

Package ‘casper’

October 9, 2013

Version 1.1.2

Date 2013-06-26

Title Characterization of Alternative Splicing based on Paired-End Reads

Author David Rossell, Camille Stephan-Otto, Manuel Kroiss

Maintainer David Rossell <rosselldavid@gmail.org>

Depends R (>= 2.14.1), Biobase, IRanges, methods, gtools, GenomicRanges, Rsamtools, plyr, gaga

Imports VGAM, mgcv, GenomicFeatures, survival, sqldf

Enhances parallel

Description Infer alternative splicing from paired-end RNA-seq data. The model is based on counting paths across exons, rather than pairwise exon connections, and estimates the fragment size and start distributions non-parametrically, which improves estimation precision.

License GPL (>=2)

LazyLoad yes

Collate GenericDefs.R ClassDefinitions.R calcDenovo.R calcExp.R
createDenovoGenome.R genePlot.R getDistrs.R modelPrior.R
pathCounts.R procBam.R procGenome.R rmShortInserts.R
simMultSamples.R simPost.R simReads.R splitGenomeByLength.R
truncnorm.R mergeExp.R wrapKnown.R

biocViews Bioinformatics, GeneExpression, DifferentialExpression, Transcription, RNASeq, High-ThroughputSequencing

R topics documented:

annotatedGenome-class	2
calcDenovo	4
calcExp	6
denovoGeneExpr-class	8

denovoGenomeExpr-class	9
genePlot	10
getDistrs	11
getIsland	13
getNreads	14
getReads	15
hg19DB	15
K562.r111	16
mergeExp	17
modelPrior	18
modelPriorAS-class	19
pathCounts	20
pathCounts-class	21
plot	22
plotExpr	23
plotPriorAS	24
procBam	25
procBam-class	26
procGenome	27
rmShortInserts	29
simMultSamples	30
simReads	31
splitGenomeByLength	32
subsetGenome	33
transcripts	34
txLength	35
wrapKnown	36
Index	38

annotatedGenome-class *Class "annotatedGenome"*

Description

Stores an annotated genome, either known or build de novo by combining a known genome with observed RNA-seq data.

Objects from the Class

Objects are typically created with a call to `procGenome` (for known genomes) or to `createDenovoGenome` (for de novo genomes).

Slots

islands GRangesList object with elements corresponding to gene islands. It indicates the start/end/name of each exon contained in the island

transcripts Each element in the list corresponds to a gene island. It indicates the exons contained in each known variant.

exon2island data.frame indicating the chromosome, start and end of each exon, and its corresponding gene island.

exonsNI GRanges indicating the chromosome, start/end and id of each exon

aliases data.frame indicating the aliases for each known transcript, i.e. transcripts having the exact same sequence of exons (as indicated in the UCSC database).

genomeVersion Character indicating the genome version from which the object was build, e.g. "hg19"

dateCreated Character indicating the date when the object was created. UCSC genomes change from time to time, so that an "hg19" genome from Jan 2012 may not be exactly the same as in Dec 2012.

denovo Logical variable. FALSE indicates that the object was created using available annotation only. TRUE indicates that new exons/islands were added based on the data observed in a particular RNA-seq experiment.

txLength Numeric vector storing transcript lengths.

Methods

show signature(object = "annotatedGenome"): Displays general information about the object.

Author(s)

Camille Stephan-Otto Attolini

See Also

[procGenome](#) and [createDenovoGenome](#) to create annotatedGenome objects.

Examples

```
showClass("annotatedGenome")
```

 calcDenovo

Estimate expression of gene splicing variants de novo.

Description

calcDenovo estimates expression of gene splicing variants, considering both known variants and variants that have not been previously described.

relativeExpr computes relative expression estimates from calcDenovo output via Bayesian model selection or averaging.

Usage

```
calcDenovo(distrs, genomeDB, pc, readLength, islandid, priorq=3, mprior,
minpp=0.001, selectBest=FALSE, method="auto", niter, exactMarginal=TRUE,
verbose=TRUE, mc.cores=1)
```

```
relativeExpr(expr, summarize="modelAvg", minProbExpr=0.5, minExpr=0.05)
```

Arguments

distrs	List of fragment distributions as generated by the getDistrs function
genomeDB	denovoGenome object containing annotated genome, as returned by the createDenovoGenome function.
pc	Named vector of exon path counts as returned by pathCounts
readLength	Read length in bp, e.g. in a paired-end experiment where 75bp are sequenced on each end one would set readLength=75.
islandid	Name of the gene island to be analyzed. If not specified, all gene islands are analyzed.
priorq	Parameter of the prior distribution on the proportion of reads coming from each variant. The prior is Dirichlet with prior sample size for each variant equal to priorq. We recommend priorq=3 as this defines a non-local prior that penalizes including variants with low estimated expression unless they are warranted by the data.
mprior	Prior on the model space, as returned by modelPrior. This argument is used to favor splicing variant configurations that are consistent with the annotated genome. If left missing a uniform prior on the model space is used, i.e. all splicing variants configurations are equally likely a priori. We strongly recommend setting mprior via modelPrior. See details.
minpp	Models (i.e. splicing configurations) with posterior probability less than minpp are not reported. This argument can help reduce substantially the amount of required memory to store the results.
selectBest	If set to TRUE only the model with highest posterior probability is reported. While this can save memory, we do not recommend this option as it may ignore a substantial amount of uncertainty.

method	Method used to perform the model search. "exact" exhaustively enumerates all possible models (warning: this is not feasible for genes over 5 exons). "rwmc" uses a random-walk MCMC scheme to focus on models with high posterior probability. "auto" uses "exact" for genes with up to 5 exons and "rwmc" for longer genes. See details.
niter	Number of MCMC iterations.
exactMarginal	Set to FALSE to estimate posterior model probabilities as the proportion of MCMC visits. Set to TRUE to use the integrated likelihoods (default). See details.
verbose	Set to TRUE to display progress information.
mc.cores	Number of processors to be used for parallel computation. Can only be used if the package multicore is available for your system. Warning: using multiple processors substantially increases the memory requirements, so set this value carefully.
expr	Object of class denovoGenomeExpr, as returned by calcDenovo.
summarize	Set to "modelAvg" to average estimates across models according to the posterior probability of each model. Set to "bestModel" to report the posterior mode, i.e. the estimate under the model with highest posterior probability.
minProbExpr	Variants with marginal posterior probability of being expressed below minProbExpr are not reported.
minExpr	Variants with estimated expression below minExpr are not reported.

Details

calcDenovo attempts to consider all possible splicing variant configurations for each gene. For example, a gene with 5 exons has 31 possible variants ($2^5 - 1$), each of which may be either expressed or not expressed. This gives rise to $2^{31} - 1$ possible configurations of expressed variants.

The routines computes the posterior probability of each model (i.e. configuration of expressed variants) via Bayes theorem.

$$P(\text{model}) \propto m(y|\text{model}) P(\text{model})$$

where $m(y|\text{model})$ is the integrated likelihood and $P(\text{model})$ is the prior probability of the model. Importantly, $P(\text{model})$ can be set by analyzing available genome annotations. For instance, genes with 20 exons have variants that tend to have close to 20 exons. They also tend to have more expressed variants than genes with 5 exons. The function modelPrior analyzes an annotated genome to set reasonable values for $P(\text{model})$.

Given that the number of models grows super-exponentially with the number of exons, an exhaustive enumeration of all possible models is not feasible unless the gene is very short (in practice, we set the limit at 5 exons). For longer genes, we use a random walk MCMC scheme that incorporates Shotgun Stochastic Search. This random exploration scheme focuses on models with large posterior probability. The longer the number of iterations (niter) the more stable the results, but this can come at a substantial computational cost.

In order to compute $P(\text{model})$ one can either use the computed $m(y|\text{model}) P(\text{model})$ (option exactMarginal==TRUE) or the proportion of MCMC visits (option exactMarginal==FALSE). Unless niter is large the former option typically provides more precise estimates.

