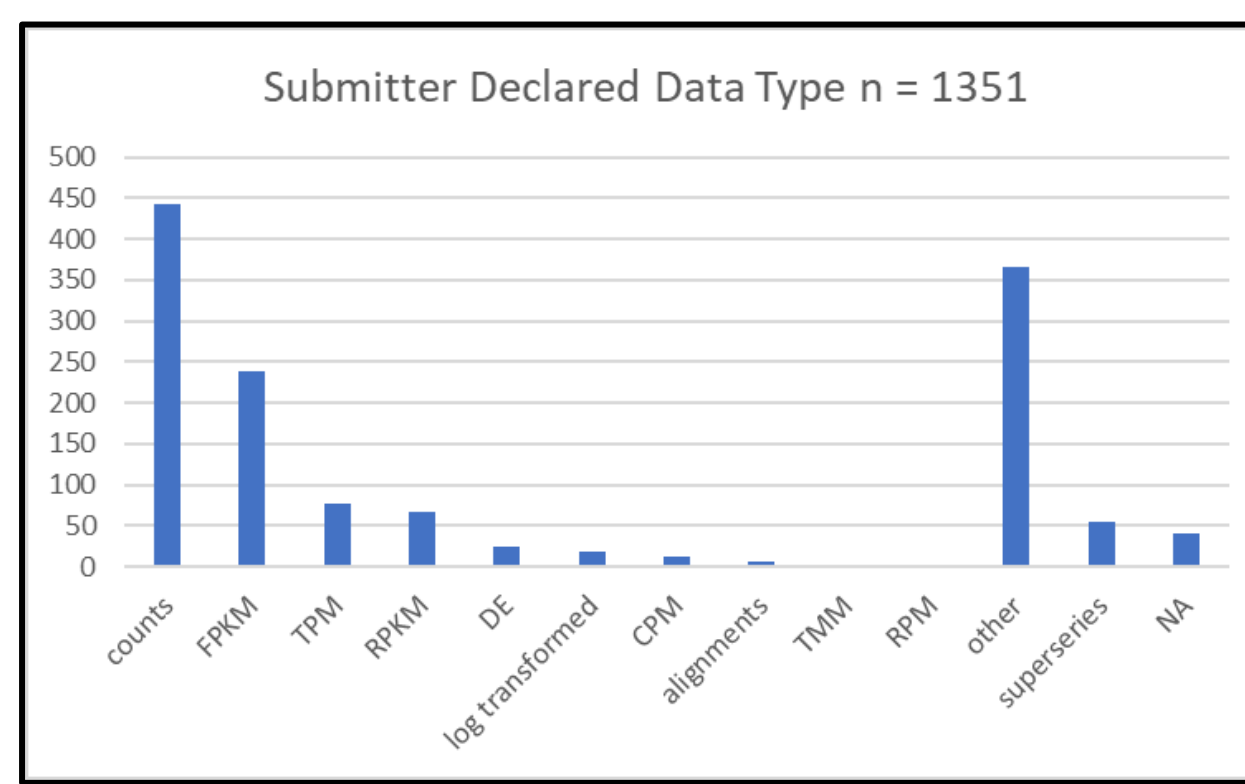
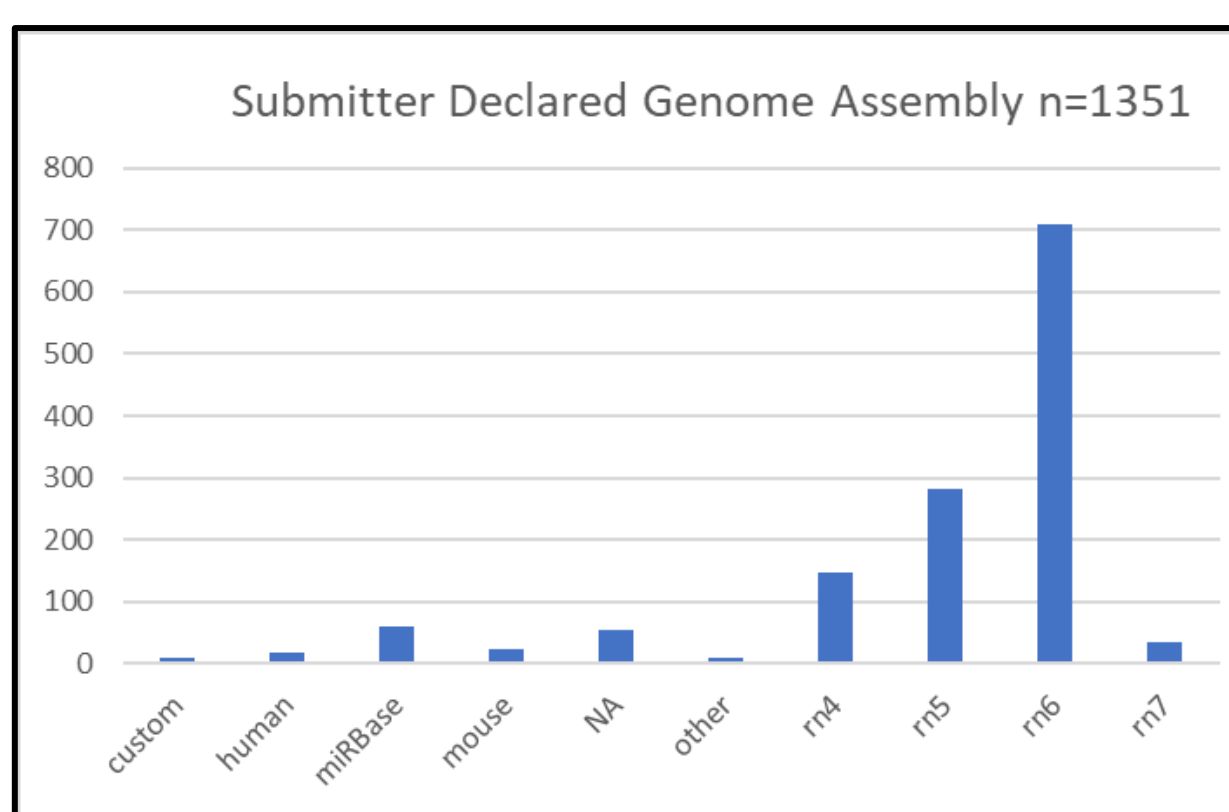


