

Package ‘ATACseqQC’

October 17, 2017

Type Package

Title ATAC-seq Quality Control

Version 1.0.5

Author Jianhong Ou, Jun Yu, Michelle Kelliher, Lucio Castilla, Nathan Lawson, Lihua Julie Zhu

Maintainer Jianhong Ou <jianhong.ou@umassmed.edu>

Description ATAC-seq, an assay for Transposase-Accessible Chromatin using sequencing, is a rapid and sensitive method for chromatin accessibility analysis. It was developed as an alternative method to MNase-seq, FAIRE-seq and DNase-seq. Comparing to the other methods, ATAC-seq requires less amount of the biological samples and time to process. In the process of analyzing several ATAC-seq dataset produced in our labs, we learned some of the unique aspects of the quality assessment for ATAC-seq data. To help users to quickly assess whether their ATAC-seq experiment is successful, we developed ATACseqQC package partially following the guideline published in Nature Method 2013 (Greenleaf et al.), including diagnostic plot of fragment size distribution, proportion of mitochondria reads, nucleosome positioning pattern, and CTCF or other Transcript Factor footprints.

Depends R (>= 3.4), BiocGenerics, S4Vectors

Imports BSgenome, Biostrings, ChIPpeakAnno, IRanges, GenomicRanges, GenomicAlignments, GenomeInfoDb, GenomicScores, graphics, grid, limma, Rsamtools, randomForest, rtracklayer, stats, stringr

Suggests RUnit, BiocStyle, knitr, BSgenome.Hsapiens.UCSC.hg19, TxDb.Hsapiens.UCSC.hg19.knownGene, phastCons100way.UCSC.hg19, motifStack, MotifDb, trackViewer

License GPL (>= 2)

LazyData TRUE

VignetteBuilder knitr

RoxygenNote 6.0.1

biocViews Sequencing, DNASeq, GeneRegulation, QualityControl, Coverage, NucleosomePositioning

NeedsCompilation no

R topics documented:

ATACseqQC-package	2
enrichedFragments	2
factorFootprints	4
fragSizeDist	5
plotFootprints	6
pwmsscores	7
readBamFile	7
shiftGAlignmentsList	8
shiftReads	9
splitBam	9
splitGAlignmentsByCut	11
writeListOfGAlignments	12

Index	14
--------------	-----------

ATACseqQC-package	<i>ATAC-seq Quality Control</i>
-------------------	---------------------------------

Description

ATAC-seq, an assay for Transposase-Accessible Chromatin using sequencing, is a rapid and sensitive method for chromatin accessibility analysis. It was developed as an alternative method to MNase-seq, FAIRE-seq and DNase-seq. Comparing to the other methods, ATAC-seq requires less amount of the biological samples and time to process. In the process of analyzing several ATAC-seq dataset produced in our labs, we learned some of the unique aspects of the quality assessment for ATAC-seq data. To help users to quickly assess whether their ATAC-seq experiment is successful, we developed ATACseqQC package partially following the guideline published in Nature Method 2013 (Greenleaf et al.), including diagnostic plot of fragment size distribution, proportion of mitochondria reads, nucleosome positioning pattern, and CTCF or other Transcript Factor footprints.

enrichedFragments	<i>enrichment for nucleosome-free fragments and nucleosome signals</i>
-------------------	--

Description

Get the enrichment signals for nucleosome-free fragments and nucleosomes.

Usage

```
enrichedFragments(bamfiles, index = bamfiles, TSS, librarySize,
  upstream = 1010L, downstream = 1010L, n.tile = 101L,
  normal.method = "quantile", adjustFragmentLength = 80L,
  TSS.filter = 0.5, seqlev = paste0("chr", c(1:22, "X", "Y")))
```

Arguments

bamfiles	A vector of characters indicates the file names of bams.
index	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
TSS	an object of GRanges indicates the transcript start sites. All the width of TSS should equal to 1. Otherwise, TSS will be reset to the center of input TSS.
librarySize upstream, downstream	A vector of numeric indicates the library size. Output of estLibSize numeric(1) or integer(1). Upstream and downstream size from each TSS.
n.tile	numeric(1) or integer(1). The number of tiles to generate for each element of TSS.
normal.method	character(1). Normalization methods, could be "none" or "quantile". See normalizeBetweenArrays .
adjustFragmentLength	numeric(1) or integer(1). The size of fragment to be adjusted to. Default is set to half of the nucleosome size (80)
TSS.filter	numeric(1). The filter for signal strength of each TSS. Default 0.5 indicates the average signal strength for the TSS from upstream to downstream bins should be greater than 0.5.
seqlev	A vector of character indicates the sequence names to be considered.

Value

A list of matrixes. In each matrix, each row record the signals for corresponding feature.

