

Package ‘GREP2’

July 25, 2018

Type Package

Title GEO RNA-Seq Experiments Processing Pipeline

Version 1.0.1

Maintainer Naim Al Mahi <mahina@mail.uc.edu>

Description

An R based pipeline to download and process Gene Expression Omnibus (GEO) RNA-seq data. For a given GEO series accession ID, this pipeline generates metadata, both gene and transcript level counts, and quality control (QC) report. This package is mainly developed to process the GEO RNA-seq datasets for the web platform GREIN. More details about GREP2 and GREIN can be found here <[doi:10.1101/326223](https://doi.org/10.1101/326223)>.

Depends R (>= 3.4.0)

Imports XML, rentrez, RCurl, GEOquery, Biobase, parallel, tximport, EnsDb.Hsapiens.v86, EnsDb.Rnorvegicus.v79, EnsDb.Mmusculus.v79, AnnotationDbi, org.Hs.eg.db, org.Mm.eg.db, org.Rn.eg.db, GenomicFeatures, utils

Suggests knitr, rmarkdown

SystemRequirements SRAtoolkit, Salmon, Java, FastQC, MultiQC

URL <https://github.com/uc-bd2k/GREP2>

BugReports <https://github.com/uc-bd2k/GREP2/issues>

License GPL-3

VignetteBuilder knitr

LazyData true

RoxygenNote 6.0.1

NeedsCompilation no

Author Naim Al Mahi [cre, aut],
Mario Medvedovic [aut, ctb]

Repository CRAN

Date/Publication 2018-07-25 09:20:03 UTC

R topics documented:

build_index	2
get_fastq	3
get_metadata	4
get_srr	4
process_geo_rnaseq	5
run_fastqc	7
run_multiqc	8
run_salmon	9
run_tximport	10
trim_fastq	11
Index	13

build_index	<i>Build index for mapping using Salmon</i>
--------------------	---------------------------------------------

Description

`build_index` for mapping reads using Salmon.

Usage

```
build_index(species = c("human", "mouse", "rat"), kmer = 31,
            ens_release = 92, destdir)
```

Arguments

<code>species</code>	name of the species. Only 'human', 'mouse', and 'rat' are allowed to use.
<code>kmer</code>	k-mer size for indexing. default is 31. See 'Salmon' for details.
<code>ens_release</code>	version of Ensembl release.
<code>destdir</code>	directory where all the files will be saved.

Value

directory of index files

References

Rob Patro, Geet Duggal, Michael I. Love, Rafael A. Irizarry, and Carl Kingsford (2017): Salmon provides fast and bias-aware <https://www.nature.com/articles/nmeth.4197>

Examples

```
#Running this function will take some time.
```

```
build_index(species="human",kmer=31,  
ens_release=92,destdir=tempdir())
```

get_fastq

Download fastq files

Description

get_fastq downloads fastq files using SRA toolkit. We recommend using Aspera for fast downloading. You need to install Aspera for using ascp option.

Usage

```
get_fastq(srr_id, library_layout = c("SINGLE", "PAIRED"),  
use_sra_file = FALSE, sra_files_dir = NULL, n_thread, destdir)
```

Arguments

srr_id SRA run accession ID.
library_layout layout of the library used. Either 'SINGLE' or 'PAIRED'.
use_sra_file logical, whether to use downloaded SRA files to get fastq files or directly download fastq files.
sra_files_dir directory where SRA files are saved. If you use use_sra_file=FALSE then sra_files_dir=NULL.
n_thread number of cores to use.
destdir directory where all the results will be saved.

Value

A single fastq file will be generated for SINGLE end reads and two files for PAIRED end reads.

Examples

```
get_fastq(srr_id="SRR5890521",library_layout="SINGLE",  
use_sra_file=FALSE,sra_files_dir=NULL,n_thread=2,  
destdir=tempdir())
```

<code>get_metadata</code>	<i>Download metadata from GEO and SRA</i>
---------------------------	-------------------------------------------

Description

Download metadata from GEO and SRA

Usage

```
get_metadata(geo_series_acc, destdir, geo_only = FALSE,
            download_method = "auto")
```

Arguments

<code>geo_series_acc</code>	GEO series accession ID.
<code>destdir</code>	directory where the metadata files will be saved.
<code>geo_only</code>	logical, whether to download GEO metadata only. Default is FALSE. If TRUE, then SRA metadata will not be downloaded.
<code>download_method</code>	download method for GEOquery. See 'download.file' from R package utils for details. Default is 'libcurl'.

Value

a list of GEO and SRA metadata.

Examples

```
get_metadata(geo_series_acc="GSE102170", destdir=tempdir(),
            geo_only=TRUE, download_method="auto")
```

<code>get_srr</code>	<i>Download SRA run files</i>
----------------------	-------------------------------

Description

`get_srr` downloads SRA files using Aspera or FTP. We recommend using Aspera for fast downloading. You need to install Aspera for using `ascp` option.

Usage

```
get_srr(srr_id, destdir, ascp, prefetch_workspace, ascp_path)
```

Arguments

srr_id	SRA run accession ID.
destdir	directory where all the results will be saved.
ascp	logical, whether to use Aspera for downloading SRA files.
prefetch_workspace	directory where SRA run files will be downloaded. This parameter is needed if ascp=TRUE. The location of this directory can be found by going to the aspera directory (/aspera/connect/bin/) and typing 'vdb-config -i'. A new window will pop-up and under the 'Workspace Name', you will find the location. Usually the default is '/home/username/ncbi/public'.
ascp_path	path to the Aspera software.

Value

SRA run accession file with extension ".sra". If you use ascp=TRUE, then downloaded files will be saved under '/prefetch_workspace/sra' directory. If ascp=FALSE, then files will be saved in the 'destdir'

Examples

```
get_srr(srr_id="SRR5890521", destdir=tempdir(), ascp=FALSE,  
prefetch_workspace=NULL, ascp_path=NULL)
```

process_geo_rnaseq *A complete pipeline to process GEO RNA-seq data*

Description

process_geo_rnaseq downloads and processes GEO RNA-seq data for a given GEO series accession ID. It filters metadata for RNA-seq samples only. We use SRA toolkit for downloading SRA data, Trimmomatic for read trimming (optional), and Salmon for read mapping.

Usage

```
process_geo_rnaseq(geo_series_acc, destdir, download_method = "auto",  
ascp = TRUE, prefetch_workspace, ascp_path, use_sra_file = FALSE,  
trim_fastq = FALSE, index_dir, other_opts = NULL, species = c("human",  
"mouse", "rat"), countsFromAbundance = c("no", "scaledTPM",  
"lengthScaledTPM"), n_thread)
```

Arguments

geo_series_acc	GEO series accession ID.
destdir	directory where all the results will be saved.
download_method	download method for GEOquery.
ascp	logical, whether to use Aspera connect to download SRA run files. If FALSE, then wget will be used to download files which might be slower than 'ascp' download.
prefetch_workspace	directory where SRA run files will be downloaded. This parameter is needed when ascp=TRUE. The location of this directory can be found by going to the aspera directory (<i>/.aspera/connect/bin/</i>) and typing 'vdb-config -i'. A new window will pop-up and under the 'Workspace Name', you will find the location. Usually the default is '/home/username/ncbi/public'.
ascp_path	path to the Aspera software.
use_sra_file	logical, whether to download SRA file first and get fastq files afterwards.
trim_fastq	logical, whether to trim fastq file.
index_dir	directory of the indexing files needed for read mapping using Salmon. See function 'build_index'.
other_opts	options other than default to use for read mapping. See Salmon documentation for the available options.
species	name of the species. Only 'human', 'mouse', and 'rat' are allowed to use.
countsFromAbundance	whether to generate counts based on abundance. Available options are: 'no', 'scaledTPM' (abundance based estimated counts scaled up to library size), 'lengthScaledTPM' (default, scaled using the average transcript length over samples and library size). See Bioconductor package tximport for further details.
n_thread	number of cores to use.

Value

a list of metadata from GEO and SRA saved in the destdir. Another list of gene and transcript level estimated counts summarized by Bioconductor package '[tximport](#)' is also saved in the destdir.

References

- Rob Patro, Geet Duggal, Michael I. Love, Rafael A. Irizarry, and Carl Kingsford (2017): Salmon provides fast and bias-aware quantification of transcript expression. *Nature methods*, 14(4), 417. <https://www.nature.com/articles/nmeth.4197>
- Charlotte Soneson, Michael I. Love, Mark D. Robinson (2015): Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Research*. <http://dx.doi.org/10.12688/f1000research.7563.1>
- Philip Ewels, Mans Magnusson, Sverker Lundin, and Max Kaller (2016): MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*, 32(19), 3047-3048. <https://doi.org/10.1093/bioinformatics/btw354>

Examples

```
geo_series_acc="GSE102170"
#You will have to build index first before running this function.

build_index(species="human",kmer=31,ens_release=92,
destdir=tempdir())
process_geo_rnaseq (geo_series_acc=geo_series_acc,destdir=tempdir(),
download_method="auto",
ascp=FALSE,prefetch_workspace=NULL,
ascp_path=NULL,use_sra_file=FALSE,trim_fastq=FALSE,
index_dir=tempdir(),species="human",
countsFromAbundance="lengthScaledTPM",n_thread=1)
```

run_fastqc

QC report for each fastq files using FastQC

Description

run_fastqc HTML report of each fastq files using FastQC. You need to install FastQC from <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Usage

```
run_fastqc(destdir, fastq_dir, n_thread)
```

Arguments

destdir	directory where all the results will be saved.
fastq_dir	directory of the fastq files.
n_thread	number of cores to use.

Value

HTML report of the fastq files under fastqc directory.

Examples

```
fastq_dir=system.file("extdata","", package="GREP2")

run_fastqc(destdir=tempdir(),fastq_dir=fastq_dir,
n_thread=2)
```

run_multiqc

Generate combined QC report for Salmon and FastQC

Description

`run_fastqc` generates a single HTML report from the fastQC reports and salmon read mapping results using MultiQC.

Usage

```
run_multiqc(fastqc_dir, salmon_dir, destdir)
```

Arguments

<code>fastqc_dir</code>	directory where all the FastQC files are saved.
<code>salmon_dir</code>	directory of the salmon files.
<code>destdir</code>	directory where you want to save the combined QC report.

Value

HTML report.

References

Philip Ewels, Mans Magnusson, Sverker Lundin, and Max Kaller (2016): MultiQC: summarize analysis results for multiple tools and samples in a single report. Bioinformatics, 32(19), 3047-3048. <https://doi.org/10.1093/bioinformatics/btw354>

Examples

```
## Not run:  
run_multiqc(fastqc_dir=tempdir(), salmon_dir=tempdir(),  
destdir=tempdir())  
  
## End(Not run)
```

run_salmon*Quantify transcript abundances using Salmon*

Description

`run_salmon` is a wrapper function for mapping reads to quantify transcript abundances using Salmon. You need to install Salmon and build index to run this function. For index building see function '`build_index`'.

Usage

```
run_salmon(srr_id, library_layout = c("SINGLE", "PAIRED"), index_dir, destdir,
           fastq_dir, use_trimmed_fastq = FALSE, other_opts = NULL, n_thread)
```

Arguments

<code>srr_id</code>	SRA run accession ID.
<code>library_layout</code>	layout of the library used. Either 'SINGLE' or 'PAIRED'.
<code>index_dir</code>	directory of the indexing files needed for read mapping using Salmon. See function ' <code>build_index</code> '.
<code>destdir</code>	directory where all the results will be saved.
<code>fastq_dir</code>	directory of the fastq files.
<code>use_trimmed_fastq</code>	logical, whether to use trimmed fastq files.
<code>other_opts</code>	other options to use. See Salmon documentation for the available options.
<code>n_thread</code>	number of cores to use.

Details

`run_salmon` We use default options of Salmon. This function works for a single sample. You can use this function in a loop for multiple samples. For other options from Salmon use '`other_opts`'.

Value

The following items will be returned and saved in the salmon directory:

1. `quant_new.sf`: plain-text, tab-separated quantification file that contains 5 column: Name,Length,EffectiveLength,TPM, and NumReads.
2. `cmd_info.json`: A JSON format file that records the main command line parameters with which Salmon was invoked for the run that produced the output in this directory.
3. `aux_info`: This directory will have a number of files (and subfolders) depending on how salmon was invoked.
4. `meta_info.json`: A JSON file that contains meta information about the run, including stats such as the number of observed and mapped fragments, details of the bias modeling etc.

5. ambig_info.tsv: This file contains information about the number of uniquely-mapping reads as well as the total number of ambiguously-mapping reads for each transcript.
6. lib_format_counts.json: This JSON file reports the number of fragments that had at least one mapping compatible with the designated library format, as well as the number that didn't.
7. libParams: The auxiliary directory will contain a text file called flenDist.txt. This file contains an approximation of the observed fragment length distribution.

References

Rob Patro, Geet Duggal, Michael I. Love, Rafael A. Irizarry, and Carl Kingsford (2017): Salmon provides fast and bias-aware quantification of transcript expression. *Nature methods*, 14(4), 417. <https://www.nature.com/articles/nmeth.4197>

Examples

```
#You will have to build index first to run this function
fastq_dir=system.file("extdata","", package="GREP2")

build_index(species="human",kmer=31,ens_release=92,
destdir=tempdir())
run_salmon(srr_id="SRR5890521",library_layout="SINGLE",
index_dir=tempdir(),destdir=tempdir(),
fastq_dir=fastq_dir,use_trimmed_fastq=FALSE,
other_opts=NULL,n_thread=2)
```

run_tximport

Wrapper function to run tximport

Description

`run_tximport` function runs `tximport` on transcript level abundances from `Salmon` to summarize to gene level. See Bioconductor package [tximport](#) for details.

Usage

```
run_tximport(srr_id, species = c("human", "mouse", "rat"), salmon_dir,
countsFromAbundance = c("no", "scaledTPM", "lengthScaledTPM"))
```

Arguments

<code>srr_id</code>	SRA run accession ID.
<code>species</code>	name of the species. Only 'human', 'mouse', and 'rat' are allowed to use.
<code>salmon_dir</code>	directory where salmon files are saved. This should be the folder created by <code>Salmon</code> and is called "salmon".

```
countsFromAbundance
    whether to generate counts based on abundance. Available options are: 'no',
    'scaledTPM' (abundance based estimated counts scaled up to library size),
    'lengthScaledTPM' (default, scaled using the average transcript length over
    samples and library size).
```

Details

We use Ensembl annotation for both genes and transcripts.

Value

a list of gene and transcript level estimated counts.

References

Charlotte Soneson, Michael I. Love, Mark D. Robinson (2015): Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. F1000Research. <http://dx.doi.org/10.12688/f1000research.7563.1>

Examples

```
#Run this function after running Salmon.
salmon_dir=system.file("extdata","", package="GREP2")
run_tximport(srr_id="SRR5890521", species="human",
salmon_dir=salmon_dir,countsFromAbundance="lengthScaledTPM")
```

trim_fastq

Trim fastq files using Trimmomatic

Description

trim_fastq trim fastq files based on the illumina instruments using Trimmomatic.

Usage

```
trim_fastq(srr_id, fastq_dir, instrument, library_layout = c("SINGLE",
    "PAIRED"), destdir, n_thread)
```

Arguments

srr_id	SRA run accession ID.
fastq_dir	directory of the fastq files.
instrument	name of the illumina sequencing platform. For example, 'HiSeq'.
library_layout	layout of the library used. Either 'SINGLE' or 'PAIRED'.
destdir	directory where the trimmed fastq files will be saved.
n_thread	number of cores.

Details

The following parameters are used as default in the trimmoatic function:

1. Remove leading low quality or N bases (below quality 3) (LEADING:3)
2. Remove trailing low quality or N bases (below quality 3) (TRAILING:3)
3. Scan the read with a 4-base wide sliding window, cutting when the average quality per base drops below 15 (SLIDINGWINDOW:4:15)
4. Drop reads below the 36 bases long (MINLEN:36)

Value

trimmed fastq files.

References

Anthony M. Bolger, Marc Lohse, and Bjoern Usadel (2014): Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics, 30(15), 2114-2120. <https://doi.org/10.1093/bioinformatics/btu170>

Examples

```
fastq_dir=system.file("extdata","", package="GREP2")
trim_fastq(srr_id="SRR5890521",fastq_dir=fastq_dir,
instrument="MiSeq",library_layout="SINGLE",
destdir=tempdir(),n_thread=2)
```

Index

build_index, 2
get_fastq, 3
get_metadata, 4
get_srr, 4
process_geo_rnaseq, 5
run_fastqc, 7
run_multiqc, 8
run_salmon, 9
run_tximport, 10
trim_fastq, 11
tximport, 6, 10
utils, 4