

# Small and long RNA-Seq pipelines in the Genboree Workbench

NIH Extracellular RNA Communication Consortium  
2<sup>nd</sup> Investigators' Meeting  
May 19<sup>th</sup>, 2014

**Sai Lakshmi Subramanian** – [sailakss@bcm.edu](mailto:sailakss@bcm.edu) (Primary Contact)  
**Kevin Riehle** – [riehe@bcm.edu](mailto:riehe@bcm.edu)

Bioinformatics Research Laboratory, Baylor College of Medicine, Houston, TX  
Data Coordination Component (DCC) of exRNA Communication Consortium

**Robert Kitchen** - [rob.kitchen@yale.edu](mailto:rob.kitchen@yale.edu)

Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT  
Data Integration and Analysis Component (DCC) of exRNA Communication Consortium

# Outline

- Introduction to Genboree
- Genboree Workbench Basics
- Data analysis tools
  - RNA-Seq data – RSEQtools
  - Setting up long RNA-Seq data analysis in the Genboree Workbench
  - Small RNA-Seq data – smallRNA-seq Pipeline
  - Setting up small RNA-Seq data analysis in the Genboree Workbench

# Genboree Workbench Basics

- Genboree URL - <http://www.genboree.org/>
- Watch Video Tutorial – Introduction to Genboree Workbench  
[https://docs.google.com/file/d/0Bz3\\_YiJBA\\_j3Tk1uOFllazdMbkk/edit](https://docs.google.com/file/d/0Bz3_YiJBA_j3Tk1uOFllazdMbkk/edit)
- How do I create a new account in Genboree?  
[http://genboree.org/theCommons/ezfaq/show/public-commons?faq\\_id=493](http://genboree.org/theCommons/ezfaq/show/public-commons?faq_id=493)
- How to find Help within the Genboree Workbench?  
[http://genboree.org/theCommons/ezfaq/show/public-commons?faq\\_id=497](http://genboree.org/theCommons/ezfaq/show/public-commons?faq_id=497)
  - Help >> Getting Started
  - Help >> Tool Map
  - Help >> Tool Help Resources

## Genboree Core Services

- Authorization & access
- User data repository management

## Genboree User Data Services

### Genome-Based Tables

- Genomic annotations
- High-density score data

### Generic Entity Tables

- Loosely typed metadata
- AVP based

System/Network ▾ Data ▾ Genome ▾ Transcriptome ▾ Cistrome ▾ Epigenome ▾ Metagenome ▾ >>

Welcome to the Genboree Workbench! [Watch Intro Video](#)

**Toolset Menus**

Data Selector

Refresh Data Filter: Select a filter... ▾

Examples and Test Data

Databases

Bowtie - Example Data

Tracks

- Class: BT474
  - BT474Read:BT474Density
  - BT474ReadChrMXY:BT474DensityChrMXY
- Class: Gene
- Class: High Density Score Data
- Class: Marker
- Class: Sequence

Lists & Selections

Sample Sets

Samples

Files

- Bowtie
  - Bowtie Sample Output Custom Index
  - Bowtie Sample Output WholeGenome
  - BT474.subset.1.fastq.gz
  - BT474.subset.2.fastq.gz
- indexFiles
  - bowtie
    - hg19 custom chrMXY
    - hg19%20custom%20chrMXY\_bowtie2.tar.gz

### Details

Attribute	Value
Download	<a href="#">Click to Download File</a>
Group	Examples and Test Data
Database	Bowtie - Example Data
Description	
Name	BT474.subset.2.fastq.gz
CreatedDate	2014/05/02 16:32:52
LastModified	2014/05/02 16:32:53

### Input Data

↑ ↓ ✕ 🎨

- BT474.subset.1.fastq.gz
- BT474.subset.2.fastq.gz

**Drag & drop INPUT ENTITIES to Input Data Panel**

### Output Targets

↑ ↓ ✕ 🎨

- Bowtie - Example Data

**Drag & drop OUTPUT DESTINATIONS to Output Targets Panel**

# Genboree Workbench Basics

➤ Create Group

[http://genboree.org/theCommons/ezfaq/show/public-commons?faq\\_id=489](http://genboree.org/theCommons/ezfaq/show/public-commons?faq_id=489)