Value

Expression set with expression estimates. `featureNames` identify each transcript via RefSeq ids, and the `featureData` contains Entrez identifiers and the length of each transcript (in bp). If `ci` type was set to a value other than "none", the `featureData` also contains the 95% credibility intervals (i.e. intervals that contain the true parameter value with 95% posterior probability).

Author(s)

Camille Stephan-Otto Attolini, Manuel Kroiss, David Rossell

References

Rossell D. QASPER: Quantifying Alternative Splicing from Paired End Reads. Technical report, IRB Barcelona (2010)

See Also

`plotExpr` to produce a plot with splicing variants and estimated expression.

Examples

```
##---- Should be DIRECTLY executable !! ----
##-- ==> Define data, use random,
##--or do help(data=index) for the standard data sets.
```

calcExp

Estimate expression of a known set of gene splicing variants.

Description

Estimate expression of gene splicing variants, assuming that the set of variants is known. When `relativeExpr` is set to FALSE, reads per kilobase per million are returned. Otherwise relative expression estimates are returned.

Usage

```
calcExp(distrs, genomeDB, pc, readLength, islandid, rpkm=TRUE, priorq=2,
priorqGeneExpr=2, ciType="none", niter=10^3, burnin=100, mc.cores=1, verbose=FALSE)
```

Arguments

<code>distrs</code>	List of fragment distributions as generated by the <code>getDistrs</code> function
<code>genomeDB</code>	knownGenome object containing annotated genome, as returned by the <code>procGenome</code> function.
<code>pc</code>	Named vector of exon path counts as returned by <code>pathCounts</code>
<code>readLength</code>	Read length in bp, e.g. in a paired-end experiment where 75bp are sequenced on each end one would set <code>readLength=75</code> .

islandid	Name of the gene island to be analyzed. If not specified, all gene islands are analyzed.
rpkkm	Set to FALSE to return relative expression levels, i.e. the proportion of reads generated from each variant per gene. These proportions add up to 1 for each gene. Set to TRUE to return reads per kilobase per million (RPKM).
priorq	Parameter of the prior distribution on the proportion of reads coming from each variant. The prior is Dirichlet with prior sample size for each variant equal to priorq. We recommend priorq=2 for estimation, as it pools the estimated expression away from 0 and 1 and returned lower estimation errors than priorq=1 in our simulated experiments.
priorqGeneExpr	Parameter for prior distribution on overall gene expression. Defaults to 2, which ensures non-zero estimates for all genes
citype	Set to "none" to return no credibility intervals. Set to "asympt" to return approximate 95% CIs (obtained via the delta method). Set to "exact" to obtain exact CIs via Monte Carlo simulation. Options "asympt" and especially "exact" can increase the computation time substantially.
niter	Number of Monte Carlo iterations. Only used when citype=="exact".
burnin	Number of burnin Monte Carlo iterations. Only used when citype=="exact".
mc.cores	Number of processors to be used for parallel computation. Can only be used if the package multicore is available for your system.
verbose	Set to TRUE to display progress information.

Value

Expression set with expression estimates. featureNames identify each transcript via RefSeq ids, and the featureData contains further information. If citype was set to a value other than "none", the featureData also contains the 95% credibility intervals (i.e. intervals that contain the true parameter value with 95% posterior probability).

Author(s)

Camille Stephan-Otto Attolini, Manuel Kroiss, David Rossell

References

Rossell D. QASPER: Quantifying Alternative Splicing from Paired End Reads. Technical report, IRB Barcelona (2010)

Examples

```
data(K562.r111)
data(hg19DB)

#Pre-process
bam0 <- rmShortInserts(K562.r111, isizeMin=100)
pbam0 <- procBam(bam0)
head(getReads(pbam0))
```

```

#Estimate distributions, get path counts
distrs <- getDistrs(hg19DB,bam=bam0,readLength=75)
pc <- pathCounts(pbam0, DB=hg19DB)

#Get estimates
eset <- calcExp(distrs=distrs, genomeDB=hg19DB, pc=pc, readLength=75, rpkm=FALSE)
head(exprs(eset))
head(fData(eset))

#Add fake sample by permuting and combine
eset2 <- eset[sample(1:nrow(eset),replace=FALSE),]
esetall <- mergeExp(eset,eset2)

#After merge samples are correctly matched
head(exprs(esetall))

```

```
denovoGeneExpr-class  Class "denovoGeneExpr"
```

Description

denovoGeneExpr stores inferred expression for de novo splicing variants for a single gene. denovoGenomeExpr stores the information for several genes (typically, the whole genome).

Objects from the Class

Objects are returned by calcDenovo. When running calcDenovo on multiple genes results are returned in a denovoGenomeExpr object. Results for a single gene can be retrieved using the `[[` operator as usual, which returns a denovoGeneExpr object.

Slots

posprob data.frame containing the posterior probability of each model

expression data.frame with the estimated expression of each variant under each model

variants RangedData indicating the exons contained in each variant.

integralSum Sum of the $\log(\text{integrated likelihood}) + \log(\text{model prior probability})$ across all considered models.

npathDeleted Number of paths that had 0 probability under all considered variants and had to be excluded for model fitting purposes.

Methods

show signature(object = "denovoGeneExpr"): Displays general information about the object.

"[" Selects a subset of genes

"["[" Selects a single gene

posprob Accesses the posterior probabilities of each model (slot posprob)

variants Accesses the variants and their respective exons (slot variants)

variants<- Replaces the value of the slot variants (can be useful for renaming variants, for instance)

Author(s)

David Rossell

See Also[procGenome](#) and [createDenovoGenome](#) to create denovoGeneExpr objects.**Examples**

```
showClass("denovoGeneExpr")
```

```
denovoGenomeExpr-class
```

```
Class "denovoGenomeExpr"
```

Description

denovoGeneExpr stores inferred expression for de novo splicing variants for a single gene. denovoGenomeExpr stores the information for several genes (typically, the whole genome).

Objects from the Class

Objects are returned by calcDenovo.

Slots

islands A list of denovoGeneExpr objects, with each element containing results for an individual gene.

Methods

show signature(object = "denovoGenomeExpr"): Displays general information about the object.

as.list Coerces the object to a list

"[" Selects a subset of genes

"[]" Selects a single gene

Author(s)

Camille Stephan-Otto Attolini

See Also[procGenome](#) and [createDenovoGenome](#) to create denovoGenomeExpr objects.**Examples**

```
showClass("denovoGeneExpr")
showClass("denovoGenomeExpr")
```

genePlot *Plot exon structure for each transcript of a given gene.*

Description

Plot exon structure for each transcript of a given gene. Optionally, aligned reads can be added to the plot.

Usage

```
genePlot(generanges, islandid, genomeDB, reads, exp, names.arg, xlab='', ylab='', xlim, cex=1, yaxt='n')
```

Arguments

generanges	Object containing the ranges with start/end of each exon.
islandid	If generanges is not specified, transcripts are obtained from island islandid from the annotated genome genomeDB.
genomeDB	Annotated genome produced with the "procGenome" function
reads	pbam object with aligned reads. This is an optional argument.
exp	ExpressionSet object with expression values, as returned by calcExp. This is an optional argument.
names.arg	Optionally, indicate the names of each transcript.
xlab	x-axis label
ylab	y-axis label
xlim	x-axis limits, defaults to start of 1st exon and end of last exon
cex	Character expansion
yaxt	The y-axis in the plot has no interpretation, hence by default it is not displayed.
col	Either single color or vector of colors to be used to draw each transcript. Defaults to rainbow colors.
...	Other arguments to be passed on to plot.

Value

A plot is produced.

