

Category	finished/draft	Total
Bacteria	700/455	1155
Archaea	53/3	56
Eukarya	19/21	40
Plasmids	932/0	932
Viruses	2387/0	2387
All Genomes	4091/479	4570

**IMG 2.7** is the **15th** release of the Integrated Microbial Genomes (IMG) genomic data management and analysis system. **IMG 2.7** was released on **December 10<sup>th</sup>, 2008**.

## IMG 2.7 Content

### Genomes

The content of **IMG 2.7** has been updated with new microbial genomes available in **RefSeq version 31** (August 30, 2008).

**IMG 2.7** contains a total of **4,570** genomes consisting of **1,155** bacterial, **56** archaeal, **40** eukaryotic genomes, **2,387** viruses (including bacterial phages), and **932** plasmids that did not come from a specific microbial genome sequencing project. Among these genomes, **4,091** are finished genomes, and **479** are draft genomes; **328** are **JGI** sequenced genomes: **243** are finished and **63** are draft genomes.

Note that **31** microbial genomes from **IMG 2.6** were **replaced** in **IMG 2.7** because (1) a "Draft" genome has been replaced by its "Finished" version or (2) the composition of the genome has changed through the addition of new replicons (plasmids, chromosomes). For replaced genomes, the gene object identifiers (gene OIDs) for the protein-coding genes (CDS) were mapped to their new version in IMG. 2.7. See IMG [Data Evolution History](#) for details.

**Plasmid names** were curated by adding strain names to organism name when available from publications or other sources.

**tRNA** and **rRNA** genes (23S, 16S and 5S) missing from the original RefSeq genome files are added using tRNAscan-SE v1.23 for tRNA genes and similarity comparisons to existing RNA genes. In IMG 2.7 **928** tRNA and **202** rRNA genes were added in **167** genomes. Furthermore, the existing RNA genes were checked and the orientation of 24 rRNA genes and 50 tRNA genes in public genomes was corrected.

IMG 2.7 includes pathways from **MetaCyc** (<http://www.metacyc.org/>) version 12.5 (Oct 27, 2008).

# IMG Statistics

Various statistics are provided via the **IMG Statistics** link on the home page of IMG, as shown below, including: (i) **IMG Total Gene Count** which consists of counting all the genes (protein coding genes, RNA genes) in IMG, except obsolete genes, and (ii) **Protein Product Names** which consists of counting all distinct protein product names associated with (predicted for) protein coding genes (CDSs); note that this count does not include RNA or obsolete genes.

Compared to **IMG 2.6**, **IMG 2.7** contains **4,930,346 genes**, an increase of **236,083 genes**.

The screenshot displays the IMG Statistics page with the following components:

- Header:** IMG logo and "INTEGRATED MICROBIAL GENOMES".
- Navigation:** "IMG Home", "Find Genomes", "Analysis Carts", "MyIMG", "Using IMG".
- Domain Statistics Table:**

Domain	Genome Count	Gene Count	% of Total
Bacteria	1155	4207336	85.34%
Archaea	56	127731	2.59%
Eukaryota	40	501349	10.17%
Plasmid	932	25409	0.52%
Viruses	2387	68521	1.39%
<b>Total</b>	<b>4570</b>	<b>4930346<sup>1</sup></b>	<b>100.00%</b>
- Function Statistics Table:**

	Total Count	Genes with	% of Total
COG	3623	3157646	64.05%
Pfam	10340	3378710	68.53%
Enzyme	5145	561985	11.40%
TIGRfam	3418	1312207	26.61%
IMG Term	3872	867245	17.59%
GO-Molecular Function	8834	576387	11.69%
GO-Cellular Component	2182	279876	5.68%
<b>Protein Product Name<sup>2</sup></b>	<b>426451</b>	<b>2894549</b>	<b>58.71%</b>
- IMG Cluster Statistics Table:**

	Total Count	Genes with	% of Total
COG	3623	3157646	64.05%
Pfam	10340	3378710	68.53%
TIGRfam	3418	1312207	26.61%
IMG Ortholog Clusters	253345	4286544	86.94%
IMG Chromosomal Cassettes	835881	4011865	81.37%
- Conserved IMG Chromosomal Cassettes by Table:**

	Total Count	Genes with	% of Total
COG Clusters	239829	2080891	42.21%
Pfam Clusters	504025	2325893	47.18%
- Pathway Statistics Table:**

	Total Count	Genes with	% of Total
COG Pathway	77	3157646	64.05%
Kegg Pathway	213	473295	9.60%
TIGRfam Roles	105	1161106	23.55%
IMG Pathway	638	3111111	6.31%
IMG Parts List	52	325178	6.60%
MetaCyc	1395	399173	8.10%
GO-Biological Process	15064	480935	9.75%
- Project Map:** A world map showing the geographic locations of genome isolation sites. A red arrow points from the "Project Map" link in the left sidebar to this map. A tooltip is visible over the North Atlantic region, showing coordinates: "Chic: water column of the Gulf of Mexico in the Baltic: Sea at a depth of 1562 59.32222, -20.05056" and "Vibrio cholerae O157".

The **Project Map** link on the home page of IMG leads to a Google Map, as illustrated above, displaying the location of isolation sites for genomes that are associated with longitude/latitude coordinates in GOLD (<http://www.genomesonline.org/>).

