

# Getting Data and Databases - HW

Kelsey Martin

2024-03-07

---

## Assignment 6

is due on March 8th, 2024

---

### Instructions:

Please turn in the answers to this assignment as a knitted R markdown document.

- For this, don't forget to comment out all package installations after you install them! You don't want to keep installing things over and over!

Answers will be graded for correctness, completeness, and how well the instructions are followed.

Turn in your assignment on canvas

@ These are the learning objectives associated with each question

---

### QUESTION 1 (5 pts)

@ Students will display that they understand how to obtain data from the SRA.

Go to the SRA website.

This time, search for an organism of your interest, along with a tissue or condition of your interest. You can further narrow it down to a type of NGS platform (i.e., RNAseq) if you'd like.

Find something that you are interested in and that seems like a good candidate. Don't get anything too big either.

In your search results, **click on your sample of interest**. You should be pulled to a new page. Here, you will see a bunch of information on your sample. Under **Runs**, go ahead and click on your sample accession number. This will pull you to the SRA browser where you can see **Metadata, Analysis, Reads, Data Access, and FASTA/FASTQ download tabs**.

Go under FASTA/FASTQ and download the FASTQ file.

Go ahead and put your FASTQ files in a new directory called seq\_data. Make sure all file paths are correct!

Now read your file in.

- info about my files: This is ChIP-seq data that was used in a study on LSR2, which is a pleiotropic transcription factor that is differentially regulated between smooth and rough morphotypes of [*Mycobacterium abscessus*]. The one ending in 57 is the WT Rough morphotype, and the one ending in 52 is the WT Smooth morphotype.

```
#if (!require("BiocManager", quietly = TRUE))
#   install.packages("BiocManager")

#BiocManager::install("ShortRead")
library(ShortRead)
```

```
## Loading required package: BiocGenerics
```

```
##
```

```
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      IQR, mad, sd, var, xtabs
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      anyDuplicated, aperm, append, as.data.frame, basename, cbind,
##      colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,
##      get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,
##      match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
##      Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,
##      table, tapply, union, unique, unsplit, which.max, which.min
```

```
## Loading required package: BiocParallel
```

```
## Loading required package: Biostrings
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
##
```

```
## Attaching package: 'S4Vectors'
```

```
## The following object is masked from 'package:utils':
```

```
##
```

```
##      findMatches
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      expand.grid, I, unname
```

```
## Loading required package: IRanges
```

```

##
## Attaching package: 'IRanges'

## The following object is masked from 'package:grDevices':
##
##     windows

## Loading required package: XVector

## Loading required package: GenomeInfoDb

##
## Attaching package: 'Biostrings'

## The following object is masked from 'package:base':
##
##     strsplit

## Loading required package: Rsamtools

## Loading required package: GenomicRanges

## Loading required package: GenomicAlignments

## Loading required package: SummarizedExperiment

## Loading required package: MatrixGenerics

## Loading required package: matrixStats

##
## Attaching package: 'MatrixGenerics'

## The following objects are masked from 'package:matrixStats':
##
##     colAlls, colAnyNAs, colAnys, colAvgsPerRowSet, colCollapse,
##     colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
##     colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
##     colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
##     colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
##     colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
##     colWeightedMeans, colWeightedMedians, colWeightedSds,
##     colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgsPerColSet,
##     rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
##     rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
##     rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
##     rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
##     rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
##     rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
##     rowWeightedSds, rowWeightedVars

```

```
## Welcome to Bioconductor
##
```

```
##
## Attaching package: 'Biobase'
```

```
## The following objects are masked from 'package:matrixStats':
##
##   anyMissing, rowMedians
```

```
# returns sequence information
sread(seq1)
```

```
##          width seq
##      [1]   151 ATGGACGGTGTTTACGCCACTGGTTGGTCGA...GAACTCCAGTCACACTTGAATCACGTATGC
##      [2]   151 TACTACTCTACTCAGTGACTACTTGTCAAAC...GGGAAAGAGTGTAGATCTCGGTGGTCGCCG
##      [3]   151 ATCTGCACGATCTGCTGGTCACCGAGAAGGT...CCCGCCGAGGTGGAGATCGGAAGAGCACA
##      [4]   151 CCACCTCCGGCGGGCACGGCTCACCCGCCAC...CAGATCATGCAGATAGATCGGAAGAGCGTC
##      [5]   151 GATGCGGGAAGACCGTTCTGCGTTACGAGGA...GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
##      ...   ...   ...
## [3861612]  151 AGTCCGGCCCAGCTGCCCCCTCCCAGGAGGA...TTAAAAAGGGGGGGGGGGGGGGGGGGGGGG
## [3861613]  151 GATCGGAAGAGCACACGTCTGAATCCAGTC...GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
## [3861614]  151 GGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG...GGGGGGGGGGGGGGGGGGGGGGGGGGGGGG
## [3861615]  151 GAACATTTCTTCTTGATCACCTGGAATGCG...TCGGAAGAGCACACGTCTGAATCCAGTCT
## [3861616]  151 CCACGCGTTGCCGCACGCCGCATCCTGATG...TCGGAAGAGCGTCGTGTAGGGAAGAGTGT
```

