

Package ‘BSgenome’

April 9, 2015

Title Infrastructure for Biostrings-based genome data packages

Description Infrastructure shared by all the Biostrings-based genome data packages

Version 1.34.1

Author Herve Pages

Maintainer H. Pages <hpages@fhcrc.org>

biocViews Genetics, Infrastructure, DataRepresentation, SequenceMatching, Annotation, SNP

Depends R (>= 2.8.0), methods, BiocGenerics (>= 0.1.2), S4Vectors (>= 0.0.7), IRanges (>= 1.99.1), GenomeInfoDb (>= 1.1.4), GenomicRanges (>= 1.17.15), Biostrings (>= 2.33.3), rtracklayer (>= 1.25.8)

Imports methods, stats, BiocGenerics, S4Vectors, IRanges, XVector, GenomeInfoDb, GenomicRanges, Biostrings, Rsamtools, rtracklayer

Suggests BiocInstaller, BSgenome.Celegans.UCSC.ce2 (>= 1.3.11), BSgenome.Hsapiens.UCSC.hg19 (>= 1.3.11), BSgenome.Hsapiens.UCSC.hg19.masked, BSgenome.Rnorvegicus.UCSC.rm5, SNPlocs.Hsapiens.dbSNP.20100427, hgu95av2probe, Biobase, RUnit

License Artistic-2.0

LazyLoad yes

Collate utils.R OnDiskNamedSequences-class.R SNPlocs-class.R InjectSNPsHandler-class.R BSgenome-class.R available.genomes.R injectSNPs.R getSeq-methods.R bsapply.R BSgenome-utils.R export-methods.R GenomeData-class.R GenomeDataList-class.R gdapply.R gdReduce.R BSgenomeForge.R

R topics documented:

available.genomes	2
bsapply	4
BSgenome-class	6

BSgenome-utils	10
BSgenomeForge	12
BSParams-class	14
export-methods	15
gdapply	17
gdReduce	18
GenomeData-class	19
GenomeDataList-class	21
getSeq-methods	22
injectSNPs	27
SNPlocs-class	29

Index	32
--------------	-----------

available.genomes	<i>Find available/installed genomes</i>
-------------------	---

Description

available.genomes gets the list of BSgenome data packages that are available in the Bioconductor repositories for your version of R/Bioconductor.

installed.genomes gets the list of BSgenome data packages that are currently installed on your system.

getBSgenome searches the installed BSgenome data packages for the specified genome and returns it as a [BSgenome](#) object.

Usage

```
available.genomes(splitNameParts=FALSE, type=getOption("pkgType"))
```

```
installed.genomes(splitNameParts=FALSE)
```

```
getBSgenome(genome, masked=FALSE)
```

Arguments

splitNameParts	Whether to split or not the package names in parts. In that case the result is returned in a data frame with 5 columns.
type	Character string indicating the type of package ("source", "mac.binary" or "win.binary") to look for.
genome	A BSgenome object, or the full name of an installed BSgenome data package, or a short string specifying a genome assembly (a.k.a. provider version) that refers unambiguously to an installed BSgenome data package.
masked	TRUE or FALSE. Whether to search for the <i>masked</i> BSgenome object (i.e. the object that contains the masked sequences) or not (the default).

Details

A BSgenome data package contains the full genome sequences for a given organism.

Its name typically has 4 parts (5 parts if it's a *masked* BSgenome data package i.e. if it contains masked sequences) separated by a dot e.g. BSgenome.Mmusculus.UCSC.mm10 or BSgenome.Mmusculus.UCSC.mm10.masked

1. The 1st part is always BSgenome.
2. The 2nd part is the name of the organism in abbreviated form e.g. Mmusculus, Hsapiens, Celegans, Scerevisiae, Ecoli, etc...
3. The 3rd part is the name of the organisation who provided the genome sequences. We formally refer to it as the *provider* of the genome. E.g. UCSC, NCBI, TAIR, etc...
4. The 4th part is the release string or number used by this organisation for this particular genome assembly. We formally refer to it as the *provider version* of the genome. E.g. mm9, mm10, hg18, hg19, GRCh38, susScr3, etc...
5. If the package contains masked sequences, its name has the .masked suffix added to it, which is typically the 5th part.

A BSgenome data package contains a single top-level object (a [BSgenome](#) object) named like the package itself that can be used to access the genome sequences.

Value

For available.genomes and installed.genomes: by default (i.e. if splitNameParts=FALSE), a character vector containing the names of the BSgenome data packages that are available (for available.genomes) or currently installed (for installed.genomes). If splitNameParts=TRUE, the list of packages is returned in a data frame with one row per package and the following columns: pkgname (character), organism (factor), provider (factor), provider_version (character), and masked (logical).

For getBSgenome: the [BSgenome](#) object containing the sequences for the specified genome. Or an error if the object cannot be found in the BSgenome data packages currently installed.

Author(s)

H. Pages

See Also

- [BSgenome](#) objects.
- [available.packages](#).

Examples

```
## -----
## available.genomes() and installed.genomes()
## -----

# What genomes are currently installed:
installed.genomes()
```

```

# What genomes are available:
available.genomes()

# Split the package names in parts:
av_gen <- available.genomes(splitNameParts=TRUE)
table(av_gen$organism)
table(av_gen$provider)

# Make your choice and install with:
library(BiocInstaller)
biocLite("BSgenome.Scerevisiae.UCSC.sacCer1")

# Have a coffee 8-)

# Load the package and display the index of sequences for this genome:
library(BSgenome.Scerevisiae.UCSC.sacCer1)
Scerevisiae # same as BSgenome.Scerevisiae.UCSC.sacCer1

## -----
## getBSgenome()
## -----

## Specify the full name of an installed BSgenome data package:
genome <- getBSgenome("BSgenome.Celegans.UCSC.ce2")
genome

## Specify a genome assembly (a.k.a. provider version):
genome <- getBSgenome("hg19")
class(genome) # BSgenome object
providerVersion(genome)
genome$chrM

genome <- getBSgenome("hg19", masked=TRUE)
class(genome) # MaskedBSgenome object
providerVersion(genome)
genome$chr22

```

bsapply

bsapply

Description

Apply a function to each chromosome in a genome.

Usage

```
bsapply(BSParams, ...)
```

Arguments

BSParams a BSParams object that holds the various parameters needed to configure the bsapply function
... optional arguments to 'FUN'.

Details

By default the exclude parameter is set to not exclude anything. A popular option will probably be to set this to "rand" so that random bits of unassigned contigs are filtered out.

Value

If BSParams sets simplify=FALSE, an ordinary list is returned containing the results generated using the remaining BSParams specifications. If BSParams sets simplify=TRUE, an sapply-like simplification is performed on the results.

Author(s)

Marc Carlson

See Also

[BSParams-class](#), [BSgenome-class](#), [BSgenome-utils](#)

Examples

```
## Load the Worm genome:
library("BSgenome.Celegans.UCSC.ce2")

## Count the alphabet frequencies for every chromosome but exclude
## mitochondrial ones:
params <- new("BSParams", X = Celegans, FUN = alphabetFrequency,
exclude = "M")
bsapply(params)

## Or we can do this same function with simplify = TRUE:
params <- new("BSParams", X = Celegans, FUN = alphabetFrequency,
exclude = "M", simplify = TRUE)
bsapply(params)

## Examples to show how we might look for a string (in this case an
## ebox motif) across the whole genome.
Ebox <- DNASTringSet("CACGTG")
pdict0 <- PDict(Ebox)

params <- new("BSParams", X = Celegans, FUN = countPDict, simplify = TRUE)
bsapply(params, pdict = pdict0)

params@FUN <- matchPDict
bsapply(params, pdict = pdict0)
```

```

## And since its really overkill to use matchPDict to find a single pattern:
params@FUN <- matchPattern
bsapply(params, pattern = "CACGTG")

## Examples on how to use the masks
library("BSgenome.Hsapiens.UCSC.hg19.masked")
genome <- BSgenome.Hsapiens.UCSC.hg19.masked
## I can make things verbose if I want to see the chromosomes getting processed.
options(verbose=TRUE)
## For the 1st example, lets use default masks
params <- new("BSPParams", X = genome, FUN = alphabetFrequency,
exclude = c(1:8,"M","X","random","hap"), simplify = TRUE)
bsapply(params)

if (interactive()) {
  ## Set up the motifList to filter out all double Ts and all double Cs
  params@motifList <-c("TT","CC")
  bsapply(params)

  ## Get rid of the motifList
  params@motifList=as.character()
}

##Enable all standard masks
params@maskList <- c(RM=TRUE,TRF=TRUE)
bsapply(params)

##Disable all standard masks
params@maskList <- c(AGAPS=FALSE,AMB=FALSE)
bsapply(params)

```

BSgenome-class

BSgenome objects

Description

The BSgenome class is a container for storing the full genome sequences of a given organism. BSgenome objects are usually made in advance by a volunteer and made available to the Bioconductor community as "BSgenome data packages". See [?available.genomes](#) for how to get the list of "BSgenome data packages" currently available.