Methods

```
signature(generanges="CompressedIRangesList", islandid="ANY", genomeDB="ANY", reads="ANY", exp="ANY")
Plots a set of transcripts. Each element in the generanges corresponds to a transcript. Each transcript should contain exon start/end positions.
```

```
signature(generanges="IRanges", islandid="ANY", genomeDB="ANY", reads="ANY", exp="ANY")
Plots a single transcript. Each range indicates the start/end of a single exon.
```

signature(generanges="IRangesList", islandid="ANY", genomeDB="ANY", reads="ANY", exp="ANY")
Plots a set of transcripts. Each element in the generanges corresponds to a transcript. Each transcript should contain exon start/end positions.

signature(generanges="GRangesList", islandid="ANY", genomeDB="ANY", reads="ANY", exp="ANY")
Plots a set of transcripts. Each element in the generanges corresponds to a transcript. Each transcript should contain exon start/end positions.

signature(generanges="GRanges", islandid="ANY", genomeDB="ANY", reads="ANY", exp="ANY")
Plots a set of transcripts. Each space in generanges corresponds to a transcript. Each transcript should contain exon start/end positions.

signature(generanges="missing", islandid="character", genomeDB="annotatedGenome", reads="GRanges", exp="ANY")
Plots all transcripts stored in genomeDB for island with identifier islandid. Individual reads are added to the plot (reads contains start/end of individual read fragments).

signature(generanges="missing", islandid="character", genomeDB="annotatedGenome", reads="missing", exp="ANY")
Plots all transcripts stored in genomeDB for island with identifier islandid.

signature(generanges="missing", islandid="character", genomeDB="annotatedGenome", reads="procBam", exp="ANY")
Plots all transcripts stored in genomeDB for island with identifier islandid. Individual reads are added to the plot (reads contains start/end of individual read fragments).

Author(s)

Camille Stephan-Otto Attolini, David Rossell

Examples

```
data(hg19DB)

#Plot an IRangesList
txs <- transcripts(entrezid="27", genomeDB=hg19DB)
genePlot(txs)

#Equivalently, indicate islandid
islandid <- getIsland(entrezid="27", genomeDB=hg19DB)
genePlot(islandid=islandid, genomeDB=hg19DB)
```

getDistrs

Compute fragment start and fragment length distributions

Description

Compute fragment start distributions by using reads aligned to genes with only one annotated variant. Estimate fragment length distribution using fragments aligned to long exons (>1000nt). Fragment length is defined as the distance between the start of the left-end read and the end of the right-end read.

Usage

```
getDistrs(DB, bam, pbam, islandid, verbose, nreads=4*10^6, readLength)
```

Arguments

DB	Annotated genome. Object of class knownGenome as returned by procGenome.
bam	Aligned reads, as returned by scanBam. It must be a list with elements 'qname', 'rname', 'pos' and 'mpos'. Ignored when argument pbam is specified.
pbam	Processed BAM object of class procBam, as returned by function procBam. Arguments bam and readLength are ignored when pbam is specified.
islandid	Island IDs of islands to be used in the read start distribution calculations (defaults to genes with only one annotated variant)
verbose	Set to TRUE to print progress information.
nreads	To speed up computations, only the first nreads are used to obtain the estimates. The default value of 4 millions usually gives highly precise estimates.
readLength	Read length in bp, e.g. in a paired-end experiment where 75bp are sequenced on each end one would set readLength=75.

Value

An object of class readDistrs with slots:

lenDis	Table with number of fragments with a given length
stDis	Cumulative distribution function (object of type closure) for relative start position

Author(s)

Camille Stephan-Otto Attolini, David Rossell

Examples

```
data(K562.r111)
data(hg19DB)
bam0 <- rmShortInserts(K562.r111, isizeMin=100)

distrs <- getDistrs(hg19DB,bam=bam0,readLength=75)

#Fragment length distribution
plot(distrs,'fragLength')

#Fragment start distribution (relative to transcript length)
plot(distrs,'readSt')
```

getIsland	<i>getIsland returns the island id associated to a given entrez or transcript id in an annotatedGenome object. getChr indicates the chromosome for a given Entrez, transcript or island id.</i>
-----------	---

Description

annotatedGenome objects store information regarding genes and transcripts. When there's an overlap in exons between several genes, these genes are grouped into gene islands. getIsland retrieves the island to which each gene or transcript was assigned, while getChr indicates the chromosome.

Usage

```
getIsland(entrezid, txid, genomeDB)
getChr(entrezid, txid, islandid, genomeDB)
```

Arguments

entrezid	Character indicating single Entrez identifier. Can be left missing and specify another identifier instead.
txid	Character indicating a single RefSeq transcript identifier. Can be left missing and specify another identifier instead.
islandid	Character indicating the gene island identifier. Can be left missing and specify another identifier instead.
genomeDB	Object of class annotatedGenome

Value

Character with island identifier

Methods

```
signature(entrezid='character', txid='missing', genomeDB='annotatedGenome') Return
  island id for given Entrez identifier
signature(entrezid='missing', txid='character', genomeDB='annotatedGenome') Return
  island id for given transcript identifier (RefSeq)
signature(entrezid='character', txid='missing', islandid='missing', genomeDB='annotatedGenome')
  Return chromosome for given Entrez identifier (RefSeq)
signature(entrezid='missing', txid='character', islandid='missing', genomeDB='annotatedGenome')
  Return chromosome for given transcript identifier (RefSeq)
signature(entrezid='missing', txid='missing', islandid='character', genomeDB='annotatedGenome')
  Return chromosome for given island identifier

signature(entrezid='character', txid='missing', islandid='missing') Return chromo-
  some for given Entrez identifier
```

signature(entrezid='missing',txid='character',islandid='missing') Return chromosome for given transcript identifier (RefSeq)

signature(entrezid='missing',txid='character',islandid='missing') Return chromosome for given island identifier

Examples

```
data(hg19DB)
getIsland(entrezid="27",genomeDB=hg19DB)
getIsland(txid="NM_005158",genomeDB=hg19DB)

getChr(entrezid="27",genomeDB=hg19DB)
getChr(txid="NM_005158",genomeDB=hg19DB)
```

getNreads

Get total number of paths in each island from a pathCounts object.

Description

getNreads returns a numeric vector with the total number of path counts in each island from a pathCounts object.

Usage

```
getNreads(pc)
```

Arguments

pc pathCounts object generated by pathCounts()

Value

Numeric vector with total number of path counts in each island of pc.

Methods

signature(pathCounts='pathCounts') Returns numeric vector with total number of path counts for each island in the pathCounts object.

Author(s)

Camille Stephan-Otto Attolini

Examples

```
##---- Should be DIRECTLY executable !! ----
##-- ==> Define data, use random,
##--or do help(data=index) for the standard data sets.
```

getReads	<i>getReads returns the reads stored in a procBam object.</i>
----------	---

Description

procBam objects store reads that have been split according to their CIGAR codes. getReads accesses these reads.

Usage

```
getReads(x)
```

Arguments

x Object of class procBam

Value

RangedData object with reads stored in x.

Methods

signature(x='procBam') Return reads stored in x.

Examples

```
#See example in calcExp
```

hg19DB	<i>Subset of human genome (UCSC hg19 version)</i>
--------	---

Description

We downloaded the human genome hg19 via procGenome and selected a few genes from chromosome 1 to use as a toy data for the vignette and examples.

Usage

```
data(hg19DB)
```

Format

An annotatedGenome object. See help(procGenome) and help(annotatedGenome-class) for details.

Examples

```
data(hg19DB)
hg19DB
slotNames(hg19DB)
```

K562.r111

Toy RNA-seq data from RGASP project.

Description

The paired-end RNA-seq data is from the RGASP project sample K562_2x75 (replicate 1, lane 1) and was obtained at ftp://ftp.sanger.ac.uk/pub/gencode/rgasp/RGASP1/inputdata/human_fastq. Reads were aligned against hg19 with tophat 2.0.2 and bowtie 0.12.5, setting the insert size at -r 200, and imported into R using scanBam from package Rsamtools. For illustration purposes, we selected reads mapping to a few genes only (namely, the genes that were also selected for the toy genome annotation in data(hg19DB).

Usage

```
data(K562.r111)
```

Format

A list indicating read id, chromosome, start and end locations and the position of the pair, as returned by scanBam.