Author(s)

Jianhong Ou

Examples

factorFootprints *plot ATAC-seq footprints infer factor occupancy genome wide*

Description

Aggregate ATAC-seq footprint for a given motif generated over binding sites within the genome.

Usage

```
factorFootprints(bamfiles, index = bamfiles, pfm, genome,
  min.score = "95%", bindingSites, seqlev = paste0("chr", c(1:22, "X",
  "Y")), upstream = 100, downstream = 100)
```

Arguments

<code>bamfiles</code>	A vector of characters indicates the file names of bams.
<code>index</code>	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
<code>pfm</code>	A Position frequency Matrix represented as a numeric matrix with row names A, C, G and T.
<code>genome</code>	An object of BSgenome .
<code>min.score</code>	The minimum score for counting a match. Can be given as a character string containing a percentage (e.g. "95 score or as a single number. See matchPWM .
<code>bindingSites</code>	A object of GRanges indicates candidate binding sites (eg. the output of fimo).
<code>seqlev</code>	A vector of characters indicates the sequence levels.
<code>upstream, downstream</code>	numeric(1) or integer(1). Upstream and downstream of the binding region for aggregate ATAC-seq footprint.

Value

an invisible list of matrixes with the signals for plot.

Author(s)

Jianhong Ou

References

Chen, K., Xi, Y., Pan, X., Li, Z., Kaestner, K., Tyler, J., Dent, S., He, X. and Li, W., 2013. DANPOS: dynamic analysis of nucleosome position and occupancy by sequencing. *Genome research*, 23(2), pp.341-351.

Examples

```

shiftedBamfile <- system.file("extdata", "GL1.bam",
                             package="ATACseqQC")
library(MotifDb)
CTCF <- query(MotifDb, c("CTCF"))
CTCF <- as.list(CTCF)
library(BSgenome.Hsapiens.UCSC.hg19)
factorFootprints(shiftedBamfile, pfm=CTCF[[1]],
                  genome=Hsapiens,
                  min.score="95%", seqlev="chr1",
                  upstream=100, downstream=100)

```

fragSizeDist *fragment size distribution*

Description

estimate the fragment size of bams

Usage

```
fragSizeDist(bamFiles, bamFiles.labels, ylim = NULL, logYlim = NULL)
```

Arguments

bamFiles	A vector of characters indicates the file names of bams.
bamFiles.labels	labels of the bam files, used for pdf file naming.
ylim	numeric(2). ylim of the histogram.
logYlim	numeric(2). ylim of log-transformed histogram for the insert.

Value

Invisible fragment length distribution list.

Author(s)

Jianhong Ou

Examples

```

bamFiles <- system.file("extdata", "GL1.bam", package="ATACseqQC")
bamFiles.labels <- "GL1"
fragSizeDist(bamFiles, bamFiles.labels,
            ylim=c(0, 1e4), logYlim=log10(c(5e-3, 2)))

```

`plotFootprints` *Plots a footprint estimated by Centipede*

Description

Visualizing the footprint profile

Usage

```
plotFootprints(Profile, Mlen = 0, xlab = "Dist. to motif (bp)",  
    ylab = "Cut-site probability", legTitle, newpage = TRUE, motif)
```

Arguments

Profile	A vector with the profile estimated by CENTIPEDE
Mlen	Length of the motif for drawing vertical lines delimiting it
xlab	Label of the x axis
ylab	Label for the y axis
legTitle	Title for one of the plot corners
newpage	Plot the figure in a new page?
motif	a pfm object.

Value

Null.

Author(s)

Jianhong Ou

Examples

pwmscores	<i>max PWM scores for sequences</i>
-----------	-------------------------------------

Description

calculate the maximal PWM scores for each given sequences

Usage

```
pwmscores(pwm, subject)
```

Arguments

pwm	A Position Weight Matrix represented as a numeric matrix with row names A, C, G and T.
subject	Typically a DNAString object. A Views object on a DNAString subject, a MaskedDNAString object, or a single character string, are also supported. IUPAC ambiguity letters in subject are ignored (i.e. assigned weight 0) with a warning.

Value

a numeric vector

Author(s)

Jianhong

readBamFile	<i>read in bam files</i>
-------------	--------------------------

Description

wrapper for readGAlignments/readGAlignmentsList to read in bam files.

Usage

```
readBamFile(bamFile, which, tag = character(0), what = c("qname", "flag",
  "mapq", "isize", "seq", "qual", "mrnm"),
  flag = scanBamFlag(isSecondaryAlignment = FALSE, isUnmappedQuery = FALSE,
  isNotPassingQualityControls = FALSE, asMates = FALSE, ...)
```

Arguments

bamFile	character(1). Bam file name.
which	A GRanges , RangesList , or any object that can be coerced to a RangesList , or missing object, from which a IRangesList instance will be constructed. See ScanBamParam .
tag	A vector of characters indicates the tag names to be read. See ScanBamParam .