➤ Create Database

[http://genboree.org/theCommons/ezfaq/show/public-commons?faq\\_id=491](http://genboree.org/theCommons/ezfaq/show/public-commons?faq_id=491)

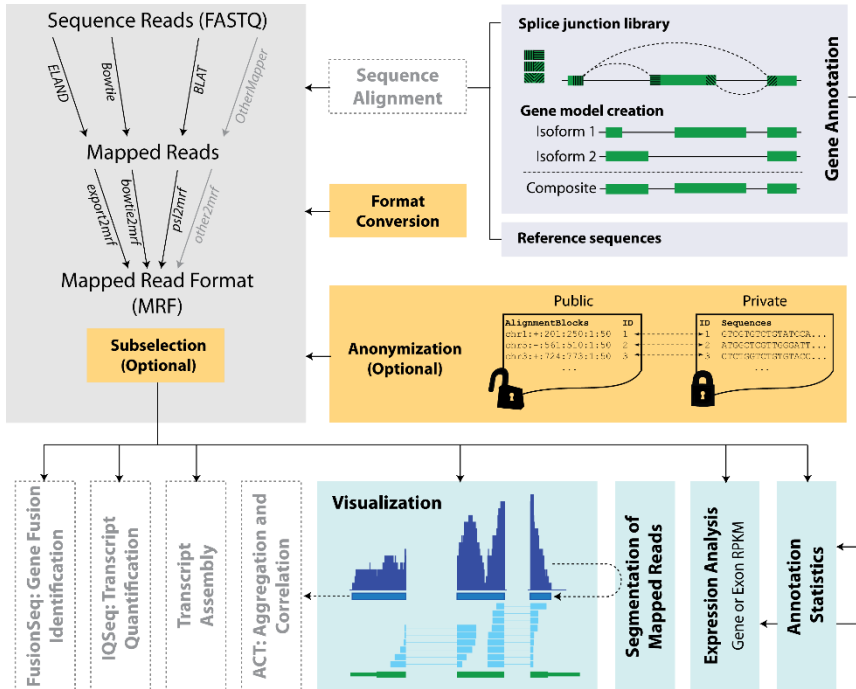
➤ Create Project

[http://genboree.org/theCommons/ezfaq/show/public-commons?faq\\_id=492](http://genboree.org/theCommons/ezfaq/show/public-commons?faq_id=492)

➤ Upload Data Files

[http://genboree.org/theCommons/ezfaq/show/public-commons?faq\\_id=495](http://genboree.org/theCommons/ezfaq/show/public-commons?faq_id=495)

# exRNA Data analysis tools



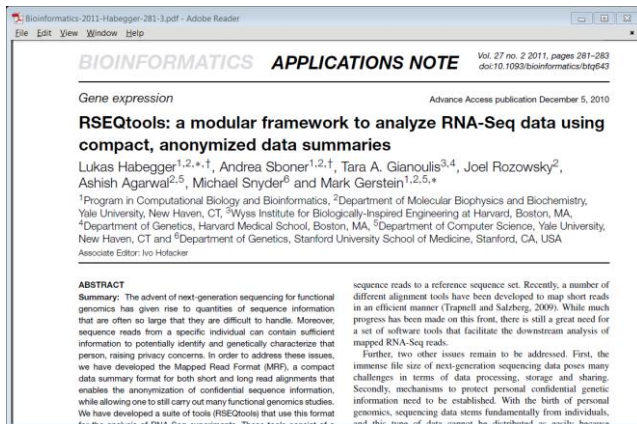
## RSEQtools

A modular and flexible framework to perform common RNA-Seq analysis tasks

### Current implemented pipeline

Generates genome-wide signal maps (hg18/19) as well as gene level expression quantifications (RPKM) for UCSC known genes.

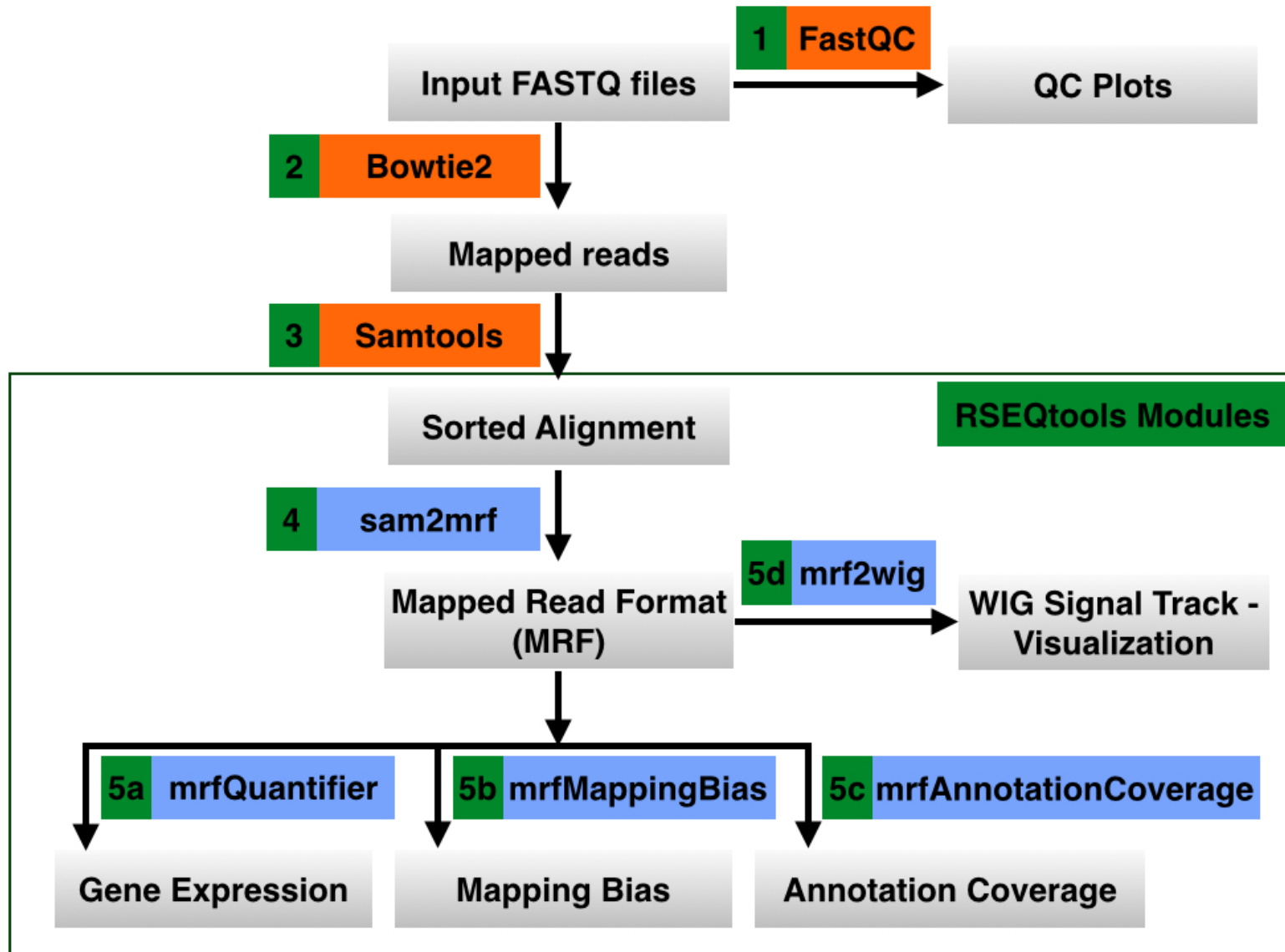
Developed by Gerstein Lab – DIAC, exRNA consortium



