

# Package ‘GreyListChIP’

October 9, 2015

**Type** Package

**Version** 1.0.1

**Title** Grey Lists -- Mask Artefact Regions Based on ChIP Inputs

**Date** 2015-05-30

**Author** Gord Brown <gdbzork@gmail.com>

**Maintainer** Gordon Brown <gordon.brown@cruk.cam.ac.uk>

**Description** Identify regions of ChIP experiments with high signal in the input, that lead to spurious peaks during peak calling. Remove reads aligning to these regions prior to peak calling, for cleaner ChIP analysis.

**License** Artistic-2.0

**LazyLoad** yes

**Depends** R (>= 3.1), methods, GenomicRanges

**Imports** GenomicAlignments, BSgenome, Rsamtools, rtracklayer, MASS, parallel, GenomeInfoDb

**Suggests** BiocStyle, BiocGenerics, RUnit

**Enhances** BSgenome.Hsapiens.UCSC.hg19

**biocViews** ChIPSeq, Alignment, Preprocessing, DifferentialPeakCalling, Sequencing, GenomeAnnotation, Coverage

**NeedsCompilation** no

## R topics documented:

calcThreshold-methods . . . . .	2
countReads-methods . . . . .	3
getKaryotype-methods . . . . .	4
greyList . . . . .	5
GreyList-class . . . . .	6
greyListBS . . . . .	8
loadKaryotype-methods . . . . .	9
makeGreyList-methods . . . . .	10

<b>Index</b>	<b>11</b>
--------------	-----------

---

`calcThreshold-methods` *Calculate Read Count Threshold*

---

**Description**

Based on the counts from `countReads`, sample counts from the set several times, estimate the parameters of the negative binomial distribution for each sample, then calculate the mean of the parameters (*size* and *mu*). Use these values to calculate the read count threshold, given the specified p-value threshold.

**Usage**

```
calcThreshold(obj, reps=100, sampleSize=30000, p=0.99, cores=1)
```

**Arguments**

<code>obj</code>	A <code>GreyList</code> object for which to calculate the threshold.
<code>reps</code>	The number of times to sample bins and estimate the parameters of the negative binomial distribution.
<code>sampleSize</code>	The number of bins to sample on each repetition.
<code>p</code>	The p-value threshold for marking bins as “grey”.
<code>cores</code>	The number of CPU cores (parallel threads) to use when sampling repeatedly from the set of counts

**Details**

This method samples from the set of counts generated during the `countReads` step. Each sample is fitted to the negative binomial distribution, and the parameters estimated. The means of the *mu* and *size* parameters is calculated, then used to choose a read count threshold, given the p-value cutoff provided. If `cores` is given, the process will use that many cores to parallelize the parameter estimation.

**Value**

The modified `GreyList` object, with the threshold added.

**Author(s)**

Gord Brown

**References**

Venables, W. N. and Ripley, B. D. (2002) *Modern Applied Statistics with S*. Fourth edition. Springer.

## Examples

```
# Load a pre-built R object with counts.
data(greyList)

# Calculate the threshold:
gl <- calcThreshold(gl, reps=10, sampleSize=1000, p=0.99, cores=1)
```

---

countReads-methods      *Count reads from a [BamFile](#)*

---

## Description

Given the tiling of the genome created when the [GreyList](#) object was created (or replaced via `getKaryotype`), count reads overlapping the bins, in preparation for estimating the threshold for grey-listing bins.

## Usage

```
countReads(obj, bamFile)
```

## Arguments

obj	A <a href="#">GreyList</a> object on which to count reads.
bamFile	A <a href="#">BamFile</a> from which to count reads.

## Details

This method counts reads contained within the bins that make up the genome tiling. Bins are overlapping (by default 1Kb bins at 512b intervals) so reads are counted once for each bin that wholly contains them.

## Value

The modified [GreyList](#) object, with added counts.