Abstract

The Rat Genome Database (RGD) is expanding and incorporating expression data content into the larger ecosystem of RGD so users can seamlessly query for coherent gene information across portals. Researchers will be able to access expression metadata and values that were submitted to public resources such as the Gene Expression Omnibus (GEO) repository, with all expression values converted to transcripts per million (TPM). 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. The expression curation tool relies on a pipeline that imports metadata from the GEO Accession Display and utilizes NLP to match ontology terms to GEO series attributes. Curators can enter missing terms, confirm the predicted term, or provide a more specific term when appropriate. Fields for descriptors such as tissue type, vertebrate trait, clinical measurement, strain, cell type, experimental condition, etc. are built into the user interface. To enhance operational efficiency, the curation process begins at the project level interface, where ontology terms are entered a single time and then propagated across all applicable samples. Curators conduct a sample level review and edit terms on a per-sample basis as needed. When the metadata are correct and as complete as possible, they are loaded into the appropriate tables in RGD's relational database. The expression values submitted for the curated GEO series will be loaded for all genes and transcripts in the corresponding files. Currently, RGD has imported 1,859 GEO series related to *Rattus norvegicus* expression studies. Of those, 1,351 have been reviewed and prioritized for curation. To date, metadata for 165 GEO series have been uploaded. Expression value types submitted to the repository represent a wide range of analysis outputs (i.e., FPKM, counts, log₂FC). The submitter declared genome assemblies in the reviewed GEO series include versions RGSC3.4-mRatBN7.2 as well as custom and non-rat references. The lack of standardization in the repository makes it difficult to identify rat data and furthermore, correlate expression values across studies. The goal of Phase Two is to standardize the expression values by developing and evaluating 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. This pipeline integrates quality control measures, alignment with the STAR¹ aligner, and abundance estimation with the RSEM² software package. Phase Three will focus on enhanced visualization of expression values. The current tabular-based view will be updated and new graphical visualizations at the gene and transcript levels are planned.

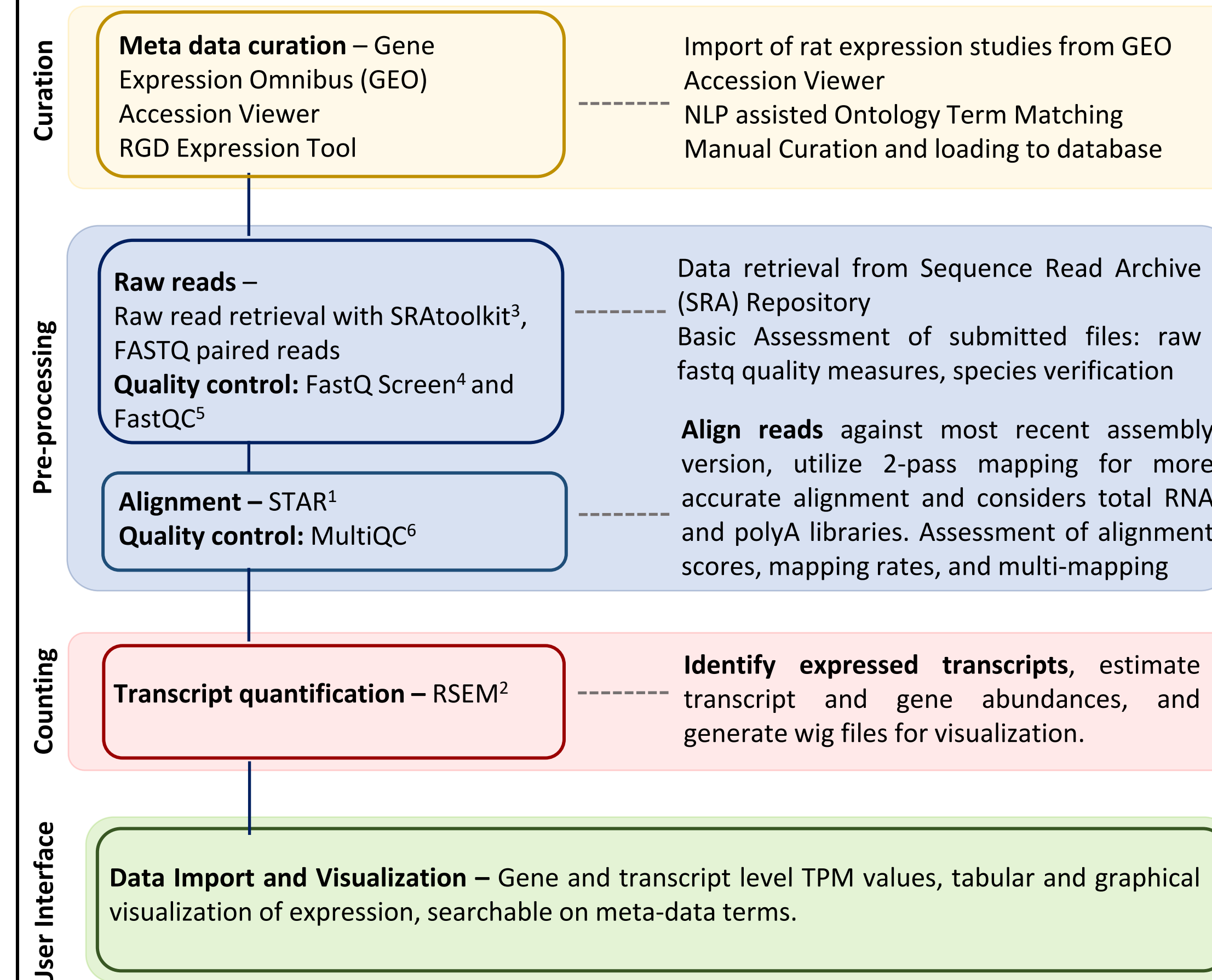
Background

1,351 studies were reviewed at the time of writing the abstract for prioritization of meta-data curation. 58% are tagged for future curation due to stringent first pass curation criteria (data type = TPM, FPKM, source = non-culture, strategy = bulk RNA, illumina sequencing). 22% of the reviewed studies are not associated with published articles. A variety of genome assemblies were reported by GEO Submitters for each GEO Accession. (Assembly names are simplified for ease of viewing in the graph and are not indicative of Ensembl or NCBI resources).

