

## GenMAPP Gene Database for *Escherichia coli* K12

Ec-K12-Std\_External\_20090529.gdb

### ReadMe

Last revised: 6/4/09

This document contains the following:

1. Overview of GenMAPP application and accessory programs
2. System Requirements and Compatibility
3. Installation Instructions
4. Gene Database Specifications
  - a. Gene ID Systems
  - b. Species
  - c. Data Sources and Versions
  - d. Database Report
5. Contact Information for support, bug reports, feature requests
6. Release notes
  - a. Current version: Ec-K12-Std\_External\_20090529.gdb
  - b. Previous version: Ec-K12-Std\_External\_20060731x.gdb
  - c. Previous version: Ec-K12-Std\_External\_20060731.gdb
7. Database Schema Diagram

#### 1. Overview of the GenMAPP application and accessory programs

GenMAPP (Gene Map Annotator and Pathway Profiler) is a free computer application for viewing and analyzing DNA microarray and other genomic and proteomic data on biological pathways. MAPPFinder is an accessory program that works with GenMAPP and Gene Ontology to identify global biological trends in gene expression data. The GenMAPP Gene Database (file with the extension *.gdb*) is used to relate gene IDs on MAPPs (*.mapp*, representations of pathways and other functional groupings of genes) to data in Expression Datasets (*.gex*, DNA microarray or other high-throughput data). GenMAPP is a stand-alone application that requires the Gene Database, MAPPs, and Expression Dataset files to be stored on the user's computer. GenMAPP and its accessory programs and files may be downloaded from <<http://www.GenMAPP.org>>. GenMAPP requires a separate Gene Database for each species. This ReadMe describes a Gene Database for *Escherichia coli* K12 that was built by the Loyola Marymount University (LMU) Bioinformatics Group using the program GenMAPP Builder 2.0, part of the open source XMLPipeDB project <<http://xmlpipedb.cs.lmu.edu/>>.

#### 2. System Requirements and Compatibility:

- This Gene Database is compatible with GenMAPP 2.0 and 2.1 and MAPPFinder 2.0. These programs can be downloaded from <<http://www.genmapp.org>>.
- System Requirements for GenMAPP 2.0/2.1 and MAPPFinder 2.0:  
Operating System: Windows 98 or higher, Windows NT 4.0 or higher (2000, XP, etc)  
Monitor Resolution: 800 X 600 screen or greater (SVGA)  
Internet Browser: Microsoft Internet Explorer 5.0 or later  
Minimum hardware configuration:  
Memory: 128 MB (512 MB or more recommended)  
Processor: Pentium III  
Disk Space: 300 MB disk (more recommended if multiple databases will be used)

### 3. Installation Instructions

- Extract the zipped archive and place the file “Ec-K12-Std\_External\_20090529.gdb” in the folder you use to store Gene Databases for GenMAPP. If you accept the default folder during the GenMAPP installation process, this folder will be C:\GenMAPP 2 Data\Gene Databases.
- To use the Gene Database, launch GenMAPP and go to the menu item *Data > Choose Gene Database*. Alternatively, you can launch MAPPFinder and go to the menu item *File > Choose Gene Database*.

### 4. Gene Database Specifications

#### a. Gene ID Systems

This *Escherichia coli* K12 Gene Database is “UniProt-centric” in that the main data source (primary ID system) for gene IDs and annotations is the UniProt complete proteome set for *Escherichia coli* K12, made available as an XML download by the Integr8 resource. In addition to UniProt IDs, this database provides the following proper gene ID systems that were cross-referenced by the UniProt data: Blattner, EchoBASE, EcoGene, GeneId (NCBI), RefSeq (protein IDs of the form AP\_#####, NP\_#####, or YP\_#####), and W3110. It also supplies UniProt-derived annotation links from the following systems: EMBL, InterPro, PDB, and Pfam. The Gene Ontology data has been acquired directly from the Gene Ontology Project. The GOA project was used to link Gene Ontology terms to UniProt IDs. Links to data sources are listed in the section below. The following SystemCodes are used:

Proper ID System	SystemCode
UniProt	S
Blattner	Ln
EchoBASE	Ec
EcoGene	Eg
GeneId	L
RefSeq	Q
W3110	W3
Affy	X

#### b. Species

This Gene Database is based on the UniProt proteome set for *Escherichia coli* K12, taxon ID 83333. Two substrains of *E. coli* K12 have had their genomes sequenced, MG1655 (taxon ID 511145) and W3110 (taxon ID 316407). Both of these strains are represented in the UniProt proteome set and are thus both represented in this Gene Database. Blattner IDs are the primary ID system for the MG1655 strain and W3110 IDs are the primary ID system for the W3110 strain. Each of these ID systems has been cross-referenced to all other proper and improper ID systems in the Gene Database.

