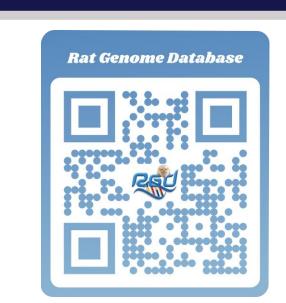# Rat gene expression data expansion at the Rat Genome Database, an update

Wendy M. Demos[1], Logan Lamers[1], Jeff L. De Pons[1], Adam C. Gibson[1], Varun Reddy Gollapally[1], G. Thomas Hayman[1], Mary L. Kaldunski[1], Akhilanand Kundurthi[1], Stanley J.F. Laulederkind[1], Jennifer R. Smith[1], Jyothi Thota[1], Marek A. Tutaj[1], Monika Tutaj[1], Mahima Vedi[1], Shur-Jen Wang[1], Kent C. Brodie[2], Stacy M. Zacher[3], Melinda R. Dwinell[1], Anne E. Kwitek[1].

*Rat Genome Database, Dept. of Physiology[1]; Clinical and Translational Science Institute[2]; Finance and Administration[3] Medical College of Wisconsin, Milwaukee, WI 53226 USA*

https://rgd.mcw.edu

## Abstract

The Rat Genome Database (RGD, https://rgd.mcw.edu ) has been expanding and incorporating gene expression data content into the larger ecosystem of RGD. Researchers will soon be able to access gene expression values and sample metadata that were submitted to public resources such as the Gene Expression Omnibus (GEO, https://www.ncbi.nlm.nih.gov/geo/ ) repository, with all data values converted to transcripts per million (TPM) data type. Integration of data submitted to GEO was designed in three phases as outlined below. Additional expression data, such as the baseline expression of strains related to the Hybrid Rat Diversity Program (https://rgd.mcw.edu/rgdweb/hrdp_panel.html ) have been integrated into JBrowse2.

In Phase One of the project, an expression curation tool was developed to aid in comprehensive Natural Language Processing (NLP) assisted manual curation of public datasets. To date, 2,109 GEO expression projects with focus on the *Rattus norvegicus* model have been reviewed and prioritized for curation. Of those projects, 307 met the primary prioritization criteria and metadata for 3,629 project samples have been loaded to the database.

The goal of Phase Two is to standardize the publicly available gene expression data by reprocessing samples in the curated GEO projects. RGD developed and is optimizing a bioinformatic pipeline that downloads and converts fastq files from the Sequence Read Archive, aligns to the most current and correct *R. norvegicus* genome assembly, and outputs TPM data type. This pipeline integrates quality control measures, verifies, or calculates sample sex, produces alignments with the STAR aligner, and estimates gene and transcript level abundance with the RSEM software package.

Phase Three focuses on enhanced visualization of gene expression values. The tabular-based display of gene expression values has been updated on the RGD gene pages. Users can view data by anatomical system and download sample metadata, TPM values, and data sources. Evaluation of JBrowse2 visualizations for expression profiles and TPM values at the gene and transcript levels is underway. These curated and reprocessed expression data, download options and visualizations will provide standardization of publicly available data and continue to add value for RGD users.

