# New Data Mining Interfaces at the Bovine Genome Database

**Colin Diesh[1], Aditi Tayal[1], Deepak Unni[1], Darren Hagen[1], Christine G. Elsik[1]**

[1] *Division of Animal Sciences, University of Missouri, Columbia, MO 65211, USA*

United States Department of Agriculture
National Institute of Food and Agriculture

BovineGenome.org
The Bovine Genome Database

## Abstract

The Bovine Genome Database (BGD, http://BovineGenome.org) is an informatics resource for Bos taurus. We have used InterMine to deploy a new data warehouse called BovineMine to provide a fast and flexible query interface. BovineMine integrates information from many data sources including RefSeq, Ensembl, UniProt, InterPro, OrthoDB, Homologene, Pubmed, Gene Ontology and BioGRID. BovineMine also includes data that we provide on our genome browsers, such as the Bovine Official Gene Set, RNAseq data, SNP and QTL. Users may perform a "Quick Search" or use the "Query Builder" for specialized searches. The "Genome Region Search" and "Overlapping Feature Search" allow users to download annotations with a specified genomic context. Users may download query results in various formats, such as tab-delimited files, GFF, Fasta, BED, JSON and XML. In addition to BovineMine, users may access data through genome browsers (GBrowse and JBrowse) and       .

We have also created specialized search interfaces that focus on differences in annotations and assembly versions, and provide easy navigation to novel genes on genome browsers. The Annotation Assembly Comparison Tool allows users to lookup locations of genes on two bovine assemblies (UMD3.1 and Btau_4.6.1). The Ensembl-NCBI Comparison Tool allows users to investigate disagreements in gene models across  gene sets. The Predicted Transcript RNAseq Read Count Tool provides counts of spliced RNAseq read alignments to predicted transcripts to aid selection of RNAseq tracks in JBrowse. The Candidate Novel Protein Coding Gene Search Tool allows users to easily navigate to locations candidate novel genes in JBrowse.
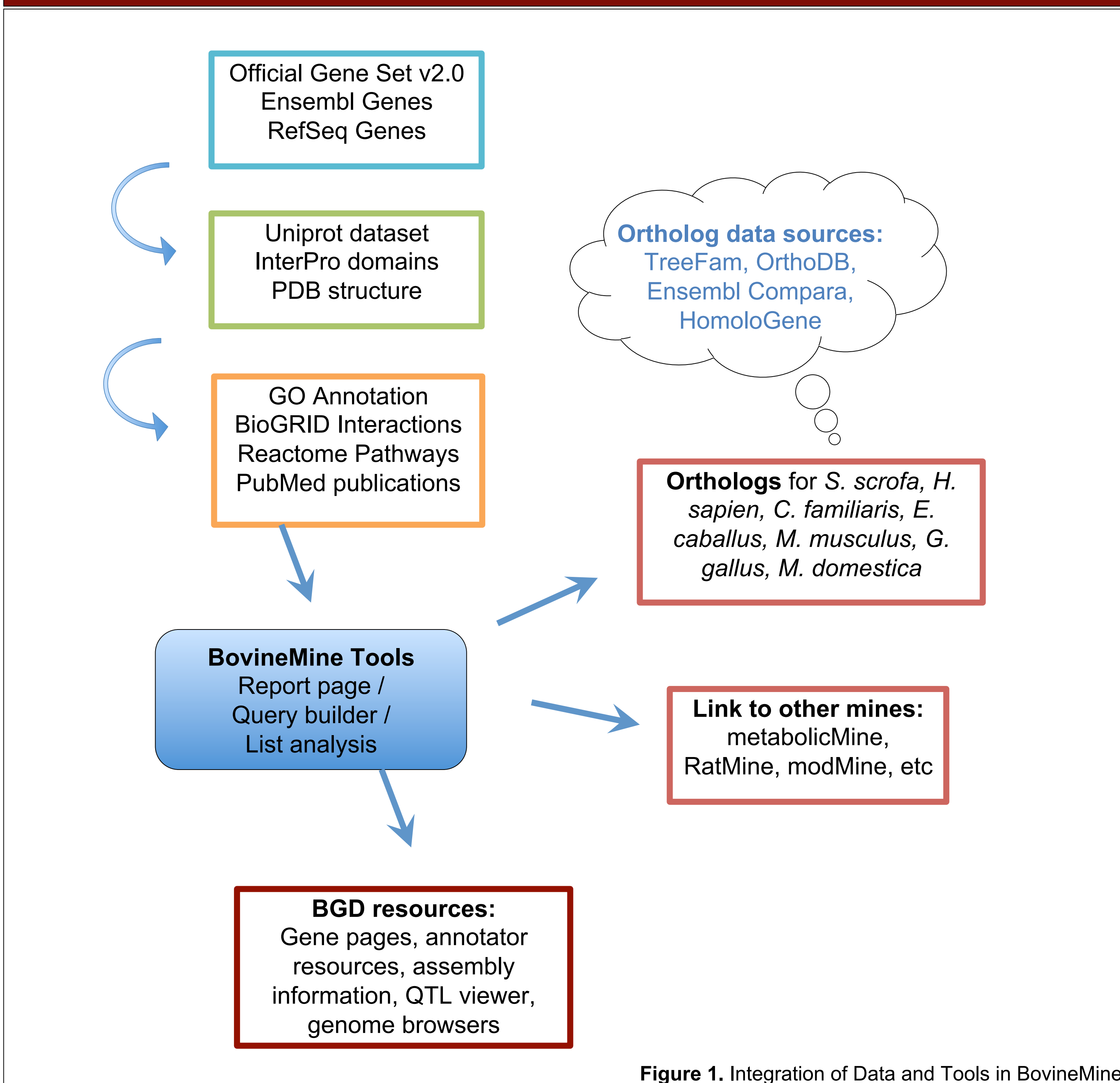
## BovineMine Datasets