[illegible]

```
# returns the read ID and length
id(seq1)
```

go ahead and read in another sample file or two from the same BioProject. You can add your additional code below.

**QUESTION 2 (5 pts)**

Using the samples you read in above, lets try out some more of shortread's functions to interrogate the data.

```
## class: ShortReadQ
## length: 6 reads; width: 151 cycles
```

```
## class: ShortReadQ
## length: 6 reads; width: 151 cycles
```





[illegible]



```

## 22      285 read SRR25491352.fastq.gz
## 23      284 read SRR25491352.fastq.gz
## 24      274 read SRR25491352.fastq.gz
## 25      273 read SRR25491352.fastq.gz
## 26      267 read SRR25491352.fastq.gz
## 27      253 read SRR25491352.fastq.gz
## 28      252 read SRR25491352.fastq.gz
## 29      224 read SRR25491352.fastq.gz
## 30      215 read SRR25491352.fastq.gz
## 31      212 read SRR25491352.fastq.gz
## 32      201 read SRR25491352.fastq.gz
## 33      200 read SRR25491352.fastq.gz
## 34      190 read SRR25491352.fastq.gz
## 35      186 read SRR25491352.fastq.gz
## 36      175 read SRR25491352.fastq.gz
## 37      166 read SRR25491352.fastq.gz
## 38      159 read SRR25491352.fastq.gz
## 39      159 read SRR25491352.fastq.gz
## 40      150 read SRR25491352.fastq.gz
## 41      149 read SRR25491352.fastq.gz
## 42      148 read SRR25491352.fastq.gz
## 43      146 read SRR25491352.fastq.gz
## 44      135 read SRR25491352.fastq.gz
## 45      130 read SRR25491352.fastq.gz
## 46      128 read SRR25491352.fastq.gz
## 47      128 read SRR25491352.fastq.gz
## 48      127 read SRR25491352.fastq.gz
## 49      127 read SRR25491352.fastq.gz
## 50      123 read SRR25491352.fastq.gz
## 51    6062 read SRR25491357.fastq.gz
## 52      356 read SRR25491357.fastq.gz
## 53      353 read SRR25491357.fastq.gz
## 54      172 read SRR25491357.fastq.gz
## 55      143 read SRR25491357.fastq.gz
## 56      129 read SRR25491357.fastq.gz
## 57      105 read SRR25491357.fastq.gz
## 58      103 read SRR25491357.fastq.gz
## 59       88 read SRR25491357.fastq.gz
## 60       87 read SRR25491357.fastq.gz
## 61       83 read SRR25491357.fastq.gz
## 62       80 read SRR25491357.fastq.gz
## 63       65 read SRR25491357.fastq.gz
## 64       65 read SRR25491357.fastq.gz
## 65       62 read SRR25491357.fastq.gz
## 66       61 read SRR25491357.fastq.gz
## 67       56 read SRR25491357.fastq.gz
## 68       54 read SRR25491357.fastq.gz
## 69       46 read SRR25491357.fastq.gz
## 70       44 read SRR25491357.fastq.gz
## 71       44 read SRR25491357.fastq.gz
## 72       40 read SRR25491357.fastq.gz
## 73       39 read SRR25491357.fastq.gz
## 74       39 read SRR25491357.fastq.gz
## 75       37 read SRR25491357.fastq.gz