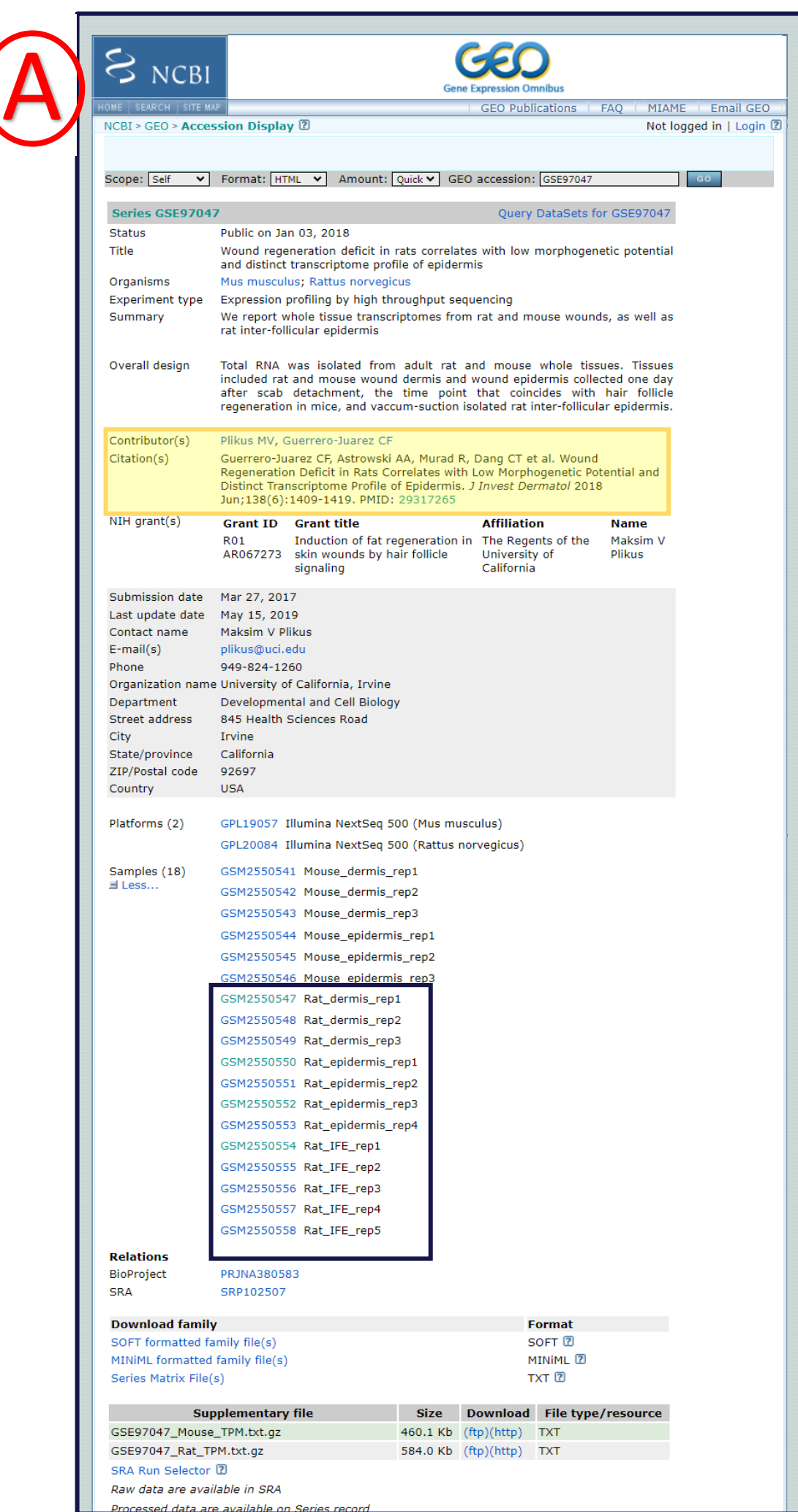


Submitters to GEO declare the data format type for the provided results files. The most common data type in the 1,351 reviewed studies fall into the *Other* category. The *Other* category is non-FPKM or TPM data such as peak data, alignment files, UMI counts, etc. *SuperSeries* category may have multiple data types as a *SuperSeries* accession may have many *SubSeries* each with a different data type. At times, a data type will be provided, but upon review, the submitted data do not reflect the type, or there is a discrepancy between the filename and declared type (For example RPKM data are provided but the filename is ACCESSION.fpkm.txt).