**Figure 1:** GEO Accession Display Interface. Inset shows sample specific details submitted to GEO.

## Phase 1: GEO Metadata Curation

Curators leverage information submitted to the GEO repository (**Figure 1**) and associated publications, if available, (**Figure 2**) and prioritize curation based on data type FPKM / TPM (for data import) and bulk RNA sequencing strategy. Most projects earmarked for future curation are due to the data type submitted to GEO followed by the sequencing strategy (**Figure 3**). The Expression Curation Tool (**Figure 4**) relies on a pipeline that imports data from the GEO and utilizes Natural Language Processing to match ontology terms to GEO accession attributes. This interface provides fields for 7 ontologies, descriptors, and experimental conditions. Curators utilize built-in ontology browsers to enter missing terms, confirm the predicted terms, or provide a more specific term when appropriate.
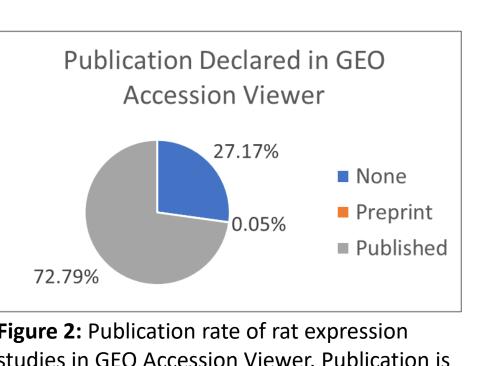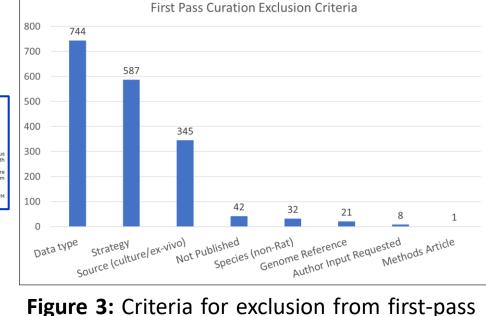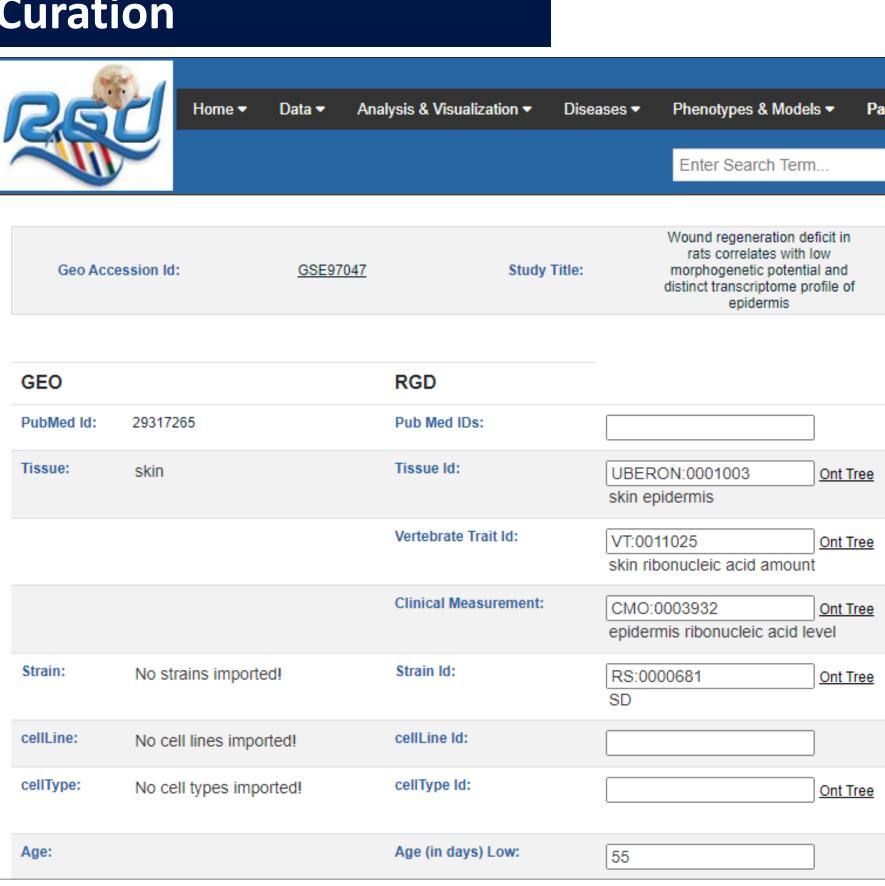


**Figure 2:** Publication rate of rat expression studies in GEO Accession Viewer. Publication is not necessary for curation if enough details are provided.



**Figure 3:** Criteria for exclusion from first-pass curation that are grounds to mark projects for future or no curation.



