

# GeneLab

Open Science for Exploration

# GeneLab Phase 2: Integrated Search & Data Federation of Space Biology Experiment Data

P.B. Tran<sup>1</sup>, D.C. Berrios<sup>1</sup>, M.M. Gurrum<sup>1</sup>, J.C.M. Hashim<sup>1</sup>, S. Raghunandan<sup>1</sup>, S.Y. Lin<sup>1</sup>, T.Q. Le<sup>1</sup>, D.M. Heher<sup>1</sup>, H.T. Thai<sup>1</sup>, J.D. Welch<sup>2</sup>, S.M. Caldwell<sup>3</sup>, O.G. Stotzky<sup>4</sup>, M.G. Skidmore<sup>4</sup>  
<sup>1</sup>Intelligent Systems Division, <sup>2</sup>Partnerships Division, <sup>3</sup>Space Biosciences Division, <sup>4</sup>Project Management and Management Division, NASA Ames Research Center, Moffett Field, CA, USA

The GeneLab project is a **science initiative** to maximize the scientific return of omics data collected from spaceflight and from ground simulations of microgravity and radiation experiments, supported by a **data system** for a public bioinformatics repository and collaborative analysis tools for these data. The mission of GeneLab is to maximize the utilization of the valuable biological research resources aboard the ISS by collecting genomic, transcriptomic, proteomic and metabolomic (so-called “omics”) data to enable the exploration of the molecular network responses of terrestrial biology to space environments using a systems biology approach. All GeneLab data are made available to researchers worldwide through its open-access data system.

GeneLab is currently being developed by NASA to support worldwide accessibility of biomedical research in order to enable the human exploration of space and improve life on earth. Open access to Phase 1 of the GeneLab Data Systems (GLDS) was implemented in April 2015. Download volumes have grown steadily, mirroring the growth in curated space biology research data sets (80 as of Sept. 2016), now exceeding 10 TB/month, with over 10,000 file downloads since the start of Phase 1. For the period April 2015 to May 2016, most frequently downloaded were data from studies of *Mus musculus* (39%) followed closely by *Arabidopsis thaliana* (30%), with the remaining downloads roughly equally split across 12 other organisms (each < 10% of total downloads). GLDS Phase 2 is focusing on interoperability, supporting data federation, including integrated search capabilities, of GLDS-housed data sets with external data sources, such as gene expression data from NIH/NCBI’s Gene Expression Omnibus (GEO), proteomic data from EBI’s PRIDE system, and metagenomic data from Argonne National Laboratory’s MG-RAST. GEO and MG-RAST employ specifications for investigation metadata that are different from those used by the GLDS and PRIDE (e.g., ISA-Tab). The GLDS Phase 2 system will implement a Google-like, full-text search engine using a Service-Oriented Architecture by utilizing publicly available RESTful web services Application Programming Interfaces (e.g., GEO Entrez Programming Utilities) and a Common Metadata Model (CMM) in order to accommodate the different metadata formats between the heterogeneous bioinformatics databases. GLDS Phase 2 completion with fully implemented capabilities will be made available to the general public in September 2017.

For more information please visit the GeneLab project website at:

<http://genelab.nasa.gov>

Explore the GeneLab Data System repository at:

<https://genelab-data.ndc.nasa.gov/genelab/>

## GLDS Phase 1 (aka “C-Gene”) Operational System

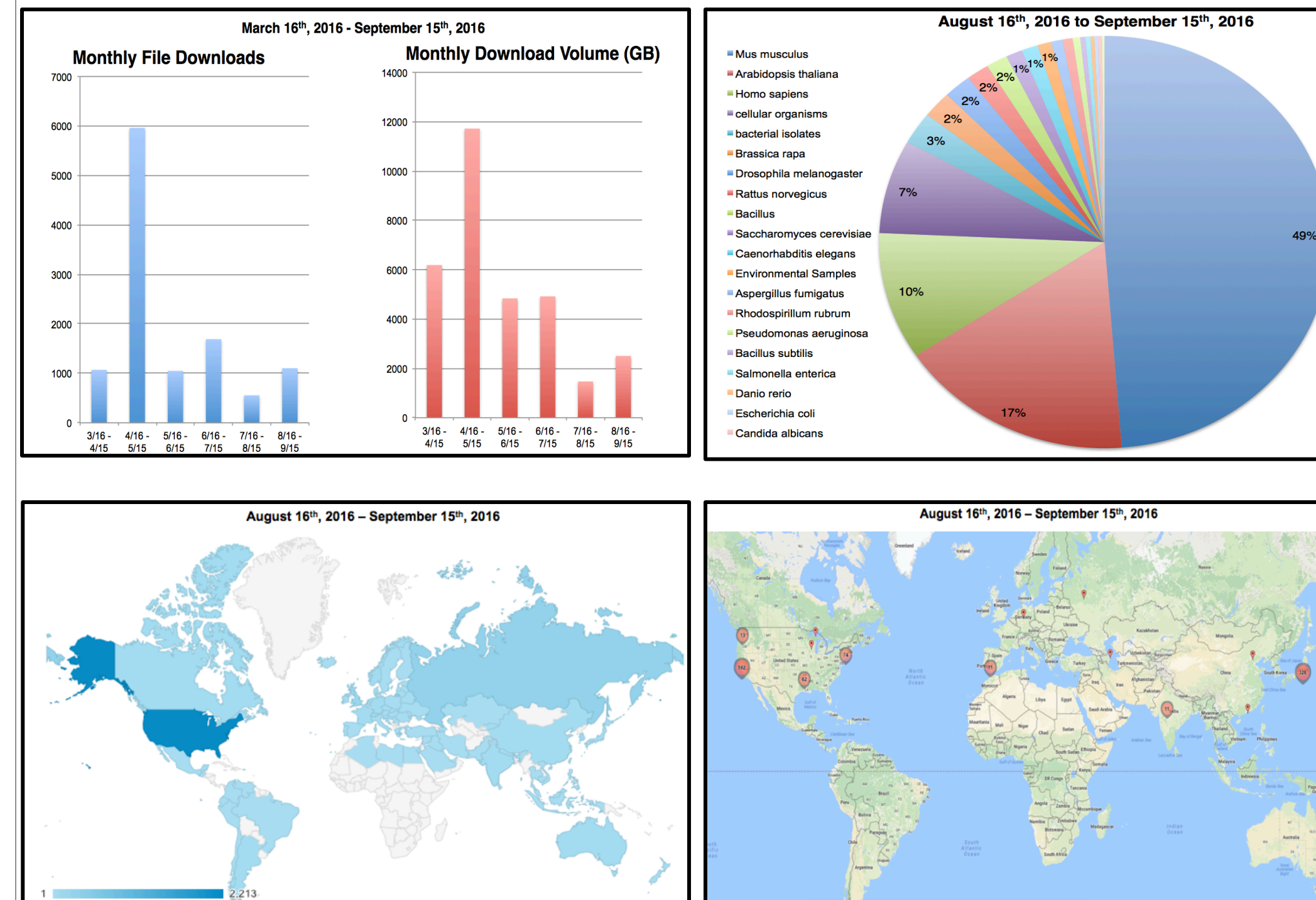
The screenshot shows the GeneLab Phase 1 web interface. At the top, there is a search bar with the text 'mouse liver STS'. Below the search bar, there are navigation links: Home, Repository, Data, Data Mining Tools, Submit Data, Contact Us. The search results page shows 'Page 1 of 4 (Total Studies: 80)'. The first result is titled 'Comparative Transcriptomic Analysis of Adult Medaka Tissues Sampled after Adaptation to a Space Environment'. Below the title, there is a table with columns: Organisms, Factors, Assay Types, Release Date, and Description. The table lists two studies: GLDS-83 and GLDS-81.