Official Gene Set v2.0
Ensembl Genes
RefSeq Genes

Uniprot dataset
InterPro domains
PDB structure

GO Annotation
BioGRID Interactions
Reactome Pathways
PubMed publications

Ortholog data sources:
TreeFam, OrthoDB, Ensembl Compara, HomoloGene

Orthologs for *S. scrofa, H. sapien, C. familiaris, E. caballus, M. musculus, G. gallus, M. domestica*

BovineMine Tools
Report page /
Query builder /
List analysis

Link to other mines:
metabolicMine,
RatMine, modMine, etc

BGD resources:
Gene pages, annotator resources, assembly information, QTL viewer, genome browsers

**Figure 1.** Integration of Data and Tools in BovineMine

## BovineMine Search Tools



**Figure 2.** Search Tools

**A.** The **Quick Search** tool can be used to find keywords from multiple datasets at once. With the **List Analyse** tool, lists of IDs can be saved and analyzed using multiple queries.

**B.** The **Query Builder** is a powerful GUI for making complex structured queries.

**C. Template queries** are prepared queries to accomplish common tasks (genes matching GO term, homologues of given gene, etc).

GO term –> Genes
Show all genes annotated with a given GO term (and its children terms)

GO Term > Name   DNA binding
Organism > Name   Bos taurus
GO Evidence Code > Code   IEA

Show Results

## Data Mining Example

The following is an example using QTL regions that are associated to milk fat yield. To begin, we use the "Region search" and paste regions corresponding to QTL locations on base pair coordinates in the format "chr:start..end":



**Figure 3.** (left) BovineMine's Region search page and (right) results page containing genes that overlap the regions of interest.

The results from the region search can then be augmented with GO Annotations.



**Figure 4.** (left) The "GO Annotation" attribute is added to the results using "Manage Columns" tool. (middle) The results of region search with GO terms is shown. (right) The table download page for the region search is used to download results to plain text.

## New Tools for Annotators at BGD



New tools help annotators navigate across alternate bovine assemblies, convert identifiers, compare gene sets and select tracks to view in the browser. For example, the RNAseq Read Count tool provides raw read counts from spliced RNAseq alignments to locations of predicted transcripts (Ensembl, RefSeq and Bovine OGSv2). The read counts help users determine which of 91 RNAseq data sets to view in Jbrowse/Web Apollo. The spliced read alignments (BAM tracks) are particularly useful in determining whether splice junctions in predicted genes are correct.

**Figure 5.** Using the RNAseq read count tool to select a BAM track that would be useful for evaluating splice sites of a particular gene model.

## Acknowledgement