

Package ‘ATACseqQC’

October 15, 2018

Type Package

Title ATAC-seq Quality Control

Version 1.4.3

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

Maintainer Jianhong Ou <jianhong.ou@duke.com>

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 (>= 1.31.2), randomForest, rtracklayer, stats, motifStack, preseqR, utils, KernSmooth

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

License GPL (>= 2)

LazyData TRUE

VignetteBuilder knitr

RoxygenNote 6.1.0

biocViews Sequencing, DNaseSeq, ATACSeq, GeneRegulation, QualityControl, Coverage, NucleosomePositioning

git_url <https://git.bioconductor.org/packages/ATACseqQC>

git_branch RELEASE_3_7
git_last_commit 29b48ac
git_last_commit_date 2018-09-27
Date/Publication 2018-10-15

R topics documented:

ATACseqQC-package	2
bamQC	3
distanceDyad	3
enrichedFragments	4
estimateLibComplexity	6
factorFootprints	7
footprintsScanner	8
fragSizeDist	9
NFRscore	10
peakdet	11
plotCorrelation	11
plotFootprints	12
prepareBindingSitesList	13
PTscore	14
pwmscores	15
readBamFile	15
readsDupFreq	16
saturationPlot	17
shiftGAlignmentsList	18
shiftReads	19
splitBam	19
splitGAlignmentsByCut	21
vPlot	22
writeListOfGAlignments	23
Index	25

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.

bamQC	<i>Mapping quality control</i>
-------	--------------------------------

Description

Check the mapping rate, PCR duplication rate, and mitochondria reads contamination.

Usage

```
bamQC(bamfile, index = bamfile, mitochondria = "chrM",
      outPath = sub(".bam", ".clean.bam", basename(bamfile)),
      doubleCheckDup = FALSE)
```

Arguments

bamfile	character(1). File name of bam.
index	character(1). File name of index file.
mitochondria	character(1). Sequence name of mitochondria.
outPath	character(1). File name of cleaned bam.
doubleCheckDup	logical(1). Double check duplicates or not if there is no tags for that.

Value

A list of quality summary.

Author(s)

Jianhong Ou

Examples

```
bamfile <- system.file("extdata", "GL1.bam", package="ATACseqQC")
bamQC(bamfile, outPath=NULL)
```

distanceDyad	<i>Distance of potential nucleosome dyad</i>
--------------	--

Description

Calculate the distance of potential nucleosome dyad and the linear model for V.

Usage

```
distanceDyad(vPlotOut, fragLenRanges = c(60, 180, 250), draw = TRUE,
            ...)
```

Arguments

vPlotOut	The output of vPlot .
fragLenRanges	A numeric vector (length=3) for fragment size of nucleosome free and mono-nucleosome. Default c(60, 180, 250).
draw	Plot the results or not. Default TRUE.
...	Parameters could be passed to plot.

Value

an invisible list with distance of nucleosome and the linear model.

Author(s)

Jianhong Ou

See Also

[vPlot](#)

Examples

```
bamfile <- system.file("extdata", "GL1.bam",
                      package="ATACseqQC")
library(MotifDb)
CTCF <- query(MotifDb, c("CTCF"))
CTCF <- as.list(CTCF)
library(BSgenome.Hsapiens.UCSC.hg19)
vp <- vPlot(bamfile, pfm=CTCF[[1]],
            genome=Hsapiens,
            min.score="95%", seqlev="chr1",
            draw=FALSE)
distanceDyad(vp)
```

enrichedFragments *enrichment for nucleosome-free fragments and nucleosome signals*

Description

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

Usage

```
enrichedFragments(bamfiles, index = bamfiles, gal, 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")))
```

estimateLibComplexity *Library complexity estimation*

Description

Estimating the library complexity.

Usage

```
estimateLibComplexity(histFile, times = 100,  
  interpolate.sample.sizes = seq(0.1, 1, by = 0.1),  
  extrapolate.sample.sizes = seq(5, 20, by = 5))
```

Arguments

histFile	A two-column matrix of integers. The 1st column is the frequency $j = 1, 2, 3, \dots$. The 2nd column is the number of genomic regions with the same frequency (j) of duplication. This file should be sorted by the first column in ascending order. For example, one row of a histogram file: 10 20 means there are 10 genomic regions, each of which is covered by 20 identical fragments at a given sequencing depth of a sequencing library.
times	An positive integer representing the minimum required number of successful estimation. Default is 100.
interpolate.sample.sizes	A numeric vector with values between (0, 1].
extrapolate.sample.sizes	A numeric vector with values greater than 1.

