# Toolkit for quick and low-cost differential expression analysis of RNA-seq data

**Alexis Norris[1], Mamatha Garige[2], Susmita Ghosh[2], Carole Sourbier[2], Heather Lombardi[1], and Mayumi Miller[3]**
[1]Office of New Animal Drug Evaluation (CVM); [2]Office of Biotechnology Products (CDER); [3]Office of Research (CVM)

## Abstract

Bulk RNA sequencing (RNA-seq) is a powerful and widely used tool to identify biomarkers. The most common RNA-seq analysis is differential expression (DE) analysis, for which there are well-benchmarked analysis tools available. However, these tools are relatively inaccessible to scientists without bioinformatics expertise. Thus, we sought to create an end-to-end toolkit for DE analysis that could be quickly employed by all FDA scientists, using existing resources and with minimal cost to the user.

Our toolkit makes use of user-friendly, open-source interfaces: Galaxy, precisionFDA, and RStudio. First, precisionFDA or CDRH's Galaxy is used to trim the sequencing reads, align, and generate gene expression counts. Then, RStudio software is used to identify differentially expressed genes and pathways. Finally, RStudio's Shiny interactive dashboards are used to explore and visualize the results.
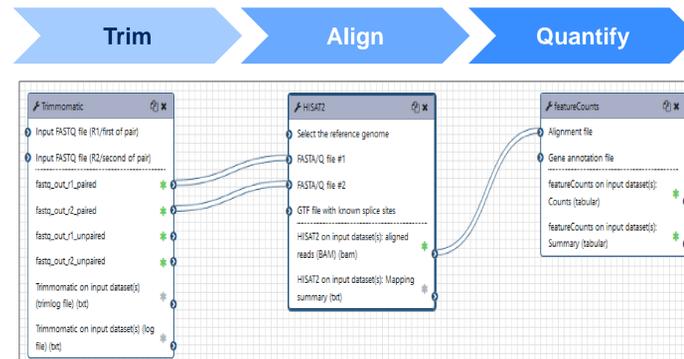
The entire analysis can be performed on a laptop, and typical analyses (6-12 samples) can be completed within one day. CDRH usage is low cost ($5 per CPU day) and Galaxy runs in a web browser without additional cost. RStudio software is free and can be run on PC, Mac, and Linux computers. The Shiny dashboards allows users to generate custom, publication quality figures and tables without costly bioinformatics support.

We have created a toolkit to perform RNA-seq analysis without the need for bioinformatics expertise. The analysis is low cost, quick, reproducible, and available to FDA scientists to support regulatory research and evaluation.

## Introduction

**Biomarker discovery with RNA-seq**
Differential expression analysis of RNA-seq data identifies genes which are dysregulated between groups. Pathways with differentially expressed genes indicate dysregulated pathways. These dysregulated genes and pathways are potential biomarkers for susceptibility, diagnosis, prognosis, treatment response, etc.

**Dysregulated genes** / **Disrupted pathways**

**Overview of toolkit analysis**
The toolkit has two parts and makes use of user-friendly, open-source interfaces that are publicly available: Galaxy, precisionFDA, and RStudio. Part I is performed using either Galaxy (**Figure 1**) or precisionFDA. Part II is performed using RStudio (**Figure 2**). Code and instructions are available at https://git.fda.gov/alexis.norris/rnaseq_de_toolkit

## Materials and Methods

**Dataset used to demonstrate the toolkit**
We used polyA-captured Illumina paired-end (75bp x 2) reads from Pertea, *et al.* Nature Methods 2016 (https://www.nature.com/articles/nprot.2016.095). The files are subsets of Geuvadis data, for chrX of six female and 6 male samples, available to download at ftp://ftp.ccb.jhu.edu/pub/RNAseq_protocol.

**Part I: Quantify gene expression**
1. Trim reads (trimmomatic)
2. Align reads to the reference genome (hisat2)
3. Quantify gene expression (featureCounts)

**Part II: Identify differentially expressed genes and pathways; visualize results**
1. Identify differentially expressed genes (edgeR)
2. Identify differentially expressed pathways (clusterProfiler)
3. Visualize results (RShiny)

Trim → Align → Quantify → Load Part I data → Identify genes → Identify pathways → Visualize genes → Visualize pathways