Accessor methods

In the code snippets below, x is a BSgenome object. Note that, because the BSgenome class contains the [GenomeDescription](#) class, then all the accessor methods for [GenomeDescription](#) objects can also be used on x.

- `sourceUrl(x)` Returns the source URL i.e. the permanent URL to the place where the FASTA files used to produce the sequences contained in `x` can be found (and downloaded).
- `seqnames(x)`, `seqnames(x) <- value` Gets or sets the names of the single sequences contained in `x`. Each single sequence is stored in a [DNAStrng](#) or [MaskedDNAStrng](#) object and typically comes from a source file (FASTA) with a single record. The names returned by `seqnames(x)` usually reflect the names of those source files but a common prefix or suffix was eventually removed in order to keep them as short as possible.
- `seqlengths(x)` Returns the lengths of the single sequences contained in `x`.
See [?length,XVector-method](#) and [?length,MaskedXString-method](#) for the definition of the length of a [DNAStrng](#) or [MaskedDNAStrng](#) object. Note that the length of a masked sequence ([MaskedXString](#) object) is not affected by the current set of active masks but the `nchar` method for [MaskedXString](#) objects is.
`names(seqlengths(x))` is guaranteed to be identical to `seqnames(x)`.
- `mseqnames(x)` Returns the index of the multiple sequences contained in `x`. Each multiple sequence is stored in a [DNAStrngSet](#) object and typically comes from a source file (FASTA) with multiple records. The names returned by `mseqnames(x)` usually reflect the names of those source files but a common prefix or suffix was eventually removed in order to keep them as short as possible.
- `names(x)` Returns the index of all sequences contained in `x`. This is the same as `c(seqnames(x), mseqnames(x))`.
- `length(x)` Returns the length of `x`, i.e., the total number of sequences in it (single and multiple sequences). This is the same as `length(names(x))`.
- `x[[name]]` Returns the sequence (single or multiple) in `x` named `name` (`name` must be a single string). No sequence is actually loaded into memory until this is explicitly requested with a call to `x[[name]]` or `x$name`. When loaded, a sequence is kept in a cache. It will be automatically removed from the cache at garbage collection if it's not in use anymore i.e. if there are no reference to it (other than the reference stored in the cache). With `options(verbose=TRUE)`, a message is printed each time a sequence is removed from the cache.
- `x$name` Same as `x[[name]]` but `name` is not evaluated and therefore must be a literal character string or a name (possibly backtick quoted).
- `masknames(x)` The names of the built-in masks that are defined for all the single sequences. There can be up to 4 built-in masks per sequence. These will always be (in this order): (1) the mask of assembly gaps, aka "the AGAPS mask"; (2) the mask of intra-contig ambiguities, aka "the AMB mask"; (3) the mask of repeat regions that were determined by the RepeatMasker software, aka "the RM mask"; (4) the mask of repeat regions that were determined by the Tandem Repeats Finder software (where only repeats with period less than or equal to 12 were kept), aka "the TRF mask". All the single sequences in a given package are guaranteed to have the same collection of built-in masks (same number of masks and in the same order).
`masknames(x)` gives the names of the masks in this collection. Therefore the value returned by `masknames(x)` is a character vector made of the first `N` elements of `c("AGAPS", "AMB", "RM", "TRF")`, where `N` depends only on the BSgenome data package being looked at ($0 \leq N \leq 4$). The man page for most BSgenome data packages should provide the exact list and permanent URLs of the source data files that were used to extract the built-in masks. For example, if you've installed the `BSgenome.Hsapiens.UCSC.hg19` package, load it and see the Note section in [?BSgenome.Hsapiens.UCSC.hg19](#).

Author(s)

H. Pages

See Also

[available.genomes](#), [GenomeDescription-class](#), [BSgenome-utils](#), [DNAString-class](#), [DNAStringSet-class](#), [MaskedDNAString-class](#), [getSeq](#), [BSgenome-method](#), [injectSNPs](#), [subseq](#), [XVector-method](#), [rm](#), [gc](#)

Examples

```
## Loading a BSgenome data package doesnt load its sequences
## into memory:
library(BSgenome.Celegans.UCSC.ce2)

## Number of sequences in this genome:
length(Celegans)

## Display a summary of the sequences:
Celegans

## Index of single sequences:
seqnames(Celegans)

## Lengths (i.e. number of nucleotides) of the single sequences:
seqlengths(Celegans)

## Load chromosome I from disk to memory (hence takes some time)
## and keep a reference to it:
chrI <- Celegans[["chrI"]] # equivalent to Celegans$chrI

chrI

class(chrI) # a DNAString instance
length(chrI) # with 15080483 nucleotides

## Single sequence can be renamed:
seqnames(Celegans) <- sub("^chr", "", seqnames(Celegans))
seqlengths(Celegans)
Celegans$I
seqnames(Celegans) <- paste0("chr", seqnames(Celegans))

## Multiple sequences:
library(BSgenome.Rnorvegicus.UCSC.rn5)
rn5 <- BSgenome.Rnorvegicus.UCSC.rn5
rn5
seqnames(rn5)
rn5_chr1 <- rn5$chr1
mseqnames(rn5)
rn5_random <- Rnorvegicus$random
rn5_random
class(rn5_random) # a DNAStringSet instance
```

```

## Character vector containing the description lines of the first
## 4 sequences in the original FASTA file:
names(rn5_random)[1:4]

## -----
## PASS-BY-ADDRESS SEMANTIC, CACHING AND MEMORY USAGE
## -----

## We want a message to be printed each time a sequence is removed
## from the cache:
options(verbose=TRUE)

gc() # nothing seems to be removed from the cache
rm(rn5_chr1, rn5_random)
gc() # rn5_chr1 and rn5_random are removed from the cache (they are
     # not in use anymore)

options(verbose=FALSE)

## Get the current amount of data in memory (in Mb):
mem0 <- gc()["Vcells", "(Mb)"]

system.time(rn5_chr2 <- rn5$chr2) # read from disk

gc()["Vcells", "(Mb)"] - mem0 # rn5_chr2 occupies 20Mb in memory

system.time(tmp <- rn5$chr2) # much faster! (sequence
                             # is in the cache)

gc()["Vcells", "(Mb)"] - mem0 # were still using 20Mb (sequences
                             # have a pass-by-address semantic
                             # i.e. the sequence data are not
                             # duplicated)

## subseq() doesnt copy the sequence data either, hence it is very
## fast and memory efficient (but the returned object will hold a
## reference to rn5_chr2):
y <- subseq(rn5_chr2, 10, 8000000)
gc()["Vcells", "(Mb)"] - mem0

## We must remove all references to rn5_chr2 before it can be
## removed from the cache (so the 20Mb of memory used by this
## sequence are freed):
options(verbose=TRUE)
rm(rn5_chr2, tmp)
gc()

## Remember that y holds a reference to rn5_chr2 too:
rm(y)
gc()

options(verbose=FALSE)
gc()["Vcells", "(Mb)"] - mem0

```

Description

Utilities for BSgenome objects.