Source

```
ftp://ftp.sanger.ac.uk/pub/gencode/rgasp/RGASP1/inputdata/human_fastq
```

References

C Trapnell, L Pachter, SL Salzberg. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics*, 2009, 25, 1105-1111. doi=10.1093/bioinformatics/btp120.

B Langmead, C Trapnell, M Pop, SL Salzberg. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology*, 2009, 10:R25.

Examples

```
data(K562.r111)
names(K562.r111)
```

mergeExp	<i>Merge splicing variant expression from multiple samples</i>
----------	--

Description

mergeExp combines the output of calcExp from multiple samples, i.e. multiple ExpressionSet objects, into a single ExpressionSet

Usage

```
mergeExp(..., sampleNames, keep=c('transcript','gene'))
```

Arguments

...	ExpressionSet objects to be combined.
sampleNames	Character vector indicating the name of each sample. Defaults to 'Sample1', 'Sample2', etc.
keep	Variables in the featureData of each individual ExpressionSet to keep in the merged output.

Details

mergeExp runs some checks to ensure that object can be combined (e.g. making sure that measurements are obtained on same set of genes), then sorts and formats each input ExpressionSet.

A label with the sample name is appended to variables in the featureData that appear in multiple samples, e.g. variable 'se' reporting standard errors (obtained by setting citype='asympt' in calcExp).

Value

Object of class ExpressionSet combining the input ExpressionSets. Its featureData contains the columns indicated in the keep argument, plus a column readCount with the total number of reads mapped to each gene (or gene island, when multiple genes have overlapping exons).

Author(s)

David Rossell

See Also

calcExp to obtain an ExpressionSet for an individual sample.

Examples

```
#See example in calcExp
```

 modelPrior

Set prior distribution on expressed splicing variants.

Description

Set prior on expressed splicing variants using the genome annotation contained in a knownGenome object.

The prior probability of variants V_1, \dots, V_n being expressed depends on n , on the number of exons in each variant V_1, \dots, V_n and the number of exons in the gene. See the details section.

Usage

```
modelPrior(genomeDB, maxExons=40, smooth=TRUE, verbose=TRUE)
```

Arguments

genomeDB	Object of class knownGenome
maxExons	The prior distribution is estimated for genes with 1 up to maxExons exons. As there are fewer genes with many exons, the prior parameters are estimated poorly. To avoid this common estimate is used for all genes with more than maxExons exons
smooth	If set to TRUE the estimated prior distribution parameters for the number of exons in a gene are smoothed using Generalized Additive Models. This step typically improves the precision of the estimates, and is only applied to genes with 10 or more exons.
verbose	Set to TRUE to print progress information.

Details

The goal is to set a prior that takes into account the number of annotated variants for genes with E exons, as well as the number of exons in each variant.

Suppose we have a gene with E exons. Let V_1, \dots, V_n be n variants of interest and let $|V_1|, \dots, |V_n|$ be the corresponding number of exons in each variant. The prior probability of variants V_1, \dots, V_n being expressed is modeled as

$$P(V_1, \dots, V_n | E) = P(n | E) P(|V_1| | E) \dots P(|V_n| | E)$$

where $P(n | E) = \text{NegBinom}(n; k_E, r_E) I(0 < n < 2^E)$ and $P(|V_i| | E) = \text{BetaBinomial}(|V_i| - 1; E - 1, \alpha_E, \beta_E)$.

The parameters $k_E, r_E, \alpha_E, \beta_E$ depend on E (the number of exons in the gene) and are estimated from the available annotation via maximum likelihood. Parameters are estimated jointly for all genes with $E \geq \text{maxExons}$ in order to improve the precision.

For `smooth==TRUE`, α_E and β_E are modeled as a smooth function of E by calling `gam` and setting the smoothing parameter via cross-validation. Estimates for genes with $E \geq 10$ are substituted by their smooth versions, which typically helps improve stability in the estimates.

Value

List with 2 components.

nvarPrior	List with prior distribution on the number of expressed variants for genes with 1,2,3... exons. Each element contains the truncated Negative Binomial parameters, observed and predicted frequencies (counting the number of genes with a given number of variants).
nexonPrior	List with prior distribution on the number of exons in a variant for genes with 1,2,3... exons. Each element contains the Beta-Binomial parameters, observed and predicted frequencies (counting the number of variants with a given number of exons)

Author(s)

David Rossell, Camille Stephan-Otto Attolini

Examples

```
data(hg19DB)
mprior <- modelPrior(hg19DB, maxExons=10)

##Prior on number of expressed variants
##Genes with 2 exons
##mprior$nvarPrior[['2']]
##Genes with 3 exons
##mprior$nvarPrior[['3']]

##Prior on the number of exons in an expressed variant
##Genes with 2 exons
##mprior$nexonPrior[['2']]
##Genes with 3 exons
##mprior$nexonPrior[['3']]
```

modelPriorAS-class *Class "modelPriorAS"*

Description

modelPriorAS stores parameters for the prior distribution on all possible alternative splicing configuration (i.e. prior on model space). This information is used for de novo reconstruction of splicing variants.

Objects from the Class

Objects are created by function modelPrior.

Slots

nvarPrior Prior on the number of variants per gene. A list with components "nbpar" containing the parameters of the Negative Binomial distribution, "obs" containing the observed counts and "pred" the Negative Binomial predicted counts.

nexonPrior Prior on the number of exons in an expressed variant. A list with components "bbpar" containing Beta-Binomial parameters, "obs" containing the observed counts and "pred" the Beta-Binomial predicted counts.

Methods

show signature(object = "modelPriorAS"): Displays general information about the object.

"[" Selects prior parameters for genes with the specified number of exons

coef Selects a single gene

Author(s)

David Rossell

See Also

[procGenome](#) and [createDenovoGenome](#) to create modelPriorAS objects.

Examples

```
showClass("modelPriorAS")
```

pathCounts

Compute exon path counts

Description

Compute counts for exon paths visited by aligned reads

Usage

```
pathCounts(reads, DB, mc.cores = 1, verbose=FALSE)
```

Arguments

reads	Object of class procBam containing aligned reads, as returned by procBam.
DB	Object of class annotatedGenome containing either a known or de novo annotated genome.
mc.cores	Number of processors to be used for parallel computing. Requires having package multicore installed and loaded.
verbose	Set to TRUE to print progress information.

Value

Named integer vector with counts of exon paths. Names are character strings built as ".exon1.exon2-exon3.exon4.", with dashes making the split between exons visited by left and right-end reads correspondingly.

Methods

`signature(reads='list')` Computes counts for exon paths from a list of `procBam` objects (usually reads processed and split by chromosome).

`signature(reads='procBam')` Compute counts for exon paths from a `procBam` object of processed reads.

Author(s)

Camille Stephan-Otto Attolini

See Also

[procGenome](#) to create an annotated genome object, [createDenovoGenome](#) to create a de novo annotated genome. See `help(getNreads)` to get number of fragments mapping to each island.

Examples

```
##---- Should be DIRECTLY executable !! ----
##-- ==> Define data, use random,
##--or do help(data=index) for the standard data sets.
```

pathCounts-class *Class "pathCounts"*

Description

Stores exon path counts.

Objects from the Class

Objects are created with a call to `pathCounts`.

Slots

counts List with one element per gene island. For each island, it contains a named vector with exon path counts. The names indicate the visited exons.

For instance, consider that for gene '1' with 2 exons we observe 10 reads in which the left end falls completely in exon 1 and the right end in exon 2. Suppose that for 5 reads the left end bridges exons 1-2 and the right end falls in exon 2. Then `pc[['1']]` would contain `c(10,5)` and `names(pc[['1']])` would contain `c(".1-2.", ".1.2-2.")`

denovo Logical variable. FALSE indicates that the counts correspond to a known genome (i.e. created with `procGenome`), and TRUE to a de novo annotated genome (i.e. created with `createDenovoGenome`).

stranded Logical variable. TRUE indicates that the path counts were obtained from an RNA-seq experiment where strand information was preserved.