<code>what</code>	A character vector naming the fields to return. Fields are described on the Rsamtools [<code>scanBam</code>] help page.
<code>flag</code>	An integer(2) vector used to filter reads based on their 'flag' entry. This is most easily created with the Rsamtools [<code>scanBamFlag</code>] helper function.
<code>asMates</code>	logical(1). Paired ends or not
<code>...</code>	parameters used by readGAlignmentsList or readGAlignments

Value

A GAlignmentsList object when `asMats=TRUE`, otherwise A GAlignments object.

Author(s)

Jianhong Ou

Examples

```
library(BSgenome.Hsapiens.UCSC.hg19)
which <- as(seqinfo(Hsapiens)["chr1"], "GRanges")
bamfile <- system.file("extdata", "GL1.bam",
                       package="ATACseqQC", mustWork=TRUE)
readBamFile(bamfile, which=which, asMates=TRUE)
```

`shiftGAlignmentsList` *shift 5' ends*

Description

shift the GAlignmentsLists by 5' ends. All reads aligning to the positive strand will be offset by +4bp, and all reads aligning to the negative strand will be offset -5bp by default.

Usage

```
shiftGAlignmentsList(gal, positive = 4L, negative = 5L)
```

Arguments

<code>gal</code>	An object of GAlignmentsList .
<code>positive</code>	integer(1). the size to be shift for positive strand
<code>negative</code>	integer(1). the size to be shift for negative strand

Value

An object of [GAlignments](#) with 5' end shifted reads.

Author(s)

Jianhong Ou

Examples

```
bamfile <- system.file("extdata", "GL1.bam", package="ATACseqQC")
tags <- c("AS", "XN", "XM", "XO", "XG", "NM", "MD", "YS", "YT")
library(BSgenome.Hsapiens.UCSC.hg19)
which <- as(seqinfo(Hsapiens)[["chr1"]], "GRanges")
gal <- readBamFile(bamfile, tag=tags, which=which, asMates=TRUE)
objs <- shiftGAlignmentsList(gal)
export(objs, "shift.bam")
```

shiftReads

*shift read for 5'end***Description**

shift reads for 5'ends

Usage

shiftReads(x, positive = 4L, negative = 5L)

Arguments

x	an object of GAlignments
positive	integer(1). the size to be shift for positive strand
negative	integer(1). the size to be shift for negative strand

Value

an object of GAlignments

Author(s)

Jianhong Ou

splitBam

*prepare bam files for downstream analysis***Description**

shift the bam files by 5'ends and split the bam files.

Usage

```
splitBam(bamfile, tags, outPath = NULL, txs, genome, conservation,
positive = 4L, negative = 5L, breaks = c(0, 100, 180, 247, 315, 473,
558, 615, Inf), labels = c("NucleosomeFree", "inter1", "mononucleosome",
"inter2", "dinucleosome", "inter3", "trinucleosome", "others"),
seqlev = paste0("chr", c(1:22, "X", "Y")), cutoff = 0.8)
```

Arguments

<code>bamfile</code>	character(1). File name of bam.
<code>tags</code>	A vector of characters indicates the tags in bam file.
<code>outPath</code>	Output file path.
<code>txs</code>	GRanges of transcripts.
<code>genome</code>	An object of BSgenome
<code>conservation</code>	An object of GScores .
<code>positive</code>	integer(1). the size to be shift for positive strand
<code>negative</code>	integer(1). the size to be shift for negative strand
<code>breaks</code>	A numeric vector for fragment size of nucleosome freee, mononucleosome, dinucleosome and trinucleosome
<code>labels</code>	A vector of characters indicates the labels for the levels of the resulting category. The length of labels = length of breaks - 1
<code>seqlev</code>	A vector of characters indicates the sequence levels.
<code>cutoff</code>	numeric(1). Cutoff value for prediction by randomForest .

Value

an invisible list of [GAlignments](#)

Author(s)

Jianhong Ou

See Also

[shiftGAlignmentsList](#), [splitGAlignmentsByCut](#), and [writeListOfGAlignments](#)

Examples

```
bamfile <- system.file("extdata", "GL1.bam", package="ATACseqQC")
tags <- c("AS", "XN", "XM", "XO", "XG", "NM", "MD", "YS", "YT")
library(BSgenome.Hsapiens.UCSC.hg19)
library(TxDb.Hsapiens.UCSC.hg19.knownGene)
txs <- transcripts(TxDb.Hsapiens.UCSC.hg19.knownGene)
library(phastCons100way.UCSC.hg19)
objs <- splitBam(bamfile, tags,
                  txs=txs, genome=Hsapiens,
                  conservation=phastCons100way.UCSC.hg19,
                  seqlev="chr1")
```

`splitGAlignmentsByCut` *split bams into nucleosome free, mononucleosome, dinucleosome and trinucleosome*

Description

use random forest to split the reads into nucleosome free, mononucleosome, dinucleosome and trinucleosome. The features used in random forest including fragment length, GC content, and UCSC phastCons conservation scores.

