

# Package ‘R453Plus1Toolbox’

March 26, 2013

**Type** Package

**Title** A package for importing and analyzing data from Roche’s Genome Sequencer System.

**Version** 1.8.1

**Date** 2012-08-16

**Author** Hans-Ulrich Klein, Christoph Bartenhagen, Christian Ruckert

**Maintainer** Hans-Ulrich Klein <h.klein@uni-muenster.de>

## Imports

BiocGenerics (>= 0.1.3), Biobase (>= 2.15.1), biomaRt, Biostrings, BSgenome, IRanges (>= 1.13.10), methods, R2HTML, Rsamtools, rtracklayer, ShortRead, VariantAnnotation

**Depends** R (>= 2.12.0), BiocGener-

ics, Biobase, Biostrings, BSgenome.Scerevisiae, UCSC.sacCer2, TeachingDemos

**Suggests** rtracklayer, ShortRead, Rsamtools, BSgenome.Hsapiens, UCSC.hg19

**Description** The R453Plus1 Toolbox comprises useful functions for the analysis of data generated by Roche’s 454 sequencing platform. It adds functions for quality assurance as well as for annotation and visualization of detected variants, complementing the software tools shipped by Roche with their product. Further, a pipeline for the detection of structural variants is provided.

**License** LGPL-3

**biocViews** HighThroughputSequencing, Infrastructure, DataImport, DataRepresentation, Visualization, QualityControl, ReportWriting

## R topics documented:

alignShortReads . . . . .	3
AnnotatedVariants-class . . . . .	4
annotateVariants . . . . .	5
assayDataAmp . . . . .	6
ava2vcf . . . . .	7
AVASet . . . . .	7
AVASet-class . . . . .	9
avaSetExample . . . . .	12
avaSetFiltered . . . . .	13

avaSetFiltered_annot	13
baseFrequency	14
baseQualityHist	15
baseQualityStats	15
breakpoints	16
Breakpoints-class	16
calculateTiTv	18
captureArray	19
complexity.dust	20
complexity.entropy	21
convertCigar	22
coverageOnTarget	22
demultiplexReads	23
detectBreakpoints	24
dinucleotideOddsRatio	27
fDataAmp	28
featureDataAmp	28
filterChimericReads	29
flowgramBarplot	31
gcContent	31
gcContentHist	32
gcPerPosition	32
genomeSequencerMIDs	33
getAlignedReads	34
getAminoAbbr	35
getVariantPercentages	35
homopolymerHist	36
htmlReport	37
MapperSet	38
MapperSet-class	40
mapperSetExample	41
mergeBreakpoints	42
mutationInfo	43
nucleotideCharts	44
plotAmpliconCoverage	44
plotChimericReads	45
plotVariants	47
plotVariationFrequency	49
positionQualityBoxplot	51
readLengthHist	51
readLengthStats	52
readSff	52
readsOnTarget	53
referenceSequences	54
regions	54
removeLinker	55
sequenceCaptureLinkers	56
sequenceQualityHist	57
setVariantFilter	57
sff2fastq	58
SFFContainer-class	59
SFFRead-class	61

<i>alignShortReads</i>	3
variants . . . . .	63
<b>Index</b>	<b>64</b>

*alignShortReads*                      *Exact alignment of DNA sequences against a reference*

### Description

This method aligns given sequences against a given reference genome using the `matchPDict` method. Only exact (no errors) and unique matches are returned.

### Usage

```
alignShortReads(object, bsGenome, seqNames, ensemblNotation)
```

### Arguments

<code>object</code>	The reads that should be aligned against either as a <code>DNAStringSet</code> or a <code>AVASet</code> instance. In the latter case the reference sequences are extracted and aligned.
<code>bsGenome</code>	A <code>bsGenome</code> instance providing the reference sequences.
<code>seqNames</code>	The names of the sequences in <code>bsGenome</code> that should be used. If omitted, all reference sequences are used.
<code>ensemblNotation</code>	If set to <code>TRUE</code> , "chr" is removed from the reference sequences' names in the returned alignment. Default value is <code>FALSE</code> .

### Details

All reads are aligned against the reference and its reverse complement. If the reads are not in 5' to 3' orientation, they should be reversed before. Note that only exact and unique alignments are reported. Use `matchPDict` directly for more flexibility.

### Value

An object of class `AlignedRead` or a `AVASet` instance.

### Author(s)

Hans-Ulrich Klein

### See Also

[matchPDict](#), [DNAStringSet](#), [AlignedRead](#), [AVASet](#)

### Examples

```
library("BSgenome.Scerevisiae.UCSC.sacCer2")
reads = DNAStringSet(c(
  "CCGTTCAAAGAGCCCTTGGCCCATAATCCACCGGTT",
  "ATCCTGCCACAGGAGTCCATGGAGGTTTCGCCA"))
alignShortReads(reads, Scerevisiae, seqNames="chrIII")
```

---

AnnotatedVariants-class *Class "AnnotatedVariants"*

---

## Description

A class for storing annotation about variants. An object of this class is returned by the method `annotateVariants`. The class has not been designed to be created by users directly.