## IMG 2.7 User Interface

The User Interface (UI) has been extended in order to improve its overall functionality and usability.

The main UI changes are marked on the UI Map diagram above and include:

- **New features**

- (i) Each **Organism Details** page contains a **Phylogenetic Distribution of Genes** that allows examining potential horizontally transferred genes of a genome based on the distribution of best BLAST hits of its protein-coding genes.
- (ii) **MetaCyc** pathways have been included into IMG 2.7 and are used for characterizing all genomes.
- (iii) An additional version of **Function Profile** is provided under **Compare Genomes** with a more flexible genome selection capability.
- (iv) **Missing Enzymes** can be examined within the context of a **KEGG Map** or **Function Profile** result.

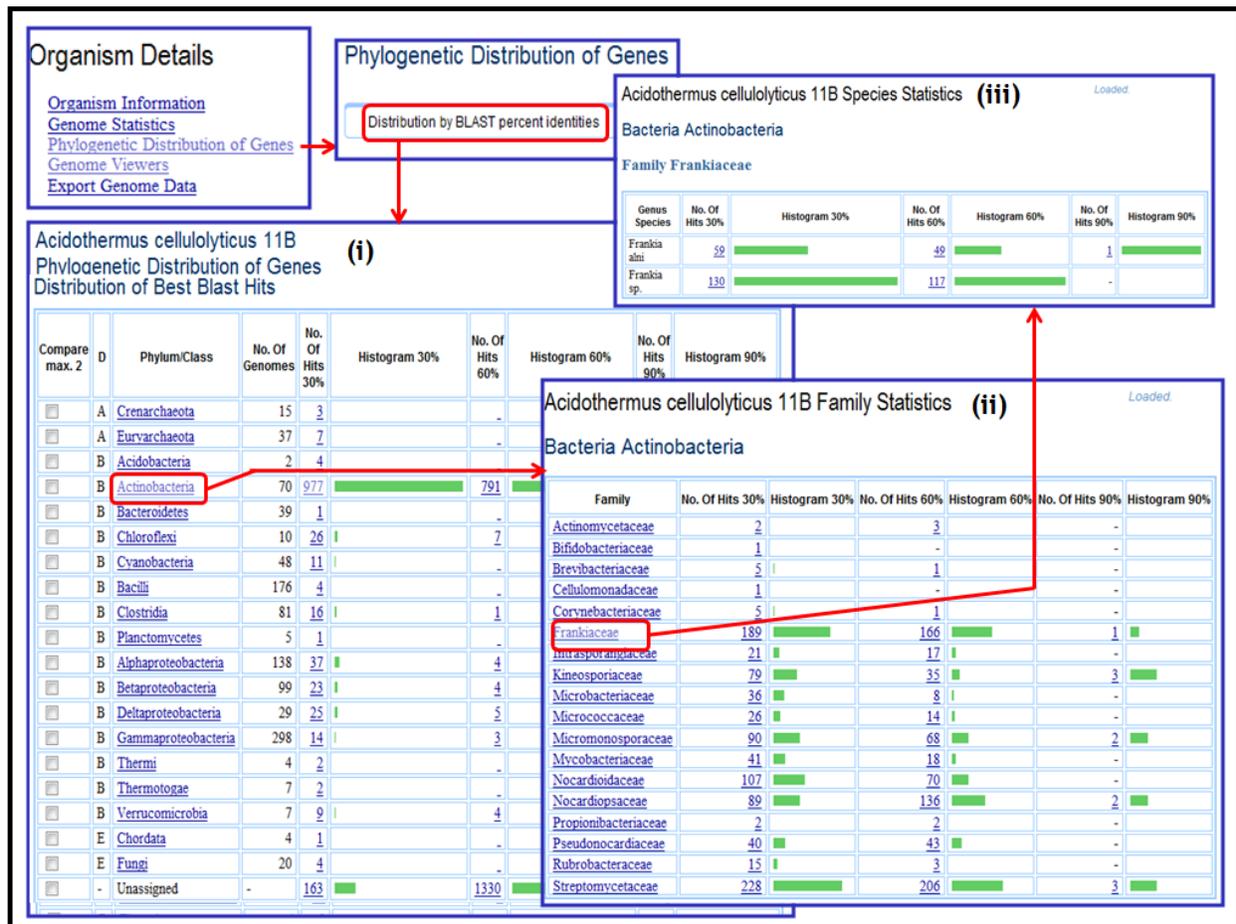
- **Extended features**

- (i) **Genome Search** and **Organism Details** extended to include additional **metadata** fields.
- (ii) **Genome Statistics** part of **Organism Details** extended with genes without enzymes but with PRIAM hits.
- (iii) **Gene Details** for pseudogenes show their individual components.
- (iv) **Function Search** extended with keyword search across all functional name filters.
- (v) **Phylogenetic Profiler** for single genes extended with **Missing Gene** search.
- (vi) **Summary Statistics** table for comparing genomes includes links to the **graphical viewers** for displaying the distribution of genes based on their COG and KEGG, Pfam, and TIGRfam associations.

