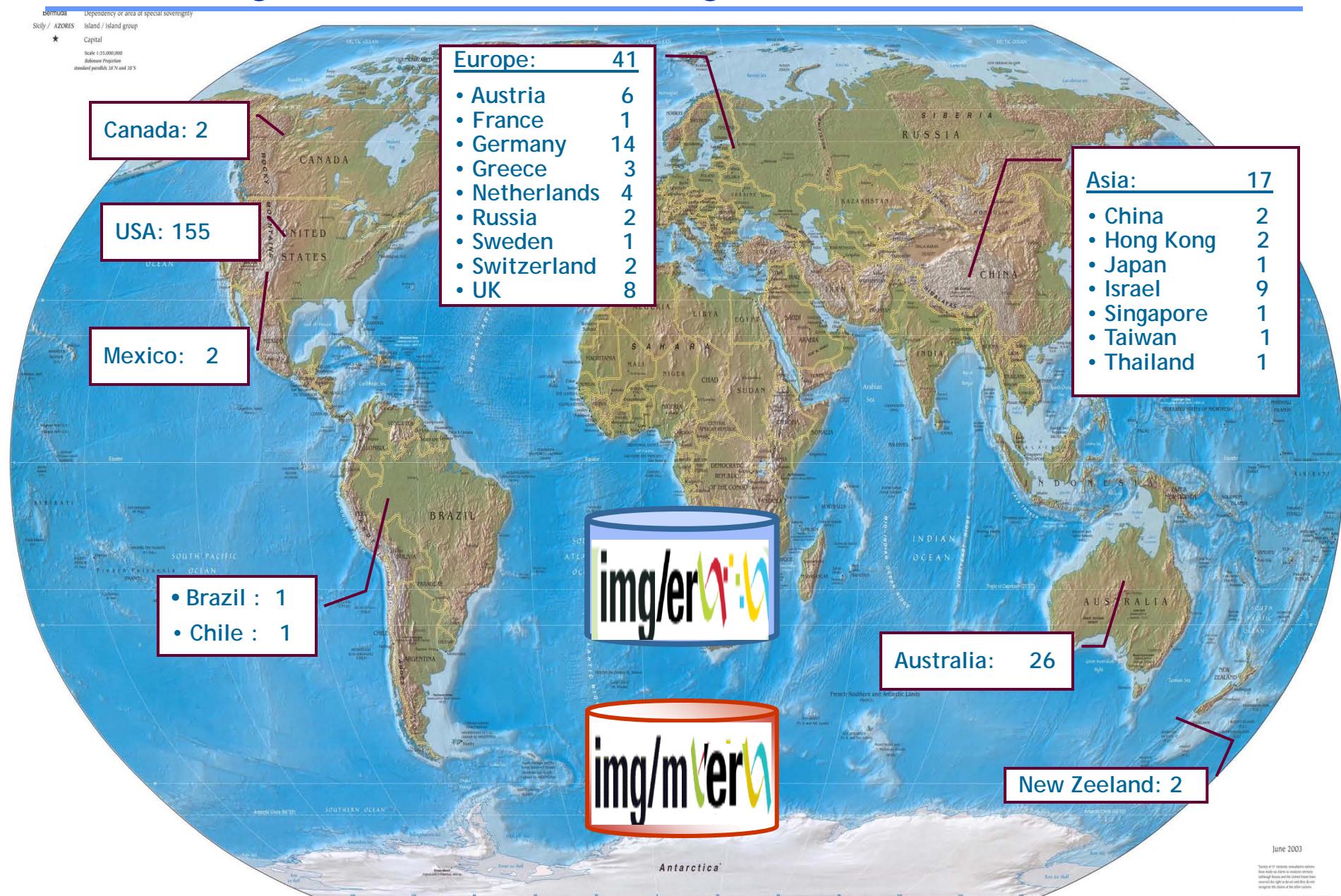
[A MetaGenome Analysis test case](#)

The methodology and steps to analyze a metagenome in IMG/M will be presented with a user case

[Athanasios Lykidis](#)



IMG ER systems: community users



Acknowledgements

Genome Biology Program



Amrita
Pati

Biological Data Management



Peter
Williams