Methods

`show` `signature(object = "pathCounts")`: Displays general information about the object.

Author(s)

Camille Stephan-Otto Attolini

Examples

```
showClass("pathCounts")
```

<code>plot</code>	<i>Plot estimated read start and fragment length distributions.</i>
-------------------	---

Description

Plots the estimated fragment length (insert size) distribution and the relative read start distribution (0 indicating transcription start, 1 transcription end). The former checks that the insert size distribution matches that described in the experimental protocol. The latter checks the extent to which reads are non-uniformly distributed (note: `casper` does NOT assume reads to be uniformly distributed, so a lack of uniformity is not a problem per se).

Usage

```
plot(x, y, ...)  
  
lines(x, ...)
```

Arguments

<code>x</code>	Object of type <code>readDistrs</code> , as returned by <code>getDistrs</code> .
<code>y</code>	Set to <code>"fragLength"</code> to plot the estimated insert size distribution. Set to <code>"readSt"</code> to plot a histogram of the estimated read start distribution.
<code>...</code>	Further arguments to be passed on to <code>plot</code> .

Methods

`signature(x = "readDistrs")` `x` is an object of type `readDistrs`, as returned by `getDistrs`.
The `plot` function allows to visualize the fragment length and read start distributions.

Examples

```
#See getDistrs examples
```

plotExpr *Plot inferred gene structure and expression.*

Description

Plots variants with sufficiently large posterior probability of being expressed along with their (marginal) estimated expression.

Usage

```
plotExpr(gene, minProbExpr = 0.5, minExpr = 0.1,
         xlab = "(kb)", ylab = "", xlim, cex = 1, yaxt = "n", col, ...)
```

Arguments

gene	denovoGeneExpr object containing results for a single gene, as returned by calcDenovo.
minProbExpr	Variants with marginal posterior probability of expression below minProbExpr are not reported
minExpr	Variants with (marginal) estimated expression below minExpr are not reported. Can be useful to remove sequence preference artifacts.
xlab	x-axis label, passed on to plot
ylab	y-axis label, passed on to plot
xlim	x-axis limits, passed on to plot
cex	Character expansion, passed on to plot
yaxt	Type of y-axis, passed on to plot
col	Colors for each variant, defaults to rainbow colors. It is possible to specify a single color.
...	Other arguments to be passed on to plot

Details

The marginal posterior probability that a variant is expressed is the sum of the posterior probabilities of all models containing that variant.

The marginal estimated expression is the average expression across all models (including those where the variant has 0 expression) weighted by the posterior probability of each model.

Methods

signature(gene = "denovoGeneExpr") gene contains the results from a de novo isoform expression analysis for a single gene, as returned by calcDenovo. When calcDenovo is run on multiple genes simultaneously, the desired gene can be selected using the "[" operator as usual.

Examples

```
#See calcDenovo examples
```

plotPriorAS	<i>Plot prior distribution on set of expressed variants (i.e. the model space).</i>
-------------	---

Description

Plots the prior distribution on the number of expressed variants and the number of exons per variant in genes with exons (as returned by function modelPrior). The prior distribution is compared to the observed frequencies to check that the assumed distributional forms are reasonable.

Usage

```
plotPriorAS(object, type = "nbVariants", exons = 1:9, xlab, ylab = "Probability", col = c("red", "blue"))
```

Arguments

object	modelPriorAS object with prior distribution on model space.
type	Set to "nbVariants" to plot the prior on the number of variants per gene. Set to "nbExons" to plot the prior on the number of exons.
exons	Vector with integers. The plot is only produced with number of exons indicated in exons.
xlab	x-axis label, passed on to plot
ylab	y-axis label, passed on to plot
col	Colors for bars showing prior probabilities and frequencies in the known genome

Methods

signature(object = "modelPriorAS") object contains the prior distribution on the model space, as returned by function modelPrior

Examples

```
#See modelPrior examples
```

procBam	<i>Process BAM object</i>
---------	---------------------------

Description

Process paired-end data stored in BAM object generated by scanBam. Outputs GRanges objects for reads and junctions.

Usage

```
procBam(bam, stranded, seed, verbose=FALSE, rname, keep.junx=FALSE,
keep.flag=FALSE, ...)
```

Arguments

bam	BAM object generated by scanBam
stranded	Set to TRUE to indicate that the RNA-seq experiment preserved the strand information.
seed	Seed for random number generator
verbose	Set to TRUE to print progress information.
rname	Chromosome to process be combined with the <code>'which'</code> argument in the scanBam function)
keep.junx	Option to store junction information. Only useful for finding denovo exons and transcripts.
keep.flag	Option to store alignment flag information.
...	Other arguments to be passed on to procBam.

Details

In case of multihits with same start position for both reads but different insertions/deletions patterns only one alignment is chosen at random.

Value

An object of class procBam containing reads with both ends correctly aligned and split according to the corresponding CIGAR. Unique identifiers by fragment are stored. Junctions spanned by reads are also stored in GRanges object if the argument `'keep.junx'` is set to TRUE.

Methods

```
signature(bam='list',stranded='logical',seed='integer',verbose='logical', rname='character',keep.ju
Process paired-end data stored in BAM object generated by scanBam. Outputs GRanges ob-
jects for reads and (optionally) junctions.
```

Author(s)

Camille Stephan-Otto Attolini

See Also

scanBam from package Rsamtools

Examples

```
##See example in calcExp
```

procBam-class	<i>Class "procBam"</i>
---------------	------------------------

Description

Stores processed bam files in a RangedData format. Each read is split into disjoint ranges according to its cigar code.

Objects from the Class

Objects are created with a call to procBam.

Slots

pbam GRanges indicating chromosome, start and end of each disjoint range. The pair id and read id within the pair are also stored.

junx GRanges indicating chromosome, start and end of junctions spanned by reads.

stranded Logical variable. TRUE indicates that the reads were obtained from an RNA-seq experiment where strand information was preserved.

In the case of stranded experiments:

plus GRanges indicating chromosome, start and end of each disjoint range for fragments originated from the positive strand. The pair id and read id within the pair are also stored.

minus GRanges indicating chromosome, start and end of each disjoint range for fragments originated from the negative strand. The pair id and read id within the pair are also stored.

pjunx GRanges indicating chromosome, start and end of junctions spanned by reads originated from the positive strand.

mjunx GRanges indicating chromosome, start and end of junctions spanned by reads originated from the negative strand.

Methods

show signature(object = "procBam"): Displays general information about the object.

Author(s)

Camille Stephan-Otto Attolini, David Rossell

Examples

```
showClass("procBam")
```

procGenome	<i>Create an annotatedGenome object that stores information about genes and transcripts</i>
------------	---

Description

procGenome processes annotations for a given transcriptome, either from a TranscriptDb object created by GenomicFeatures package (e.g. from UCSC) or from a user-provided GRanges object (e.g. by importing a gtf file output by Cufflinks RABT module with import from rtracklayer package).

createDenovoGenome creates a de novo annotated genome by combining UCSC annotations and observed RNA-seq data.

Usage

```
procGenome(genDB, genome, mc.cores = 1)
```

```
createDenovoGenome(reads, DB, minLinks=2,
maxLinkDist=1e+05, maxDist=1000, minConn=2, minJunx=3, minLen=12, mc.cores=1)
```

Arguments

genDB	Either annotations downloaded from UCSC (e.g. genDB<-makeTranscriptDbFromUCSC(genome="hg19", tablename="refGene")) or GRanges object (e.g. genDB <- import('transcripts.gtf')). See details.
genome	Character indicating genome version (e.g. "hg19", "dm3")
mc.cores	Number of cores to use in parallel processing (multicore package required)
DB	annotatedGenome object, as returned by procGenome
minLinks	Minimum number of reads joining two exons to merge their corresponding genes
maxLinkDist	Maximum distance between two exons to merge their corresponding genes. A value of 0 disables this option.
maxDist	Maximum distance between two exons with reads joining them to merge their corresponding genes.
minConn	Minimum number of fragments connecting a new exon to an annotated one to add to denovo genome.

minJunx	Minimum number of junctions needed to redefine an annotated exon's end or start.
minLen	Minimum length of a junction to consider as a putative intron.
reads	Processed reads stored in a RangedData, as returned by procBam

Details

These functions create the annotation objects that are needed for subsequent functions. Typically these objects are created only once for a set of samples.

