```
###############################################################################
# Mapping and counting RNA-seq reads with Rsubread R Bioconductor package #
###############################################################################

#####
# Toy example
#####
# small sample of reads (56,168 read-pairs) from a single replicate
# sequencing: strand specific, paired-end, Illumina, 100-100nts
# small chromosome fragment (1-839990 nt part of the 3L chromosome of Drosophila
simulans contains 96 genes)

#####
# Mapping
#####
# Call Rsubread library
library(Rsubread)

# Indexing reference genome
buildindex(basename="Genome_sequence_ToyExample_index",
reference="../Data/Genome_sequence_ToyExample.fa", memory=1500)
# p7 at Rsubread.pdf: An index needs to be built before read mapping can be
performed. This function creates a hash
# table for the reference genome, which can then be used by Subread and Subjunc
aligners for read
# alignment.
# Highly repetitive subreads (or uninformative subreads) are excluded from the
hash table so as to
# reduce mapping ambiguity.

#####?
# Check how many uninformative subreads were found!
# Open Genome_sequence_ToyExample.fa fasta file with "less" unix command in
another terminal.
# Can you see potential uninformative regions in the genome?
#####?

# Mapping reads
subjunc(index="Genome_sequence_ToyExample_index",
readfile1="../Data/Trimmed_reads_ToyExample_1.fq",
readfile2="../Data/Trimmed_reads_ToyExample_2.fq",
input_format="FASTQ",
output_format="BAM", output_file="Mapped_reads_ToyExample.BAM",
nthreads=1, phredOffset=64, unique=TRUE,
minFragLength=50, maxFragLength=10000, PE_orientation="fr")
# p2 at Rsubread.pdf: Subjunc perform global alignments. The seed-and-vote
paradigm
# enables efficient and accurate alignments to be carried out."
## phredOffset: sanger 33; illumina 64
## PE_orientation: character string giving the orientation of the two reads from
the same pair. It
# has three possible values including fr, ff and rf. Letter f denotes the
forward
# strand and letter r the reverse strand. fr by default (ie. the first read in
the pair
# is on the forward strand and the second read on the reverse strand).

#####?
# Copy the "Summary" here!
#####?
```

```
52 #####
53 # Counting
54 #####
55 Counts_ToyExample=featureCounts(files="Mapped_reads_ToyExample.BAM",
   annot.ext="../Data/Genome_annotation_ToyExample.gtf", isGTFAnnotationFile=T,
   GTF.featureType="exon", GTF.attrType="gene_id", useMetaFeatures=T,
   isPairedEnd=T, requireBothEndsMapped=T,   checkFragLength=F, nthreads=1,
   strandSpecific=2, reportReads=T)
56 # p12 at Rsubread.pdf: This function assigns mapped sequencing reads to genomic
   features
57 ## GTF.featureType: a character string giving the feature type used to select
   rows in the
58 # GTF annotation which will be used for read summarization. exon by default.
59 ## GTF.attrType: a character string giving the attribute type in the GTF
   annotation which will be
60 # used to group features (eg. exons) into meta-features (eg. genes). gene_id by
61 # default.
62 ## useMetaFeatures: logical indicating whether the read summarization should be
   performed at the
63 # feature level (eg. exons) or meta-feature level (eg genes). If TRUE, features
   in
64 # the annotation (each row is a feature) will be grouped into meta-features
   using
65 # their values using the "gene_id" attribute in the GTF-format annotation file,
   and reads will assiged
66 # to the meta-features instead of the features.
67 ## requireBothEndsMapped: logical indicating if both ends from the same fragment
   are required to be
68 # successfully aligned before the fragment can be assigned to a feature or
   metafeature.
69 ## checkFragLength: logical indicating if the two ends from the same fragment
   are required to satisify
70 # the fragment length criteria before the fragment can be assigned to a feature
   or
71 # meta-feature. The fragment length criteria are specified via minFragLength and
   maxFragLength.
72 ## strandSpecific: integer indicating if strand-specific read counting should be
   performed. It has
73 # three possible values: 0 (unstranded), 1 (stranded) and 2 (reversely
   stranded).
74 ## reportReads: logical indicating if read counting result for each
   read/fragment is saved to a
75 # file. If TRUE, read counting results for reads/fragments will be saved to a
   tab-
76 # delimited file that contains four columns including name of read/fragment,
   sta-
77 # tus(assigned or the reason if not assigned), name of target feature/meta-
   feature
78 # and number of hits if the read/fragment is counted multiple times. Name of the
79 # file is the same as name of the input read file except a suffix
   '.featureCounts' is
80 # added. Multiple files will be generated if there is more than one input read
   file.
81
82 #####?
83 # What pecentage of the reads were counted in total?
84 #####?
85
86 # p16 at Rsubread.pdf: Description of featureCounts variable
87
88 # Names of objects of the featureCounts variable
```

```r
 89 names(Counts_ToyExample)
 90
 91 # Counts of the first genes
 92 head(Counts_ToyExample$counts)
 93
 94 # Histogram of the counts
 95 hist(Counts_ToyExample$counts)
 96 hist(log10(Counts_ToyExample$counts))
 97
 98 # Write the count table to a file
 99 write.table(Counts_ToyExample$counts, "Counts_ToyExample.tsv", quote=F,
    sep="\t")
100
101
102 ##############################################################################
    ###
103 # Differential expression analysis of real data with edgeR R Bioconductor
    package #
104 ##############################################################################
    ###
105 library("edgeR")
106
107 ######
108 # Preprocessing
109 ######
110
111 # Read file with read counts
112 Counts=read.table(file="../Data/Dsim_count_table.tsv", header=T, row.names=1)
113 head(Counts)
114 # Number of genes
115 dim(Counts)
116
117 # Keep those genes that were expressed in at a reasonable level (25 pairs) in
    all samples
118 KeptCounts=Counts[rowSums(Counts>=25)==6, ]
119 # Number of kept genes
120 dim(KeptCounts)
121
122 #########################
123 # Pair wise DE analysis #
124 #########################
125 # an example
126
127 ## "Treatment" groups
128 Group_PW=factor(c("C15","C15","C15","C23","C23","C23"))
129 # Tell R that a variable is nominal by making it a factor. The factor stores the
    nominal values as a vector of integers in the range [ 1... k ] (where k is the
    number of unique values in the nominal variable), and an internal vector of
    character strings (the original values) mapped to these integers.
130
131 ## Differential Expression list
132 DE_list_PW=DGEList(KeptCounts, group=Group_PW)
133 ## DGElist data class
134 # edgeR stores data in a simple list-based data object called a DGEList. This
    type of object is
135 # easy to use because it can be manipulated like any list in R.
136
137 ## TMM --  Trimmed Mean of M-values -- normalization of read counts
138 DE_list_PW=calcNormFactors(DE_list_PW, method=c("TMM"))
139 # The calcNormFactors function normalizes for RNA composition by finding a set
    of scaling
```

```r
140 # factors for the library sizes that minimize the log-fold changes between the
    samples for most
141 # genes. The default method for computing these scale factors uses a trimmed
    mean of M-
142 # values (TMM) between each pair of samples [*]
143 # [*]: Robinson, M.D. and Oshlack, A. (2010). A scaling normalization method for
    differential
144 # expression analysis of RNA-seq data. Genome Biology 11, R25.
145
146 ## Estimating dispersions for the pair wise DE analysis
147 DE_list_PW=estimateCommonDisp(DE_list_PW)
148 DE_list_PW=estimateTrendedDisp(DE_list_PW)
149 DE_list_PW=estimateTagwiseDisp(DE_list_PW)
150 ## Pseudo counts
151 # In general, edgeR functions work directly on the raw counts. For the most
    part, edgeR does
152 # not produce any quantity that could be called a "normalized count".
153 # An exception is the internal use of pseudo-counts by the classic edgeR
    functions estimateCommonDisp
154 # and exactTest. The exact negative binomial test [*] computed by exactTest and
    the con-
155 # ditional likelihood [*] used by estimateCommonDisp and estimateTagwiseDisp
    require the
156 # library sizes to be equal for all samples. These functions therefore compute
    normalized counts
157 # called pseudo-counts by the method of Robinson and Smyth [*]. The pseudo-
    counts are
158 # computed for a specific purpose, and their computation depends on the
    experimental design
159 # as well as the library sizese. Users are therefore disuaded from interpreting
    the psuedo-counts
160 # as general purpose normalized counts.
161 # [*]: Robinson, M.D. and Smyth, G.K. (2008). Small-sample estimation of
    negative binomial
162 # dispersion, with applications to SAGE data. Biostatistics 9, 321–332.
163 ## Average log2 CPM (Counts per million)
164 # log-CPM value for each count:
165 # log2((rgi + 0.5)/Ri+1)×10^6)
166 # rgi: read(pair) count of gene g for sample i
167 # Ri: the total number of mapped read(pair)s for sample i (i.e. the library size
    of sample i)
168
169 ## DE testing with tagwise dispersion
170 DE_list_PW.tgw=exactTest(DE_list_PW, dispersion="tagwise", pair=c("C15","C23"))
171 # Once negative binomial models are fitted and dispersion estimates are
    obtained, we can proceed with testing
172 # procedures for determining differential expression using the exact test.
173 # The exact test is only applicable to experiments with a single factor.
174 # edgeR uses the quantile-adjusted conditional maximum likelihood (qCML) method
    for ex-
175 # periments with single factor.
176 # Compared against several other estimators (e.g. maximum likelihood estimator,
    Quasi-
177 # likelihood estimator etc.) using an extensive simulation study, qCML is the
    most reliable in
178 # terms of bias on a wide range of conditions and specifically performs best in
    the situation
179 # of many small samples with a common dispersion, the model which is applicable
    to Next-
180 # Gen sequencing data.
181 # The qCML method calculates the likelihood by conditioning on the total counts
```

```
181 for each
182 # tag, and uses pseudo counts after adjusting for library sizes.
183
184 ## Which genes were differentially expressed according to the Benjamini-Hochberg
    corrected p-values?
185 Result=DE_list_PW.tgw$table
186 Result$adj.PValue=p.adjust(Result$PValue, method="BH")
187 Up=Result[Result$adj.PValue<0.05 & Result$logFC>0 ,]
188 Down=Result[Result$adj.PValue<0.05 & Result$logFC<0 ,]
189
190 # Write the DE genes and the Result table to files
191 write(rownames(Up), "DE_Up_genes.txt")
192 write(rownames(Down), "DE_Down_genes.txt")
193 write.table(Result, "Result.tsv", sep="\t")
194
```