**Figure 4:** Expression curation tool interface. Terms are first added at the project level and propagated across all applicable samples. Curators can then modify individual sample terms and attributes as needed based on the reference, within the next window of the interface prior to loading metadata to the underlying data tables.

## Phase 2: Data Analysis

Submission to GEO does not require a standard reference (**Figure 5**) or data type (**Figure 6**). The RGD RNAseq pipeline (**Figure 7**) uses a sample accession list generated from information provided in the GEO series matrix file and the SRA run table. The workflow is launched from a bash menu and requires the accession list as input. Data are retrieved from the cloud as .sra files with the SRAtoolkit to a local server and transformed to fastq files. Multiple quality control steps occur throughout the process to evaluate sequencing quality measures, confirm species and sex (**Table 1, Figure 8**). Alignment to the rat genome reference is completed by the STAR aligner and abundance measures are calculated by RSEM. Verification of the workflow included comparing data submitted to GEO with generated data output. Expression thresholds as defined by Expression Atlas were used to verify similarity between results (**Figure 9**).



**Figure 5:** *R. norvegicus* genome reference versions declared in GEO for rat expression studies. This graph does not discern between RefSeq or Ensembl versions.
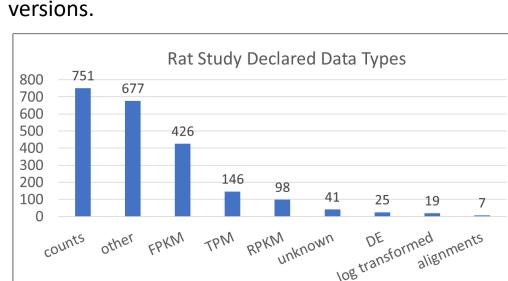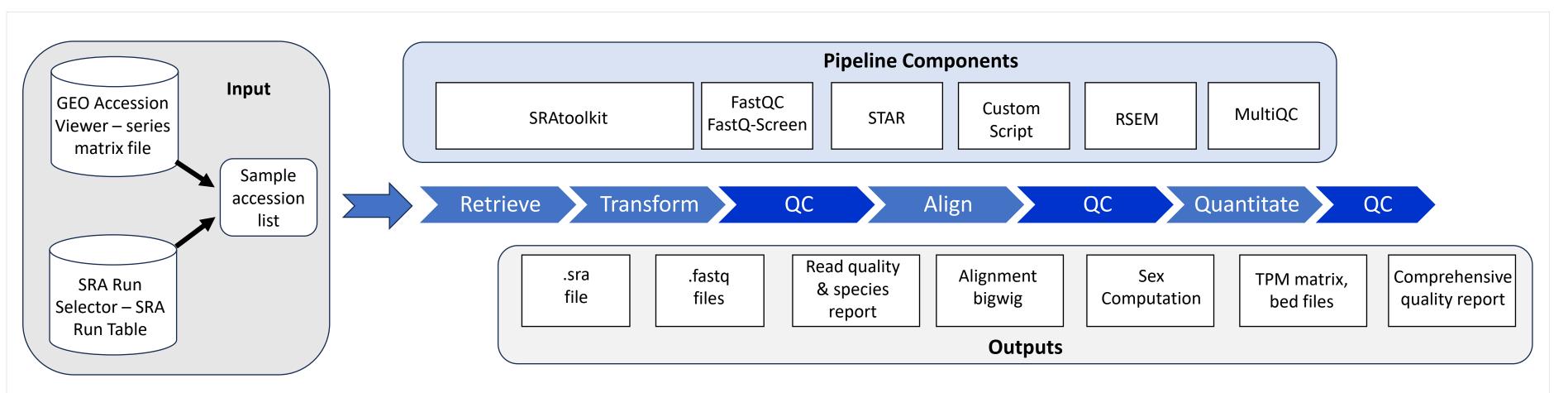


**Figure 6:** Processed or raw data are submitted as supplemental files to GEO. Submitters may declare the data type in the data processing field at the time of submission.



**Figure 7:** Schematic of the RGD RNAseq workflow highlighting input, software components, and outputs. Sample-specific attributes are gathered from the GEO Accession Viewer and the SRA Run Selector and compiled into a text file. The attributes are leveraged in various stages of the process. The workflow pauses prior to each QC step, ensuring curator review prior moving forward to the next stage. BigWig files are intended for visualization in JBrowse2, matrix files will be imported into the database for display on the gene pages.
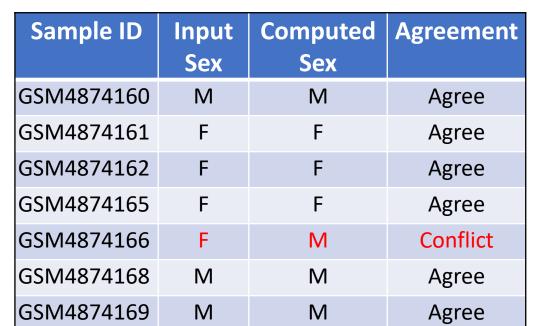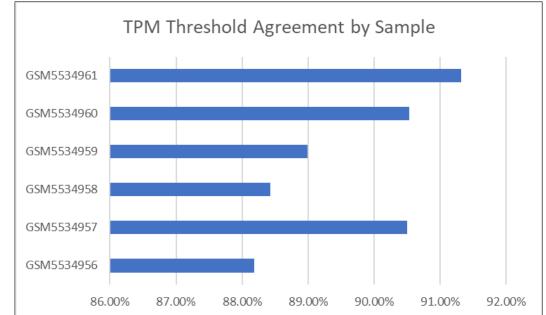
| Sample ID | Input Sex | Computed Sex | Agreement |
|---|---|---|---|
| GSM4874160 | M | M | Agree |
| GSM4874161 | F | F | Agree |
| GSM4874162 | F | F | Agree |
| GSM4874165 | F | F | Agree |
| GSM4874166 | F | M | Conflict |
| GSM4874168 | M | M | Agree |
| GSM4874169 | M | M | Agree |

**Table 1:** Example of sample sex quality control output. In this case, the sex declared of sample GSM4874166 in the GEO accession viewer does not agree with the computed sex.



**Figure 8:** Sample sex can be verified by viewing the bigwig files in JBrowse2. In this example, the alignment supports the computed sex of male is accurate.



**Figure 9:** Agreement of expression levels between sample transcript TPM data submitted for GSE182706 and RGD generated data. Data were normalized to common transcripts.

## Phase 3: GEO Data Visualization

To date, a total of 31 studies encompassing about 22,000 records with an average of approximately 32,000 individual expression values per record are viewable in the database. Initial loads were imported from Expression Atlas, leveraging their standardization of the metadata fields, and the mapping of the sequences to the Rnor6.0 assembly. Additionally, 3 rat studies and one study with data for both rat and dog were imported from GEO. In all cases, samples are untreated controls and expression values are in TPM. As curated GEO projects are subjected to reprocessing by the RGD RNAseq workflow, those data will be loaded to the database. Data are viewable and downloadable from any gene page within the RGD (**Figure 10**). Publication information for each data point is available from the 'Reference' link on the expanded expression table. (**Figure 11**).



**Figure 10:** RNA-seq expression value display from a RGD gene page. Users can navigate to the expression section by using the navigation tools on the left-hand side of the browser window (highlighted in blue). Data are presented by organ system within the interface. Users can download all data points as a table or limit the table data to a single system.
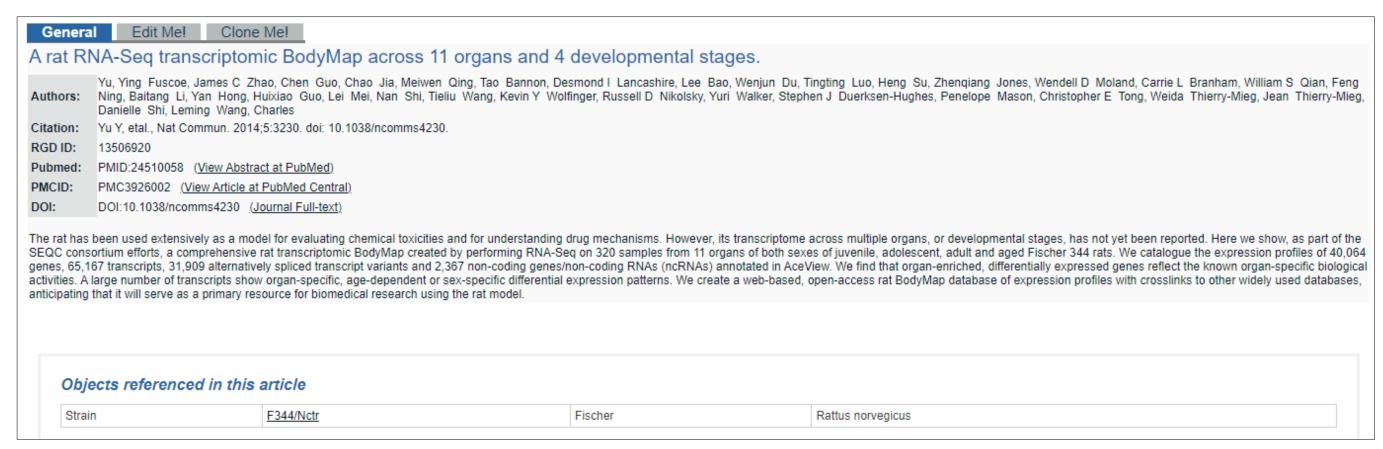


**Figure 11:** Reference information. Links located under the 'Reference' column heading of the expression table will bring users to a new window that provides information regarding the origin of the data and related objects.

## Expanded Content

In addition to providing standardized publicly available data, RGD recently released baseline expression profiles of classic inbred strains of the of the Hybrid Rat Diversity Panel. The first release includes gastrocnemius, left retroperitoneal fat pad and thymus of the F344/DuCrl and LE/StmMcwi strains (**Figure 12**). Further development is underway to provide gene and transcript level expression visualizations. TPM values as displayed in **Figure 13** may be an available option for data submitted to the Community Projects Portal.
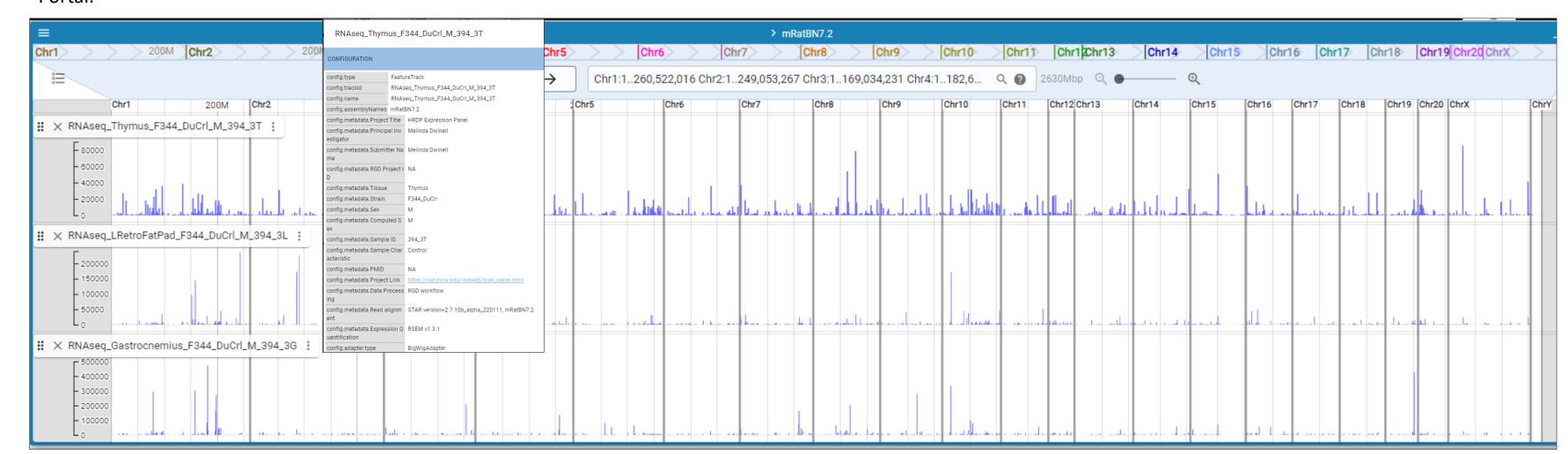