Value

invisible estimates, a data frame of 3 columns: relative sequence depth, number of distinct fragments, number of putative sequenced reads.

Author(s)

Haibo Liu, Feng Yan

See Also

[readsDupFreq](#)

Examples

```
library(preseqR)  
data(FisherButterfly)  
estimateLibComplexity(histFile=FisherButterfly, times=100)
```

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, maxSiteNum = 1e+06,
  anchor = "cut site")
```

Arguments

bamfiles	A vector of characters indicates the file names of bams. All the bamfiles will be pulled together.
index	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
pfm	A Position frequency Matrix represented as a numeric matrix with row names A, C, G and T.
genome	An object of BSgenome .
min.score	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 .
bindingSites	A object of GRanges indicates candidate binding sites (eg. the output of fimo).
seqlev	A vector of characters indicates the sequence levels.
upstream, downstream	numeric(1) or integer(1). Upstream and downstream of the binding region for aggregate ATAC-seq footprint.
maxSiteNum	numeric(1). Maximal number of predicted binding sites. if predicted binding sites is more than this number, top maxSiteNum binding sites will be used.
anchor	"cut site" or "fragment center". If "fragment center" is used, the input bamfiles must be paired-end.

Value

an invisible list of matrixes with the signals for plot. It includes: - signal mean values of coverage for positive strand and negative strand in feature regions - spearman.correlation spearman correlations of cleavage counts in the highest 10-nucleotide-window and binding prediction score. - bindingSites predicted binding sites.

Author(s)

Jianhong Ou, Julie Zhu

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

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

footprintsScanner	<i>scan ATAC-seq footprints infer factor occupancy genome wide</i>
-------------------	--

Description

Aggregate ATAC-seq footprint for a bunch of motifs generated over binding sites within the genome.

Usage

```
footprintsScanner(bamfiles, index = bamfiles, bindingSitesList,
                 seqlev = "chr1", upstream = 100, downstream = 100)
```

Arguments

bamfiles	A vector of characters indicates the file names of bams. All the bamfiles will be pulled together.
index	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
bindingSitesList	A object of GRangesList indicates candidate binding sites (eg. the output of fimo).
seqlev	A vector of characters indicates the sequence levels.
upstream, downstream	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. It includes: - signal mean values of coverage for positive strand and negative strand in feature regions - spearman.correlation spearman correlations of cleavage counts in the highest 10-nucleotide-window and binding prediction score. - bindingSites predicted binding sites.

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

```
bamfile <- system.file("extdata", "GL1.bam",
                      package="ATACseqQC")
bsl <- system.file("extdata", "jolma2013.motifs.bindingList.95.rds",
                  package="ATACseqQC")
bindingSitesList <- readRDS(bsl)
footprintsScanner(bamfile, bindingSitesList=bindingSitesList)
```

fragSizeDist

fragment size distribution

Description

estimate the fragment size of bams

Usage

```
fragSizeDist(bamFiles, bamFiles.labels, index = bamFiles, 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.
index	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
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 <- dir(system.file("extdata", package="ATACseqQC"), "GL.*.bam$", full.names=TRUE)
bamFiles.labels <- sub(".bam", "", basename(bamFiles))
fragSizeDist(bamFiles, bamFiles.labels)
```

NFRscore

*Nucleosome Free Regions (NFR) score***Description**

NFR score is a ratio between cut signal adjacent to TSS and that flanking the corresponding TSS. Each TSS window of 400 bp is first separated into 3 sub-regions: the most upstream 150 bp (n1), the most downstream of 150 bp (n2), and the middle 100 bp (nf). Then the number of fragments with 5' ends overlapping each region are calculated for each TSS. The NFR score for each TSS is calculated as $NFR\text{-score} = \log_2(nf) - \log_2((n1+n2)/2)$. A plot can be generated with the NFR scores as Y-axis and the average signals of 400 bp window as X-axis, very much like a MA plot for gene expression data.

Usage

```
NFRscore(obj, txs, seqlev = intersect(seqlevels(obj), seqlevels(txs)),
         nucleosomeSize = 150, nucleosomeFreeSize = 100)
```

Arguments

obj an object of [GAlignments](#)
txs GRanges of transcripts
seqlev A vector of characters indicates the sequence levels.
nucleosomeSize numeric(1) or integer(1). Default is 150
nucleosomeFreeSize numeric(1) or integer(1). Default is 100

Value

A object of [GRanges](#) with NFR scores

Author(s)

Jianhong Ou

Examples