# New Features

## Organism Details – Phylogenetic Distribution of Genes

**Phylogenetic Distribution of Genes** provides a glimpse into the evolutionary history of the genes in the genome based on the distribution of best BLAST hits of its protein-coding genes. The genes that were likely vertically inherited are expected to have higher sequence similarity to the genes in the genomes within the same taxonomic group, while those horizontally transferred may have their best BLAST hits to the genes in distantly related organisms. Since this tool considers best BLAST hits and does not perform phylogenetic tree reconstruction and analysis, the results can be used as a first approximation of the evolutionary history of the genes and require manual analysis to establish whether the genes of interest were indeed horizontally transferred.



**FIGURE 1. Organism Details – Phylogenetic Distribution of Genes.**

The phylogenetic distribution of best BLAST hits of protein-coding genes in a selected genome is displayed as a histogram, as shown in Figure 1(i); counts correspond to the number of genes that have best BLASTp hits to proteins of other genomes in a specific phylum or class with more than 90% identity (right column), 60-90% identity (middle column) and 30-60% identity (left column). Gene counts in the histogram are linked to the lists of genes in the selected genome that have best BLAST hit in a certain phylum/class with specified percent identity. The genes in the list can be sorted either by their oids ("Table View") or by their assignment to COGs, which

in turn can be classified according to COG Functional Categories (“COG Functional Cat.”) or COG Pathways (“COG Pathways”). The genes in the table can be selected and added to **Gene Cart** or analyzed through the corresponding **Gene Details**.

For genes that have best BLASTp hits with 30%, 60-90%, and 90% identity, the tool displays summary statistics of COG functional categories and pathways associated with these genes either across all phyla/classes or across two selected phyla/classes. The statistics show the number of genes associated with a specific COG category or pathway, as well as the percentage of these genes out of the total number of genes with affiliation to this particular class or phylum in the certain interval of percent identity of the best BLAST hits (shown in parenthesis). The phylogenetic distribution of best BLAST hits can be **projected** onto the **families** in a phylum/class (see Figure 1(ii)), and then further onto species in a family (see Figure 1(iii)).

### Find Functions – MetaCyc Pathways

IMG includes the **MetaCyc** collection of pathways available from <http://www.metacyc.org/>. IMG genomes are associated with MetaCyc pathways via enzymes predicted using PRIAM.

The screenshot displays the 'Find Functions' section of the IMG web interface, specifically the 'MetaCyc' tab. It is divided into several panels:

- MetaCyc Pathways (i):** A hierarchical tree of pathways. The '2,3-dihydroxybenzoate biosynthesis (588)' pathway is selected and highlighted with a red box. A red arrow points from this box to the 'Function List (v)' panel.
- MetaCyc Pathway Details (ii):** A table showing enzyme details for the selected pathway. A red box highlights the 'Add Selected to Function Cart' button, with a red arrow pointing to the 'Function List (iv)' panel.
- Function List (v):** A table listing selected MetaCyc pathway IDs and names:
 

Selection	Function ID	Name
<input checked="" type="checkbox"/>	MetaCyc:PWY-5787	oligomeric urushiol biosynthesis
<input checked="" type="checkbox"/>	MetaCyc:PWY-5886	4-hydroxyphenylpyruvate biosynthesis
<input checked="" type="checkbox"/>	MetaCyc:PWY-5901	2,3-dihydroxybenzoate biosynthesis
<input checked="" type="checkbox"/>	MetaCyc:PWY-981	salicylate biosynthesis
- Function List (iv):** A table listing selected enzyme EC numbers and names:
 

Selection	Function ID	Name
<input checked="" type="checkbox"/>	EC:1.3.1.28	2,3-dihydro-2,3-dihydroxybenzoate dehydrogenase.
<input checked="" type="checkbox"/>	EC:3.3.2.1	Isochorismatase.
<input checked="" type="checkbox"/>	EC:5.4.4.2	Isochorismate synthase.
- MetaCyc Pathway Diagram (iii):** A metabolic map showing the conversion of chorismate to isochorismate and then to 2,3-dihydro-2,3-dihydroxybenzoate, involving enzymes like isochorismatase and isochorismate synthase.