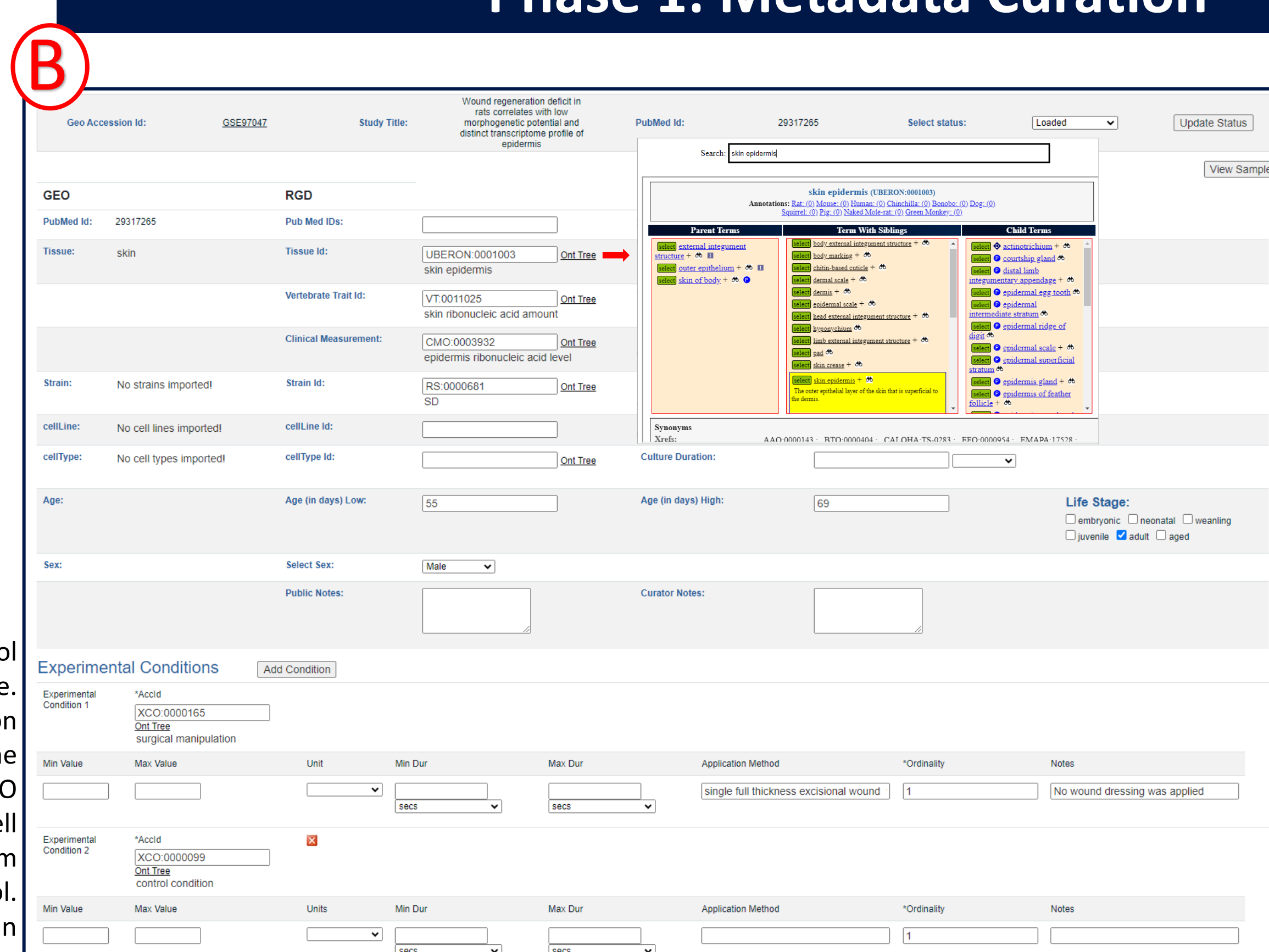
Methods



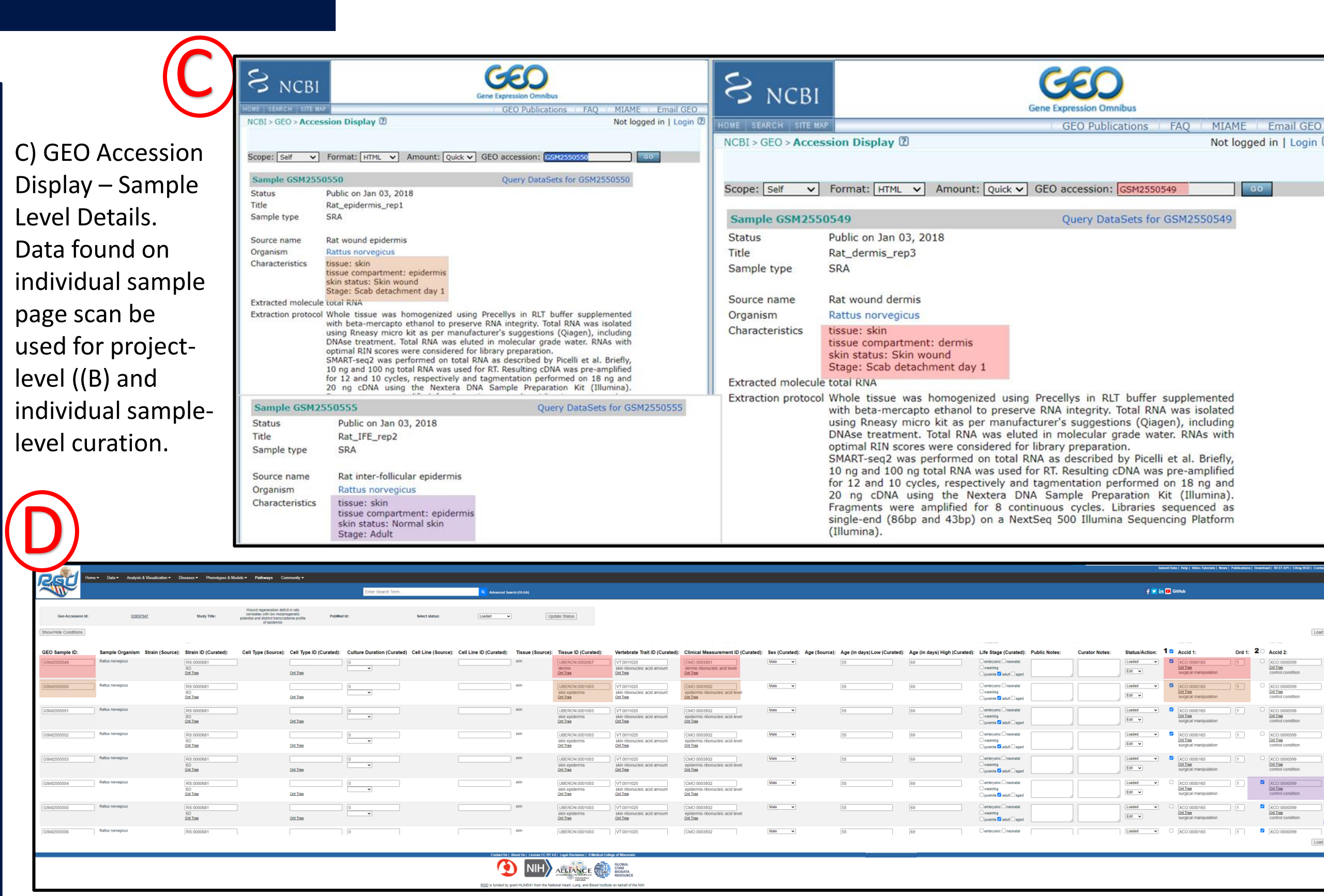
Phase 1: Metadata Curation



A) GEO Accession Display – Landing page for a given Gene Series. Experiment summary, submitter contact information, and publication ID (if available) are provided. Links to the individual sample detail web pages are provided (boxed in blue). Near the bottom of the page, download links for the entire series may be available, or data may be provided on a per-sample basis (<https://www.ncbi.nlm.nih.gov/geo/query/>).

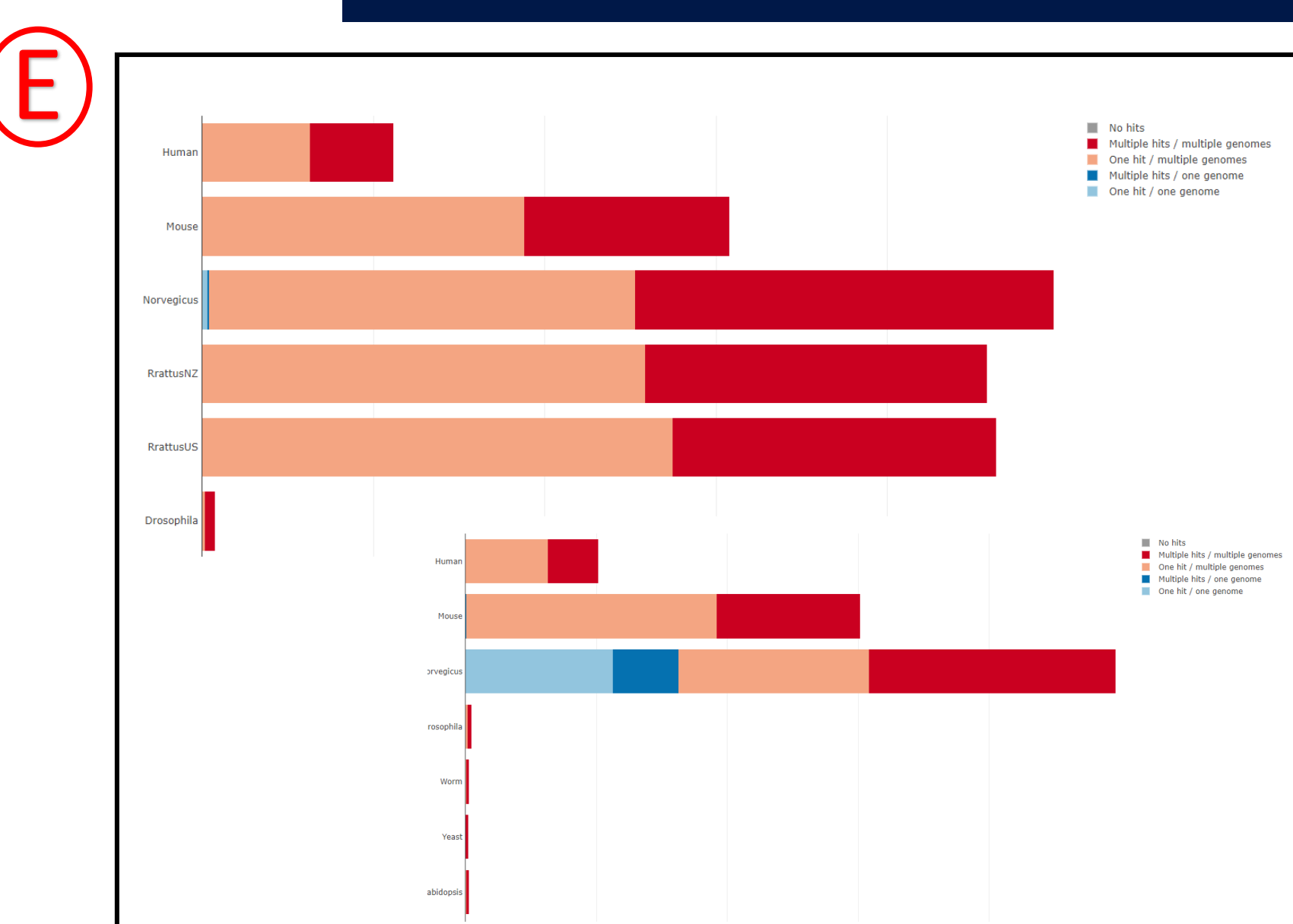


B) RGD Expression Curation tool – project level curation page. Where the NLP based prediction software finds a match, the original text from the GEO record that was matched as well as the suggested ontology term are shown in the curation tool. Curators can utilize the built-in ontology browsers to enter missing terms, confirm the predicted terms or provide a more specific term when appropriate. The prepopulated data are verified by curators and missing information is manually gathered from the GEO Accession Viewer (highlighted in yellow in Figure A) as well as the publication and its references when available. If necessary, curators will reach out to corresponding authors to obtain clarifying information. Terms are entered a single time for a GEO Accession and propagated across all applicable samples. If a submitter provides more than 1 citation, additional PMIDs can be attached to the data.