# RNA-Seq Analysis using RSEQtools

- **INPUT:** Takes 1 or 2 FASTQ sequences
- **Gene Annotations:** UCSC KnownGene for hg18, hg19
- **OUTPUT:**
  - Alignments in [SAM and BAM](#) formats - compressed in *AnalysisName\_alignments.tar.gz*
  - All result files compressed in *AnalysisName\_results.tar.gz*
  - Key result files for quick access (also found in *AnalysisName\_results.tar.gz*):
    - File with gene expression values
    - Mapping bias file
    - Annotation Coverage File
  - Signal Tracks of mapped reads (for visualization in genome browsers) - Uploaded as signal track in user db (if specified in the UI)
  - FastQC Output Plots
  - Custom Bowtie2 index, if generated

# RNA-Seq Analysis pipeline in the Genboree Workbench





# Example Dataset for RNA-Seq analysis using RSEQtools

**Data Selector**

Refresh Data Filter:

- genboree.org
  - Examples and Test Data
    - Databases
      - Bowtie - Example Data
      - RSEQtools - Example Data
        - RSEQtools hg18 - Example Data
          - Tracks **Signal Tracks for Visualization**
          - Lists & Selections
          - SampleSets
          - Samples
          - Files
            - FastQC **FASTQ reads QC outputs**
              - RSEQtools hg18 Sample Output Custom Index
              - RSEQtools hg18 Sample Output
            - indexFiles **Custom Bowtie2 indexes**
            - RSEQtools **RSEQtools Pipeline Result Files**
              - RSEQtools hg18 Sample Output Custom Index
              - RSEQtools hg18 Sample Output
                - jobFile.json
                - RSEQtools%20hg18%20Sample%20Output\_alignments.zip
                - RSEQtools%20hg18%20Sample%20Output\_coverage.txt
                - RSEQtools%20hg18%20Sample%20Output\_geneExpression.txt
                - RSEQtools%20hg18%20Sample%20Output\_mappingBias.txt
                - RSEQtools%20hg18%20Sample%20Output\_results.zip
              - sample.fastq.gz **RNA-seq Input FASTQ file**
          - Projects **Project page with FastQC output plots**
            - RSEQtools Example Project