```
library(GenomicRanges)
bamfile <- system.file("extdata", "GL1.bam",
                      package="ATACseqQC", mustWork=TRUE)
gal1 <- readBamFile(bamfile=bamfile, tag=character(0),
                  which=GRanges("chr1", IRanges(1, 1e6)),
                  asMates=FALSE)
library(TxDb.Hsapiens.UCSC.hg19.knownGene)
txs <- transcripts(TxDb.Hsapiens.UCSC.hg19.knownGene)
nfr <- NFRscore(gal1, txs)
```

peakdet *Detect peak positions*

Description

Detect the peaks positions and valley positions. The algorithm is modified from `github::dgroner/peakdet`

Usage

```
peakdet(y, delta = 0, silence = TRUE)
```

Arguments

y	A vector of numeric where to search peaks
delta	A numeric of length 1, defining the local threshold for peak detection. If it is set to 0, the delta will be set to 1/10 of the range of y.
silence	logical(1). If false, echo the delta value when delta is set as 0.

Value

A list with peakpos and valleypos. Both peakpos and valleypos are vectors of numeric which indicate the positions of peak or valley.

plotCorrelation *plot Correlations of multiple samples*

Description

plot PCA or heatmap for multiple bamfiles. The correlation is calculated by the counts in promoter regions.

Usage

```
plotCorrelation(objs, txs, seqlev = intersect(seqlevels(objs[[1]]),
  seqlevels(txs)), upstream = 2000, downstream = 500,
  type = c("heatmap", "PCA"), ...)
```

Arguments

objs	an object of GAlignmentsList
txs	GRanges of transcripts
seqlev	A vector of characters indicates the sequence levels.
upstream	numeric(1) or integer(1). Start position of promoter. Default is 2000
downstream	numeric(1) or integer(1). End position of promoter. Default is 500
type	Figure type, heatmap or PCA plot.
...	parameters could be passed to downstream functions such as plot for pca or heatmap for heatmap.

Value

A invisible object of [GRanges](#) with counts

Author(s)

Jianhong Ou

Examples

```
library(GenomicRanges)
library(GenomicAlignments)
path <- system.file("extdata", package="ATACseqQC", mustWork=TRUE)
bamfiles <- dir(path, "*.bam$", full.name=TRUE)
gals <- lapply(bamfiles, function(bamfile){
  readBamFile(bamFile=bamfile, tag=character(0),
              which=GRanges("chr1", IRanges(1, 1e6)),
              asMates=FALSE)
})
library(TxDb.Hsapiens.UCSC.hg19.knownGene)
txs <- transcripts(TxDb.Hsapiens.UCSC.hg19.knownGene)
plotCorrelation(GAlignmentsList(gals), txs, seqlev="chr1")
```

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, xlim, ylim, newpage = TRUE,
  motif, segmentation)
```

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
xlim	xlim
ylim	ylim
newpage	Plot the figure in a new page?
motif	a pfm object.
segmentation	the segmentation position and abundance

Value

Null.

Author(s)

Jianhong Ou

Examples

```
library(MotifDb)
CTCF <- query(MotifDb, c("CTCF"))
CTCF <- as.list(CTCF)
motif <- new("pfm", mat=CTCF[[1]], name="CTCF")
ATACseqQC:::plotFootprints(Profile=sample.int(500),
                           Mlen=ncol(CTCF[[1]]), motif=motif)
```

prepareBindingSitesList

helper function for preparing the binding list

Description

helper function for preparing the binding list

Usage

```
prepareBindingSitesList(pfms, genome, seqlev = paste0("chr", c(1:22, "X",
"Y")), expSiteNum = 5000)
```

Arguments

pfms	A list of Position frequency Matrix represented as a numeric matrix with row names A, C, G and T.
genome	An object of BSgenome .
seqlev	A vector of characters indicates the sequence levels.
expSiteNum	numeric(1). Expect number of predicted binding sites. if predicted binding sites is more than this number, top expSiteNum binding sites will be used.

Examples

```
library(MotifDb)
motifs <- query(MotifDb, c("Hsapiens"))
motifs <- as.list(motifs)
library(BSgenome.Hsapiens.UCSC.hg19)
#bindingSitesList <- prepareBindingSitesList(motifs, genome=Hsapiens)
```

PTscore	<i>Promoter/Transcript body (PT) score</i>
---------	--

Description

PT score is calculated for coverage of promoter divided by the coverage of transcripts body. PT score will show if the signal is enriched in promoters.