Usage

```
splitGAlignmentsByCut(obj, txs, genome, conservation, breaks = c(0, 100, 180,
247, 315, 473, 558, 615, Inf), labels = c("NucleosomeFree", "inter1",
"mononucleosome", "inter2", "dinucleosome", "inter3", "trinucleosome",
"others"), labelsOfNucleosomeFree = "NucleosomeFree",
labelsOfMononucleosome = "mononucleosome", trainningSetPercentage = 0.15,
cutoff = 0.8, halfSizeOfNucleosome = 80L)
```

Arguments

<code>obj</code>	an object of GAlignments
<code>txs</code>	GRanges of transcripts
<code>genome</code>	an object of BSgenome
<code>conservation</code>	an object of GScores .
<code>breaks</code>	a numeric vector for fragment size of nucleosome freee, mononucleosome, dinucleosome and trinucleosome. The breaks pre-defined here is following the description of Greenleaf's paper (see reference).
<code>labels</code>	a character vector for labels of the levels of the resulting category.
<code>labelsOfNucleosomeFree</code> , <code>labelsOfMononucleosome</code>	character(1). The label for nucleosome free and mononucleosome.
<code>trainningSetPercentage</code>	numeric(1) between 0 and 1. Percentage of trainning set from top coverage.
<code>cutoff</code>	numeric(1) between 0 and 1. cutoff value for prediction.
<code>halfSizeOfNucleosome</code>	numeric(1) or integer(1). Thre read length will be adjusted to half of the nucleosome size to enhance the signal-to-noise ratio.

Value

a list of GAlignments

Author(s)

Jianhong Ou

References

- Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. and Greenleaf, W.J., 2013. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nature methods*, 10(12), pp.1213-1218.
- Chen, K., Xi, Y., Pan, X., Li, Z., Kaestner, K., Tyler, J., Dent, S., He, X. and Li, W., 2013. DANPOS: dynamic analysis of nucleosome position and occupancy by sequencing. *Genome research*, 23(2), pp.341-351.

Examples

```
library(GenomicRanges)
bamfile <- system.file("extdata", "GL1.bam",
                      package="ATACseqQC", mustWork=TRUE)
tags <- c("AS", "XN", "XM", "XO", "XG", "NM", "MD", "YS", "YT")
gal1 <- readBamFile(bamFile=bamfile, tag=tags,
                     which=GRanges("chr1", IRanges(1, 1e6)),
                     asMates=FALSE)
names(gal1) <- mcols(gal1)$qname
library(BSgenome.Hsapiens.UCSC.hg19)
library(TxDb.Hsapiens.UCSC.hg19.knownGene)
txs <- transcripts(TxDb.Hsapiens.UCSC.hg19.knownGene)
library(phastCons100way.UCSC.hg19)
splitGAlignmentsByCut(gal1, txs=txs, genome=Hsapiens,
                      conservation=phastCons100way.UCSC.hg19)
```

writeListOfGAlignments

export list of GAlignments into bam files

Description

wrapper for [export](#) to export list of GAlignment into bam files.

Usage

```
writeListOfGAlignments(objs, outPath = ".")
```

Arguments

objs	A list of GAlignments .
outPath	character(1). Output file path.

Value

status of export.

Author(s)

Jianhong Ou

Examples

```
library(GenomicAlignments)
gal1 <- GAlignments(seqnames=Rle("chr1"), pos=1L, cigar="10M",
                     strand=Rle(strand(c("+"))), names="a", score=1)
galist <- GAlignmentsList(a=gal1)
writeListOfGAlignments(galist)
```

Index

ATACseqQC (ATACseqQC-package), 2
ATACseqQC-package, 2

BSgenome, 4, 10

DNAString, 7

enrichedFragments, 2
estLibSize, 3
export, 12

factorFootprints, 4
fragSizeDist, 5

GAlignments, 8, 10–12
GAlignmentsList, 8
GRanges, 3, 4, 7, 10
GScores, 10, 11

MaskedDNAString, 7
matchPWM, 4

normalizeBetweenArrays, 3

plotFootprints, 6
pwmscores, 7

randomForest, 10
RangesList, 7
readBamFile, 7
readGAlignments, 8
readGAlignmentsList, 8
Rsamtools, 8

ScanBamParam, 7
shiftGAlignmentsList, 8, 10
shiftReads, 9
splitBam, 9
splitGAlignmentsByCut, 10, 11

Views, 7

writeListOfGAlignments, 10, 12