

Analysis of RNA-Seq Data with R/Bioconductor

...

Thomas Girke

December 8, 2012

Overview

RNA-Seq Analysis

Aligning Short Reads

Viewing Results in IGV Genome Browser

Outline

Overview

RNA-Seq Analysis

Aligning Short Reads

Viewing Results in IGV Genome Browser

Packages for RNA-Seq Analysis in R

- GenomicRanges [Link](#): high-level infrastructure for range data
- Rsamtools [Link](#): BAM support
- rtracklayer [Link](#): Annotation imports, interface to online genome browsers
- DESeq [Link](#): RNA-Seq DEG analysis
- edgeR [Link](#): RNA-Seq DEG analysis
- DEXSeq [Link](#): RNA-Seq Exon analysis

RNA-Seq versus DGE

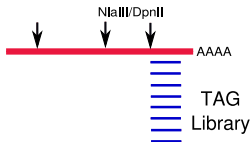
RNA-seq



Sequencing ↓

1. Alternative splicing
2. Limited expression profiling
3. SNP detection
4. Many other applications

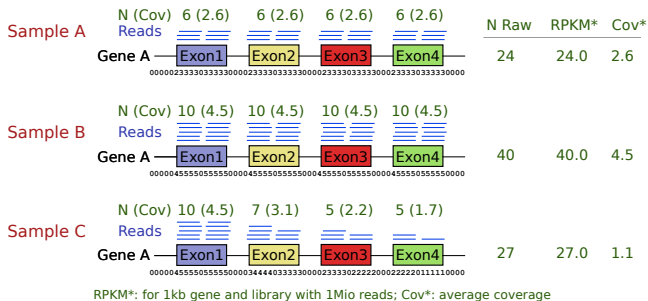
DGE



Sequencing ↓

1. Expression profiling
→ more appropriate for many biosamples

Identification of Differentially Expressed Genes



Normalization often by library size.

RNA-Seq Analysis Workflow

- Read mapping
- Counting reads overlapping with genes
- Analysis of differentially expressed genes (DEGs)
- Clustering of co-expressed genes
- Gene set/GO term enrichment analysis

Outline

Overview

RNA-Seq Analysis
Aligning Short Reads

Viewing Results in IGV Genome Browser

Data Sets and Experimental Variables

- To make the following sample code work, please download and unpack the sample data [Link](#) in the directory of your current R session.
- It contains four simplified alignment files from RNA-Seq experiment SRA023501 [Link](#) and a shortened GFF to allow fast analysis on a laptop.
- The alignments were created by aligning the reads with Bowtie against the Arabidopsis reference genome.
- **Note:** usually, the aligned reads would be stored in BAM format and then imported into R with the `readBamGappedAlignments` function (see below)!

This information could be imported from an external targets file

```
> targets <- read.delim("./data/targets.txt")  
> targets
```

| | Samples | Factor | Fastq |
|---|---------|--------|-----------------|
| 1 | AP3_f14 | AP3 | SRR064154.fastq |
| 2 | AP3_f14 | AP3 | SRR064155.fastq |
| 3 | T1_f14 | TRL | SRR064166.fastq |
| 4 | T1_f14 | TRL | SRR064167.fastq |

Outline

Overview

RNA-Seq Analysis
Aligning Short Reads

Viewing Results in IGV Genome Browser

Align Reads and Output Indexed Bam Files

Note: this step requires the command-line tool `bowtie2` [Link](#). If it is not available on a system then one can skip this mapping step and download the pre-generated Bam files from here: [Link](#)

```
> library(Rsamtools)
> dir.create("results") # Note: all output data will be written to directory 'results'
> system("bowtie2-build ./data/tair10chr.fasta ./data/tair10chr.fasta") # Build indexed reference genome
> targets <- read.delim("./data/targets.txt") # Import experiment design information
> targets
> input <- paste("./data/", targets[,3], sep="")
> output <- paste("./results/", targets[,3], ".sam", sep="")
> reference <- "./results/tair10chr.fasta"
> for(i in seq(along=targets[,3])) {
+   command <- paste("bowtie2 -x ./data/tair10chr.fasta -U", input[i], "-S", output[i])
+   system(command)
+   asBam(file=output[i], destination=gsub(".sam", "", output[i]), overwrite=TRUE, indexDestination=1)
+   unlink(output[i])
+ }
```

Import Annotation Data from GFF

Annotation data from GFF

```
> library(rtracklayer); library(GenomicRanges); library(Rsamtools)
> gff <- import.gff("./data/TAIR10_GFF3_trunc.gff", asRangedData=FALSE)
> seqlengths(gff) <- end(ranges(gff[which(elementMetadata(gff)[,"type"]=="chromosome"),]))
> subgene_index <- which(elementMetadata(gff)[,"type"] == "exon")
> gffsub <- gff[subgene_index,] # Returns only gene ranges
> gffsub[1:4, c(2,5)]
```

GRanges with 4 ranges and 2 metadata columns:

| | seqnames | ranges | strand | type | group |
|-----|----------|--------------|--------|----------|--------------------|
| | <Rle> | <IRanges> | <Rle> | <factor> | <factor> |
| [1] | Chr1 | [3631, 3913] | + | exon | Parent=AT1G01010.1 |
| [2] | Chr1 | [3996, 4276] | + | exon | Parent=AT1G01010.1 |
| [3] | Chr1 | [4486, 4605] | + | exon | Parent=AT1G01010.1 |
| [4] | Chr1 | [4706, 5095] | + | exon | Parent=AT1G01010.1 |

seqlengths:

| Chr1 | Chr2 | Chr3 | Chr4 | Chr5 | ChrC | ChrM |
|--------|--------|--------|--------|--------|--------|--------|
| 100000 | 100000 | 100000 | 100000 | 100000 | 100000 | 100000 |

```
> ids <- gsub("Parent=|\\.\\.\\.*", "", elementMetadata(gffsub)$group)
> gffsub <- split(gffsub, ids) # Coerce to GRangesList
```

Read Counting for Exonic Gene Ranges - Old

Number of reads overlapping gene ranges

```
> samples <- as.character(targets$Fastq)
> samplespath <- paste("./results/", samples, ".bam", sep="")
> names(samplespath) <- samples
> countDF <- data.frame(row.names=names(gffsub))
> for(i in samplespath) {
+   aligns <- readBamGappedAlignments(i) # Substitute next two lines with this one.
+   counts <- countOverlaps(gffsub, aligns, ignore.strand=TRUE)
+   countDF <- cbind(countDF, counts)
+ }
> colnames(countDF) <- samples
> countDF[1:4,]
```

| | SRR064154.fastq | SRR064155.fastq | SRR064166.fastq | SRR064167.fastq |
|-----------|-----------------|-----------------|-----------------|-----------------|
| AT1G01010 | 52 | 26 | 60 | 75 |
| AT1G01020 | 146 | 77 | 82 | 64 |
| AT1G01030 | 5 | 1 | 13 | 14 |
| AT1G01040 | 483 | 347 | 302 | 358 |

```
> write.table(countDF, "./results/countDF", quote=FALSE, sep="\t", col.names = NA)
> countDF <- read.table("./results/countDF")
```

Read Counting for Exonic Gene Ranges - New

The `summarizeOverlaps` function from the `GenomicRanges` is easier to use and provides more options. See [here](#) [Link](#) for details.

```
> bfl <- BamFileList(samplespath, index=character())
> countDF2 <- summarizeOverlaps(gffsub, bfl, mode="Union", ignore.strand=TRUE)
> countDF2 <- assays(countDF2)$counts
> countDF2[1:4,]
```

| | SRR064154.fastq | SRR064155.fastq | SRR064166.fastq | SRR064167.fastq |
|-----------|-----------------|-----------------|-----------------|-----------------|
| AT1G01010 | 52 | 26 | 60 | 75 |
| AT1G01020 | 146 | 77 | 82 | 64 |
| AT1G01030 | 5 | 1 | 13 | 14 |
| AT1G01040 | 480 | 346 | 282 | 335 |

Simple RPKM Normalization

RPKM: reads per kilobase of exon model per million mapped reads

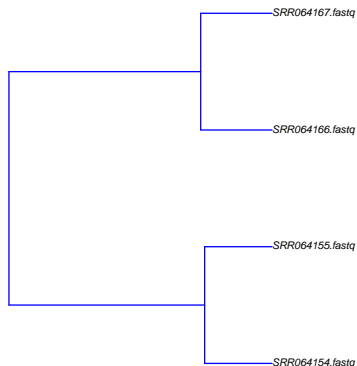
```
> returnRPKM <- function(counts, gffsub) {  
+   geneLengthsInKB <- sum(width(reduce(gffsub)))/1000 # Length of exon union per gene in kbp  
+   millionsMapped <- sum(counts)/1e+06 # Factor for converting to million of mapped reads.  
+   rpm <- counts/millionsMapped # RPK: reads per kilobase of exon model.  
+   rpkm <- rpm/geneLengthsInKB # RPKM: reads per kilobase of exon model per million mapped reads.  
+   return(rpkm)  
+ }  
> countDFrpkm <- apply(countDF, 2, function(x) returnRPKM(counts=x, gffsub=gffsub))  
> countDFrpkm[1:4,]
```

| | SRR064154.fastq | SRR064155.fastq | SRR064166.fastq | SRR064167.fastq |
|-----------|-----------------|-----------------|-----------------|-----------------|
| AT1G01010 | 52.478177 | 24.641790 | 387.88096 | 365.35877 |
| AT1G01020 | 140.199699 | 69.439799 | 504.40560 | 296.65869 |
| AT1G01030 | 4.471187 | 0.839801 | 74.46772 | 60.43155 |
| AT1G01040 | 131.564008 | 88.765242 | 526.94917 | 470.71264 |

QC Check

QC check by computing a sample correlating matrix and plotting it as a tree

```
> library(ape)
> d <- cor(countDF, method="spearman")
> hc <- hclust(dist(1-d))
> plot.phylo(as.phylo(hc), type="p", edge.col=4, edge.width=2, show.node.label=TRUE, no.margin=TRUE)
```



Identify DEGs with Simple Fold Change Method

Compute mean values for replicates

```
> source("http://faculty.ucr.edu/~tgirke/Documents/R_BioCond/My_R_Scripts/colAg.R")
> countDFrpk_mean <- colAg(myMA=countDFrpk, group=c(1,1,2,2), myfct=mean)
> countDFrpk_mean[1:4,]
```

| | SRR064154.fastq_SRR064155.fastq | SRR064166.fastq_SRR064167.fastq |
|-----------|---------------------------------|---------------------------------|
| AT1G01010 | 38.559984 | 376.61987 |
| AT1G01020 | 104.819749 | 400.53214 |
| AT1G01030 | 2.655494 | 67.44964 |
| AT1G01040 | 110.164625 | 498.83091 |

Log2 fold changes

```
> countDFrpk_mean <- cbind(countDFrpk_mean, log2ratio=log2(countDFrpk_mean[,2]/countDFrpk_mean[,1]))
> countDFrpk_mean <- countDFrpk_mean[is.finite(countDFrpk_mean[,3]), ]
> degs2fold <- countDFrpk_mean[countDFrpk_mean[,3] >= 1 | countDFrpk_mean[,3] <= -1,]
> degs2fold[1:4,]
```

| | SRR064154.fastq_SRR064155.fastq | SRR064166.fastq_SRR064167.fastq | log2ratio |
|-----------|---------------------------------|---------------------------------|-----------|
| AT1G01010 | 38.559984 | 376.61987 | 3.287933 |
| AT1G01020 | 104.819749 | 400.53214 | 1.934007 |
| AT1G01030 | 2.655494 | 67.44964 | 4.666758 |
| AT1G01040 | 110.164625 | 498.83091 | 2.178890 |

```
> write.table(degs2fold, "./results/degs2fold", quote=FALSE, sep="\t", col.names = NA)
> degs2fold <- read.table("./results/degs2fold")
```

Identify DEGs with DESeq Library

Raw count data are expected here!

```
> library(DESeq)
> countDF <- read.table("./results/countDF")
> conds <- targets$Factor
> cds <- newCountDataSet(countDF, conds) # Creates object of class CountDataSet derived from eSet class
> counts(cds)[1:4, ] # CountDataSet has similar accessor methods as eSet class.
```

| | SRR064154.fastq | SRR064155.fastq | SRR064166.fastq | SRR064167.fastq |
|-----------|-----------------|-----------------|-----------------|-----------------|
| AT1G01010 | 52 | 26 | 60 | 75 |
| AT1G01020 | 146 | 77 | 82 | 64 |
| AT1G01030 | 5 | 1 | 13 | 14 |
| AT1G01040 | 483 | 347 | 302 | 358 |

```
> cds <- estimateSizeFactors(cds) # Estimates library size factors from count data. Alternatively, one can
> cds <- estimateDispersions(cds) # Estimates the variance within replicates
> res <- nbinomTest(cds, "AP3", "TRL") # Calls DEGs with nbinomTest
> res <- na.omit(res)
> res2fold <- res[res$log2FoldChange >= 1 | res$log2FoldChange <= -1,]
> res2foldpadj <- res2fold[res2fold$padj <= 0.05, ]
> res2foldpadj[1:4,1:8]
```

| | id | baseMean | baseMeanA | baseMeanB | foldChange | log2FoldChange | pval | padj |
|----|-----------|-----------|-----------|------------|------------|----------------|--------------|--------------|
| 6 | AT1G01050 | 600.91989 | 273.11390 | 928.725867 | 3.40050744 | 1.765750 | 5.341934e-12 | 5.876128e-11 |
| 7 | AT1G01060 | 302.03514 | 169.61662 | 434.453652 | 2.56138611 | 1.356925 | 4.032316e-06 | 2.710612e-05 |
| 8 | AT1G01070 | 29.86593 | 5.66738 | 54.064485 | 9.53959078 | 3.253927 | 2.421810e-05 | 1.274082e-04 |
| 15 | AT2G01008 | 18.59928 | 34.71273 | 2.485829 | 0.07161145 | -3.803666 | 2.417262e-04 | 1.083291e-03 |

Identify DEGs with edgeR Library

Raw count data are expected here!

```
> library(edgeR)
> countDF <- read.table("./results/countDF")
> y <- DGEList(counts=countDF, group=conds) # Constructs DGEList object
> y <- estimateCommonDisp(y) # Estimates common dispersion
> y <- estimateTagwiseDisp(y) # Estimates tagwise dispersion
> et <- exactTest(y, pair=c("AP3", "TRL")) # Computes exact test for the negative binomial distribution.
> topTags(et, n=4)
```

Comparison of groups: TRL-AP3

| | logFC | logCPM | PValue | FDR |
|-----------|----------|----------|--------------|--------------|
| AT4G00050 | 5.568721 | 10.34789 | 5.473636e-54 | 8.265190e-52 |
| AT1G01050 | 4.203665 | 12.09034 | 4.331645e-45 | 3.270392e-43 |
| AT3G01120 | 3.865628 | 14.07787 | 1.674890e-41 | 8.430278e-40 |
| AT1G01060 | 3.757064 | 10.98503 | 3.360358e-34 | 1.268535e-32 |

```
> edge <- as.data.frame(topTags(et, n=50000))
> edge2fold <- edge[edge$logFC >= 1 | edge$logFC <= -1,]
> edge2foldpadj <- edge2fold[edge2fold$FDR <= 0.01, ]
```

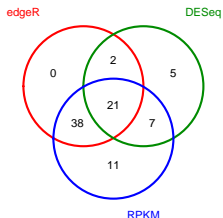
Merge Results and Compute Overlaps Among Methods

```
> bothDF <- merge(res, countDFrpkm_mean, by.x=1, by.y=0, all=TRUE); bothDF <- na.omit(bothDF)
> cor(bothDF[, "log2FoldChange"], bothDF[, "log2ratio"], method="spearman")
```

```
[1] 0.9985348
```

```
> source("http://faculty.ucr.edu/~tgirke/Documents/R_BioCond/My_R_Scripts/overLapper.R")
> setlist <- list(edgeR=rownames(edge2foldpadj), DESeq=as.character(res2foldpadj[,1]), RPKM=rownames(degs2f
> OList <- overLapper(setlist=setlist, sep="_", type="vennsets")
> counts <- sapply(OList$Venn_List, length)
> vennPlot(counts=counts)
```

Venn Diagram



Unique objects: All = 84; S1 = 61; S2 = 35; S3 = 77

Enrichment of GO Terms in DEG Sets

GO Term Enrichment Analysis

```
> library(GOstats); library(GO.db); library(ath1121501.db)
> geneUniverse <- rownames(countDF)
> geneSample <- res2foldpadj[,1]
> params <- new("GOHyperGParams", geneIds = geneSample, universeGeneIds = geneUniverse,
+             annotation="ath1121501", ontology = "MF", pvalueCutoff = 0.5,
+             conditional = FALSE, testDirection = "over")
> hgOver <- hyperGTest(params)
> summary(hgOver)[1:4,]
```

| | GOMFID | Pvalue | OddsRatio | ExpCount | Count | Size | |
|---|------------|-------------|-----------|----------|-------|------|---|
| 1 | GO:0008324 | 0.002673178 | 18 | 2.126582 | 6 | 7 | cation transmembrane transport |
| 2 | GO:0015075 | 0.002673178 | 18 | 2.126582 | 6 | 7 | ion transmembrane transport |
| 3 | GO:0015077 | 0.002673178 | 18 | 2.126582 | 6 | 7 | monovalent inorganic cation transmembrane transport |
| 4 | GO:0015078 | 0.002673178 | 18 | 2.126582 | 6 | 7 | hydrogen ion transmembrane transport |

```
> htmlReport(hgOver, file = "data/MyhyperGresult.html")
```

Outline

Overview

RNA-Seq Analysis
Aligning Short Reads

Viewing Results in IGV Genome Browser

Inspect Results in IGV

View results in IGV

- 1 Download and open IGV [Link](#)
- 2 Select in menu in top left corner A. thaliana (TAIR10)
- 3 Upload the following indexed/sorted Bam files with File -> Load from URL...

http://faculty.ucr.edu/~tgirke/HTML_Presentations/Manuals/Workshop_Dec_6_10_2012/Rrnaseq/results/SRR064154.fastq

http://faculty.ucr.edu/~tgirke/HTML_Presentations/Manuals/Workshop_Dec_6_10_2012/Rrnaseq/results/SRR064155.fastq

http://faculty.ucr.edu/~tgirke/HTML_Presentations/Manuals/Workshop_Dec_6_10_2012/Rrnaseq/results/SRR064166.fastq

http://faculty.ucr.edu/~tgirke/HTML_Presentations/Manuals/Workshop_Dec_6_10_2012/Rrnaseq/results/SRR064167.fastq

- 4 To view area of interest, enter its coordinates Chr1:49,457-51,457 in position menu on top.



Analysis of Differential Exon Usage with DEXSeq

Number of reads overlapping gene ranges

```
> source("data/Fct/gffexonDEXSeq.R")
> gffexonDEXSeq <- exons2DEXSeq(gff=gff)
> ids <- as.character(elementMetadata(gffexonDEXSeq)[, "ids"])
> countDFdex <- data.frame(row.names=ids)
> for(i in samplespath) {
+   aligns <- readBamGappedAlignments(i) # Substitute next two lines with this one.
+   counts <- countOverlaps(gffexonDEXSeq, aligns)
+   countDFdex <- cbind(countDFdex, counts)
+ }
> colnames(countDFdex) <- samples
> countDFdex[1:4,1:2]
```

| | SRR064154.fastq | SRR064155.fastq |
|---|-----------------|-----------------|
| Parent=AT1G01010:E001__Chr1_3631_3913+_Parent=AT1G01010.1 | 2 | 4 |
| Parent=AT1G01010:E002__Chr1_3996_4276+_Parent=AT1G01010.1 | 2 | 1 |
| Parent=AT1G01010:E003__Chr1_4486_4605+_Parent=AT1G01010.1 | 3 | 3 |
| Parent=AT1G01010:E004__Chr1_4706_5095+_Parent=AT1G01010.1 | 6 | 1 |

```
> write.table(countDFdex, "./results/countDFdex", quote=FALSE, sep="\t", col.names = NA)
> countDFdex <- read.table("./results/countDFdex")
```


Analysis of Differential Exon Usage with DEXSeq

Identify genes with differential exon usage

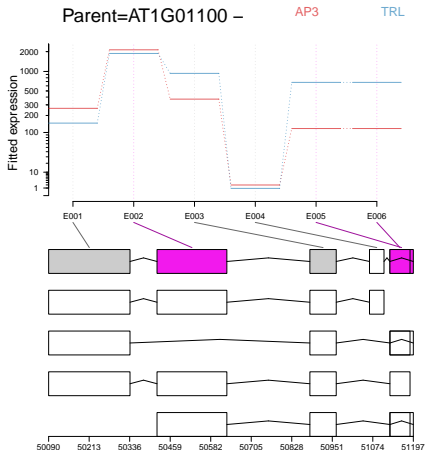
```
> library(DEXSeq)
> samples <- as.character(targets$Factor); names(samples) <- targets$Fastq
> countDFdex[is.na(countDFdex)] <- 0
> ## Construct ExonCountSet from scratch
> exset <- newExonCountSet2(countDF=countDFdex) # fData(exset)[1:4,]
> ## Performs normalization
> exset <- estimateSizeFactors(exset)
> ## Evaluate variance of the data by estimating dispersion using Cox-Reid (CR) likelihood estimation
> exset <- estimateDispersions(exset)
> ## Fits dispersion-mean relation to the individual CR dispersion values
> exset <- fitDispersionFunction(exset)
> ## Performs Chi-squared test on each exon and Benjmini-Hochberg p-value adjustment for mutliple testing
> exset <- testForDEU(exset)
> ## Estimates fold changes of exons
> exset <- estimateLog2FoldChanges(exset)
> ## Obtain results in data frame
> deuDF <- DEUresultTable(exset)
> ## Count number of genes with differential exon usage
> table(tapply(deuDF$padjust < 0.01, geneIDs(exset), any))

FALSE TRUE
    20    1
```

DEXSeq Plots

Sample plot showing fitted expression of exons

```
> plotDEXSeq(exset, "Parent=AT1G01100", displayTranscripts=TRUE, expression=TRUE, legend=TRUE)
> ## Generate many plots and write them to results directory
> mygeneIDs <- unique(as.character(na.omit(deuDF[deuDF$geneID %in% unique(deuDF$geneID,)]["geneID"])))
> DEXSeqHTML(exset, geneIDs=mygeneIDs, path="results", file="DEU.html")
```



Session Information

```
> sessionInfo()
```

```
R version 2.15.1 (2012-06-22)
```

```
Platform: x86_64-unknown-linux-gnu (64-bit)
```

```
locale:
```

```
[1] C
```

```
attached base packages:
```

```
[1] stats      graphics  utils      datasets  grDevices  methods    base
```

```
other attached packages:
```

```
[1] DEXSeq_1.4.0           xtable_1.7-0           ath1121501.db_2.8.0   org.At.tair.db_2.8.0  GO.db_2.8.0  
[11] AnnotationDbi_1.20.1  edgeR_3.0.0            limma_3.14.1          DESeq_1.10.1          lattice_0.20-10  
[21] rtracklayer_1.18.0   GenomicRanges_1.10.2  IRanges_1.16.2        BiocGenerics_0.4.0
```

```
loaded via a namespace (and not attached):
```

```
[1] AnnotationForge_1.0.2 BSgenome_1.26.1       GSEABase_1.20.0       RBGL_1.34.0           RColorBrewer_1.1-2  
[10] bitops_1.0-4.1       gee_4.13-18           genefilter_1.40.0     geneplotter_1.36.0    grid_2.15.1  
[19] splines_2.15.1       statmod_1.4.16        stats4_2.15.1         stringr_0.6.1         survival_2.36-1
```