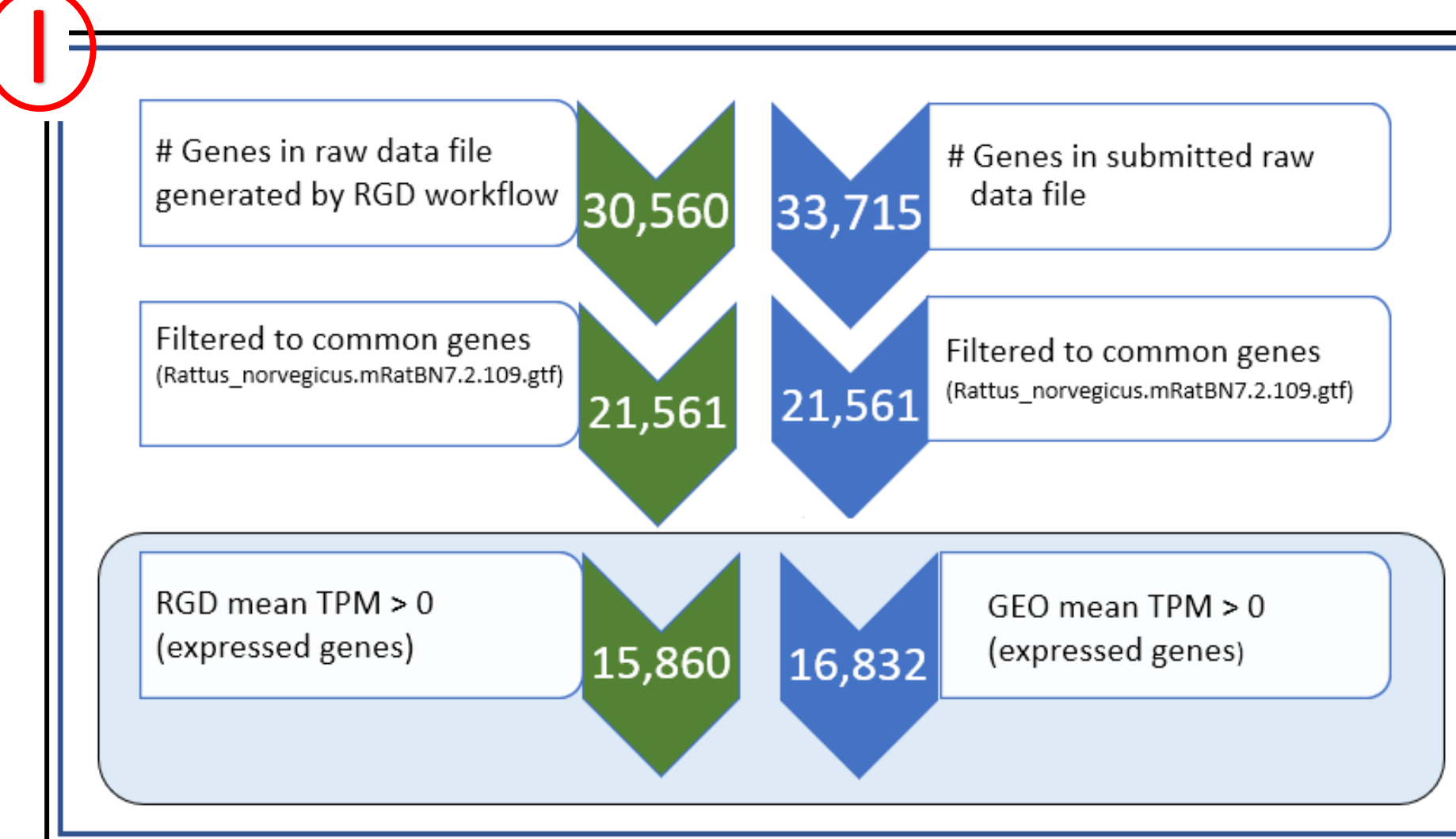
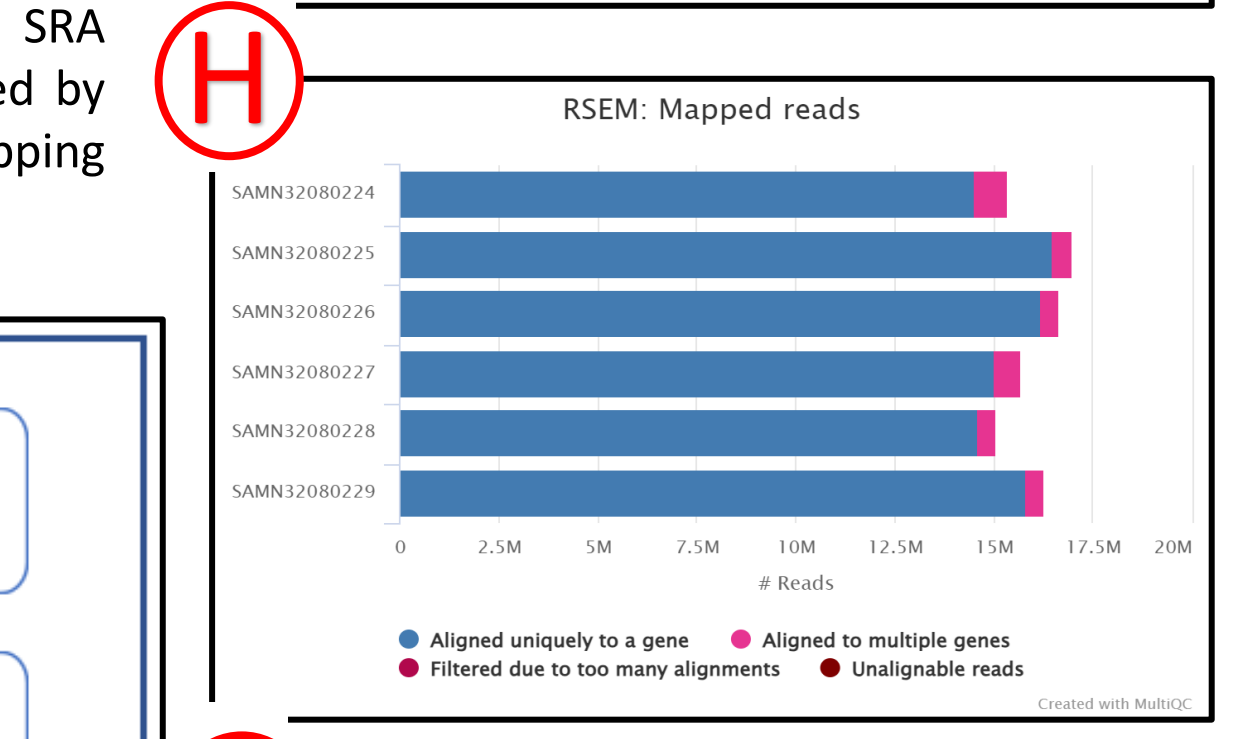
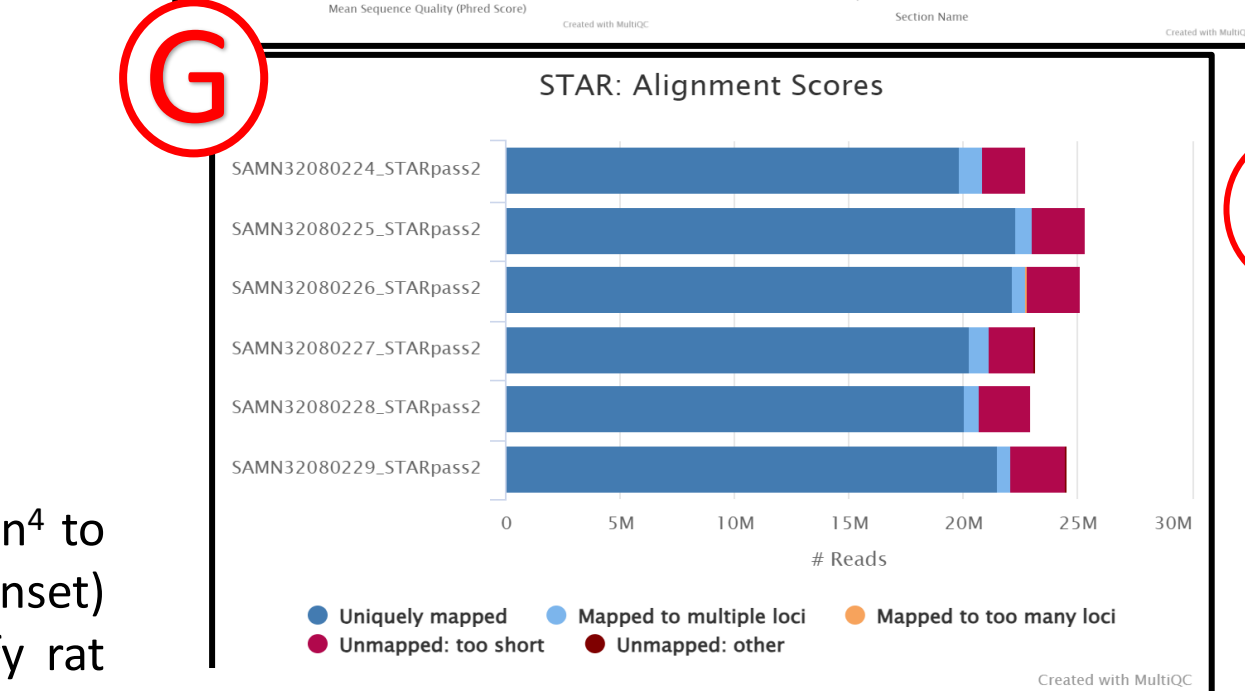


C) GEO Accession Display – Sample Level Details. Data found on individual sample page scan be used for project-level (B) and individual sample-level curation.