**FIGURE 2. Find Functions – MetaCyc Pathways.**

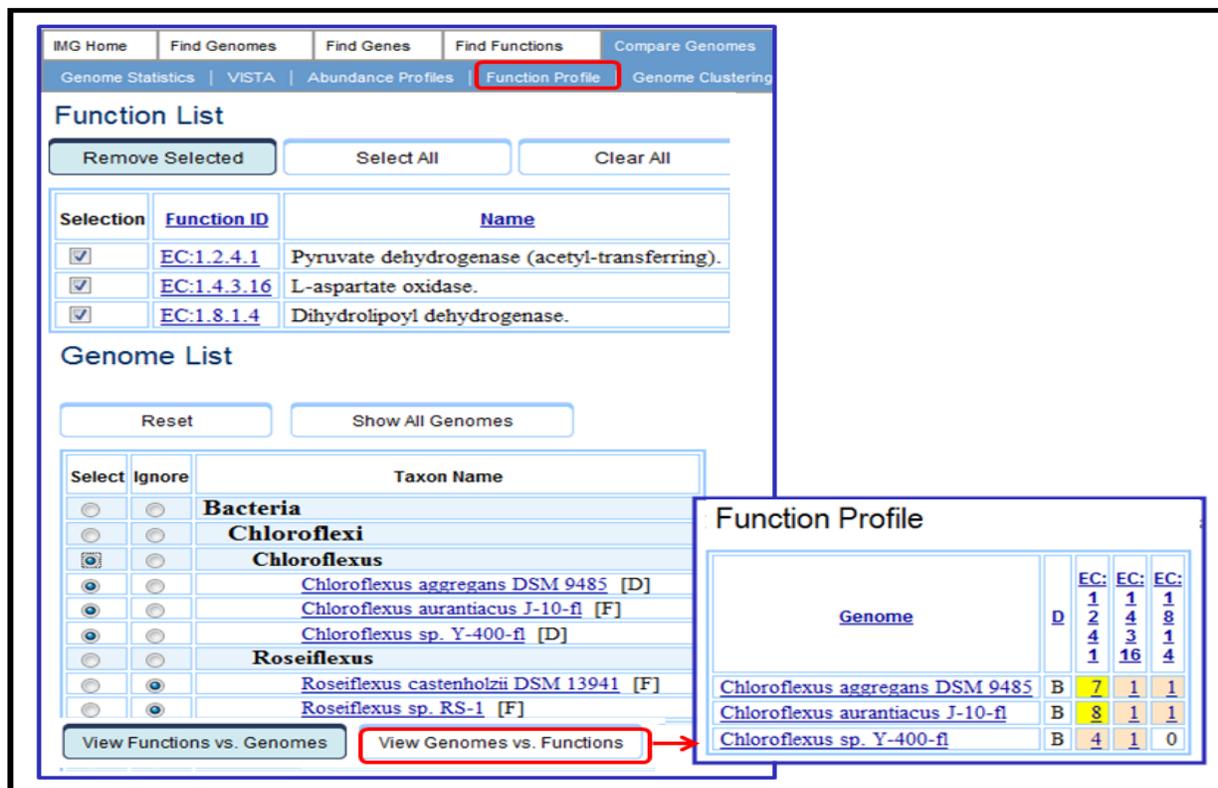
From the **Find Function** top-level menu, the **MetaCyc** option on the second-level menu leads to the **MetaCyc Pathways** browser, as shown in Figure 2(i). MetaCyc pathways associated with IMG genomes are marked and associated with the number of related IMG genes. MetaCyc

pathways are organized hierarchically within four broad categories: biosynthesis, degradation/utilization/assimilation, generation of precursor metabolites, and energy and detoxification. Each category is organized in sub-categories (classes) of pathways.

MetaCyc pathways can be examined using the **MetaCyc Pathway Details** pages, as shown in Figure 2(ii). The **MetaCyc Pathway Details** provides a link to the specification of the pathway at the MetaCyc site, as illustrated in Figure 2(iii), together with a list of the enzymes associated with a specific reaction in the pathway. For each enzyme, the number of genes associated with this enzyme is also provided, together with a link that leads to the list of these genes. By clicking on the left-column checkbox for an enzyme entry in the MetaCyc Pathway Details page, enzymes can be added to the Function Cart for further analysis, as illustrated in Figure 2(iv). Pathways can be also added directly to the Function Cart from the MetaCyc Pathways browser, as illustrated in Figure 2(v).

### Compare Genomes – Function Profile

A new version of **Function Profile** has been added under the **Compare Genomes** main menu tab. Similar to the version available under **Analysis Carts/ Functions**, the **Function Profile** displays the number (abundance) of genes associated with a function selected from the **Function List** across selected genomes, as shown in Figure 3.



**FIGURE 3. Compare Genomes – Function Profile.**

Genome selection is provided via a phylogenetically organized list of genomes, whereby genomes can be selected either individually or as phylogenetic groups. The list of genomes displays the genomes that have been selected and saved using the **Genome Browser**. The entire list of genomes can be displayed using the “Show All Genomes” button.

The **Function Profile** results are displayed in a tabular format with each column displaying the profile of a specific function across the genomes, as illustrated in the right side pane of Figure 3. Each cell in the profile-result table displays the count (abundance) of genes in a specific genome associated with a specific function, and contains a link to the associated list of genes.

### Missing Enzymes – KEGG Maps & Function Profile

Genomes may have potentially “missing” associations between enzymes and their genes. We call such associations **missing enzymes**. Missing enzymes can be examined using either a **KEGG Pathway Map** for a genome of interest or a **Functional Profile** involving genomes and enzymes of interest, as illustrated in Figure 3.

Once a KEGG pathway is selected using the **KEGG Browser** under **Find Functions**, you can view its map for a selected genome using the “Find missing enzymes” option, as illustrated in Figure 4(ii). On the **KEGG Map**, such as that shown in Figure 4(iii), enzymes that are associated with genes of the target genome are colored blue, while so called “missing” enzyme are colored either **green**, for enzymes that have a PRIAM hit to genes of the target genome, or **white** for enzymes without PRIAM hits. Clicking on a missing enzyme will lead to a **Find Candidate Genes for Missing Function** page, as shown in Figure 4(iv). Note that selection of a (green colored) missing enzyme that has a PRIAM hit enhances the chances of finding for it good candidate genes.