If interested in quantifying expression for known transcripts only, one would typically use `procGenome` with a `TranscriptDb` from the usual Bioconductor annotations. It is also possible to provide `procGenome` with any desired set of transcripts, so that one can use other genome annotations or even de novo transcripts found by some third-party software. These annotations should be provided as a `GRanges` object. For instance, one could use the transcripts found by the Cufflinks RABT module. These transcripts are saved in `gtf` files, which can be imported into Bioconductor as `GRanges` with function `import` (package `rtracklayer`). `procGenome` is intended to be used in combination with `calcExp`, which quantifies expression for a set of given transcripts.

`createDenovoGenome` performs a de novo transcript discovery following our own strategy, and is intended to be used in combination with `calcDenovo`, which performs a model search additional to quantifying expression. It is convenient to create a single genome by combining information from all samples, as this increases the power to detect new exons and guarantees that the same exons and gene islands are used across samples.

Value

Object of class `annotatedGenome`.

Methods

`signature(genDB = "transcriptDb")` `genDB` is usually obtained with a call to `makeTranscriptDbFromUCSC` (package `GenomicFeatures`), e.g. `genDB<-makeTranscriptDbFromUCSC(genome="hg19", tablename="refGene")`

`signature(genDB = "GRanges")` `genDB` stores information about all transcripts and their respective exons. Chromosome, start, end and strand are stored as usual in `GRanges` objects. `genDB` must have a column named "type" taking the value "transcript" for rows corresponding to transcript and "exon" for rows corresponding to exons. It must also store transcript and gene ids. For instance, Cufflinks RABT module creates a `gtf` file with information formatted in this manner for known and de novo predicted isoforms.

Author(s)

Camille Stephan-Otto Attolini

See Also

See `annotatedGenome-class` for a description of the class. To obtain the island id for an entrez gene id or a transcript id see `help(getIsland)`. To extract the transcripts for a entrez or island id see `help(transcripts)`. See `splitGenomeByLength` for splitting an `annotatedGenome` according to gene length.

Examples

```
##Known transcripts from UCSC
## genDB<-makeTranscriptDbFromUCSC(genome="hg19", tablename="refGene")
## hg19DB <- procGenome(genDB, "hg19")

##Transcripts inferred by Cufflinks-RABT module
##genDB.Cuff <- import('transcripts.gtf')
##hg19DB.Cuff <- procGenome(genDB.Cuff, genome='hg19')
```

rmShortInserts	<i>Remove reads with short insert sizes from imported BAM files.</i>
----------------	--

Description

In paired-end experiments short inserts (i.e. the 2 ends being very close to each other), may indicate RNA degradation or that a short RNA (e.g. miRNA) is being sequenced. Typically the goal is not to study alternative splicing for such short/degraded RNA; in this case it is recommendable to remove such short inserts to avoid biasing the insert size distribution. Requiring a minimum insert size can also result in significantly faster computations when quantifying alternative splicing via `calc` or `calcDenovo`.

Usage

```
rmShortInserts(bam, isizeMin=100)
```

Arguments

<code>bam</code>	Object with aligned reads, as returned by <code>scanBam</code>
<code>isizeMin</code>	Reads with insert size smaller than <code>isizeMin</code> will be removed.

Value

Named list, in the same format as that returned by `scanBam`.

Note

The insert size is stored in objects imported with `scanBam` in the element named `isize`.

Author(s)

David Rossell

Examples

```
##---- Should be DIRECTLY executable !! ----
##-- ==> Define data, use random,
##--or do help(data=index) for the standard data sets.
```

simMultSamples *Simulate paired end reads for multiple samples based on pilot data*

Description

Draw true expression levels for each group and future individual samples within each group, and simulate future observed data (in the form of Casper expression estimates).

These simulations can serve as the basis for sample size calculation: if one were to sequence `nsamples` new RNA-seq samples, what data would we expect to see? The simulation is posterior predictive, i.e. based on the current available data `x`.

Usage

```
simMultSamples(B, nsamples, nreads, readLength, x, groups='group', distrs, genomeDB, verbose=TRUE, mc.
```

Arguments

<code>B</code>	Number of simulations to obtain
<code>nsamples</code>	Vector indicating number of future samples per group, e.g. <code>nsamples=c(5,5)</code> to simulate 5 new samples for 2 groups.
<code>nreads</code>	number of RNA-seq reads per sample. <code>simMultSamples</code> assumes that exactly these many reads will be aligned to the genome.
<code>readLength</code>	Read length, i.e. in an experiment with paired reads at 100bp each, <code>readLength=100</code> .
<code>x</code>	ExpressionSet containing pilot data. <code>x[[group]]</code> indicates groups to be compared
<code>groups</code>	Name of column in <code>pData(x)</code> indicating the groups
<code>distrs</code>	Fragment start and length distributions. It can be either an object of class <code>readDistrs</code> , or a list where each element is of class <code>readDistrs</code> . In the latter case, an element is chosen at random for each individual sample (so that uncertainty in these distributions is taken into account)
<code>genomeDB</code>	annotatedGenome object
<code>verbose</code>	Set to TRUE to print progress
<code>mc.cores</code>	Number of cores to use in function. <code>mc.cores>1</code> requires package <code>multicore</code>

Value

List of length `B`, where each element is itself a list containing an individual simulation.

<code>simTruth</code>	<code>data.frame</code> indicating the mean and standard deviation of the Normal distribution used to generate data from each group
<code>simExpr</code>	ExpressionSet with Casper expression estimates, as returned by <code>calcExp</code> . <code>pData(simExpr)</code> indicates group information, and <code>fData(simExpr)</code> the number of observed reads for each gene

Author(s)

Victor Pena

simReads	<i>Function to simulate paired end reads following given read start and fragment length distributions and gene and variant expressions.</i>
----------	---

Description

This function generates path counts and bam files with simulated paired end reads according to given read start distribution, fragment length distribution and gene and variant expressions.

Usage

```
simReads(islandid, nSimReads, pis, rl, seed, writeBam, distrs, genomeDB,
repSims=FALSE, bamFile=NULL, stranded=FALSE, verbose=TRUE, chr=NULL, mc.cores=1)
```

Arguments

islandid	Island ID's from the genomeDB object to simulate reads
nSimReads	Named numeric vector with number of fragments to simulate in each island.
pis	Named numeric vector with relative expression of transcripts. Expressions add up to one for each island to simulate.
rl	Read length
seed	Seed of the random numbers generator
writeBam	Set to 1 to generate bam files with the simulated reads
distrs	Object of class 'readDistrs' with read start and fragment length distributions
genomeDB	Object of class 'annotatedGenome' with the genome to generate reads from
repSims	Set to TRUE to return relative read starts and fragment lengths from the simulation
bamFile	Name of the bam file to write reads to. Must end with '.bam'
stranded	Set to TRUE to preserve gene strand when generating reads. The 'XS' tag will be added to reads in the bam file and the returned 'pc' object will be stranded
verbose	Set to TRUE to print progress
chr	Characters vector with chromosomes to simulate. Defaults to whole genome simulations.
mc.cores	Number of cores to use in function

Value

Nsim	Numerical vector with the number of reads simulated for each island.
pc	Object of class 'pathCounts' with simulated path counts
sims	Only if 'repSims' is set to TRUE. List with vectors of length 'n' with the following elements: -'varl': Length of variant for corresponding read -'st' Start of fragment relative to variant start (not in genomic coordinates) -len:Fragment length -'strand':Strand of gene for simulated read

Author(s)

Camille Stephan-Otto Attolini

Examples

```

data(hg19DB)
data(K562.r111)
distrs <- getDistrs(hg19DB,bam=K562.r111,readLength=75)

islandid <- c('10319','463')
txs <- unlist(lapply(hg19DB@transcripts[islandid], names))
pis <- vector(mode='numeric', length=length(txs))
npis <- sapply(hg19DB@transcripts[islandid],length)
pis[1:npis[1]] <- rep(1/npis[1],npis[1])
pis[-1:-npis[1]] <- rep(1/npis[2],npis[2])
names(pis) <- txs
nSimReads <- c(100, 100)
names(nSimReads) <- islandid

simpc <- simReads(islandid=islandid, nSimReads=nSimReads, pis=pis,
r1=75, repSims=TRUE, seed=1, writeBam=FALSE, distrs=distrs,genomeDB=hg19DB)

```

splitGenomeByLength *Split an annotatedGenome object into subsets according to gene length*

Description

splitGenomeByLength splits an annotatedGenome according to gene length (bp), which allows estimating the fragment start and length distribution for each subset separately.