Usage

```
PTscore(obj, txs, seqlev = intersect(seqlevels(obj), seqlevels(txs)),  
        upstream = 2000, downstream = 500)
```

Arguments

obj	an object of GAlignments
txs	GRanges of transcripts
seqlev	A vector of characters indicates the sequence levels.
upstream	numeric(1) or integer(1). Start position of promoter. Default is 2000
downstream	numeric(1) or integer(1). End position of promoter. Default is 500

Value

A object of [GRanges](#) with PT scores

Author(s)

Jianhong Ou

Examples

```
library(GenomicRanges)  
bamfile <- system.file("extdata", "GL1.bam",  
                      package="ATACseqQC", mustWork=TRUE)  
gal1 <- readBamFile(bamfile=bamfile, tag=character(0),  
                  which=GRanges("chr1", IRanges(1, 1e6)),  
                  asMates=FALSE)  
library(TxDb.Hsapiens.UCSC.hg19.knownGene)  
txs <- transcripts(TxDb.Hsapiens.UCSC.hg19.knownGene)  
pt <- PTscore(gal1, txs)
```

pwmcores *max PWM scores for sequences*

Description

calculate the maximal PWM scores for each given sequences

Usage

```
pwmcores(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. IU-PAC ambiguity letters in subject are ignored (i.e. assigned weight 0) with a warning.

Value

a numeric vector

Author(s)

Jianhong

readBamFile *read in bam files*

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

readsDupFreq

Calculating duplication frequency

Description

Calculating the frequency of read duplication based on alignment status determined by rname, strand, pos, cigar, mrnm, mpos and isize.

Usage

```
readsDupFreq(bamFile, index = bamFile)
```

Arguments

bamFile	A character vector of length 1L containing the name of a BAM file. Only a BAM file with duplication reads are meaningful for estimating the library complexity. For example, a raw BAM file output by aligners, or a BAM file with mitochondrial reads removed.
index	A character vector of length 1L containing the name of a BAM index file.

Value

A two-column matrix of integers. The 1st column is the frequency $j = 1, 2, 3, \dots$. The 2nd column is the number of genomic regions with the same frequency (j) of duplication. The frequency column is in ascending order.

Author(s)

Haibo Liu

Examples

```
bamFile <- system.file("extdata", "GL1.bam", package = "ATACseqQC")
freq <- readsDupFreq(bamFile)
```

saturationPlot	<i>Plotting Saturation curves</i>
----------------	-----------------------------------

Description

Plotting the saturation curves.

Usage

```
saturationPlot(subsamplingPeakFiles, subsamplingSizes, sep = "\t",
  header = FALSE, fdr = 0.05, fdrCol = 9, startCol = 2,
  endCol = 3, skipLines = 1, peakCaller = "MACS2", outPrefix,
  span = 2, degree = 2)
```

Arguments

subsamplingPeakFiles	A character vector containing peak files from peak calling tools, such as MACS2. Currently only MACS2 output is supported.
subsamplingSizes	A named vector of integers, which are the sizes of subsamples for peak calling. The names of subsamplingPeakFiles should be identical to the basenames of subsamplingPeakFiles.
sep	A character vector of length 1L, which is the column separator used in peak files.
header	A boolean (TRUE or FALSE) vector of length 1L, showing whether there are column headers in the peak files.
fdr	A decimal between 0 and 1, a cutoff of statistical significance of peak detection.
fdrCol	An integer, column index for fdr.
startCol	An integer, column index for start positions of peak regions.
endCol	An integer, column index for end positions of peak regions.
skipLines	An integer, the number of lines (comments or instruction) to skip when peak files are read into R.

peakCaller	A character vector of length 1L containing the name of the peak caller used to generate the peak files, such as "MACS2". Currently only MACS2 output (XXX.narrowPeak or XXX.broadPeak) is support.
outPrefix	A character vector of length 1L, the file prefix for outputting saturation plots.
span	An integer, the span parameter for loess smoothing to fit a smoothed saturation curve.
degree	An integer, the degree of local polynomial used for loess.

Value

A data frame of three columns: subsamplingSizes, the number of subsampled fragments; numPeaks, the number of peaks with fdr less than a given threshold when a given number of fragmetns are subsampled; breadth, the total breadth of peaks with fdr less than a given threshold for give subsampling when a given number of fragmetns are subsampled.

Author(s)

Haibo Liu

Examples

```
if(interactive()){
}
```

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

gal	An object of GAlignmentsList .
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](#) 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, index = bamfile, 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

bamfile	character(1). File name of bam.
tags	A vector of characters indicates the tags in bam file.
index	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
outPath	Output file path.
txs	GRanges of transcripts.
genome	An object of BSgenome
conservation	An object of GScores .
positive	integer(1). the size to be shift for positive strand
negative	integer(1). the size to be shift for negative strand
breaks	A numeric vector for fragment size of nucleosome free, mononucleosome, dinucleosome and trinucleosome
labels	A vector of characters indicates the labels for the levels of the resulting category. The length of labels = length of breaks - 1
seqlev	A vector of characters indicates the sequence levels.
cutoff	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

```
if(Sys.getenv("USER")=="jianhongou"){
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 bam 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",
  trainingSetPercentage = 0.15, cutoff = 0.8,
  halfSizeOfNucleosome = 80L)
```

Arguments

obj	an object of GAlignments
txs	GRanges of transcripts
genome	an object of BSgenome
conservation	an object of GScores .
breaks	a numeric vector for fragment size of nucleosome free, mononucleosome, dinucleosome and trinucleosome. The breaks pre-defined here is following the description of Greenleaf's paper (see reference).
labels	a character vector for labels of the levels of the resulting category.
labelsOfNucleosomeFree, labelsOfMononucleosome	character(1). The label for nucleosome free and mononucleosome.
trainingSetPercentage	numeric(1) between 0 and 1. Percentage of training set from top coverage.
cutoff	numeric(1) between 0 and 1. cutoff value for prediction.
halfSizeOfNucleosome	numeric(1) or integer(1). The 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)
```

vPlot

V-plot

Description

Aggregate ATAC-seq Fragment Midpoint vs. Length for a given motif generated over binding sites within the genome.

Usage

```
vPlot(bamfiles, index = bamfiles, pfm, genome, min.score = "95%",
      bindingSites, seqlev = paste0("chr", c(1:22, "X", "Y")),
      upstream = 200, downstream = 200, maxSiteNum = 1e+06,
      draw = TRUE, ...)
```

Arguments

bamfiles	A vector of characters indicates the file names of bams. All the bamfiles will be pulled together.
index	The names of the index file of the 'BAM' file being processed; This is given without the '.bai' extension.
pfm	A Position frequency Matrix represented as a numeric matrix with row names A, C, G and T.
genome	An object of BSgenome .

min.score	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 .
bindingSites	A object of GRanges indicates candidate binding sites (eg. the output of fimo).
seqlev	A vector of characters indicates the sequence levels.
upstream, downstream	numeric(1) or integer(1). Upstream and downstream of the binding region for aggregate ATAC-seq footprint.
maxSiteNum	numeric(1). Maximal number of predicted binding sites. if predicted binding sites is more than this number, top maxSiteNum binding sites will be used.
draw	Plot or not. Default TRUE.
...	parameters could be used by smoothScatter

Value

an invisible data.frame for plot.

Author(s)

Jianhong Ou

References

Jorja G. Henikoff, Jason A. Belsky, Kristina Krassovsky, David M. MacAlpine, and Steven Henikoff. Epigenome characterization at single base-pair resolution. PNAS 2011 108 (45) 18318-18323

Examples

```
bamfile <- system.file("extdata", "GL1.bam",
                      package="ATACseqQC")
library(MotifDb)
CTCF <- query(MotifDb, c("CTCF"))
CTCF <- as.list(CTCF)
library(BSgenome.Hsapiens.UCSC.hg19)
vPlot(bamfile, pfm=CTCF[[1]],
      genome=Hsapiens,
      min.score="95%", seqlev="chr1",
      ylim=c(30, 250))
```

```
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

bamQC, 3
BSgenome, 7, 13, 20, 22

distanceDyad, 3
DNAStrng, 15

enrichedFragments, 4
estimateLibComplexity, 6
estLibSize, 5
export, 23

factorFootprints, 7
footprintsScanner, 8
fragSizeDist, 9

GAlignments, 10, 14, 18, 20, 21, 24
GAlignmentsList, 11, 18
GRanges, 5, 7, 10, 12, 14, 16, 20, 23
GRangesList, 8
GScores, 20, 21

IntegerRangesList, 16

MaskedDNAStrng, 15
matchPWM, 7, 23

NFRscore, 10
normalizeBetweenArrays, 5

peakdet, 11
plotCorrelation, 11
plotFootprints, 12
prepareBindingSitesList, 13
PTscore, 14
pwmScores, 15

randomForest, 20
readBamFile, 15
readGAlignments, 16
readGAlignmentsList, 16
readsDupFreq, 6, 16
Rsamtools, 16

saturationPlot, 17
ScanBamParam, 16
shiftGAlignmentsList, 18, 20
shiftReads, 19
smoothScatter, 23
splitBam, 19
splitGAlignmentsByCut, 20, 21

Views, 15
vPlot, 4, 22

writeListOfGAlignments, 20, 23