Usage

```
## S4 method for signature BSgenome
matchPWM(pwm, subject, min.score = "80%", exclude = "",
         maskList = logical(0))
## S4 method for signature BSgenome
countPWM(pwm, subject, min.score = "80%", exclude = "",
         maskList = logical(0))
## S4 method for signature BSgenome
vmatchPattern(pattern, subject, max.mismatch = 0, min.mismatch = 0,
             with.indels = FALSE, fixed = TRUE, algorithm = "auto",
             exclude = "", maskList = logical(0), userMask =
             RangesList(), invertUserMask = FALSE)
## S4 method for signature BSgenome
vcountPattern(pattern, subject, max.mismatch = 0, min.mismatch = 0,
             with.indels = FALSE, fixed = TRUE, algorithm = "auto",
             exclude = "", maskList = logical(0), userMask =
             RangesList(), invertUserMask = FALSE)
## S4 method for signature BSgenome
vmatchPDict(pdDict, subject, max.mismatch = 0, min.mismatch = 0,
           fixed = TRUE, algorithm = "auto", verbose = FALSE,
           exclude = "", maskList = logical(0))
## S4 method for signature BSgenome
vcountPDict(pdDict, subject, max.mismatch = 0, min.mismatch = 0,
           fixed = TRUE, algorithm = "auto", collapse = FALSE,
           weight = 1L, verbose = FALSE, exclude = "", maskList = logical(0))
```

Arguments

pwm	A numeric matrix with row names A, C, G and T representing a Position Weight Matrix.
pattern	A DNAStrng object containing the pattern sequence.
pdict	A DNAStrngSet object containing the pattern sequences.
subject	A BSgenome object containing the subject sequences.
min.score	The minimum score for counting a match. Can be given as a character string containing a percentage (e.g. "85%") of the highest possible score or as a single number.

<code>max.mismatch</code> , <code>min.mismatch</code>	The maximum and minimum number of mismatching letters allowed (see ?lowlevel-matching for the details). If non-zero, an inexact matching algorithm is used.
<code>with.indels</code>	If TRUE then indels are allowed. In that case, <code>min.mismatch</code> must be 0 and <code>max.mismatch</code> is interpreted as the maximum "edit distance" allowed between any pattern and any of its matches (see ?matchPattern for the details).
<code>fixed</code>	If FALSE then IUPAC extended letters are interpreted as ambiguities (see ?lowlevel-matching for the details).
<code>algorithm</code>	For <code>vmatchPattern</code> and <code>vcountPattern</code> one of the following: "auto", "naive-exact", "naive-inexact", "boyer-moore", "shift-or", or "indels". For <code>vmatchPDict</code> and <code>vcountPDict</code> one of the following: "auto", "naive-exact", "naive-inexact", "boyer-moore", or "shift-or".
<code>collapse</code> , <code>weight</code>	ignored arguments.
<code>verbose</code>	TRUE or FALSE.
<code>exclude</code>	A character vector with strings that will be used to filter out chromosomes whose names match these strings.
<code>maskList</code>	A named logical vector of <code>maskStates</code> preferred when used with a <code>BSGenome</code> object. When using the <code>bsapply</code> function, the masks will be set to the states in this vector.
<code>userMask</code>	A RangesList , containing a mask to be applied to each chromosome. See bsapply .
<code>invertUserMask</code>	Whether the <code>userMask</code> should be inverted.

Value

A [GRanges](#) object for `matchPWM` with two `elementMetadata` columns: "score" (numeric), and "string" (`DNAStrngSet`).

A [GRanges](#) object for `vmatchPattern`.

A [GRanges](#) object for `vmatchPDict` with one `elementMetadata` column: "index", which represents a mapping to a position in the original pattern dictionary.

A `data.frame` object for `countPWM` and `vcountPattern` with three columns: "seqname" (factor), "strand" (factor), and "count" (integer).

A [DataFrame](#) object for `vcountPDict` with four columns: "seqname" ('factor' Rle), "strand" ('factor' Rle), "index" (integer) and "count" ('integer' Rle). As with `vmatchPDict` the index column represents a mapping to a position in the original pattern dictionary.

Author(s)

P. Aboyoun

See Also

[matchPWM](#), [matchPattern](#), [matchPDict](#), [bsapply](#)

Examples

```

library(BSgenome.Celegans.UCSC.ce2)
data(HNF4alpha)

pwm <- PWM(HNF4alpha)
matchPWM(pwm, Celegans)
countPWM(pwm, Celegans)

pattern <- consensusString(HNF4alpha)
vmatchPattern(pattern, Celegans, fixed = "subject")
vcountPattern(pattern, Celegans, fixed = "subject")

vmatchPDict(HNF4alpha[1:10], Celegans)
vcountPDict(HNF4alpha[1:10], Celegans)

```

BSgenomeForge

The BSgenomeForge functions

Description

A set of functions for making a BSgenome data package.

Usage

```

## Top-level BSgenomeForge function:

forgeBSgenomeDataPkg(x, seqs_srcdir=".", destdir=".", verbose=TRUE)

## Low-level BSgenomeForge functions:

forgeSeqlengthsFile(seqnames, prefix="", suffix=".fa",
                    seqs_srcdir=".", seqs_destdir=".", verbose=TRUE)

forgeSeqFiles(seqnames, mseqnames=NULL,
              seqfile_name=NA, prefix="", suffix=".fa",
              seqs_srcdir=".", seqs_destdir=".",
              ondisk_seq_format=c("2bit", "rda", "fa.rz", "fa"),
              verbose=TRUE)

forgeMasksFiles(seqnames, nmask_per_seq,
                seqs_destdir=".",
                ondisk_seq_format=c("2bit", "rda", "fa.rz", "fa"),
                masks_srcdir=".", masks_destdir=".",
                AGAPSfiles_type="gap", AGAPSfiles_name=NA,
                AGAPSfiles_prefix="", AGAPSfiles_suffix="_gap.txt",
                RMfiles_name=NA, RMfiles_prefix="", RMfiles_suffix=".fa.out",
                TRFfiles_name=NA, TRFfiles_prefix="", TRFfiles_suffix=".bed",
                verbose=TRUE)

```

Arguments

- `x` A `BSgenomeDataPkgSeed` object or the name of a `BSgenome` data package seed file. See the `BSgenomeForge` vignette in this package for more information.
- `seqs_srcdir`, `masks_srcdir` Single strings indicating the path to the source directories i.e. to the directories containing the source data files. Only read access to these directories is needed. See the `BSgenomeForge` vignette in this package for more information.
- `destdir` A single string indicating the path to the directory where the source tree of the target package should be created. This directory must already exist. See the `BSgenomeForge` vignette in this package for more information.
- `ondisk_seq_format` Specifies how the single sequences should be stored in the forged package. Can be `"2bit"`, `"rda"`, `"fa.rz"`, or `"fa"`. If `"2bit"` (the default), then all the single sequences are stored in a single `twoBit` file. If `"rda"`, then each single sequence is stored in a separated serialized `XString` object (one per single sequence). If `"fa.rz"` or `"fa"`, then all the single sequences are stored in a single FASTA file (compressed in the RAZip format if `"fa.rz"`).
- `verbose` TRUE or FALSE.
- `seqnames`, `mseqnames` A character vector containing the names of the single (for `seqnames`) and multiple (for `mseqnames`) sequences to forge. See the `BSgenomeForge` vignette in this package for more information.
- `seqfile_name`, `prefix`, `suffix` See the `BSgenomeForge` vignette in this package for more information, in particular the description of the `seqfile_name`, `seqfiles_prefix` and `seqfiles_suffix` fields of a `BSgenome` data package seed file.
- `seqs_destdir`, `masks_destdir` During the forging process the source data files are converted into serialized `Biostrings` objects. `seqs_destdir` and `masks_destdir` must be single strings indicating the path to the directories where these serialized objects should be saved. These directories must already exist.
`forgeSeqlengthsFile` will produce a single `.rda` file. Both `forgeSeqFiles` and `forgeMasksFiles` will produce one `.rda` file per sequence.
- `nmask_per_seq` A single integer indicating the desired number of masks per sequence. See the `BSgenomeForge` vignette in this package for more information.
- `AGAPSfiles_type`, `AGAPSfiles_name`, `AGAPSfiles_prefix`, `AGAPSfiles_suffix`, `RMfiles_name`, `RMfiles_prefix` These arguments are named accordingly to the corresponding fields of a `BSgenome` data package seed file. See the `BSgenomeForge` vignette in this package for more information.