Usage

```
splitGenomeByLength(DB, breaks=c(0,3000,5000,Inf))
```

Arguments

DB Object containing annotated genome. Must be of class annotatedGenome, as returned by procGenome or createDenovoGenome.

breaks Breakpoints to define gene subgroups.

Details

By default groups are <3000bp, 3000-5000bp, >5000bp, which work well for the human genome. Further sub-divisions may result in unstable estimates of fragment start and length distributions.

Value

List where each component is of class annotatedGenome.

Author(s)

David Rossell

See Also

procGenome and createDenovoGenome for creating annotatedGenome objects. getDistrs for estimating fragment start and length distribution.

Examples

```
##Not run
## genDB<-makeTranscriptDbFromUCSC(genome="hg19", tablename="refGene")
## hg19DB <- procGenome(genDB, "hg19")
## hg19split <- splitGenomeByLength(hg19DB)
```

subsetGenome	<i>subsetGenome subsets an object of class annotatedGenome for a set of island IDs or chromosome names.</i>
--------------	---

Description

~~ Methods for function subsetGenome in package **casper** ~~ Subset an annotatedGenome object by islands or chromosomes.

Usage

```
subsetGenome(islands, chr, genomeDB)
```

Arguments

islands Vector of characters with the island IDs to retrieve from genome.

chr Vector of characters with the names of chromosomes to retrieve from genome.

genomeDB annotatedGenome object with genome to subset.

Methods

```
signature(islands = "character", chr = "missing", genomeDB = "annotatedGenome")
  Subset annotatedGenome object by a set of island IDs.
```

```
signature(islands = "missing", chr = "character", genomeDB = "annotatedGenome")
  Subset annotatedGenome object by chromosomes.
```

transcripts	<i>Returns all transcripts associated to a given island or entrez id stored in an annotatedGenome object.</i>
-------------	---

Description

annotatedGenome objects store information regarding genes and transcripts. When there's an overlap in exons between several genes, these genes are grouped into gene islands.

The function transcripts retrieves all stored transcripts for a given gene or island.

Usage

```
transcripts(entrezid, islandid, genomeDB)
```

Arguments

entrezid	Character indicating single Entrez identifier. Can be left missing and specify islandid instead.
islandid	Character indicating island identifier. Can be left missing and specify entrezid instead
genomeDB	Object of class annotatedGenome

Value

IRangesList where each element in the list corresponds to a different transcript.

Methods

```
signature(x = "readDistrs") x is an object of type readDistrs, as returned by getDistrs.
  The plot function allows to visualize the fragment length and read start distributions.
```

Examples

```
data(hg19DB)
tx <- transcripts(entrezid="27", genomeDB=hg19DB)
tx
```

txLength *~~ Methods for Function txLength in Package **casper** ~~*

Description

~~ Methods for function txLength in package **casper** ~~ Function to retrieve transcript lengths from annotated genome (class genomeDB).

Usage

```
txLength(islandid, txid, genomeDB)
```

Arguments

islandid	Retrieve length for transcripts in island islandid.
txid	Retrieve length for txid transcripts.
genomeDB	Annotated genome of class genomeDB.

Details

When called for the first time lengths are calculated and stored in the object genomeDB. Subsequent calls refer to these computed values.

Value

Named numeric vector with transcript lengths.

Methods

```
signature(islandid = "character", txid = "missing", genomeDB = "annotatedGenome")  
Retrieve lengths from genomeDB for transcripts in islandid islands.
```

```
signature(islandid = "missing", txid = "character", genomeDB = "annotatedGenome")  
Retrieve lengths from genomeDB for txid transcripts.
```

```
signature(islandid = "missing", txid = "missing", genomeDB = "annotatedGenome")  
Retrieve or calculate lengths for all transcripts in the annotated genome genomeDB.
```

 wrapKnown

Run all necessary steps to analyze a bam file with the casper pipeline.

Description

Function to analyze a bam file to generate expression estimates, read start and fragment length distributions, path counts and processed reads.

Usage

```
wrapKnown(bamFile, verbose=FALSE, seed=1, mc.cores.int=1,
mc.cores=1, genomeDB, readLength, rpkm=TRUE, priorq=2, priorqGeneExpr=2,
citype='none', niter=10^3, burnin=100, keep.pbam=FALSE,
keep.multihits=TRUE, chroms=NULL)
```

Arguments

bamFile	Bam file with the sample to analyze. Index of bam file must be in the same directory.
verbose	Set to TRUE to display progress information.
seed	Set seed of random number generator.
mc.cores.int	Number of cores to use when loading bam files. This is a memory intensive step, therefore number of cores must be chosen according to available RAM memory.
mc.cores	Number of cores to use in expression estimation.
genomeDB	knownGenome object containing annotated genome, as returned by the procGenome function.
readLength	Read length in bp, e.g. in a paired-end experiment where 75bp are sequenced on each end one would set readLength=75.
rpkm	Set to FALSE to return relative expression levels, i.e. the proportion of reads generated from each variant per gene. These proportions add up to 1 for each gene. Set to TRUE to return reads per kilobase per million (RPKM).
priorq	Parameter of the prior distribution on the proportion of reads coming from each variant. The prior is Dirichlet with prior sample size for each variant equal to priorq. We recommend priorq=2 for estimation, as it pools the estimated expression away from 0 and 1 and returned lower estimation errors than priorq=1 in our simulated experiments.
priorqGeneExpr	Parameter for prior distribution on overall gene expression. Defaults to 2, which ensures non-zero estimates for all genes
citype	Set to "none" to return no credibility intervals. Set to "asypm" to return approximate 95% CIs (obtained via the delta method). Set to "exact" to obtain exact CIs via Monte Carlo simulation. Options "asypm" and especially "exact" can increase the computation time substantially.
niter	Number of Monte Carlo iterations. Only used when citype=="exact".

burnin	Number of burnin Monte Carlo iterations. Only used when <code>citype=="exact"</code> .
keep.pbam	Set to TRUE to save processed bam object, as returned by <code>procBam</code> . This object can require substantial memory during execution and disk storage upon saving and is not needed for a default analysis.
keep.multihits	Set to FALSE to discard reads aligned to multiple positions.
chroms	Manually set chromosomes to be processed. By default only main chromosomes are considered (except 'chrM')

Details

Function executes the functions `procBam`, `getDistrs` and `pathCounts` for each chromosome. Once finished, distributions are merged and expressions computed. Distributions, processed reads, path counts and expression estimates are returned.