#### c. Data Sources and Versions

- This *Escherichia coli* K12 Gene Database was built on May 29, 2009; this build date is reflected in the filename “Ec-K12-Std\_External\_20090529.gdb.” All date fields internal to the Gene Database (and not usually seen by regular GenMAPP users) have been filled with this build date.
- UniProt complete proteome set for *Escherichia coli* K12, made available as an XML download by the Integr8 resource:  
<http://www.ebi.ac.uk/integr8/FtpSearch.do?orgProteomeId=18>  
 Filename: “18.E\_coli\_K12.xml” (downloaded as a compressed .gz file and extracted)  
 Version information for the proteome sets can be found at  
<http://www.ebi.ac.uk/integr8/HelpAction.do?action=searchById&refId=5>  
 The proteome set used for this version of the *Escherichia coli* K12 Gene Database was based

on UniProt Knowledgebase release 14.7 and InterPro Knowledgebase release 18.0 on January 20, 2009.

- Gene Ontology gene associations are provided by the GOA project: <<http://www.ebi.ac.uk/GOA/>> as a tab-delimited text file. The *Escherichia coli* K12 GOA file was accessed from the Integr8 proteome set download page: <<http://www.ebi.ac.uk/integr8/FtpSearch.do?orgProteomeId=18>> Filename: "18.E\_coli\_K12\_.goa". The GOA file for this version of the *Escherichia coli* K12 Gene Database was based on the GOA Proteome Sets 45.0 released in January 2009.
- Gene Ontology data is downloaded from <<http://www.geneontology.org/GO.downloads.ontology.shtml>> Data is released daily. For this version of the *Escherichia coli* K12 Gene Database we used the February 4, 2009 release. Filename: "go\_daily-termdb.obo-xml.gz".
- Affymetrix probe set identifiers from the Affy table in the Ec-K12-Std\_External\_20060731x.gdb Gene Database were directly copied into this database without further annotation or verification of the data. These Affymetrix probe set identifiers and associations to UniProt and Blattner were collected from <<http://www.affymetrix.com>>, in the form of annotation files (.csv). Files were downloaded on February 14, 2007 and added to the Gene Database on February 15, 2007 by GenMAPP.org. Blattner associations were extracted from the "Transcript ID(Array Design)" field and UniProt associations were extracted from the "SwissProt" field.

#### d. Database Report

- UniProt is the primary ID system for the *Escherichia coli* K12 Gene Database. The UniProt table contains all 4207 UniProt IDs contained in the UniProt proteome set for this species.
- The Blattner IDs were derived from the cross-references in the UniProt proteome set for *Escherichia coli* K12. We compared our Blattner table with the table in the supplementary material from Riley et al. (2006) *Escherichia coli* K-12: a cooperatively developed annotation snapshot—2005. *Nucleic Acids Research* 34: 1-9, "Supplementary\_Table\_1\_Annotation\_E.\_coli\_Genes.xls". Our Blattner table contains 4328 identifiers. There are 316 Blattner IDs reported in the Riley et al. table that are not in our Gene Database. Of these:
  - 157 are RNA genes (tRNA, rRNA, or misc\_RNA)
  - 1 is the origin of replication
  - 154 are protein coding sequences (CDS)
  - 4 do not have a feature designationConversely, there are 161 Blattner IDs in our table that are not in the Riley table:
  - 104 are in the form of "ECOK12Fxxx" and correspond to proteins encoded by plasmid F
  - 4 are in the form of "ECKxxxx"
  - 44 are in the form of "bxxxx"
  - 9 contain a period in the ID, such as "bxxxx.y"
- The W3110 IDs were derived from the cross-references in the UniProt proteome set for *Escherichia coli* K12. We compared our W3110 table with the table in the supplementary material from Riley et al. (2006) *Escherichia coli* K-12: a cooperatively developed annotation snapshot—2005. *Nucleic Acids Research* 34: 1-9, "Supplementary\_Table\_1\_Annotation\_E.\_coli\_Genes.xls". Our W3110 table contains 4240 identifiers. There are 374 W3110 IDs reported in the Riley et al. table that are not in our Gene Database. Of these:
  - 157 are RNA genes (tRNA, rRNA, or misc\_RNA)
  - 1 is the origin of replication
  - 195 are protein coding sequences (CDS)

21 do not have a feature designation

Conversely, there are 115 W3110 IDs in our table that are not in the Riley table:

104 are in the form of "ECOK12Fxxx" and correspond to proteins encoded by plasmid F

