

Integrated Microbial Genomes genomic data analysis system

Markowitz VM., Palaniappan K., Szeto E., Korzeniewski F., Werner G., Greiner A., Padki A., Zhao X., Taylor K., Hugenholtz P., Mavromatis K., Anderson I., Lykidis A., Ivanova N., Kyripides NC.

IMG's primary goal is to provide a time release of the JGI microbial genomes in the integrated context of publicly available microbial genome data.

In the long run IMG aims at providing a data management system that will:

- Support the analysis of all genomes
- Support community annotation of genomic sequences
- Support modeling of cellular networks
- Become a platform for understanding genomes
- Enable scientists to explore microbial community genomic data in the context of relevant environmental, geographical, geochemical and phylogenetic data
- Provide documentation and clarity of data, structural and operational semantics

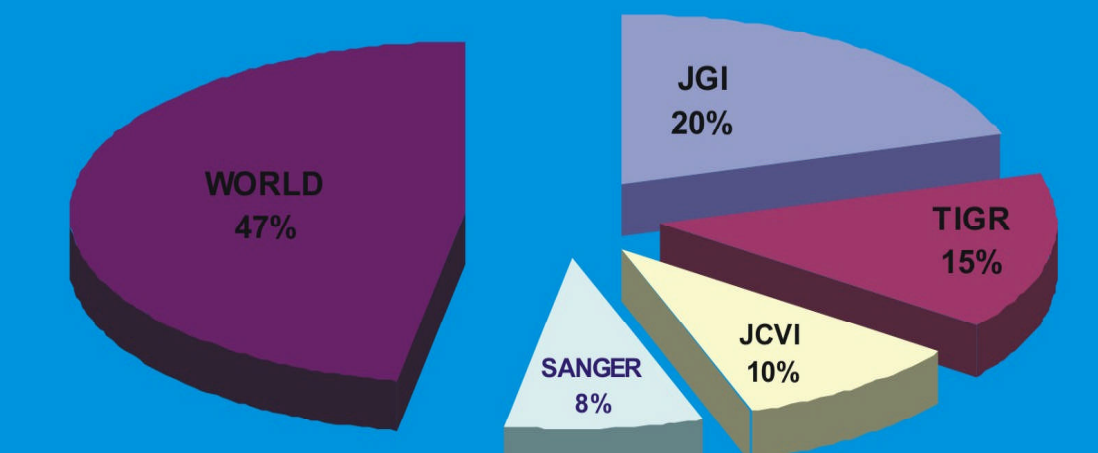
<http://img.jgi.doe.gov/>

IMG genomes

	JGI	Total
Bacteria	30 / 68	195 / 68
Archaea	0 / 4	20 / 4
Eukarya	0 / 0	9 / 0
All organisms	30 / 72	224 / 72
	Finished / draft	

Why JGI?

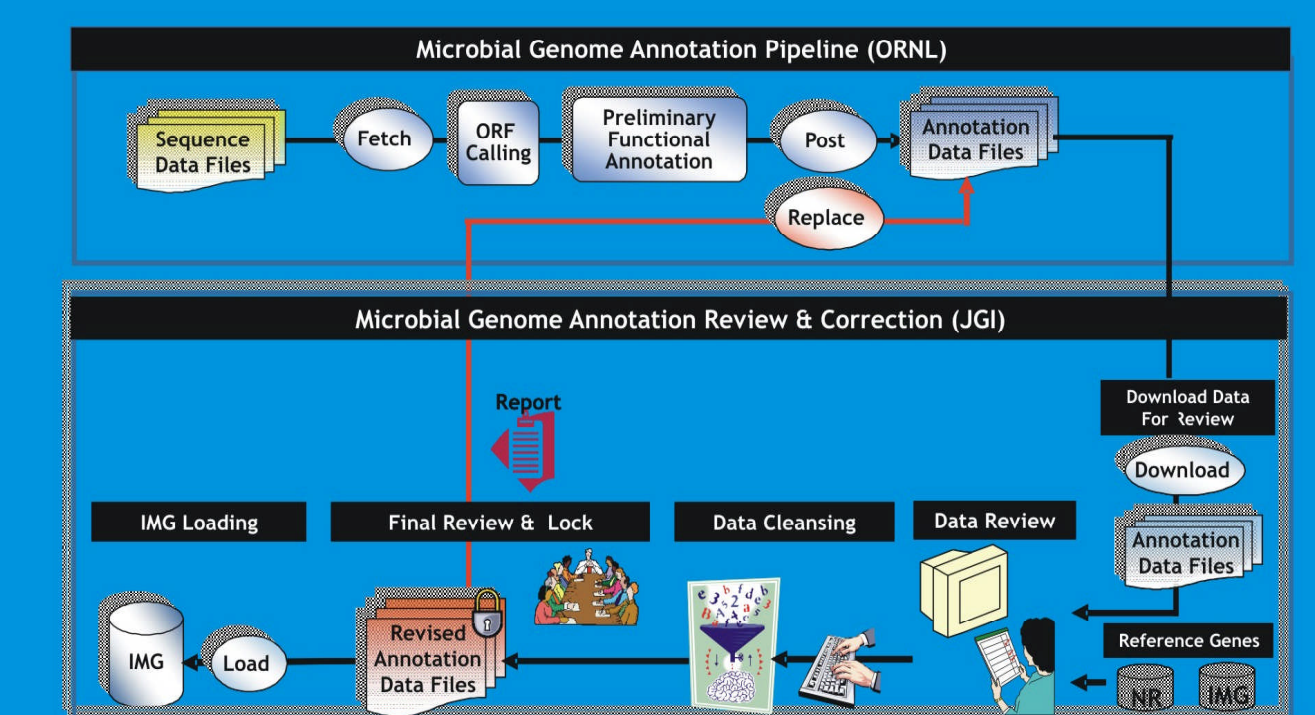
JGI is the major worldwide contributor of bacterial genomic sequences



© http://www.genomesonline.org (March 2005)

Data quality

Genomic data are manually curated before entering IMG system, providing a more consistent and clean data set



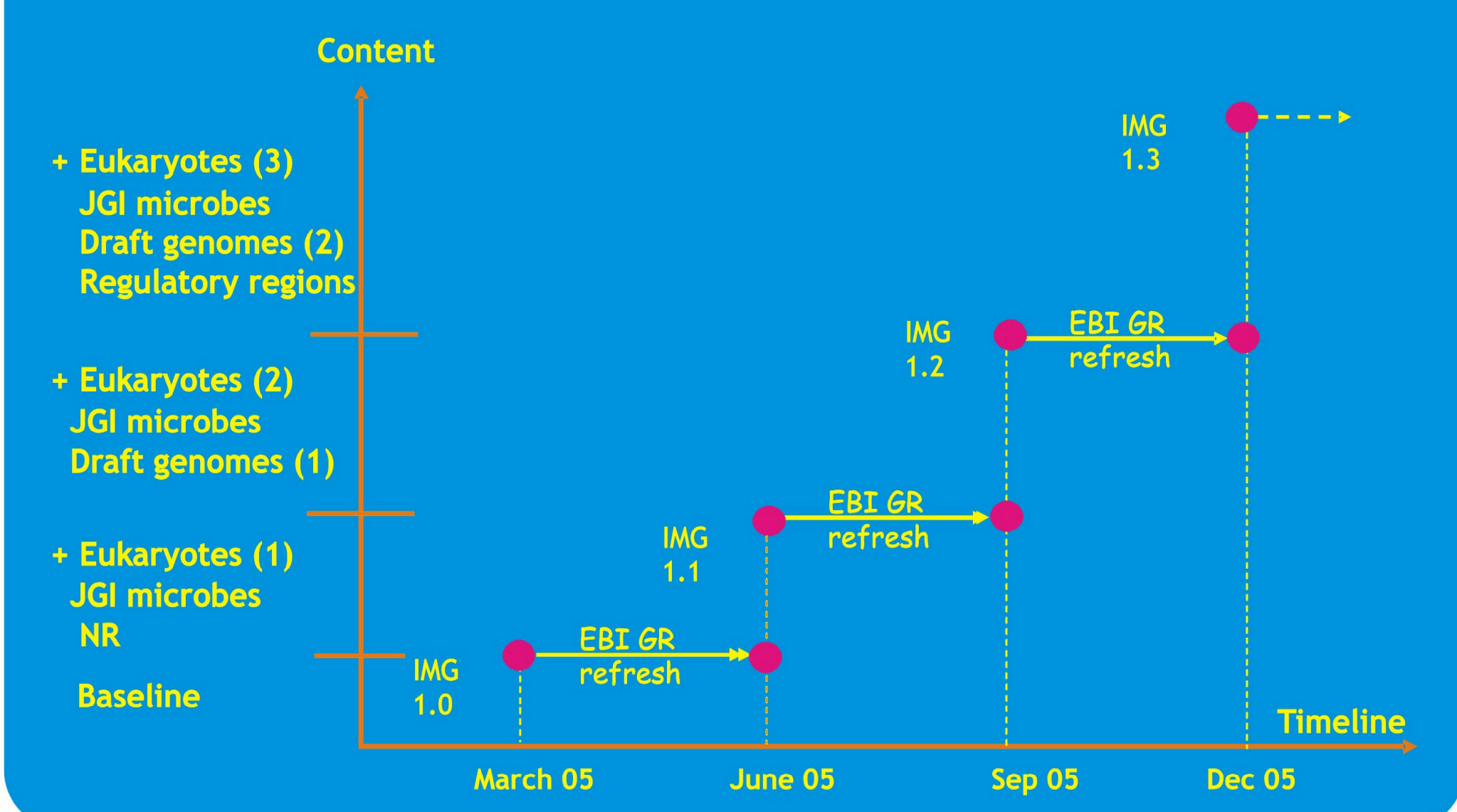
Manual curation includes identification and correction of:

- Overlapping ORFs
- ORFs with wrong 5' start sites
- Pseudogenes
- Dubious ORFs
- Missed ORFs

Future plans

Through quarterly releases, IMG aims at providing high levels of

- genome diversity, in terms of the number of genomes integrated into the system from public resources
- data quality, in terms of the coherence of annotations, the use of sound validation and correction procedures, as well as corroboration of annotations from other public microbial genome data resources.
- annotation coverage, in terms of the breadth of the functional annotations



This work was performed under the auspices of the US Department of Energy's Office of Science, Biological and Environmental Research Program and the by the University of California, Lawrence Livermore National Laboratory under Contract No. W-7405-Eng-48, Lawrence Berkeley National Laboratory under contract No. DE-AC03-76SF00098 and Los Alamos National Laboratory under contract No. W-7405-ENG-36

The Organism Browser provides multiple entries to the IMG system. You can view an alphabetical list of all organisms whose genomes are available for study, or you can view the list as a phylogenetic tree. You can then click on an individual organism to view details about its genome and follow links to useful groupings of genes within that genome (e.g., genes connected to KEGG pathways).

To select a subset of organisms, click to place a checkmark in the box next to each organism of interest, then click on "Save Selections." The selections you save will be used in multiple areas of IMG. These include the Gene Search, Gene Ortholog Neighborhoods (reached via the Gene Details page), Phylogenetic Profiler, Organism Stats, and the highlighting of orthologs in ortholog and homolog lists.

With the Phylogenetic Profiler, you can find genes in a specific organism of interest whose sequences are similar to those of genes in certain organisms and exclude those that are similar to genes in other organisms.

With the Gene Search, you can retrieve a list of genes to analyze individually or add to your gene cart. You can search by keyword or by running a BLAST query. The number of results is limited to the "Max. Gene List Results" number set in your Preferences

As you work in IMG, you can collect genes from the various IMG organisms into your Gene Cart. From any list of genes, click the checkboxes for the genes you want to add and then click "Save Selections." Once a gene is in your cart, you can use the Gene Cart page for doing alignments, comparing gene neighborhoods, and downloading sequence. You can also remove genes from the cart by selecting them and clicking "Remove Selected."

The Organism Stats reports statistics for the currently selected organisms. The data includes numbers of bases, scaffolds, and genes with various characteristics (e.g., genes connected to KEGG pathways, genes in InterPro).

Future extensions

Future viewers will show the phylogenetic distribution of one (top viewer) or multiple (bottom viewer) orthologous genes across IMG genomes. Many more features are in development.

The Gene Information table includes identifying information, locus information, chemical information about the product, and KEGG information. To see more information about the product, click the "Show All Gene Information" button.

To quickly run a BLASTp of the sequence against NR (NCBI's nonredundant protein database), SwissProt, or Pfam, click the "BLASTp against External Databases" button.

The Evidence for Function Prediction provides a view of the gene's immediate neighborhood. The gene itself is shown in red. Genes in the same positional cluster that are also in the same KEGG pathway are highlighted in green. When you move the cursor over any gene, you will see a popup box with the locus tag, gene name, and scaffold coordinates (except in Internet Explorer for the Macintosh). Click the arrow representing any gene to see the Gene Details page for that gene.

To see more of the gene of interest's chromosome, click "Show in Chromosome Viewer." To see the neighborhoods of the gene and its orthologs, click "Show ortholog neighborhood regions in user-selected organisms."

The Chromosome Viewer shows all the genes in a range of sequence, with color coding by COG group. The gene whose details you were viewing is shown in red.

Colored EC numbers on the KEGG map are links. Red links represent the current gene and point to its Gene Details page. Green links indicate genes in a positional cluster with the current gene and link to the Gene Details page for the positional cluster gene. Blue links are other genes identified in the same organism. Click a purple link to see a list of genes in the same organism with the same EC number. Orange links are "EC equivalents," or genes found in other organisms that have the same EC number assignment.

To see more of the gene of interest's chromosome, click "Show in Chromosome Viewer." To see the neighborhoods of the gene and its orthologs, click "Show ortholog neighborhood regions in user-selected organisms." You can simultaneously view neighborhoods for the gene of interest and its orthologs in organisms you've selected. The Chromosome Viewer shows all the genes in a range of sequence, with color coding by COG group. The gene whose details you were viewing is shown in red.