Details

These functions are intended for Bioconductor users who want to make a new `BSgenome` data package, not for regular users of these packages. See the `BSgenomeForge` vignette in this package (`vignette("BSgenomeForge")`) for an extensive coverage of this topic.

Author(s)

H. Pages

Examples

```

seqs_srcdir <- system.file("extdata", package="BSgenome")
seqnames <- c("chrX", "chrM")

## Forge .rda sequence files:
forgeSeqFiles(seqnames, prefix="ce2", suffix=".fa.gz",
              seqs_srcdir=seqs_srcdir,
              seqs_destdir=tempdir(), ondisk_seq_format="rda")

## Forge .2bit sequence files:
forgeSeqFiles(seqnames, prefix="ce2", suffix=".fa.gz",
              seqs_srcdir=seqs_srcdir,
              seqs_destdir=tempdir(), ondisk_seq_format="2bit")

## Sanity checks:
library(BSgenome.Celegans.UCSC.ce2)
genome <- BSgenome.Celegans.UCSC.ce2

load(file.path(tempdir(), "chrX.rda"))
stopifnot(genome$chrX == chrX)
load(file.path(tempdir(), "chrM.rda"))
stopifnot(genome$chrM == chrM)

ce2_sequences <- import(file.path(tempdir(), "single_sequences.2bit"))
ce2_sequences0 <- DNASTringSet(list(chrX=genome$chrX, chrM=genome$chrM))
stopifnot(identical(names(ce2_sequences0), names(ce2_sequences)) &&
          all(ce2_sequences0 == ce2_sequences))

```

BSPARAMS-class

*Class "BSPARAMS"***Description**

A parameter class for representing all parameters needed for running the bsapply method.

Objects from the Class

Objects can be created by calls of the form `new("BSPARAMS", ...)`.

Slots

X: a BSgenome object that contains chromosomes that you wish to apply FUN on

FUN: the function to apply to each chromosome in the BSgenome object 'X'

exclude: this is a character vector with strings that will be used to filter out chromosomes whose names match these strings.

simplify: TRUE/FALSE value to indicate whether or not the function should try to simplify the output for you.

maskList: A named logical vector of maskStates preferred when used with a BSgenome object. When using the bsapply function, the masks will be set to the states in this vector.

motifList: A character vector which should contain motifs that the user wishes to mask from the sequence.

userMask: A [RangesList](#) object, where each element masks the corresponding chromosome in X. This allows the user to conveniently apply masks besides those included in X.

invertUserMask: A logical indicating whether to invert each mask in userMask.

Methods

`bsapply(p)` Performs the function FUN using the parameters contained within BSParams.

Author(s)

Marc Carlson

See Also

[bsapply](#)

export-methods

Export a BSgenome object as a FASTA or twoBit file

Description

`export` methods for [BSgenome](#) objects.

NOTE: The `export` generic function and most of its methods are defined and documented in the `rtracklayer` package. This man page only documents the 2 `export` methods define in the `BSgenome` package.

Usage

```
## S4 method for signature BSgenome,FastaFile,ANY
export(object, con, format, ...)
## S4 method for signature BSgenome,TwoBitFile,ANY
export(object, con, format, ...)
```

Arguments

`object` The [BSgenome](#) object to export.

con	A FastaFile or TwoBitFile object. Alternatively con can be a single string containing the path to a FASTA or twoBit file, in which case either the file extension or the format argument needs to be "fasta", "twoBit", or "2bit". Also note that in this case, the export method that is called is either the method with signature <code>c("ANY", "character", "missing")</code> or the method with signature <code>c("ANY", "character", "character")</code> , both defined in the rtracklayer package. If object is a BSgenome object and the file extension or the format argument is "fasta", "twoBit", or "2bit", then the flow eventually reaches one of 2 methods documented here.
format	If not missing, should be "fasta", "twoBit", or "2bit" (case insensitive for "twoBit" and "2bit").
...	Extra arguments passed down to other methods. The method for TwoBitFile objects forwards them to bsapply .

Author(s)

Michael Lawrence

See Also

- [BSgenome](#) objects.
- The [export](#) generic, and [FastaFile](#) and [TwoBitFile](#) objects in the [rtracklayer](#) package.

Examples

```
library(BSgenome.Celegans.UCSC.ce2)
genome <- BSgenome.Celegans.UCSC.ce2

## Export as FASTA file.
out1_file <- file.path(tempdir(), "Celegans.fasta")
export(genome, out1_file)

## Export as twoBit file.
out2_file <- file.path(tempdir(), "Celegans.2bit")
export(genome, out2_file)

## Sanity checks:
dna0 <- DNASTringSet(as.list(genome))

system.time(dna1 <- import(out1_file))
stopifnot(identical(names(dna0), names(dna1)) && all(dna0 == dna1))

system.time(dna2 <- import(out2_file)) # importing twoBit is 10-20x
# faster than importing non
# compressed FASTA
stopifnot(identical(names(dna0), names(dna2)) && all(dna0 == dna2))
```

`gdapply`*Applies a function to elements of a `GenomeData`*

Description

WARNING: Starting with BioC 3.0, `GenomeData` and `GenomeDataList` objects are deprecated. Note that the `GenomeData/GenomeDataList` containers predate the `GRanges/GRangesList` containers and, most of the times, the latter can be used instead of the former. Please let us know on the [bioc-devel mailing list](http://bioconductor.org/help/mailling-list/) (<http://bioconductor.org/help/mailling-list/>) if you have a use case where you think there are significant benefits in using `GenomeData/GenomeDataList` over `GRanges/GRangesList`, or if you have questions or concerns about this.

Returns a list of values obtained by applying a function to elements of a `GenomeData` or `GenomeDataList` object.

Usage

```
gdapply(X, FUN, ...)
```

Arguments

<code>X</code>	An object of class <code>GenomeData</code> or <code>GenomeDataList</code> .
<code>FUN</code>	A function to be applied to each chromosome-level sub-element of <code>X</code> .
<code>...</code>	Further arguments; passed to <code>FUN</code>

Value

Typically an object of the same class as `X`.

Author(s)

Deepayan Sarkar

See Also

[GenomeData-class](#), [GenomeDataList-class](#)

 gdReduce

Reduces arguments to a single GenomeData instance

Description

WARNING: Starting with BioC 3.0, GenomeData and GenomeDataList objects are deprecated. Note that the GenomeData/GenomeDataList containers predate the GRanges/GRangesList containers and, most of the times, the latter can be used instead of the former. Please let us know on the bioc-devel mailing list (<http://bioconductor.org/help/mailling-list/>) if you have a use case where you think there are significant benefits in using GenomeData/GenomeDataList over GRanges/GRangesList, or if you have questions or concerns about this.

This function accepts one or more objects that are reduced, with a user-specified function, to a single [GenomeData](#) instance.

Usage

```
gdReduce(f, ..., init, right = FALSE, accumulate = FALSE, gdArgs = list())
```

Arguments

f	An object of class "function", accepting two instances of classes appropriate for the ... arguments, and returning an object suitable for subsequent use in f and incorporation into GenomeData.
...	Objects to be reduced. All objects should be of the same class, as dictated by methods defined on gdReduce. A function to be applied to each chromosome-level sub-element of X.
init	An R object of the same kind as the elements of ...
right	A logical indicating whether to proceed from left to right (default) or right to left.
accumulate	A logical indicating whether the successive reduce combinations should be accumulated. By default, only the final combination is used.
gdArgs	Additional arguments passed to the GenomeData constructor used to assemble the final object.

Details

The gdReduce method for [GenomeData](#) objects successively combines [GenomeData](#) elements of ... using f; all arguments assigned to ... must be of class [GenomeData](#). f is a function accepting two objects returned by "[[" applied to the successive elements of ..., returning a single [GenomeData](#) object to be used in subsequent calls to f. init, right, and accumulate are as described for [Reduce](#). gdArgs can be used to provide metadata information to the constructor used to create the final [GenomeData](#) object.

Currently the gdReduce method for [GenomeDataList](#) objects works when a single [GenomeDataList](#) object x is provided as ... and it does `gdReduce(f, x[[1]], x[[2]] ... x[[N]], init, right, accumulate, gdArgs)` where N is the length of x i.e. the number of [GenomeData](#) objects in it.