## Author(s)

Gord Brown

## See Also

[GreyList](#), [BamFile](#)

## Examples

```
# Load a pre-built GreyList object.
data(greyList)

path <- system.file("extra", package="GreyListChIP")
## Not run: fn <- file.path(path, "sample_chr21.bam")
## Not run: gl <- countReads(gl, fn)
```

---

getKaryotype-methods *Replace the karyotype of a [GreyList](#) object*

---

## Description

Though a [BSgenome](#) object (or a karyotype file) is supplied when the [GreyList](#) object is created, it is conceivable that the user might want to replace it. This method allows that.

## Usage

```
getKaryotype(obj, genome, tileSize=1024)
```

## Arguments

obj	A <a href="#">GreyList</a> object.
genome	A <a href="#">BSgenome</a> object, from which to take the karyotype.
tileSize	The size in nucleotides of each tile. Overlapping tiles will be generated, spaced at 1/2 the width of the tiles.

## Value

Returns the [GreyList](#) object with a new genome and tiling.

## Author(s)

Gord Brown

## See Also

[GreyList](#), [BSgenome](#)

## Examples

```
# Load a pre-built GreyList object.
data(greyList)
library(BSgenome.Hsapiens.UCSC.hg19)

# Replace the karyotype, updating the genome tiling.
## Not run: gl <- getKaryotype(gl, BSgenome.Hsapiens.UCSC.hg19)
```

---

`greyList`*A sample `GreyList` object for use in examples.*

---

## Description

This is a sample `GreyList` object, covering only human chromosome 21 (from genome version hg19). The input library used to generate this grey list can be found in the European Nucleotide Archive, under accession number ERR336953.

The library was made from a culture of the MCF-7 cell line, bought from ATCC. The library was sequenced to a depth of 35,716,191 reads on an Illumina Genome Analyzer IIX. The reads were aligned to human reference genome hg19 (GRCh37) using BWA version 0.7.5a with default parameters. Approximately 96% of reads aligned to the reference genome. Reads aligning to chromosome 21 were extracted using Samtools. The chromosome 21-only karyotype file was created by deleting all lines except chromosome 21, in a file generated by `fetchChromSizes` as described in the vignette. This package was then used to create the `GreyList` sample object.

When printed, the object displays several important slots in the object (if they have been filled with calculated values). For example, this object has all its slots filled, indicating that the analysis is complete:

```
GreyList on karyotype file karyotype_chr21.txt
  tiles: 94004
  files: jc899_chr21.bam
  size (mean): 0.370332362145541
  mu (mean): 10.2330008719269
  params: reps=10, sample size=1000, p-value=0.99
  threshold: 81
  regions: 118
  coverage: 4.45%
```

The fields are described in the class's documentation, but briefly, we can see:

1. the name of the karyotype file (or `BSgenome` object),
2. the number of tiles (overlapping by 1/2 the tile width),
3. the BAM file(s) used for read counting (currently only 1 is allowed),
4. the two estimated parameters of the negative binomial distribution, `size` and `mu`,
5. the input parameters,
6. the calculated read depth cutoff (over 1kb tiles),
7. the number of distinct regions, and
8. the percentage coverage of the reference genome.

The fact that all the fields are present indicates that the regions have been generated; otherwise fields still without values would be omitted. Of course any stage can be re-run with different parameters.

**Usage**

```
data(greyList)
```

**Format**

An S4 [GreyList](#) object.

**Value**

A [GreyList](#) object named `gl`.

---

GreyList-class	"GreyList" <i>Objects</i>
----------------	---------------------------

---

**Description**

Regions of high signal in the input samples of a ChIP experiment can lead to artefacts in peak calling. This class generates "grey lists" of such regions, for use in filtering reads before peak calling (or filtering peaks after peak calling, though it is generally safer to filter first).

**Objects from the Class**

Objects can be created by calls of the form `new("GreyList", genome, ...)`, where `genome` is a "BSgenome" object describing a genome, such as `BSgenome.Hsapiens.UCSC.hg19`. Alternatively, a karyotype file can be provided explicitly: `new("GreyList", karyoFile=fn, ...)`. Either `genome` or `karyoFile` must be provided; if both are present, the [BSgenome](#) object takes precedence.