4 are in the form of "ECKxxxx"

4 are in the form of "JWxxxx" and are listed below:

JW1277	UPF0509 protein yciZ
JW1589	Uncharacterized protein ydgU
JW3648	LexA-regulated protein tisA
JW3649	LexA-regulated protein tisB

3 contain a period in the ID, such as "JWxxxx.1" and are listed below:

JW0680.1	Uncharacterized protein ybfK
JW1494.1	Two-component-system connector protein yneN
JW2623.1	Putative UPF0401 protein ypjI

- The *Escherichia coli* K12 Gene Database also contains 4022 EchoBASE IDs, 4190 EcoGene IDs, 4097 GeneIds, and 8187 RefSeq IDs that were cross-referenced by the UniProt XML.

## 5. Contact Information for support, bug reports, feature requests

- The Gene Database for *Escherichia coli* K12 was built by the Loyola Marymount University (LMU) Bioinformatics Group using the program GenMAPP Builder, part of the open source XMLPipeDB project <<http://xmlpipedb.cs.lmu.edu/>>.
- For support, bug reports, or feature requests relating to XMLPipeDB or GenMAPP Builder, please consult the XMLPipeDB Manual found at <<http://xmlpipedb.cs.lmu.edu/documentation.shtml>> or go to our SourceForge site <<http://sourceforge.net/projects/xmlpipedb/>>.
- For issues related to the *Escherichia coli* K12 Gene Database, please contact:  
Kam D. Dahlquist, PhD.  
Department of Biology  
Loyola Marymount University  
1 LMU Drive, MS 8220  
Los Angeles, CA 90045-2659  
kdahlquist@lmu.edu
- For issues related to GenMAPP 2.0/2.1 or MAPPFinder 2.0 please contact GenMAPP support directly by e-mailing [genmapp@gladstone.ucsf.edu](mailto:genmapp@gladstone.ucsf.edu) or [GenMAPP@googlegroups.com](mailto:GenMAPP@googlegroups.com).

## 6. Release Notes

### a. Current version: Ec-K12-Std\_External\_20090529.gdb

- This release is the first major update to the standard *Escherichia coli* K12 Gene Database. The data are updated according to the Data Source and Version information contained in section (4c) above. In this release, the following proper ID systems were added: GeneId (NCBI), RefSeq (protein), and W3110.
- Affymetrix probe set identifiers from the Affy table in the Ec-K12-Std\_External\_20060731x.gdb Gene Database were directly copied into this database without further annotation or verification of the data.
- Comparison of the number of IDs in the previous version of this database and the current version:

ID System	Previous version: Ec-K12-Std_External_20060731x.gdb	Current version: Ec-K12-Std_External_20090529.gdb
Blattner	4466	4328
EchoBASE	4156	4022
EcoGene	4224	4190
EMBL	2486	2532
GeneId (NCBI)	not present	4097
GeneOntology	3182	4304
Interpro	3511	3687
PDB	2902	3723
Pfam	2015	2004
RefSeq	not present	8187
UniProt	4329	4207
W3110	not present	4240
Affy	24692	24692