Value

An object of class `GenomeData`, containing elements corresponding to the intersection of all named elements of . . . (in the case of the method for `GenomeData` objects) or all elements in the single `GenomeDataList` object passed to it (in the case of the method for `GenomeDataList` objects).

Author(s)

Martin Morgan

See Also

[Reduce](#), [GenomeData-class](#), [GenomeDataList-class](#)

Examples

```
## Not run:
gdReduce
showMethods("gdReduce")

gd <- GenomeData(list(chr1 = IRanges(1, 10), chrX = IRanges(2, 5)),
                 organism = "Mmusculus", provider = "UCSC",
                 providerVersion = "mm9")

gdr <- gdReduce(function(x, y) {
  ## "[[" returns IRanges instances, construct a synthetic version
  IRanges(c(start(x), start(y)), c(end(x), end(y)))
}, GenomeDataList(list(gd, gd[2])))
gdr[["chr1"]]
gdr[["chrX"]]

## End(Not run)
```

GenomeData-class

Data on the genome

Description

WARNING: Starting with BioC 3.0, `GenomeData` and `GenomeDataList` objects are deprecated. Note that the `GenomeData/GenomeDataList` containers predate the [GRanges/GRangesList](#) containers and, most of the times, the latter can be used instead of the former. Please let us know on the `bioc-devel` mailing list (<http://bioconductor.org/help/mailling-list/>) if you have a use case where you think there are significant benefits in using `GenomeData/GenomeDataList` over [GRanges/GRangesList](#), or if you have questions or concerns about this.

`GenomeData` formally represents genomic data as a list, with one element per chromosome in the genome.

Details

This class facilitates storing data on the genome by formalizing a set of metadata fields for storing the organism (e.g. *Mmusculus*), genome build provider (e.g. UCSC), and genome build version (e.g. mm9).

The data is represented as a list, with one element per chromosome (or really any sequence, like a gene). There are no constraints as to the data type of the elements.

Note that as a [SimpleList](#), it is possible to store chromosome-level data (e.g. the lengths) in the `elementMetadata` slot. The organism, provider and providerVersion are all stored in the `SimpleList` metadata, so they may be retrieved in list form by calling `metadata(x)`.

Accessor methods

In the code snippets below, `x` is a `GenomeData` object.

`organism(x)`: Get the single string indicating the organism, if specified, otherwise NULL.

`provider(x)`: Get the single string indicating the genome build provider, if specified, otherwise NULL.

`providerVersion(x)`: Get the single string indicating the genome build version, if specified, otherwise NULL.

Constructor

```
GenomeData(listData = list(), providerVersion = metadata[["providerVersion"]],
  Creates a GenomeData with the elements from the listData parameter, a list. The other arguments correspond to the metadata fields, and, with the exception of elementMetadata, should all be either single strings or NULL (unspecified). Additional global metadata elements may be passed in metadata, in list-form, and via .... The elements in metadata are always overridden by the explicit arguments, like organism and those in .... elementMetadata should be an DataTable or NULL.
```

Coercion

`as(from, "data.frame")`: Coerces each subelement to a data frame, and binds them into a single data frame with an additional column indicating chromosome

`as(from, "RangesList")`: Coerces each subelement to a [Ranges](#) and combines them into a [RangesList](#) with the same names. The “universe” metadata property is set to the `providerVersion` of `from`.

`as(from, "RangedData")`: Coerces each subelement to a [RangedData](#) and combines them into a single `RangedData` with the same names. The “universe” metadata property is set to the `providerVersion` of `from`.

Author(s)

Michael Lawrence

See Also

The [GRanges](#) and [GRangesList](#) classes defined and documented in the **GenomicRanges** package.

[GenomeDataList-class](#), a container for storing a list of `GenomeData` objects and useful e.g. for storing data on multiple samples.

[SimpleList-class](#), the base of this class.

[gdapply](#) for applying a function to elements of a `GenomeData` object.

[gdReduce](#) for successively combining `GenomeData` objects into a single `GenomeData` objects.

Examples

```
## Not run:
gd <- GenomeData(list(chr1 = IRanges(1, 10), chrX = IRanges(2, 5)),
                 organism = "Mmusculus", provider = "UCSC",
                 providerVersion = "mm9")

organism(gd)
providerVersion(gd)
provider(gd)
gd[["chr1"]] # get data for chromosome 1

## End(Not run)
```

`GenomeDataList-class` *List of GenomeData objects*

Description

WARNING: Starting with BioC 3.0, `GenomeData` and `GenomeDataList` objects are deprecated. Note that the `GenomeData/GenomeDataList` containers predate the [GRanges/GRangesList](#) containers and, most of the times, the latter can be used instead of the former. Please let us know on the `bioc-devel` mailing list (<http://bioconductor.org/help/mailling-list/>) if you have a use case where you think there are significant benefits in using `GenomeData/GenomeDataList` over [GRanges/GRangesList](#), or if you have questions or concerns about this.

`GenomeDataList` is a list of [GenomeData](#) objects. It could be useful for storing data on multiple experiments or samples.

Details

This class inherits from [SimpleList](#) and requires that all of its elements to be instances of `GenomeData`.

One should try to take advantage of the metadata storage facilities provided by `SimpleList`. The `elementMetadata` field, for example, could be used to store the experimental design, while the `metadata` field could store the experimental platform.

Constructor

```
GenomeDataList(listData = list(), metadata = list(), elementMetadata = NULL):
```

Creates a `GenomeDataList` with the elements from the `listData` parameter, a list of `GenomeData` instances. The other arguments correspond to the optional metadata stored in [SimpleList](#).

Coercion

`as(from, "data.frame")`: Coerces each subelement to a data frame, and binds them into a single data frame with an additional column indicating chromosome

Author(s)

Michael Lawrence

See Also

The [GRanges](#) and [GRangesList](#) classes defined and documented in the **GenomicRanges** package.

[GenomeData](#), the type of elements stored in this class.

[SimpleList](#)

Examples

```
## Not run:
gd <- GenomeData(list(chr1 = IRanges(1, 10), chrX = IRanges(2, 5)),
                 organism = "Mmusculus", provider = "UCSC",
                 providerVersion = "mm9")
gdl <- GenomeDataList(list(gd), elementMetadata = DataFrame(induced = TRUE))
gdl[[1]] # get first element

## End(Not run)
```

getSeq-methods

getSeq method for BSgenome objects

Description

A [getSeq](#) method for extracting a set of sequences (or subsequences) from a [BSgenome](#) object.

Usage

```
## S4 method for signature BSgenome
getSeq(x, names, start=NA, end=NA, width=NA,
       strand="+", as.character=FALSE)
```

Arguments

x	A BSgenome object. See the available.genomes function for how to install a genome.
names	A character vector containing the names of the sequences in x where to get the subsequences from, or a GRanges object, or a GRangesList object, or a named RangesList object, or a named Ranges object. The RangesList or Ranges object must be named according to the sequences in x where to get the subsequences from.

	If names is missing, then seqnames(x) is used.
	See ?BSgenome-class for details on how to get the lists of single sequences and multiple sequences (respectively) contained in a BSgenome object.
start, end, width	Vector of integers (eventually with NAs) specifying the locations of the subsequences to extract. These are not needed (and it's an error to supply them) when names is a GRanges , GRangesList , RangesList , or Ranges object.
strand	A vector containing "+"s or/and "-"s. This is not needed (and it's an error to supply it) when names is a GRanges or GRangesList object.
as.character	TRUE or FALSE. Should the extracted sequences be returned in a standard character vector?
...	Additional arguments. (Currently ignored.)

Details

L, the number of sequences to extract, is determined as follow:

- If names is a [GRanges](#) or [Ranges](#) object then $L = \text{length}(\text{names})$.
- If names is a [GRangesList](#) or [RangesList](#) object then $L = \text{length}(\text{unlist}(\text{names}))$.
- Otherwise, L is the length of the longest of names, start, end and width and all these arguments are recycled to this length. NAs and negative values in these 3 arguments are solved according to the rules of the SEW (Start/End/Width) interface (see [?solveUserSEW](#) for the details).

If names is neither a [GRanges](#) or [GRangesList](#) object, then the strand argument is also recycled to length L.