**Slots**

`genome`: The [BSgenome](#) object corresponding with the genome the reads are aligned to

`karyotype`: The [Seqinfo](#) object from the [BSgenome](#) object, or made from the `karyo_file`

`karyo_file`: The name of a file containing chromosome sizes for the reference genome of interest, one per line, as "chromName chromLength" pairs.

`tiles`: A [GRanges](#) object with an overlapping tiling of the genome (by default 1Kb tiles every 512b).

`counts`: A numeric vector holding the counts corresponding to the tiling and the BAM file provided.

`files`: A vector of BAM filenames that were used to generate the counts (currently only accepts one).

`size_param`: The computed estimates of the "size" parameter of the negative binomial distribution, estimated by `MASS::fitdistr` from repeated sampling from the counts.

`size_stderr`: The standard errors of the "size" parameters, as estimated by `MASS::fitdistr`.

`size_mean`: The mean of the "size" estimates.

**mu\_param:** Computed estimates of the "mu" parameter of the negative binomial distribution, estimated by `MASS::fitdistr` from repeated sampling from the counts.

**mu\_stderr:** The standard errors of the "mu" parameter.

**mu\_mean:** The mean of the "mu" estimates.

**reps:** How many samples from the counts were taken.

**sample\_size:** How many values were sampled from the counts, for each estimate of "size" and "mu".

**pvalue:** The requested p-value threshold.

**threshold:** The calculated threshold, based on the p-value.

**max\_gap:** The largest gap to consider when merging nearby regions (i.e. if there are "grey" regions up to this many nucleotides apart, merge them into one long region).

**regions:** A [GRanges](#) object defining the final grey list regions.

**coverage:** The percentage of the genome covered by the grey list regions.

## Methods

**calcThreshold** signature(obj = "GreyList"): Calculate the cutoff for reads in bins, based on fitting the counts to a negative binomial distribution.

**countReads** signature(obj = "GreyList"): Count reads in bins across the genome.

**export** signature(object = "GreyList", con = "character", format = "missing"): Write the grey list to a file.

**initialize** signature(.Object = "GreyList"): Create an initial object (invoked automatically by `new("GreyList", ...)`).

**loadKaryotype** signature(obj = "GreyList"): Load a genome description from a file. The file format is one line per chromosome, with the name of the chromosome followed by white space followed by an integer indicating the length of the chromosome.

**getKaryotype** signature(obj = "GreyList"): Get the karyotype of a genome from a [BSgenome](#) object.

**makeGreyList** signature(obj = "GreyList"): Compute the actual grey list, after calculating the threshold.

**show** signature(object = "GreyList"): Display the grey list.

## Author(s)

Gord Brown (<[gdbzork@gmail.com](mailto:gdbzork@gmail.com)>)

## See Also

[BSgenome](#), [Seqinfo](#)

## Examples

```
showClass("GreyList")

# Load a karyotype file:
path <- system.file("extra", package="GreyListChIP")
fn <- file.path(path, "karyotype_chr21.txt")

# Create a GreyList object:
gl <- new("GreyList", karyoFile=fn)
```

---

greyListBS

*Construct a grey list with default arguments*

---

## Description

This function is a convenience function, wrapping several steps of grey list construction into one step. If you are content to accept the package's defaults, and are using a [BSgenome](#) object to supply the karyotype, this function might be of use to you.

## Usage

```
greyListBS(genome, bam)
```

## Arguments

genome	a BSgenome object, for the relevant genome.
bam	a BAM file to use for making the grey list.

## Value

An object of class GreyList.

## Author(s)

Gord Brown

## See Also

[GreyList](#)

## Examples

```
# If you want to accept the defaults for everything, you can create the
# GreyList in one step using a BSgenome object:
library(BSgenome.Hsapiens.UCSC.hg19)
path <- system.file("extra", package="GreyListChIP")
## Not run: fn <- file.path(path, "sample_chr21.bam")
## Not run: gl <- greyListBS(BSgenome.Hsapiens.UCSC.hg19, fn)
```



---

loadKaryotype-methods *Load a karyotype from a file*

---

### Description

Load a karyotype from a file, (re)generate the tiling, for a [GreyList](#) object.

### Usage

```
loadKaryotype(obj, karyoFile, tileSize=1024)
```

### Arguments

obj	A <a href="#">GreyList</a> object to set a new karyotype for.
karyoFile	A text file describing a genome's karyotype. The format is one line per chromosome (or contig or whatever), with the name of the chromosome, some white space, and an integer giving the length of the chromosome in nucleotides.
tileSize	The width of tiles on which to count. Tiles will be placed every $\text{tileSize}/2$ nucleotides, to catch regions of high signal that might otherwise be split across (non-overlapping) tiles and hence missed.

### Value

Returns the [GreyList](#) object with a new genome and tiling, loaded from the provided file.

### Author(s)

Gord Brown

### See Also

[GreyList](#)

### Examples

```
# load a pre-built GreyList object:
data(greyList)

# Get a karyotype file:
path <- system.file("extra", package="GreyListChIP")
fn <- file.path(path, "karyotype_chr21.txt")

# Replace the karyotype in the GreyList:
gl <- loadKaryotype(gl, fn)
```

---

makeGreyList-methods    *Generate a grey list from a [GreyList](#) object*

---

### Description

Create the actual grey list, based on the threshold calculated by `calcThreshold` and the read counts from `countReads`.

### Usage

```
makeGreyList(obj,maxGap=16384)
```

### Arguments

<code>obj</code>	The <a href="#">GreyList</a> object to create the list for.
<code>maxGap</code>	If the distance between neighbouring grey regions is less than or equal to <code>maxGap</code> , the regions will be merged into one big region.

### Details

Create the grey list as a [GRanges](#) object. Merge grey regions if they are separated by up to `maxGap` bp.

### Value

The modified [GreyList](#), with the regions added.

### Author(s)

Gord Brown

### See Also

[GreyList](#), [GRanges](#)

### Examples

```
# load a pre-built GreyList object:
data(greyList)

# calculate the actual regions:
gl <- makeGreyList(gl)
```

# Index

- \*Topic **classes**
  - GreyList-class, 6
- \*Topic **datasets**
  - greyList, 5
- \*Topic **grey list**
  - greyListBS, 8
- \*Topic **methods**
  - calcThreshold-methods, 2
  - countReads-methods, 3
  - getKaryotype-methods, 4
  - loadKaryotype-methods, 9
  - makeGreyList-methods, 10
  
- BamFile, 3
- BSgenome, 4, 6–8
  
- calcThreshold (calcThreshold-methods), 2
- calcThreshold, GreyList-method (calcThreshold-methods), 2
- calcThreshold-methods, 2
- countReads (countReads-methods), 3
- countReads, GreyList-method (countReads-methods), 3
- countReads-methods, 3
  
- export, GreyList, character, missing-method (GreyList-class), 6
  
- getKaryotype (getKaryotype-methods), 4
- getKaryotype, GreyList-method (getKaryotype-methods), 4
- getKaryotype-methods, 4
- gl (greyList), 5
- GRanges, 6, 7, 10
- GreyList, 3–6, 8–10
- GreyList (GreyList-class), 6
- greyList, 5
- GreyList-class, 6
- greyListBS, 8
  
- initialize, GreyList, BSgenome-method (GreyList-class), 6
  
- loadKaryotype (loadKaryotype-methods), 9
- loadKaryotype, GreyList-method (loadKaryotype-methods), 9
- loadKaryotype-methods, 9
  
- makeGreyList (makeGreyList-methods), 10
- makeGreyList, GreyList-method (makeGreyList-methods), 10
- makeGreyList-methods, 10
  
- Seqinfo, 6, 7
- show, GreyList-method (GreyList-class), 6