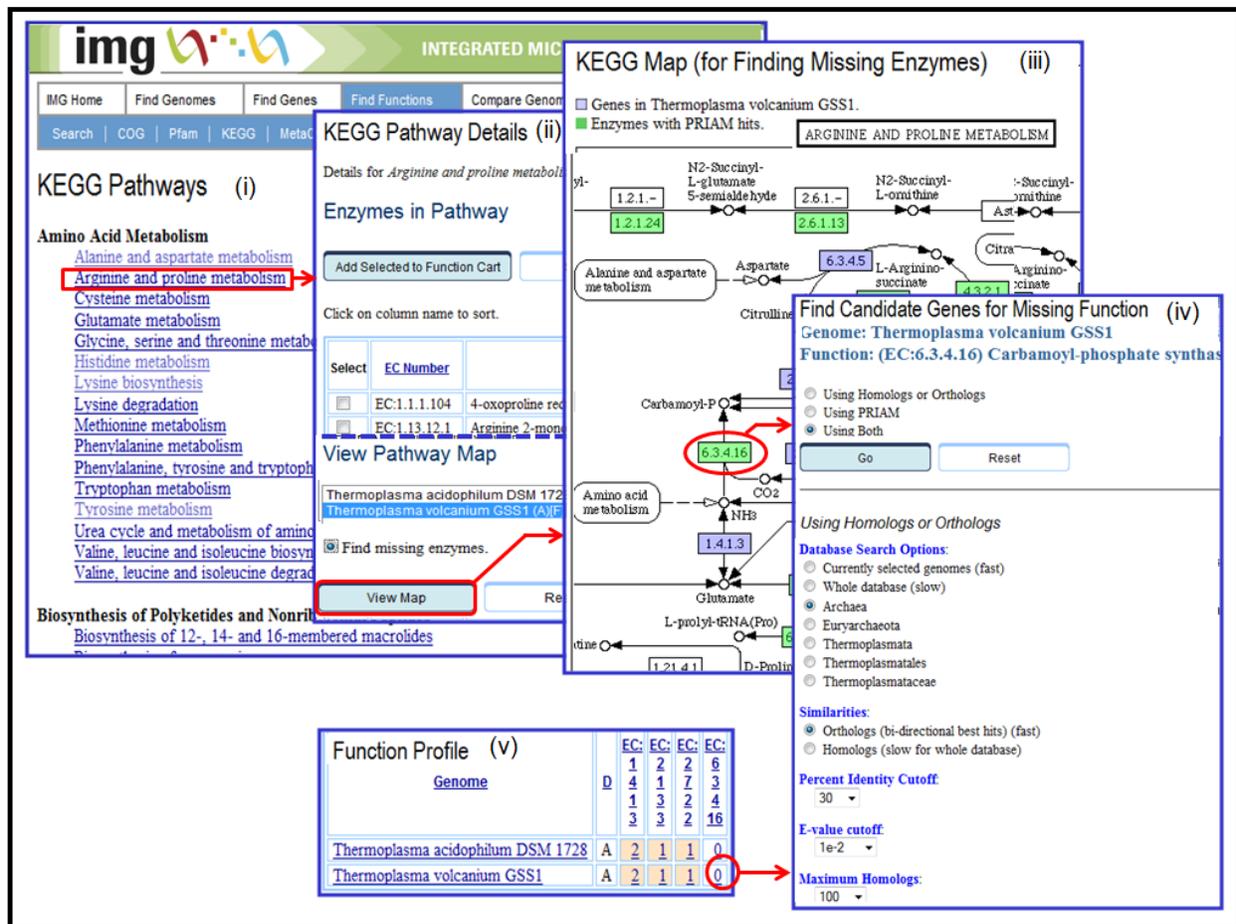


FIGURE 4. Examining Missing Enzymes via a KEGG Pathway Map or Function Profile.

You can find candidate genes of your target genome that could be associated with a missing enzyme by searching for genes that have **homologs/orthologs** associated with the missing enzyme, as illustrated in Figure 4(iv). You can search across all the genomes available in the system, across a subset of genomes within a certain domain/phyla/class, or only across the selected genomes. You can change the default values set for percent identity and e-value cutoffs and the number of retrieved homologs. Alternatively, you can employ **PRIAM** for finding genes that could be associated with the “missing” enzyme. You can change the default values set for percent identity, e-value, and percent alignment cutoffs. The result of the search for candidate genes consists of a list of genes that can be selected and included into the **Gene Cart**.

In the result for a **Function Profile** involving enzymes, missing enzymes are identified by a “0”. Clicking on the “0” identifying a missing enzyme, as shown in Figure 4(v), will also lead to a **Find Candidate Genes for Missing Function** page.

# Extended Features

## Find Genomes – Genome Search

**Genome Search** based on metadata has been extended to allow searches based on Oxygen Requirement, Motility, Sporulation, Salinity, Temperature Range, Phenotype, Disease, Relevance, Habitat, Cell Arrangement, Energy Source, and Metabolism, as shown in Figure 5.

Metadata based genome searches are carried out directly on the IMG-GOLD database that underlies the GOLD catalog of genome projects (<http://www.genomesonline.org/>). Search results include the metadata values for the fields involved in the search, as illustrated in the right pane of Figure 5.

**Metadata** in the **Organism Details** page for a genome has been also extended to include the metadata available for the corresponding project in GOLD (<http://www.genomesonline.org/>). Projects with longitude/latitude coordinates are associated with a Google Map of their location and are included into the **Project Map** of all IMG genomes, as mentioned above.

The screenshot displays the 'Genome Search by Metadata' interface. On the left, there are several filter categories with dropdown menus: Oxygen Requirement (Aerobe selected), Motility (Motile selected), Sporulation (Nonsporulating selected), Salinity (Halophile selected), Temperature Range (Hyperthermophile selected), and Phenotype (Acetic-acid selected). Each filter is followed by an 'and' connector. The main search area contains dropdown menus for Disease, Relevance, and Habitat. The 'Genome Metadata Search Results' pane on the right shows a table of search results.