Here is how the names passed to the names argument are matched to the names of the sequences in [BSgenome](#) object x. For each name in names:

- (1): If x contains a single sequence with that name then this sequence is used for extraction;
- (2): Otherwise the names of all the elements in all the multiple sequences are searched. If the names argument is a character vector then name is treated as a regular expression and [grep](#) is used for this search, otherwise (i.e. when the names are supplied via a higher level object like [GRanges](#) or [GRangesList](#)) then name must match exactly the name of the sequence. If exactly 1 sequence is found, then it is used for extraction, otherwise (i.e. if no sequence or more than 1 sequence is found) then an error is raised.

Value

Normally a [DNASTringSet](#) object (or character vector if as.character=TRUE).

With the 2 following exceptions:

1. A [DNASTringSetList](#) object (or [CharacterList](#) object if as.character=TRUE) of the same shape as names if names is a [GRangesList](#) object.
2. A [DNASTring](#) object (or single character string if as.character=TRUE) if $L = 1$ and names is not a [GRanges](#), [GRangesList](#), [RangesList](#), or [Ranges](#) object.

Note

Be aware that using `as.character=TRUE` can be very inefficient when extracting a "big" amount of DNA sequences (e.g. millions of short sequences or a small number of very long sequences).

Note that the masks in `x`, if any, are always ignored. In other words, masked regions in the genome are extracted in the same way as unmasked regions (this is achieved by dropping the masks before extraction). See [?MaskedDNAString-class](#) for more information about masked DNA sequences.

Author(s)

H. Pages; improvements suggested by Matt Settles and others

See Also

[getSeq](#), [available.genomes](#), [BSgenome-class](#), [DNAString-class](#), [DNAStringSet-class](#), [MaskedDNAString-class](#), [GRanges-class](#), [GRangesList-class](#), [RangesList-class](#), [Ranges-class](#), [grep](#)

Examples

```
## -----
## A. SIMPLE EXAMPLES
## -----

## Load the Caenorhabditis elegans genome (UCSC Release ce2):
library(BSgenome.Celegans.UCSC.ce2)

## Look at the index of sequences:
Celegans

## Get chromosome V as a DNAString object:
getSeq(Celegans, "chrV")
## which is in fact the same as doing:
Celegans$chrV

## Not run:
## Never try this:
getSeq(Celegans, "chrV", as.character=TRUE)
## or this (even worse):
getSeq(Celegans, as.character=TRUE)

## End(Not run)

## Get the first 20 bases of each chromosome:
getSeq(Celegans, end=20)

## Get the last 20 bases of each chromosome:
getSeq(Celegans, start=-20)

## -----
## B. EXTRACTING SMALL SEQUENCES FROM DIFFERENT CHROMOSOMES
## -----
```

```

myseqs <- data.frame(
  chr=c("chrI", "chrX", "chrM", "chrM", "chrX", "chrI", "chrM", "chrI"),
  start=c(NA, -40, 8510, 301, 30001, 9220500, -2804, -30),
  end=c(50, NA, 8522, 324, 30011, 9220555, -2801, -11),
  strand=c("+", "-", "+", "+", "-", "-", "+", "-")
)
getSeq(Celegans, myseqs$chr,
       start=myseqs$start, end=myseqs$end)
getSeq(Celegans, myseqs$chr,
       start=myseqs$start, end=myseqs$end, strand=myseqs$strand)

## -----
## C. USING A GRanges OBJECT
## -----

gr1 <- GRanges(seqnames=c("chrI", "chrI", "chrM"),
               ranges=IRanges(start=101:103, width=9))
gr1 # all strand values are "*"
getSeq(Celegans, gr1) # treats strand values as if they were "+"

strand(gr1)[1] <- "-"
getSeq(Celegans, gr1)

strand(gr1)[1] <- "+"
getSeq(Celegans, gr1)

strand(gr1)[2] <- "*"
if (interactive())
  getSeq(Celegans, gr1) # Error: cannot mix "*" with other strand values

gr2 <- GRanges(seqnames=c("chrM", "NM_058280_up_1000"),
               ranges=IRanges(start=103:102, width=9))
gr2
if (interactive()) {
  ## Because the sequence names are supplied via a GRanges object, they
  ## are not treated as regular expressions:
  getSeq(Celegans, gr2) # Error: sequence NM_058280_up_1000 not found
}

## -----
## D. USING A GRangesList OBJECT
## -----

gr1 <- GRanges(seqnames=c("chrI", "chrII", "chrM", "chrII"),
               ranges=IRanges(start=101:104, width=12),
               strand="+")
gr2 <- shift(gr1, 5)
gr3 <- gr2
strand(gr3) <- "-"

gr1 <- GRangesList(gr1, gr2, gr3)
getSeq(Celegans, gr1)

```

```

## -----
## E. EXTRACTING A HIGH NUMBER OF RANDOM 40-MERS FROM A GENOME
## -----

extractRandomReads <- function(x, density, readlength)
{
  if (!is.integer(readlength))
    readlength <- as.integer(readlength)
  start <- lapply(seqnames(x),
                 function(name)
                 {
                   seqlength <- seqlengths(x)[name]
                   sample(seqlength - readlength + 1L,
                        seqlength * density,
                        replace=TRUE)
                 })
  names <- rep.int(seqnames(x), elementLengths(start))
  ranges <- IRanges(start=unlist(start), width=readlength)
  strand <- strand(sample(c("+", "-"), length(names), replace=TRUE))
  gr <- GRanges(seqnames=names, ranges=ranges, strand=strand)
  getSeq(x, gr)
}

## With a density of 1 read every 100 genome bases, the total number of
## extracted 40-mers is about 1 million:
rndreads <- extractRandomReads(Celegans, 0.01, 40)

## Notes:
## - The short sequences in rndreads can be seen as the result of a
##   simulated high-throughput sequencing experiment. A non-realistic
##   one though because:
##   (a) It assumes that the underlying technology is perfect (the
##       generated reads have no technology induced errors).
##   (b) It assumes that the sequenced genome is exactly the same as
##       the reference genome.
##   (c) The simulated reads can contain IUPAC ambiguity letters only
##       because the reference genome contains them. In a real
##       high-throughput sequencing experiment, the sequenced genome
##       of course doesnt contain those letters, but the sequencer
##       can introduce them in the generated reads to indicate
##       ambiguous base-calling.
## - Those reads are coming from the plus and minus strands of the
##   chromosomes.
## - With a density of 0.01 and the reads being only 40-base long, the
##   average coverage of the genome is only 0.4 which is low. The total
##   number of reads is about 1 million and it takes less than 10 sec.
##   to generate them.
## - A higher coverage can be achieved by using a higher density and/or
##   longer reads. For example, with a density of 0.1 and 100-base reads
##   the average coverage is 10. The total number of reads is about 10
##   millions and it takes less than 1 minute to generate them.
## - Those reads could easily be mapped back to the reference by using
##   an efficient matching tool like matchPDict() for performing exact

```

```

## matching (see ?matchPDict for more information). Typically, a
## small percentage of the reads (4 to 5% in our case) will hit the
## reference at multiple locations. This is especially true for such
## short reads, and, in a lower proportion, is still true for longer
## reads, even for reads as long as 300 bases.

## -----
## F. SEE THE BSgenome CACHE IN ACTION
## -----

options(verbose=TRUE)
first20 <- getSeq(Celegans, end=20)
first20
gc()
stopifnot(length(ls(Celegans@.seqs_cache)) == 0L)
## One more gc() call is needed in order to see the amount of memory in
## used after all the chromosomes have been removed from the cache:
gc()

```

injectSNPs

SNP injection

Description

Inject SNPs from a SNPlocs data package into a genome.

Usage

```

injectSNPs(x, snps)

SNPlocs_pkgname(x)

## S4 method for signature BSgenome
snpcount(x)
## S4 method for signature BSgenome
snplocs(x, seqname, ...)

## Related utilities
available.SNPs(type=getOption("pkgType"))
installed.SNPs()

```

Arguments

x	A BSgenome object.
snps	A SNPlocs object or the name of a SNPlocs data package. This object or package must contain SNP information for the single sequences contained in x. If a package, it must be already installed (injectSNPs won't try to install it).
seqname	The name of a single sequence in x.

type	Character string indicating the type of package ("source", "mac.binary" or "win.binary") to look for.
...	Further arguments to be passed to snplocs method for SNPlocs objects.