Value

distr	Object of class <code>readDistrs</code>
pbam	List of objects of class <code>procBam</code> with one element per chromosome
pc	Object of class <code>pathCounts</code>
exp	Object of class <code>expressionSet</code>

Author(s)

Camille Stephan-Otto Attolini

Examples

```
## genDB<-makeTranscriptDbFromUCSC(genome="hg19", tablename="refGene")
## hg19DB <- procGenome(genDB, "hg19")
## bamFile="/path_to_bam/sorted.bam"
## ans <- wrapKnown(bamFile=bamFile, mc.cores.int=4, mc.cores=3, genomeDB=hg19DB, readLength=101)
## names(ans)
## head(exprs(ans$exp))
```

Index

- *Topic **textasciitildeSAM/BAM**
 - rmShortInserts, 29
- *Topic **textasciitilde\textasciitilde**
 - other possible keyword(s)**
 - textasciitilde\textasciitilde**
 - subsetGenome, 33
 - txLength, 35
- *Topic **textasciitildeannotation**
 - procGenome, 27
 - splitGenomeByLength, 32
- *Topic **textasciitildepaired-end sequencing**
 - rmShortInserts, 29
- *Topic **classes**
 - annotatedGenome-class, 2
 - denovoGeneExpr-class, 8
 - denovoGenomeExpr-class, 9
 - modelPriorAS-class, 19
 - pathCounts-class, 21
 - procBam-class, 26
- *Topic **datagen**
 - simMultSamples, 30
 - simReads, 31
- *Topic **datasets**
 - hg19DB, 15
 - K562.r111, 16
- *Topic **hplots**
 - plot, 22
 - plotExpr, 23
 - plotPriorAS, 24
- *Topic **hplot**
 - genePlot, 10
- *Topic **manip**
 - getIsland, 13
 - getReads, 15
 - mergeExp, 17
 - pathCounts, 20
 - procBam, 25
 - transcripts, 34

- *Topic **methods**
 - subsetGenome, 33
- *Topic **models**
 - calcDenovo, 4
 - calcExp, 6
- *Topic **stats**
 - getDistrs, 11
 - modelPrior, 18
- [, denovoGeneExpr, ANY, ANY, ANY-method
(denovoGeneExpr-class), 8
- [, denovoGenomeExpr, ANY, ANY, ANY-method
(denovoGenomeExpr-class), 9
- [, modelPriorAS, ANY, ANY, ANY-method
(modelPriorAS-class), 19
- [[, denovoGeneExpr, ANY, ANY, ANY-method
(denovoGeneExpr-class), 8
- [[, denovoGeneExpr, ANY, ANY-method
(denovoGeneExpr-class), 8
- [[, denovoGenomeExpr-method
(denovoGenomeExpr-class), 9
- annotatedGenome
(annotatedGenome-class), 2
- annotatedGenome-class, 2
- as.list, denovoGenomeExpr-method
(denovoGenomeExpr-class), 9
- calcDenovo, 4
- calcExp, 6
- coef, modelPriorAS, ANY, ANY, ANY-method
(modelPriorAS-class), 19
- createDenovoGenome, 3, 9, 20, 21
- createDenovoGenome (procGenome), 27
- denovoGeneExpr-class, 8
- denovoGenomeExpr-class, 9
- genePlot, 10
- genePlot, CompressedIRangesList, ANY, ANY, ANY-method
(genePlot), 10

- genePlot, GRanges, ANY, ANY, ANY, ANY-method
 (genePlot), 10
 genePlot, GRangesList, ANY, ANY, ANY, ANY-method
 (genePlot), 10
 genePlot, IRanges, ANY, ANY, ANY, ANY-method
 (genePlot), 10
 genePlot, IRangesList, ANY, ANY, ANY, ANY-method
 (genePlot), 10
 genePlot, missing, character, annotatedGenome, GRanges, GRangesList, IRanges, IRangesList-method
 (genePlot), 10
 genePlot, missing, character, annotatedGenome, missing, missing-method
 (genePlot), 10
 genePlot, missing, character, annotatedGenome, procBam, ExpressionSet-method
 (genePlot), 10
 genePlot-methods (genePlot), 10
 getChr (getIsland), 13
 getChr, character, missing, missing, annotatedGenome-method
 (getIsland), 13
 getChr, missing, character, missing, annotatedGenome-method
 (getIsland), 13
 getChr, missing, missing, character, annotatedGenome-method
 (getIsland), 13
 getChr, missing, missing, missing, annotatedGenome-method
 (getIsland), 13
 getChr-method (getIsland), 13
 getDistrs, 11
 getIsland, 13
 getIsland, character, missing, annotatedGenome-method
 (getIsland), 13
 getIsland, missing, character, annotatedGenome-method
 (getIsland), 13
 getNreads, 14
 getNreads, pathCounts-method
 (getNreads), 14
 getReads, 15
 getReads, procBam-method (getReads), 15

 hg19DB, 15

 K562.r111, 16

 lines (plot), 22
 lines, readDistrs-method (plot), 22

 mergeExp, 17
 modelPrior, 18
 modelPriorAS-class, 19

 pathCounts, 20
 pathCounts, list-method (pathCounts), 20
 pathCounts, procBam-method (pathCounts), 20
 pathCounts-class, 21
 pathCounts-method (pathCounts), 20
 plot, 22
 plot, readDistrs, ANY-method (plot), 22
 plotExpr, 23
 plotExpr, denovoGeneExpr-method
 (plotExpr), 23
 plotExpr, missing, missing-method
 (plotExpr), 23
 plotPriorAS, 24
 plotPriorAS, modelPriorAS-method
 (plotPriorAS), 24
 plotPriorAS-methods (plotPriorAS), 24
 posprob (denovoGeneExpr-class), 8
 posProb, denovoGeneExpr-method
 (denovoGeneExpr-class), 8
 procBam, 25
 procBam, list, logical, integer, logical, character-method
 (procBam), 25
 procBam, list, logical, integer, logical, missing-method
 (procBam), 25
 procBam, list, logical, integer, missing, missing-method
 (procBam), 25
 procBam, list, logical, missing, logical, missing-method
 (procBam), 25
 procBam, list, logical, missing, missing, missing-method
 (procBam), 25
 procBam, list, missing, integer, logical, missing-method
 (procBam), 25
 procBam, list, missing, integer, missing, missing-method
 (procBam), 25
 procBam, list, missing, missing, logical, missing-method
 (procBam), 25
 procBam, list, missing, missing, missing, missing-method
 (procBam), 25
 procBam, list, missing, numeric, missing, missing-method
 (procBam), 25
 procBam-class, 26
 procBam-method (procBam), 25
 procGenome, 3, 9, 20, 21, 27
 procGenome, GRanges, ANY-method
 (procGenome), 27
 procGenome, GRanges-method (procGenome), 27
 procGenome, TranscriptDb, ANY-method
 (procGenome), 27

procGenome, TranscriptDb-method
 (procGenome), 27

relativeExpr (calcDenovo), 4

rmShortInserts, 29

show, annotatedGenome-method
 (annotatedGenome-class), 2

show, denovoGeneExpr-method
 (denovoGeneExpr-class), 8

show, denovoGenomeExpr-method
 (denovoGenomeExpr-class), 9

show, modelPriorAS-method
 (modelPriorAS-class), 19

show, pathCounts-method
 (pathCounts-class), 21

show, procBam-method (procBam-class), 26

simMultSamples, 30

simReads, 31

splitGenomeByLength, 32

subsetGenome, 33

subsetGenome, character, missing, annotatedGenome-method
 (subsetGenome), 33

subsetGenome, missing, character, annotatedGenome-method
 (subsetGenome), 33

subsetGenome-methods (subsetGenome), 33

transcripts, 34

transcripts, character, missing, annotatedGenome-method
 (transcripts), 34

transcripts, missing, character, annotatedGenome-method
 (transcripts), 34

txLength, 35

txLength, character, missing, annotatedGenome-method
 (txLength), 35

txLength, missing, character, annotatedGenome-method
 (txLength), 35

txLength, missing, missing, annotatedGenome-method
 (txLength), 35

txLength-methods (txLength), 35

variants (denovoGeneExpr-class), 8

variants, denovoGeneExpr-method
 (denovoGeneExpr-class), 8

variants<- (denovoGeneExpr-class), 8

variants<- , denovoGeneExpr-method
 (denovoGeneExpr-class), 8

wrapKnown, 36