Select	D	C	Genome Name	Oxygen Requirement	Motility	Salinity
<input type="checkbox"/>	B	F	<a href="#">Halobella chejuensis KCTC 2396</a>	Aerobe	Motile	Halophile
<input type="checkbox"/>	A	F	<a href="#">Halorcula marismortui ATCC 43049</a>	Aerobe	Motile	Halophile
<input type="checkbox"/>	A	F	<a href="#">Halobacterium sp. NRC-1</a>	Aerobe	Motile	Halophile
<input type="checkbox"/>	A	D	<a href="#">Halorubrum lacusprofundi ATCC 49239</a>	Aerobe	Motile	Halophile
<input type="checkbox"/>	B	D	<a href="#">Marinobacter algicola DG893</a>	Aerobe	Motile	Halophile
<input type="checkbox"/>	B	F	<a href="#">Salinibacter ruber DSM 13855</a>	Aerobe	Motile	Halophile

FIGURE 5. Find Genomes – Genome Search.

## Organism Details – Genome Statistics

The **Genome Statistics of Organism Details** has been extended with a count of “Genes without enzymes, but with PRIAM hits” which leads to a list of genes that could be associated with enzymes predicted by PRIAM, as illustrated in Figure 6. Note that only the top PRIAM hits are provided in this list.

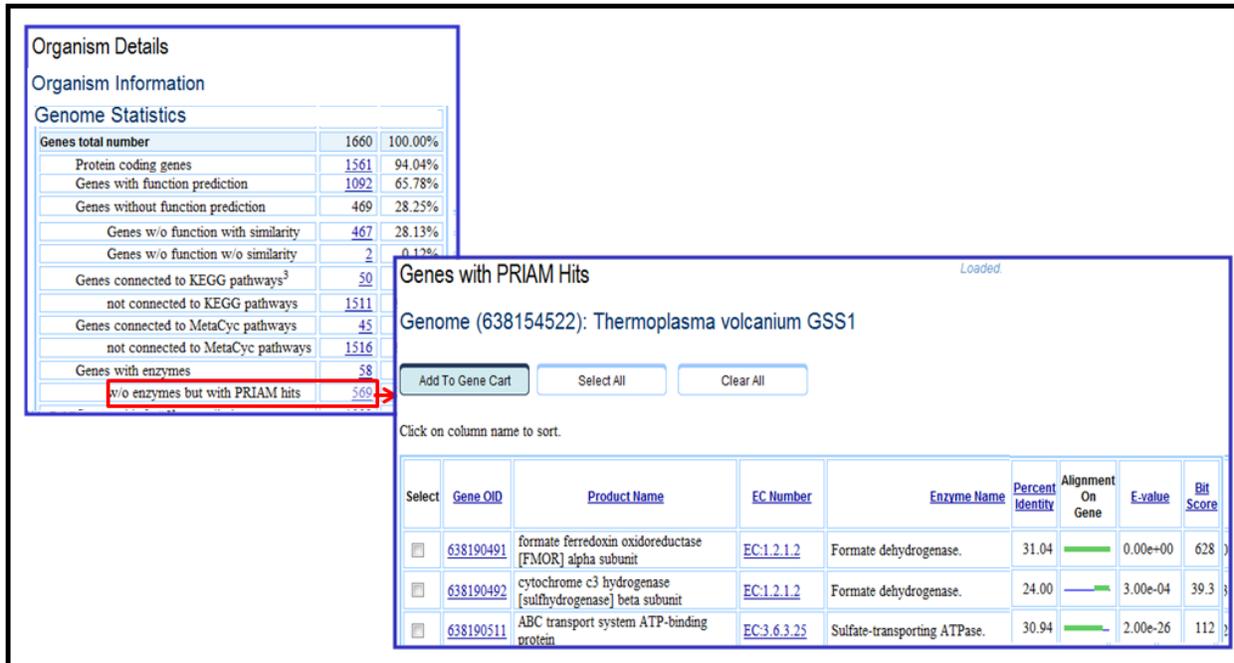


FIGURE 6. Genome Statistics- Genes without Enzymes, but with PRIAM hits.

## Gene Details – Pseudogenes

The chromosome and gene neighborhood viewers provided in the **Gene Details** page has been extended to display the components of pseudogenes, as illustrated in the right side pane of Figure 7. The previous versions of these viewers were displaying pseudogenes as one contiguous region, as shown in the left pane of Figure 7.

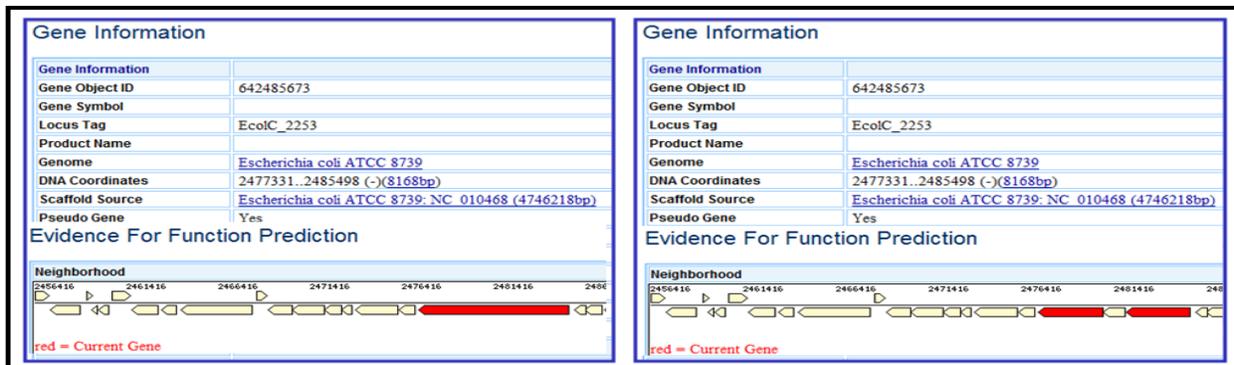


FIGURE 7. Gene Details – Pseudogenes.

## Find Functions – *Function Search*

The **Function Search** has been extended with keyword search across all functional name filters, as illustrated in Figure 8(i).

The screenshot shows the 'Find Functions' interface with the following components:

- Search Terms and Pathways (i):** A search bar containing 'kinase' and a dropdown menu for filters. The 'KEGG Pathway Enzymes' filter is selected.
- Search All Results (ii):** A table showing the number of hits for each filter.
 

Function	Hits
Gene Product Name	9213
IMG Term	1083
InterPro	423
GO	293
Enzyme Name	150
COG	129
Pfam	128
KEGG Pathway	114
TIGRfam	82
IMG Pathways	5
IMG Parts List	1
- Function Search Results (iii):** A table showing detailed results for the 'KEGG Pathway Enzymes' filter.
 

Select	Function	Name	Gene Count	Genome Count
<input type="checkbox"/>	EC:2.7.1.1	Hexokinase.	35	22
<input type="checkbox"/>	EC:2.7.1.100	S-methyl-5-thioribose kinase.	103	89
<input type="checkbox"/>	EC:2.7.1.103	Viomycin kinase.	2	2
<input type="checkbox"/>	EC:2.7.1.105	6-phosphofructo-2-kinase.	21	9
- KEGG Pathway Enzymes (iv):** A detailed view of the 'KEGG Pathway Enzymes' filter results, showing a list of enzymes under the heading 'Amino Acid Metabolism'.
  - Arginine and proline metabolism**
    - EC:2.7.2.2 Carbamate kinase. (554)
    - EC:2.7.3.2 Creatine kinase. (15)
    - EC:2.7.3.3 Arginine kinase. (19)
  - Glutamate metabolism**
    - EC:2.7.1.59 N-acetylglucosamine kinase. (40)
    - EC:2.7.2.2 Carbamate kinase. (554)
  - Glycine, serine and threonine metabolism**
    - EC:2.7.1.31 Glycerate kinase. (655)
    - EC:2.7.1.32 Choline kinase. (16)

**FIGURE 8. Find Functions – Function Search across all Function Name Filters.**

The result of this search, shown in Figure 8(ii), lists the number of hits for each function name filter, which in turns leads to detailed results for each filter, as illustrated in Figures 7(iii) and 8(iv).

## Find Genes – Phylogenetic Profiler for Single Genes

The **Phylogenetic Profiler for Single Genes** has been extended with a **Missing Gene** search capability provided on the **Phylogenetic Profiler for Single Genes Results** page, as illustrated in Figure 9(i).

**Missing Gene** search is useful for checking whether genes in a query genome (e.g., *Thermoplasma volcanium* in Figure 9) that have no homologs in other genomes (e.g., *Thermoplasma acidophilum* in Figure 9) are missing in these genomes. The search involves TBLASTn of the potentially missing gene against the sequence of these genomes, as illustrated in Figure 9(ii), where the missing gene is found in *Thermoplasma acidophilum*.

The figure shows a screenshot of the 'Phylogenetic Profiler for Single Genes' web interface. It is divided into two main sections: (i) 'Phylogenetic Profiler for Single Genes Results' and (ii) 'BLAST against Thermoplasma acidophilum DSM 1728'.

**Section (i): Phylogenetic Profiler for Single Genes Results**

At the top right, there is a table for selecting taxonomic groups:

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 type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	Thermoplasma acidophilum DSM 1728 [F]
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Thermoplasma volcanium GSS1 [F]

Below this are buttons: 'Add Selected to Gene Cart', 'Select All', and 'Clear All'.

A red box highlights the 'Missing Gene?' link. Below it, a text box says: 'Tblastn of the first selected gene in the list below against the genomes selected in Without Homologs In Genomes.'

The main results table is as follows:

Select	Result Row	Gene Object ID	Locus Tag	Gene Name	Length	COG	Enzyme	Pfam
<input type="checkbox"/>	40	638190768	TVG0275213	hypothetic protein	113aa	-	-	-
<input checked="" type="checkbox"/>	41	638190774	TVG280069	LSU ribosomal protein L40E (IMGterm)	49aa	COG1552	-	pfam01020

**Section (ii): BLAST against Thermoplasma acidophilum DSM 1728**

The BLAST results show a significant alignment with the query gene. The text includes:

```

TBLASTN 2.2.15 [Oct-15-2006]
Sequences producing significant alignments:
Score E
(bits) Value
638154521.AL139299 Thermoplasma acidophilum DSM 1728 complete ge... 103 9e-25
>638154521.AL139299 Thermoplasma acidophilum DSM 1728 complete genome.
Length = 1564906
Score = 103 bits (257), Expect = 9e-25
Identities = 48/49 (97%), Positives = 49/49 (100%)
Frame = +2
Query: 1 MAFPEAVERRLNKKICMRCYARNSIRATRCRCGYTGLRLKKKERSAGK 49
MAFPEAVERRLNKKICMRCYARNSIRATRCRCGYTGLRLKKKER+AGK
Sbjct: 1365728 MAFPEAVERRLNKKICMRCYARNSIRATRCRCGYTGLRLKKKERAAGK 1365874
    
```

FIGURE 9. Find Genes – Phylogenetic Profiler for Single Genes.

## Compare Genomes – Summary Statistics

The **Summary Statistics** table for comparing genomes selected using the **Genome Browser** includes links to **graphical viewers** for displaying the distribution of genes based on their COG and KEGG, Pfam, and TIGRfam associations, as illustrated in Figure 10.

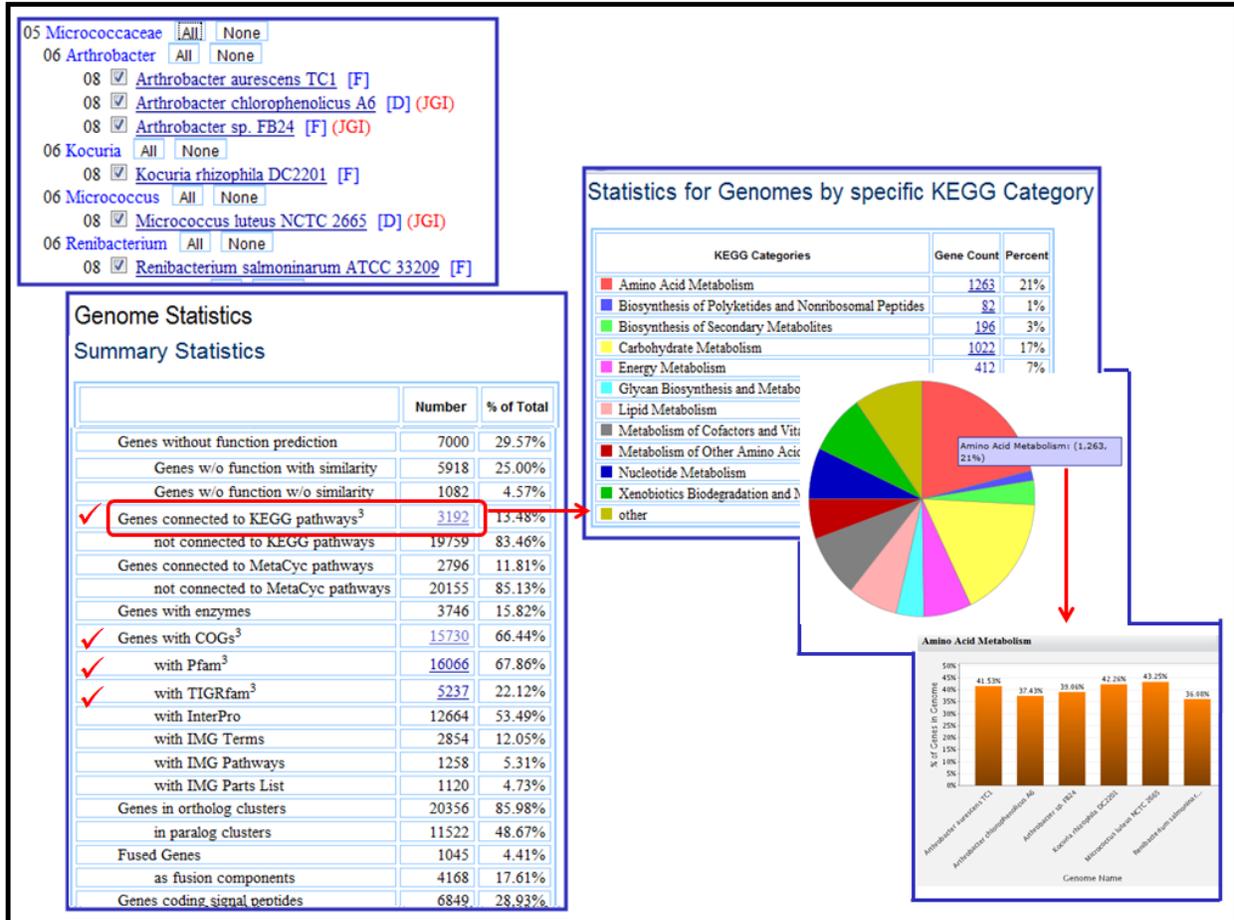


FIGURE 10. Compare Genomes – Summary Statistics.

For example, clicking on the number of genes connected to KEGG pathways, as shown in Figure 10, will display in a tabular and pie chart format the count of genes associated with each KEGG category across all selected genomes, as illustrated in the right pane of Figure 10. Clicking on a KEGG category on the *pie chart* or on the colored coded square for a KEGG category in the table will display a *bar chart* with the percent of genes for each genome associated with that KEGG category, as illustrated on the lower side left pane of Figure 10.