Value

`injectSNPs` returns a copy of the original genome `x` where some or all of the single sequences from `x` are altered by injecting the SNPs stored in `snps`. The SNPs in the altered genome are represented by an IUPAC ambiguity code at each SNP location.

`SNPlocs_pkgname`, `snpcount` and `snplocs` return NULL if no SNPs were injected in `x` (i.e. if `x` is not a [BSgenome](#) object returned by a previous call to `injectSNPs`). Otherwise `SNPlocs_pkgname` returns the name of the package from which the SNPs were injected, `snpcount` the number of SNPs for each altered sequence in `x`, and `snplocs` their locations in the sequence whose name is specified by `seqname`.

`available.SNPs` returns a character vector containing the names of the [SNPlocs](#) data packages that are currently available on the Bioconductor repositories for your version of R/Bioconductor. A [SNPlocs](#) data package contains basic SNP information (location and alleles) for a given organism.

`installed.SNPs` returns a character vector containing the names of the [SNPlocs](#) data packages that are already installed.

Note

`injectSNPs`, `SNPlocs_pkgname`, `snpcount` and `snplocs` have the side effect to try to load the [SNPlocs](#) data package that was specified thru the `snps` argument if it's not already loaded.

Author(s)

H. Pages

See Also

[BSgenome-class](#), [IUPAC_CODE_MAP](#), [injectHardMask](#), [letterFrequencyInSlidingView](#), [.inplaceReplaceLetterAt](#)

Examples

```
## What SNPlocs data packages are already installed:
installed.SNPs()

## What SNPlocs data packages are available:
available.SNPs()

if (interactive()) {
  ## Make your choice and install with:
  source("http://bioconductor.org/biocLite.R")
  biocLite("SNPlocs.Hsapiens.dbSNP.20100427")
}

## Inject SNPs from dbSNP into the Human genome:
library(BSgenome.Hsapiens.UCSC.hg19.masked)
genome <- BSgenome.Hsapiens.UCSC.hg19.masked
```

```

SNPlocs_pkgname(genome)

genome2 <- injectSNPs(genome, "SNPlocs.Hsapiens.dbSNP.20100427")
genome2 # note the extra "with SNPs injected from ..." line
SNPlocs_pkgname(genome2)
snpcount(genome2)
head(snplocs(genome2, "chr1"))

alphabetFrequency(genome$chr1)
alphabetFrequency(genome2$chr1)

## Find runs of SNPs of length at least 25 in chr1. Might require
## more memory than some platforms can handle (e.g. 32-bit Windows
## and maybe some Mac OS X machines with little memory):
is_32bit_windows <- .Platform$OS.type == "windows" &&
  .Platform$r_arch == "i386"
is_macosx <- substr(R.version$os, start=1, stop=6) == "darwin"
if (!is_32bit_windows && !is_macosx) {
  chr1 <- injectHardMask(genome2$chr1)
  ambiguous_letters <- paste(DNA_ALPHABET[5:15], collapse="")
  lf <- letterFrequencyInSlidingView(chr1, 25, ambiguous_letters)
  sl <- slice(as.integer(lf), lower=25)
  v1 <- Views(chr1, start(sl), end(sl)+24)
  v1
  max(width(v1)) # length of longest SNP run
}

```

SNPlocs-class

SNPlocs objects

Description

The SNPlocs class is a container for storing known SNP locations for a given organism. SNPlocs objects are usually made in advance by a volunteer and made available to the Bioconductor community as "SNPlocs data packages". See [?available.SNPs](#) for how to get the list of "SNPlocs data packages" currently available.

This man page's main focus is on how to extract information from a SNPlocs object.

Usage

```

snpcount(x)

snplocs(x, seqname, ...)
## S4 method for signature SNPlocs
snplocs(x, seqname, as.GRanges=FALSE, caching=TRUE)

snpid2loc(x, snpid, ...)
## S4 method for signature SNPlocs
snpid2loc(x, snpid, caching=TRUE)

```

```

snpid2alleles(x, snpid, ...)
## S4 method for signature SNPlocs
snpid2alleles(x, snpid, caching=TRUE)

snpid2grange(x, snpid, ...)
## S4 method for signature SNPlocs
snpid2grange(x, snpid, caching=TRUE)

```

Arguments

x	A SNPlocs object.
seqname	The name of the sequence for which to get the SNP locations and alleles. If as.GRanges is FALSE, only one sequence can be specified (i.e. seqname must be a single string). If as.GRanges is TRUE, an arbitrary number of sequences can be specified (i.e. seqname can be a character vector of arbitrary length).
as.GRanges	TRUE or FALSE. If TRUE, then the SNP locations and alleles are returned in a GRanges object. Otherwise (the default), they are returned in a data frame (see below).
caching	Should the loaded SNPs be cached in memory for faster further retrieval but at the cost of increased memory usage?
snpid	The SNP ids to look up (e.g. rs ids). Can be integer or character vector, with or without the "rs" prefix. NAs are not allowed.
...	Additional arguments, for use in specific methods.

Value

snpcount returns a named integer vector containing the number of SNPs for each sequence in the reference genome.

By default (i.e. when as.GRanges=FALSE), snplocs returns a data frame with 1 row per SNP and the following columns:

1. RefSNP_id: RefSNP ID (aka "rs id") with "rs" prefix removed. Character vector with no NAs and no duplicates.
2. alleles_as_ambig: A character vector with no NAs containing the alleles for each SNP represented by an IUPAC nucleotide ambiguity code. See [?IUPAC_CODE_MAP](#) in the **Biostrings** package for more information.
3. loc: The 1-based location of the SNP relative to the first base at the 5' end of the plus strand of the reference sequence.

Otherwise (i.e. when as.GRanges=TRUE), it returns a [GRanges](#) object with extra columns "RefSNP_id" and "alleles_as_ambig". Note that all the elements (genomic ranges) in this [GRanges](#) object have their strand set to "+" and that all the sequence lengths are set to NA.

snpid2loc and snpid2alleles both return a named vector (integer vector for the former, character vector for the latter) where each (name, value) pair corresponds to a supplied SNP id. For both functions the name in (name, value) is the chromosome of the SNP id. The value in (name, value) is

the position of the SNP id on the chromosome for `snpid2loc`, and a single IUPAC code representing the associated alleles for `snpid2alleles`.

`snpid2grange` returns a [GRanges](#) object similar to the one returned by `snplocs` (when used with `as.GRanges=TRUE`) and where each element corresponds to a supplied SNP id.

Author(s)

H. Pages

See Also

- [available.SNPs](#)
- [injectSNPs](#)
- [IUPAC_CODE_MAP](#) in the **Biostrings** package.

Examples

```
## COMING SOON!
```

Index

*Topic **classes**

- BSgenome-class, 6
- BSPParams-class, 14
- GenomeData-class, 19
- GenomeDataList-class, 21
- SNPlocs-class, 29

*Topic **manip**

- available.genomes, 2
- bsapply, 4
- BSgenomeForge, 12
- gdapply, 17
- gdReduce, 18
- getSeq-methods, 22
- injectSNPs, 27

*Topic **methods**

- BSgenome-class, 6
- BSgenome-utils, 10
- export-methods, 15
- GenomeData-class, 19
- GenomeDataList-class, 21
- SNPlocs-class, 29

*Topic **utilities**

- BSgenome-utils, 10
- export-methods, 15
- .inplaceReplaceLetterAt, 28
- [[, BSgenome-method (BSgenome-class), 6
- [[<-, BSgenome-method (BSgenome-class), 6
- \$, BSgenome-method (BSgenome-class), 6
- as.list, BSgenome-method (BSgenome-class), 6
- available.genomes, 2, 6, 8, 22, 24
- available.packages, 3
- available.SNPs, 29, 31
- available.SNPs (injectSNPs), 27
- bsapply, 4, 11, 15, 16
- BSgenome, 2, 3, 10, 15, 16, 22, 23, 27, 28
- BSgenome (BSgenome-class), 6
- BSgenome-class, 5, 6, 24, 28