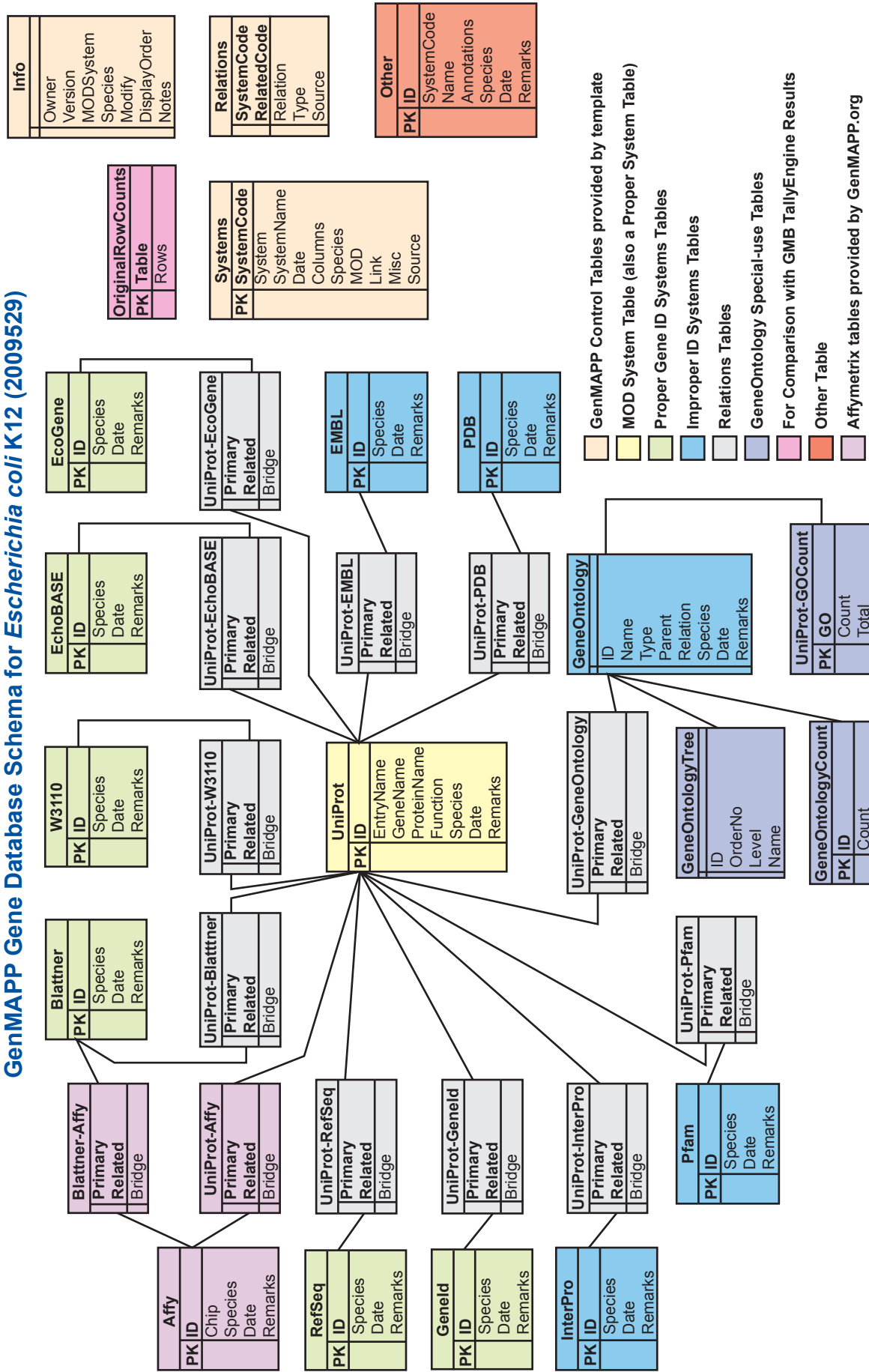
**b. Previous version: Ec-K12-Std\_External\_20060731x.gdb**

- This release is a modification of the first release of a standard *Escherichia coli* K12 Gene Database. The data in the Gene Database remains the same with the following additions.
- This Gene Database also contains Affymetrix probe set identifiers for all Affymetrix *E. coli* microarrays (E\_coli\_2 Array, E. coli Genome Sense Array, E. coli Genome Antisense Array). Affymetrix identifiers were added to the Gene Database by GenMAPP.org. These Affymetrix probe set identifiers were related to both UniProt and Blattner identifiers.
- This ReadMe document was updated on March 2, 2007.

**c. Previous version: Ec-K12-Std\_External\_20060731.gdb**

- This was the first release of a standard *Escherichia coli* K12 Gene Database.
- Unlike the official Gene Databases from GenMAPP.org, in the *Escherichia coli* K12 Gene Database, PDB has been designated as an “improper” gene ID system and cannot be used as an ID for gene objects on MAPPs. This action was taken because PDB IDs can refer to structures containing two or more different polypeptides and thus do not refer to a unique protein molecule.
- This ReadMe document was updated on September 24, 2006. The Gene Database itself has not changed. However, changes were made to this document to reflect the following:
  - The filename of the database has been changed from “Ec-Std\_20060731.gdb” to “Ec-K12-Std\_External\_20060731.gdb” to be consistent with the name of the database that can be downloaded via GenMAPP.org. “K12” was added to reduce ambiguity between *Escherichia coli* K12 and other strains of *E. coli*. “External” was added to indicate that the Gene Database was created by a group external to GenMAPP.org.
  - The URL for the XMLPipeDB Project has changed from <http://www.cs.lmu.edu/~xmlpipedb/> to <http://xmlpipedb.cs.lmu.edu/>.
  - A clarification to the taxon IDs used in the UniProt source data was made.

# GenMAPP Gene Database Schema for *Escherichia coli* K12 (2009529)



NOTE: Some Relations tables are not shown. All possible pairwise Relations tables exist between Proper ID systems and between Proper and Improper ID systems, but not between Improper ID systems (i.e., Proper-Proper, Proper-Improper, but NOT Improper-Improper).