Simple keyword search capability with browsing and navigation

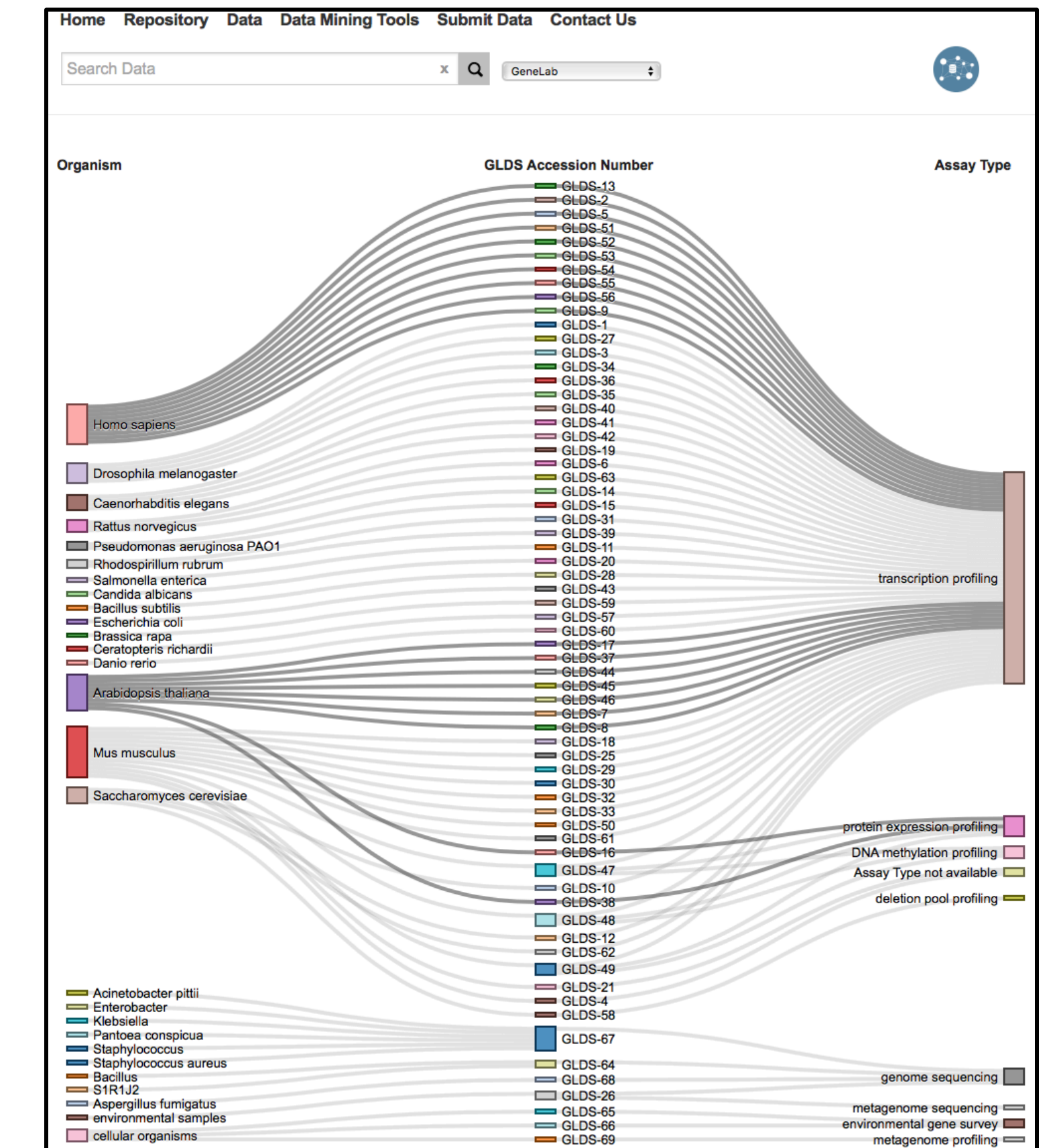
## GLDS Phase 2 (aka “X-Gene”) Work-In-Progress System

The screenshot shows the GeneLab Phase 2 web interface. At the top, there is a search bar with the text 'mouse liver STS'. Below the search bar, there is a dropdown menu with options 'GeneLab' and 'NIH GEO'. The search results page shows 'Total Search Results Found: 12440'. The first result is titled 'STS-135 Liver Transcriptomics'. Below the title, there is a table with columns: Organism, Factor, Assay Type, and Accession. The table lists one study: GLDS-25.

Google-like, full-text search engine with data federation/integration to NIH/NCBI's GEO database (over 80K GSE metadata records ingested)



GLDS Metrics (As of 9/15/2016)



Data Visualization (GLDS Organisms vs. Assay Types)