- BSgenome-utils, 5, 8, 10
- BSgenome.Hsapiens.UCSC.hg19, 7
- BSgenomeDataPkgSeed (BSgenomeForge), 12
- BSgenomeDataPkgSeed-class (BSgenomeForge), 12
- BSgenomeForge, 12
- BSPParams (BSPParams-class), 14
- BSPParams-class, 5, 14
- CharacterList, 23
- class:BSgenome (BSgenome-class), 6
- class:BSgenomeDataPkgSeed (BSgenomeForge), 12
- class:BSPParams (BSPParams-class), 14
- class:InjectSNPsHandler (injectSNPs), 27
- class:SNPlocs (SNPlocs-class), 29
- coerce, GenomeData, data.frame-method (GenomeData-class), 19
- coerce, GenomeData, RangedData-method (GenomeData-class), 19
- coerce, GenomeData, RangesList-method (GenomeData-class), 19
- coerce, GenomeDataList, data.frame-method (GenomeDataList-class), 21
- compatibleGenomes (SNPlocs-class), 29
- compatibleGenomes, SNPlocs-method (SNPlocs-class), 29
- countPWM, BSgenome-method (BSgenome-utils), 10
- DataFrame, 11
- DataTable, 20
- DNAStrng, 7, 10, 23
- DNAStrng-class, 8, 24
- DNAStrngSet, 7, 10, 23
- DNAStrngSet-class, 8, 24
- DNAStrngSetList, 23
- export, 15, 16

- export, BSgenome, FastaFile, ANY-method (export-methods), 15
- export, BSgenome, TwoBitFile, ANY-method (export-methods), 15
- export-methods, 15
- FastaFile, 16
- forgeBSgenomeDataPkg (BSgenomeForge), 12
- forgeBSgenomeDataPkg, BSgenomeDataPkgSeed-method (BSgenomeForge), 12
- forgeBSgenomeDataPkg, character-method (BSgenomeForge), 12
- forgeBSgenomeDataPkg, list-method (BSgenomeForge), 12
- forgeMasksFiles (BSgenomeForge), 12
- forgeSeqFiles (BSgenomeForge), 12
- forgeSeqlengthsFile (BSgenomeForge), 12
- gc, 8
- gdapply, 17, 21
- gdapply, GenomeData, function-method (gdapply), 17
- gdapply, GenomeDataList, function-method (gdapply), 17
- gdReduce, 18, 21
- gdReduce, GenomeData-method (gdReduce), 18
- gdReduce, GenomeDataList-method (gdReduce), 18
- GenomeData, 17–19, 21, 22
- GenomeData (GenomeData-class), 19
- GenomeData-class, 17, 19, 19
- GenomeDataList, 17–19
- GenomeDataList (GenomeDataList-class), 21
- GenomeDataList-class, 17, 19, 21, 21
- GenomeDescription, 6
- GenomeDescription-class, 8
- getBSgenome (available.genomes), 2
- getSeq, 22, 24
- getSeq, BSgenome-method (getSeq-methods), 22
- getSeq-methods, 22
- GRanges, 11, 17–19, 21–23, 30, 31
- GRanges-class, 24
- GRangesList, 17–19, 21–23
- GRangesList-class, 24
- grep, 23, 24
- injectHardMask, 28
- injectSNPs, 8, 27, 31
- injectSNPs, BSgenome-method (injectSNPs), 27
- InjectSNPsHandler (injectSNPs), 27
- InjectSNPsHandler-class (injectSNPs), 27
- installed.genomes (available.genomes), 2
- installed.SNPs (injectSNPs), 27
- LUPAC_CODE_MAP, 28, 30, 31
- length, BSgenome-method (BSgenome-class), 6
- letterFrequencyInSlidingView, 28
- MaskedDNAString, 7
- MaskedDNAString-class, 8, 24
- MaskedXString, 7
- masknames (BSgenome-class), 6
- masknames, BSgenome-method (BSgenome-class), 6
- matchPattern, 11
- matchPDict, 11
- matchPWM, 11
- matchPWM, BSgenome-method (BSgenome-utils), 10
- mseqnames (BSgenome-class), 6
- mseqnames, BSgenome-method (BSgenome-class), 6
- names, BSgenome-method (BSgenome-class), 6
- newSNPlocs (SNPlocs-class), 29
- organism, GenomeData-method (GenomeData-class), 19
- organism, SNPlocs-method (SNPlocs-class), 29
- provider, GenomeData-method (GenomeData-class), 19
- provider, SNPlocs-method (SNPlocs-class), 29
- providerVersion, GenomeData-method (GenomeData-class), 19
- providerVersion, SNPlocs-method (SNPlocs-class), 29
- RangedData, 20
- Ranges, 20, 22, 23
- Ranges-class, 24

- RangesList, [11](#), [15](#), [20](#), [22](#), [23](#)
- RangesList-class, [24](#)
- Reduce, [18](#), [19](#)
- referenceGenome (SNPlocs-class), [29](#)
- referenceGenome, SNPlocs-method (SNPlocs-class), [29](#)
- releaseDate, SNPlocs-method (SNPlocs-class), [29](#)
- releaseName, SNPlocs-method (SNPlocs-class), [29](#)
- rm, [8](#)
- seqinfo, BSgenome-method (BSgenome-class), [6](#)
- seqinfo, SNPlocs-method (SNPlocs-class), [29](#)
- seqinfo<-, BSgenome-method (BSgenome-class), [6](#)
- seqnames, SNPlocs-method (SNPlocs-class), [29](#)
- seqnames<-, BSgenome-method (BSgenome-class), [6](#)
- show, BSgenome-method (BSgenome-class), [6](#)
- show, GenomeData-method (GenomeData-class), [19](#)
- show, SNPlocs-method (SNPlocs-class), [29](#)
- SimpleList, [20](#)–[22](#)
- SimpleList-class, [21](#)
- SNPcount (injectSNPs), [27](#)
- snpcount (SNPlocs-class), [29](#)
- SNPcount, BSgenome-method (injectSNPs), [27](#)
- snpcount, BSgenome-method (injectSNPs), [27](#)
- SNPcount, InjectSNPsHandler-method (injectSNPs), [27](#)
- snpcount, InjectSNPsHandler-method (injectSNPs), [27](#)
- snpcount, SNPlocs-method (SNPlocs-class), [29](#)
- snpid2alleles (SNPlocs-class), [29](#)
- snpid2alleles, SNPlocs-method (SNPlocs-class), [29](#)
- snpid2grange (SNPlocs-class), [29](#)
- snpid2grange, SNPlocs-method (SNPlocs-class), [29](#)
- snpid2loc (SNPlocs-class), [29](#)
- snpid2loc, SNPlocs-method (SNPlocs-class), [29](#)
- SNPlocs, [27](#), [28](#)
- SNPlocs (SNPlocs-class), [29](#)
- snplocs, [28](#)
- snplocs (SNPlocs-class), [29](#)
- SNPlocs, BSgenome-method (injectSNPs), [27](#)
- snplocs, BSgenome-method (injectSNPs), [27](#)
- SNPlocs, InjectSNPsHandler-method (injectSNPs), [27](#)
- snplocs, InjectSNPsHandler-method (injectSNPs), [27](#)
- snplocs, SNPlocs-method (SNPlocs-class), [29](#)
- SNPlocs-class, [29](#)
- SNPlocs_pkgname (injectSNPs), [27](#)
- SNPlocs_pkgname, BSgenome-method (injectSNPs), [27](#)
- SNPlocs_pkgname, InjectSNPsHandler-method (injectSNPs), [27](#)
- solveUserSEW, [23](#)
- sourceUrl (BSgenome-class), [6](#)
- sourceUrl, BSgenome-method (BSgenome-class), [6](#)
- species, SNPlocs-method (SNPlocs-class), [29](#)
- subseq, XVector-method, [8](#)
- TwoBitFile, [16](#)
- vcountPattern, BSgenome-method (BSgenome-utils), [10](#)
- vcountPDict, BSgenome-method (BSgenome-utils), [10](#)
- vmatchPattern, BSgenome-method (BSgenome-utils), [10](#)
- vmatchPDict, BSgenome-method (BSgenome-utils), [10](#)
- XString, [13](#)