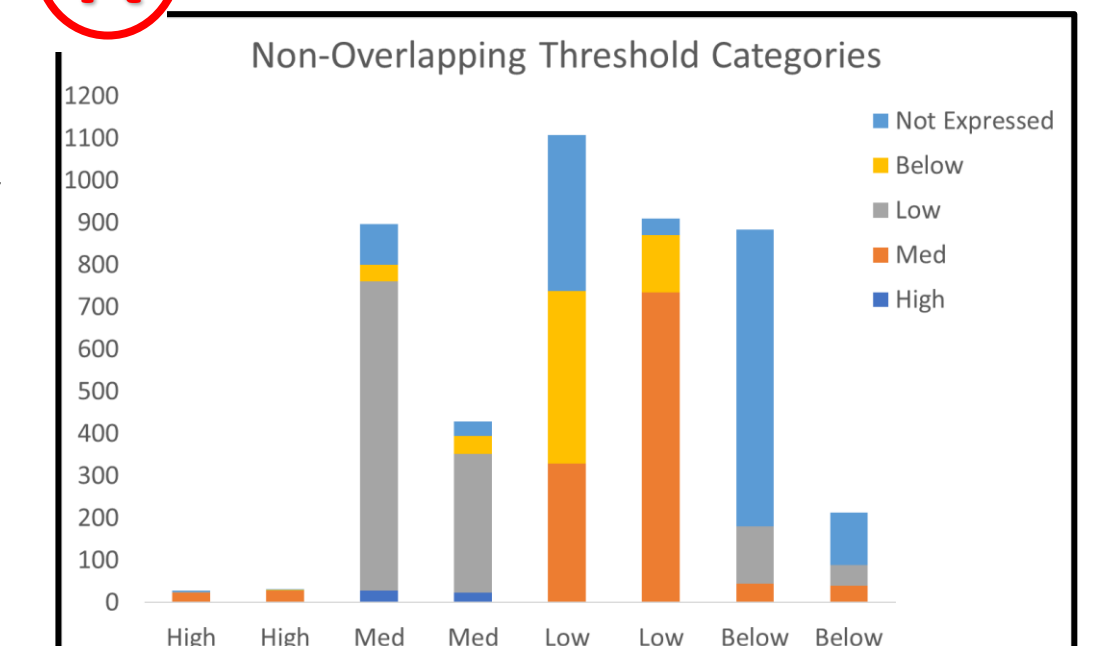
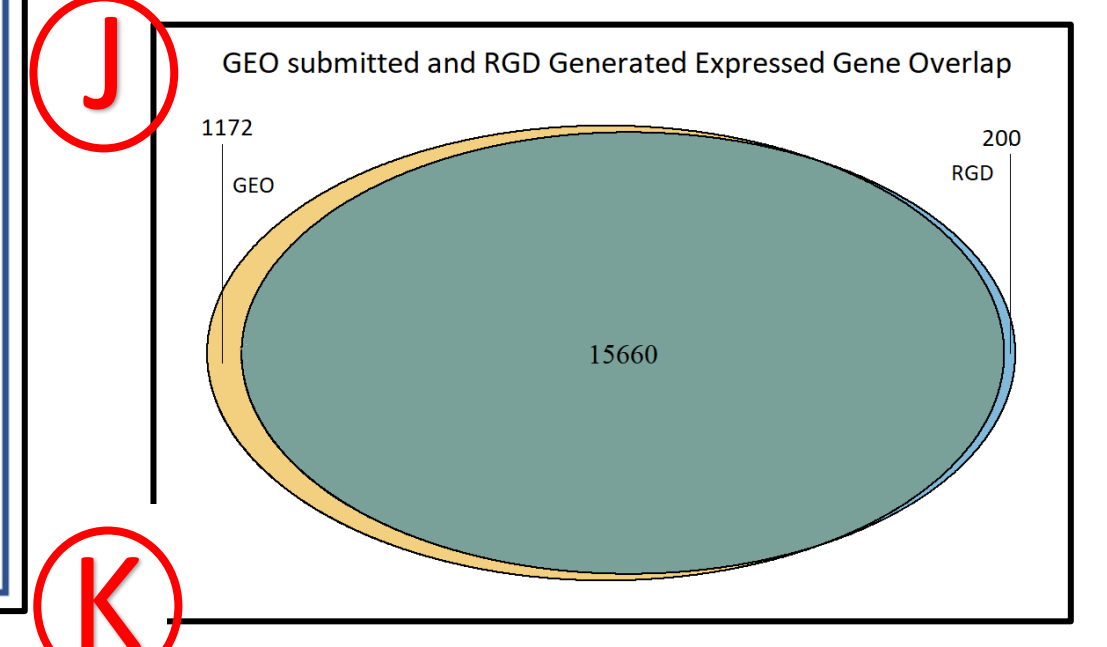
Phase 2: Data Re-analysis