```

```
## 76      35 read SRR25491357.fastq.gz
## 77      34 read SRR25491357.fastq.gz
## 78      33 read SRR25491357.fastq.gz
## 79      29 read SRR25491357.fastq.gz
## 80      29 read SRR25491357.fastq.gz
## 81      27 read SRR25491357.fastq.gz
## 82      24 read SRR25491357.fastq.gz
## 83      24 read SRR25491357.fastq.gz
## 84      24 read SRR25491357.fastq.gz
## 85      23 read SRR25491357.fastq.gz
## 86      23 read SRR25491357.fastq.gz
## 87      22 read SRR25491357.fastq.gz
## 88      22 read SRR25491357.fastq.gz
## 89      21 read SRR25491357.fastq.gz
## 90      21 read SRR25491357.fastq.gz
## 91      20 read SRR25491357.fastq.gz
## 92      20 read SRR25491357.fastq.gz
## 93      19 read SRR25491357.fastq.gz
## 94      18 read SRR25491357.fastq.gz
## 95      16 read SRR25491357.fastq.gz
## 96      16 read SRR25491357.fastq.gz
## 97      15 read SRR25491357.fastq.gz
## 98      15 read SRR25491357.fastq.gz
## 99      15 read SRR25491357.fastq.gz
## 100     14 read SRR25491357.fastq.gz
```

Are there any interesting things about your data files that you see?

I looked at the frequent sequences that occurred in the runs, and a lot of them were really long strings.

We are pretty limited by what we can do not together and within the scope of the class, so we won't delve deeper into the other shortread functions. But go ahead and check them out.

- The 'trimtails' function looks pretty useful!

There are other packages like shortread including "Rqc".

Rqc is an optimized tool designed for quality control and assessment of high-throughput sequencing data. It performs parallel processing of entire files and produces a report which contains a set of high-resolution graphics.

---

### QUESTION 3 (5 pts)

@ Students will gain an understanding of the different types of databases out there, how they differ, and

Pick one **primary database** that you think could be relevant to your research or interests and describe what data you could obtain from there that would help you moving forward.

Now do the same with one **secondary database** type.

Finally, do the same for one **specialized database** type.

Place your responses here:

Primary Database (raw data/info)

- [GenBank] (<https://www.ncbi.nlm.nih.gov/genbank/>): getting raw sequence information, comparing sequences

Secondary Database (Derived info, analysis from primary info)

- [PATRIC] (<https://www.bv-brc.org/>): This is a bioinformatics resource for bacterial human pathogens. It

Specialized Database

- [Mycobrowser] (<https://mycobrowser.epfl.ch/>): This is a repository for annotated genes in mycobacteria

- [TB Genome Annotation Portal] (<https://orca2.tamu.edu/U19/>): This is similar to Mycobrowser, but this

---

## QUESTION 4 (5 pts)

@ Students will learn about a specific R function of their choice, and implement an example of its usage

Choose an R function you think is interesting and would like to learn more about. Also mention the R package it can be found in.

*#I am interested in using BiomaRt to retrieve gene information and annotations*

```
if (!require("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
```

```
BiocManager::install("biomaRt")
```

```
## Bioconductor version 3.18 (BiocManager 1.30.22), R 4.3.2 (2023-10-31 ucrt)
```

```
## Old packages: 'GenomeInfoDb', 'S4Arrays'
```

Describe what the function does or how it can be used.

This package from BioConductor is a way to interact with gene information databases and get that information. It is most well documented for use with the [ENSEMBL database] (<https://useast.ensembl.org/index.html>). It also has pretty good documentation on [BioConductor] (<https://bioconductor.org/packages/release/bioc/>) and the [Github page] (<https://github.com/grimbough/biomaRt>) is up to date.

Give a functional example of its use. I want to try it too.

```
library(biomaRt)
# Specify we are using Ensembl dataset
ensembl <- useMart("ensembl")
#specify which species dataset to use
ensembl <- useDataset("hsapiens_gene_ensembl", mart = ensembl)
#One thing I can do with it is get the sequence of a gene I'm interested. Here I am getting the coding sequence
cbl_sequence <- getSequence(id = "ENSG00000110395", type = "ensembl_gene_id", mart = ensembl, seqType="coding")
cbl_sequence
```

```
##
## 1
## 2
## 3
## 4
## 5 ATGCCCGGCAACGTGAAGAAGAGCTCTGGGGCCGGGGCGGCAGCGGCTCCGGGGGCTCGGGTTCGGGTGGCCTGATTGGGCTCATGAAGGACGCCTT
## 6
##  ensembl_gene_id
## 1 ENSG00000110395
## 2 ENSG00000110395
## 3 ENSG00000110395
## 4 ENSG00000110395
## 5 ENSG00000110395
## 6 ENSG00000110395
```

Also provide a link to some documentation describing this function. The **NAME** notation is used for creating a hyperlink. You can provide the link below by filling in the syntax: