

Appendix 7

Descriptions of a Selected Sampling of Databases Having Relevance to the GTL Knowledgebase

Several existing databases have created effective systems for storing and analyzing genomic, metagenomic, proteomic, and other data. Having implemented successful data analysis tools, information management strategies, user capabilities, and architectures, these systems can provide viable examples of components envisioned for the GTL Knowledgebase. Moreover, many of these databases will provide important supplements to and links with other GKB capabilities. Descriptions of several such systems and their features follow.

Integrated Microbial Genomes (IMG, <http://img.jgi.doe.gov>). Developed through collaboration between the Department of Energy's (DOE) Joint Genome Institute (JGI) and the Biological Data Management and Technology Center at DOE Lawrence Berkeley National Laboratory, DOE's IMG is a data management, analysis, and annotation platform that enables the efficient comparative analysis of all complete public microbial genomes, draft or finished, produced at JGI and throughout the world. IMG currently integrates data from 4570 genomes (1155 bacteria, 56 archaea, 40 eukaryotes, 932 plasmids, and 2387 viruses), consisting of more than 4.9 million genes, with publicly available metabolic pathway collections and protein family information. IMG offer various tools for comparing genes, pathways, and functions across genomes; visualizing the physical distribution of genes within genomes; investigating the evolutionary history of genes; and developing user-defined pathways and functional categories to aid the analysis of poorly characterized genomes.

In addition to supporting the analysis of complete genomic sequences data from microbial isolates, the **Integrated Microbial Genomes with Microbiome Sampling (DOE IMG/M, <http://img.jgi.doe.gov/m>)** portal supports comparative analyses of more than 40 community sequences generated with various metagenomic sequencing technologies and data processing methods. IMG/M allows examination of profiles of functional annotations across microbial communities and isolate organisms of interest as well as analysis of strain-level heterogeneity within a species population in metagenomic data.

Phytozome (<http://www.phytozome.net>). This tool for green plant comparative genomics is a joint project of the DOE Joint Genome Institute and the Center for Integrative Genomics at the University of California, Berkeley. Phytozome provides access to nine sequenced and annotated green plant genomes, including poplar, grape, sweet sorghum, rice, soybean, green algae, moss, spikemoss, and the small flowering plant *Arabidopsis*. Clusters of orthologous and paralogous genes that represent the modern descendents of ancestral gene sets can be analyzed to explore genes associated with significant evolutionary events related to the development of green plants.

Funded by DOE, the database also offers resources for community annotation, integrates functional genomic data, and provides novel Web-based viewing and analysis tools for proteomic, gene expression microarray, and phenotype microarray data. Interactive heat maps allow users to compare microarray data for microbes under multiple stress conditions. Users also can analyze correlations between gene expressions from different experiments. Among the major new features of MicrobesOnline is the ability to search the data compendium for genes with expression profiles similar to those resulting from query profiles.

Community Cyberinfrastructure for Advanced Marine Microbial Ecology Research and Analysis (CAMERA, <http://camera.calit2.net>). The ability to explore large metagenomic datasets can enhance the research of microbial ecologists. Until recently, large-scale metagenomic analysis has been limited by the availability of computational resources that provide scientists with easy-to-use, scalable, and fully integrated Web frameworks. One such resource—CAMERA—features a rich data repository, associated bioinformatic tools, and cyberinfrastructure for conducting such analyses.

CAMERA was launched in 2007 as a collaboration between the University of California at San Diego's Calit2 division and the J. Craig Venter Institute (JCVI). With funding from the Gordon and Betty Moore Foundation, CAMERA contains 12 metagenomic datasets consisting of 14 million genomic fragment sequences. Genomic data are layered with associated geographical, temporal, and physico-chemical metadata to assist in metagenomic analyses. Additional capabilities enable homology identifications using Basic Local Alignment Search Tool (BLAST), an algorithmic resource for sequence comparisons. Furthermore, CAMERA provides graphical tools for viewing sequence regions that indicate genomic conservation and divergence and for correlating such regions with environmental factors. This new platform allows microbial researchers to begin to analyze large-scale sampling and sequencing endeavors such as JCVI's Global Ocean Sampling expedition.

Comprehensive Microbial Resource (CMR, <http://cmr.jcvi.org>). Containing more than 600 sequenced prokaryotic genomes, the CMR database provides researchers with information on inter- and intragenomic relationships for comparative genomics, genome diversity, and evolutionary studies. CMR—which is operated by JCVI—enables a wide variety of data retrievals and offers scientists numerous analytical tools for exploring the system's prokaryotic genomes. These data retrievals can be based on different gene properties that include molecular weight, hydrophobicity, guanine-cytosine (GC) content, functional-role assignments, and taxonomy. The system also has special Web-based analysis tools for precomputed homology searches, whole-genome dot plots, batch downloads, and searches across genomes using various data types.

Since consistent annotation is essential for robust genomic comparisons, CMR features primary annotations, as assigned by GenBank, and secondary annotations provided by JCVI. In addition, CMR provides comprehensive views of genes and gene annotations, genome-level structures, pathway maps, codon usage tables, GC plots, the ability to generate and visualize whole-genome alignments between two bacteria, and tabulated summary data from both individual genomes and CMR's entire genome collection.



Pathema (<http://pathema.jcvi.org>). As the Web resource for JCVI's Bioinformatics Resource Center, Pathema provides detailed curation of six target pathogens: *Bacillus anthracis*, *Clostridium botulinum*, *Burkholderia mallei*, *Burkholderia pseudomallei*, *Clostridium perfringens*, and *Entamoeba histolytica*. Funded by the National Institute of Allergy and Infectious Disease (NIAID), the center is one of eight designed to support biodefense and infectious disease research. Initially developed at JCVI, Pathema is co-maintained by the Institute of Genome Sciences (IGS) at the University of Maryland School of Medicine.

The Pathema website is separated into four main taxonomic clades: *Bacillus*, *Burkholderia*, *Clostridium*, and *Entamoeba*, allowing developers to customize clade-specific sites to each research community's needs. Pathema's dataset includes Gene Ontology assignments, metabolic pathway identification, transporter characterization, and predicted ortholog analysis and identification. The center's overarching goal is to provide a core online resource to accelerate scientific progress in understanding, detecting, diagnosing, and treating several categories of NIAID priority pathogens and other agents involved in new and re-emerging infectious diseases. Bioinformatics software with significant new capabilities, novel data types, Web resources, and analysis tools specifically geared toward biodefense are available on Pathema. Such capabilities—including intergenomic comparisons—help identify potential targets for vaccine development, therapeutics, and diagnostics. The site also serves as a focal point for the biodefense research community by disseminating data from bacterial genome-sequencing projects and by providing access to results of intergenomic comparisons.