# exRNA Data analysis tools

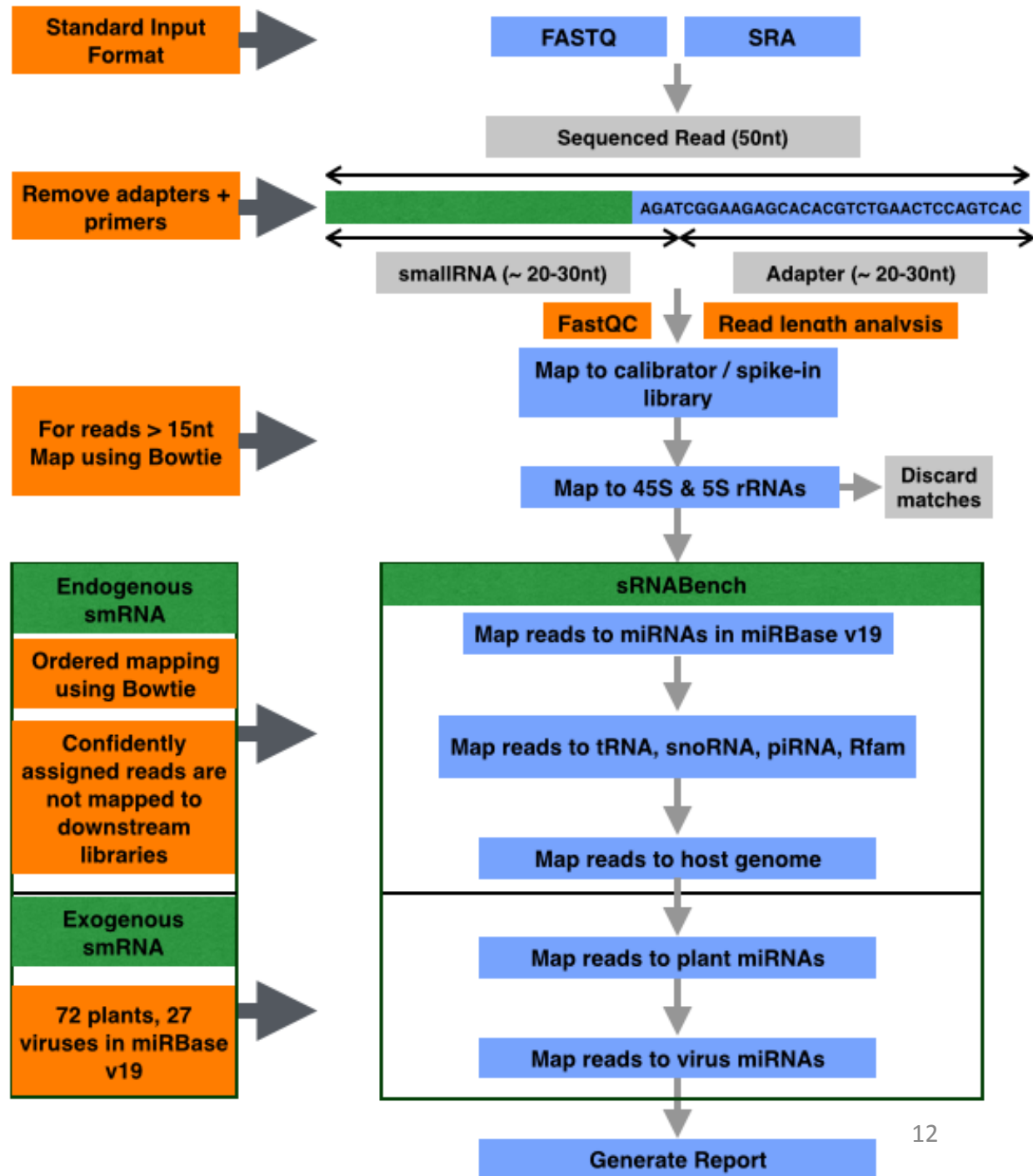
## smallRNA-seq Pipeline

- pre-processing: remove adapter & primer sequences
- contaminants: no poly-A purification leaves rRNAs, etc
- mapping: small reads tend to multi-map to large genomes
- diversity: many species of small RNA (not just miRNAs!)
- quantification: different normalization

# smallRNA-seq Pipeline

- **INPUT:** Takes a single-end FASTQ sequence file
- **Supported Genome build:** hg19
- **OUTPUT:**
  - All result files compressed in *AnalysisName\_results.tar.gz*
    - Read length distribution
    - Endogenous Mapping results – rRNA, miRNA, tRNA, snoRNA, piRNA, Rfam
    - Exogenous Mapping results – Plant and virus miRNAs in miRBase
    - 4 files containing sequence reads following
      - Adapter removal
      - rRNA removal
      - Endogenous mapping
      - Exogenous mappings

# smallRNA-seq pipeline in the Genboree Workbench



# Example Dataset for smallRNA-seq Pipeline

**Data Selector**

Refresh Data Filter:

- genboree.org
  - Examples and Test Data
    - Databases
      - Bowtie - Example Data
      - RSEQtools - Example Data
      - RSEQtools hg18 - Example Data
      - smallRNA-seq Pipeline - Example Data
        - Tracks
        - Lists & Selections
        - SampleSets
        - Samples
        - Files
          - smallRNAseqPipeline **smallRNA-seq Pipeline Result Files**
            - smallRNA-seq Pipeline Sample Output
              - jobFile.json
              - smallRNA-seq%20Pipeline%20Sample%20Output\_results.zip
            - SRR822433.fastq.gz smallRNA-seq exRNA Input FASTQ file**

# Live Demo of long and small RNA-seq pipelines in the Genboree Workbench

## Tutorial Screencasts to setup your analysis in the Genboree Workbench

- Long RNA-Seq data analysis - [https://docs.google.com/file/d/0Bz3\\_YiJBA\\_j3WHBpQmp2S0ljb0U/edit](https://docs.google.com/file/d/0Bz3_YiJBA_j3WHBpQmp2S0ljb0U/edit)
- small RNA-Seq data analysis - [https://docs.google.com/file/d/0Bz3\\_YiJBA\\_j3b1EzR3ZEWGFam8/edit](https://docs.google.com/file/d/0Bz3_YiJBA_j3b1EzR3ZEWGFam8/edit)