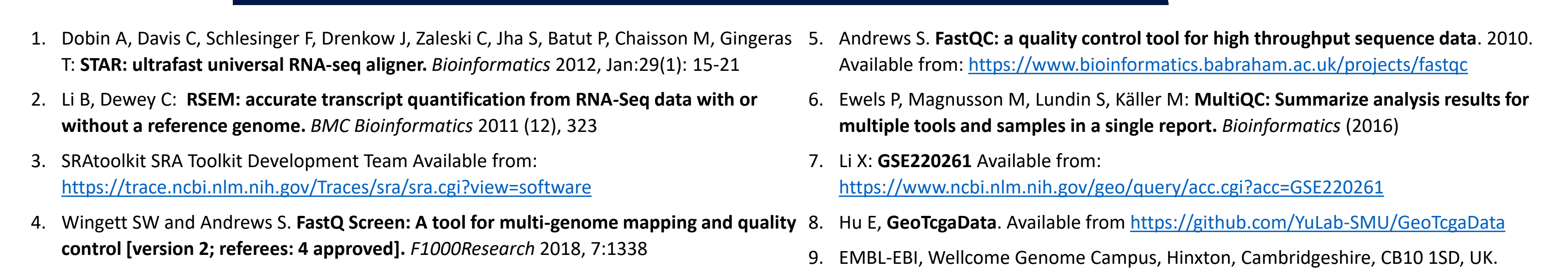
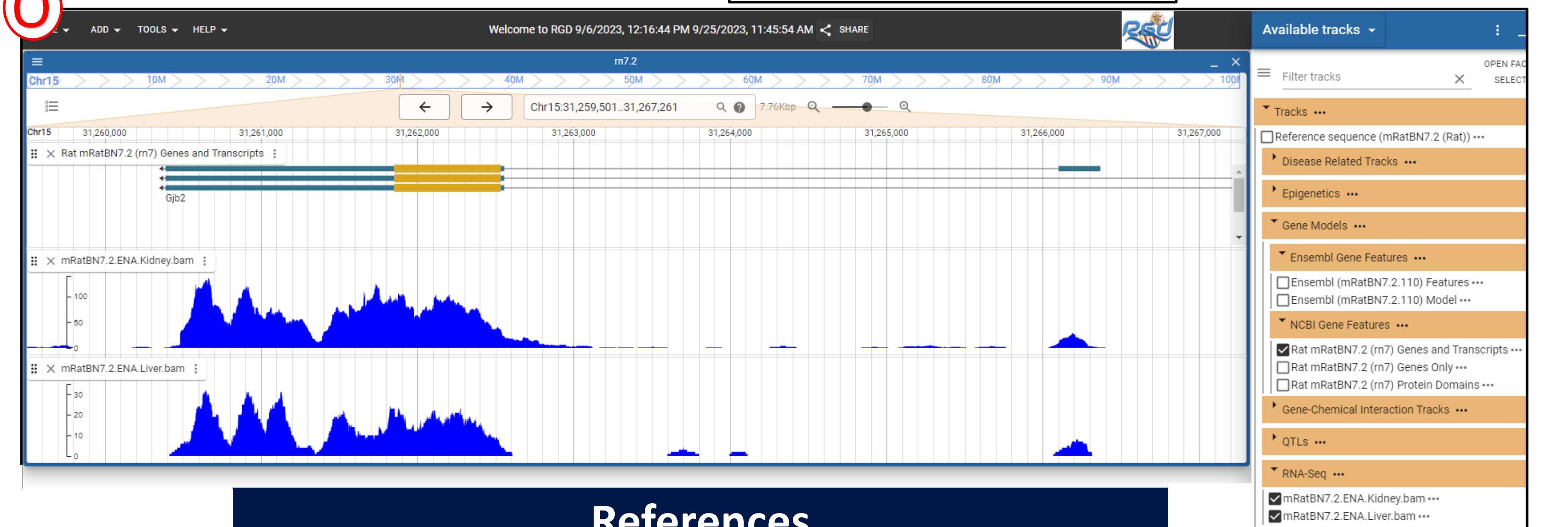
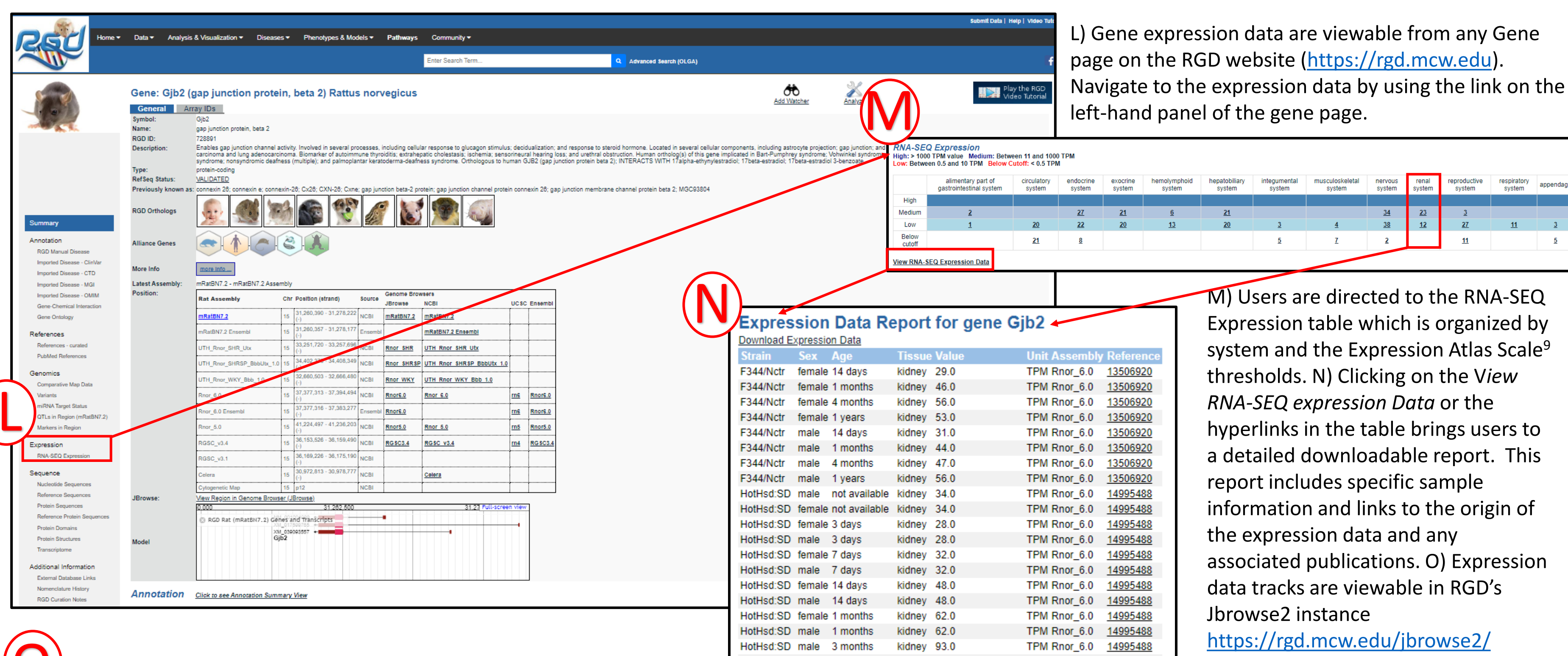
E) The Pre-processing stage of the data re-analysis pipeline utilizes FastQ Screen⁴ to verify the organism matches what is provided in the GEO Accession Viewer. (Inset) Both *R. norvegicus* and *R. rattus* assemblies are included in the tool to verify rat species. description and F) FastQC⁵ assesses the general quality of the SRA downloaded and converted FASTQ files (Data for all samples in series compiled by MultiQC⁶). MultiQC⁶ is used to assess the G) STAR¹ alignment and H) RSEM² mapping qualities.



I) Data for Accession GSE220261⁷ was used as a test to compare output by the RNAseq re-analysis pipeline. FPKM and count data were submitted to GEO and there is no associated publication. Other than knowing the data were aligned to 'rn7', very few details were provided about the data. FPKM data were converted to TPM with R library GeoTcgaData⁸. Data were normalized by limiting the analysis to gene definitions in the Ensembl v109 gtf. J) The overall overlap of expressed genes between the submitted and RGD pipeline results was almost 92% for both the FPKM and TPM analysis. K) Using the expression thresholds as defined by the Expression Atlas Scale⁹ the genes categorized as medium expression had the highest overlap rate (85%). The graph shows the count and classification for non-overlapping expression results within each threshold category.



Phase 3: Data Visualization



L) Gene expression data are viewable from any Gene page on the RGD website (<https://rgd.mcg.edu>). Navigate to the expression data by using the link on the left-hand panel of the gene page.

M) Users are directed to the RNA-SEQ Expression table which is organized by system and the Expression Atlas Scale⁹ thresholds. N) Clicking on the View RNA-SEQ expression Data or the hyperlinks in the table brings users to a detailed downloadable report. This report includes specific sample information and links to the origin of the expression data and any associated publications. O) Expression data tracks are viewable in RGD's JBrowse2 instance <https://rgd.mcg.edu/ibrowse2/>

Acknowledgements: We gratefully acknowledge our funding support from the National Institutes of Health (R01HL064541, U24HG010859, R24OD024617) and the researchers who faithfully use our website and data!

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