**Figure 12:** RNAseq profiles for gastrocnemius, left retroperitoneal fat pad and thymus of a male F344/DuCrl rat aligned to mRatBN7.2 in JBrowse2. Data have recently been released for public viewing from the track menu at https://rgd.mcw.edu/jbrowse2 . Profiles generated by alignment to GRCr8 and detailed track information (inset) will be available in a future release of JBrowse2.
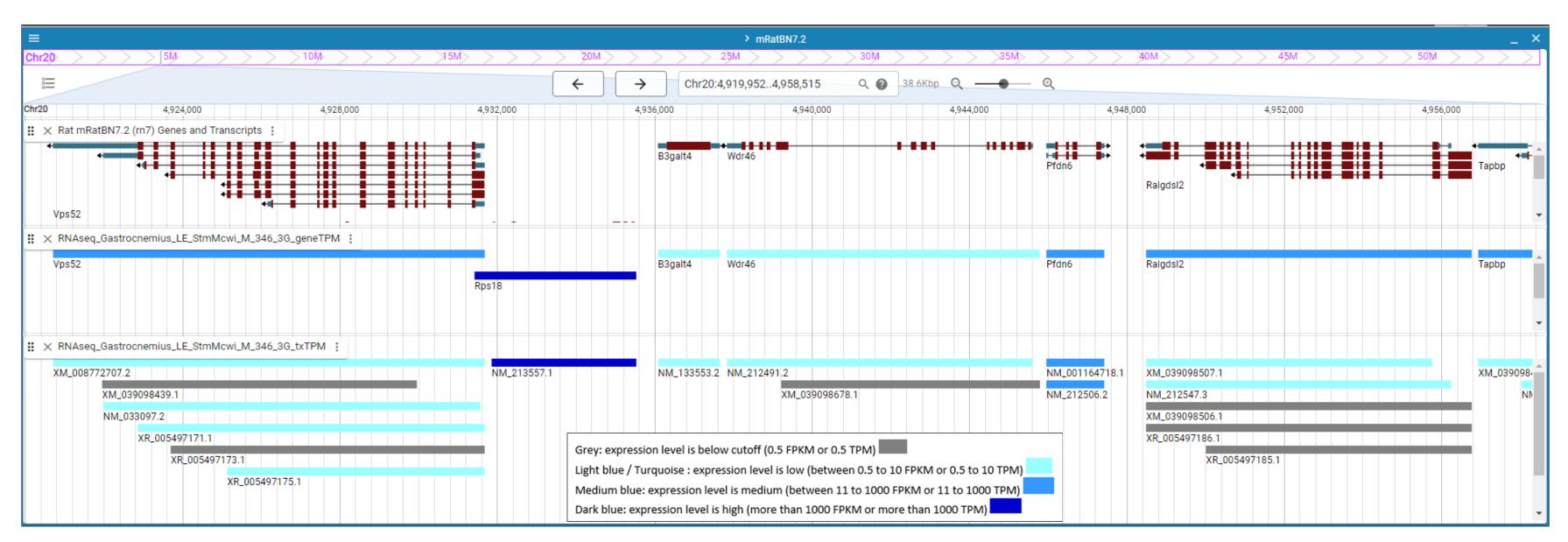


**Figure 13:** Gene (center row) and transcript (bottom row) level TPM tracks in JBrowse2. Expression levels are color-coded as below cutoff, low, medium, and high as defined by Expression Atlas (threshold ranges displayed in the inset). Transcripts or genes with a TPM value equal to zero are not displayed.

This and other recent RGD presentations will be freely available for viewing and download in RGD's Presentations Archive. https://rgd.mcw.edu/wg/com-menu/poster_archive/