**Figure 1.** Part I using Galaxy. Public instance available at www.usegalaxy.org; FDA instance available from CDRH HPC group.
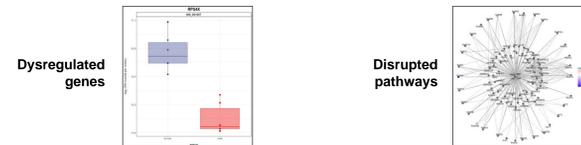
**Figure 2.** Part II. To explore the results (both tables and visually) and generate high-quality figures for publication, Shiny Apps can be deployed using RStudio. https://rstudio.com/
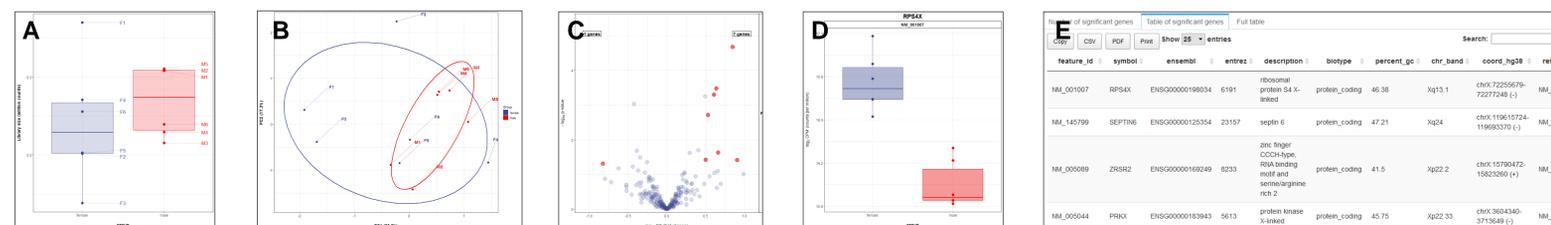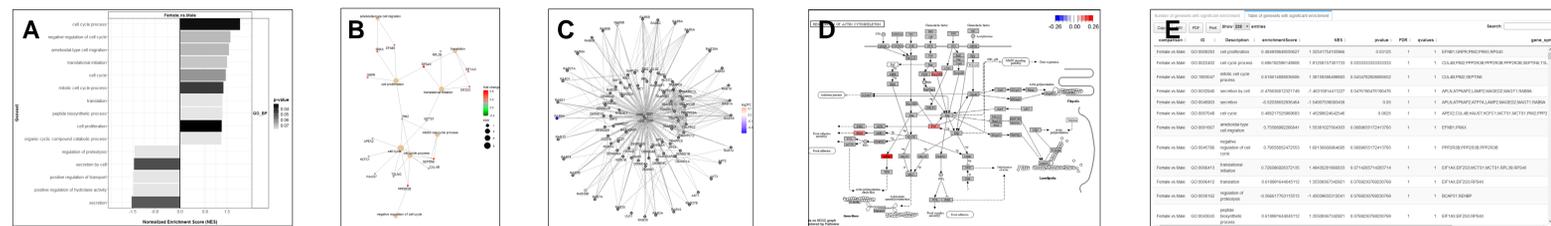
## Results and Discussion

**Figure 3.** Examples of gene visualizations. **A**. library sizes; **B**. Principal Component Analysis (PCA) clustering; **C**. volcano plot; **D**. boxplot of expression for a gene of interest; and **E**. table of full DE results that can be filtered and exported.

**Figure 4.** Examples of pathway visualizations. **A**. barplot of the top pathways; **B**. gene overlap (gene ~ concept map); **C**. Reactome pathway's network view; **D**. KEGG pathway; and **E**. table of full pathway results that can be filtered and exported.

## Conclusion

**RNA-seq DE Toolkit enables regulatory research**
*Quick*
- Typical analyses can be completed within 1 day
*Inexpensive*
- Galaxy is free; CDRH usage is only $5/CPU day
- precisionFDA is free
- RStudio is free
*Reproducible*
- Methods are documented
- Methods, data, and Shiny Apps can be easily shared
*Flexible*
- The Galaxy portion is now also available on precisionFDA
*Reduced barrier to RNA-seq analysis*
- No coding required for Galaxy, Shiny Apps, and precisionFDA
- Minimal coding for RStudio

**Current efforts**
*Expand user testing*
- The toolkit has been tested by beta users on PC laptops
- We're currently testing on Mac and Linux computers
*Eliminate the need to install RStudio software on laptop*
- Run Rmd and RShiny apps on precisionFDA directly

**For more information, contact:**
Alexis.Norris@fda.hhs.gov or Mayumi.Miller@fda.hhs.gov