## Details

The list encapsulated by this class has one element for each variant. Each element is a nested list with the elements `genes`, `transcripts`, `exons` and `snps`. All these elements are data frames listing `genes`, `transcripts`, `exons` or `snps` respectively that were affected by the variant. Use the example below to explore the data frames' contents.

## Objects from the Class

Objects can be created by calls of the form `new("AnnotatedVariants")`. The method `annotateVariants` returns `AnnotatedVariants`-objects.

## Slots

`annotatedVariants`: Object of class "list" with one entry for each variant.

## Methods

**`annotatedVariants`** signature(`object` = "AnnotatedVariants"): Get the list with variants.

**`annotatedVariants<-`** signature(`object` = "AnnotatedVariants", `value` = "list"): Set a new list with variants.

**`names`** signature(`x` = "AnnotatedVariants"): Get the names of the with variants.

**`names<-`** signature(`x` = "AnnotatedVariants", `value` = "character"): Set the names of the variants.

## Author(s)

Hans-Ulrich Klein

## See Also

[annotateVariants](#), [htmlReport](#)

## Examples

```
variants = data.frame(
  start=c(106157528, 106154991,106156184),
  end=c(106157528, 106154994,106156185),
  chromosome=c("4", "4", "4"),
  strand=c("+", "+", "+"),
  seqRef=c("A", "ATAG", "---"),
  seqMut=c("G", "----", "ATA"),
  seqSur=c("TACAGAA", "TTTATAGATA", "AGC---TCC"),
  stringsAsFactors=FALSE)
```

```
rownames(variants) = c("snp", "del", "ins")
## Not run: av = annotateVariants(variants)
## Not run: annotatedVariants(av)[["snp"]]
```

annotateVariants      *Adds genomic information to variants*

### Description

This method annotates given genomic variants (mutations). Annotation includes affected genes, exons and codons. Resulting amino acid changes are returned as well as dbSNP identifiers, if the mutation is already known. All information is fetched from Ensembl via biomaRt using the datasets `hsapiens_gene_ensembl` and `hsapiens_snp`.

### Usage

```
annotateVariants(object, bsGenome)
```

### Arguments

object	A data frame storing variants or an instance of AVASet/MapperSet or a data frame (see details).
bsGenome	An object of class BSgenome giving the genome to be used as reference sequence to calculate amino acid changes. This argument is only applicable when object is of type <a href="#">MapperSet</a> . Default is 'BSgenome.Hsapiens.UCSC.hg19'. Note that the genome should fit to the Ensembl annotation.

### Details

If a data frame is given, the following columns must be present:

start	genomic start position in the current Ensembl genome
end	genomic end position in the current Ensembl genome
chromosome	chromosome in ensembl notation (i.e. "1", "2", ..., "Y")
strand	"+" or "-" relative to the nucleotide bases given below
seqRef	reference sequence
seqMut	sequence of the observed variant
seqSur	reference sequence extended for 3 bases in both directions

The rownames of the data frame are used as mutations' names (IDs). See examples for a properly defined data frame.

### Value

An object of class AnnotatedVariants. Affected genes, transcripts and exon as well as known SNPs are stored in a list-like structure. See the documentation of class [AnnotatedVariants-class](#) for details.

### Author(s)

Hans-Ulrich Klein

**See Also**

[AnnotatedVariants-class](#), [AVASet-class](#), [MapperSet-class](#), [htmlReport](#)

**Examples**

```
variants = data.frame(
  start=c(106157528, 106154991,106156184),
  end=c(106157528, 106154994,106156185),
  chromosome=c("4", "4", "4"),
  strand=c("+", "+", "+"),
  seqRef=c("A", "ATAG", "---"),
  seqMut=c("G", "---", "ATA"),
  seqSur=c("TACAGAA", "TTTATAGATA", "AGC---TCC"),
  stringsAsFactors=FALSE)
rownames(variants) = c("snp", "del", "ins")
## Not run: annotateVariants(variants)
```

---

assayDataAmp

*Access the amplicon data of an AVASet.*

---

**Description**

Similar to `assayData` of the `Biobase ExpressionSet`, this function returns the assay data of the amplicon slot of an instance of the `AVASet`.

**Usage**

```
assayDataAmp(object)
```

**Arguments**

`object`            An `link{AVASet-class}` object.

**Value**

The assay data of the amplicon slot consists of a list of two data frames with the number of forward and reverse reads of all amplicons for each sample (see [AVASet-class](#) for details).

**Author(s)**

Christoph Bartenhagen

**See Also**

[fDataAmp](#), [featureDataAmp](#), [AVASet-class](#)

**Examples**

```
# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)

# show contents of amplicon assay data
assayDataAmp(avaSetExample)
```

---

ava2vcf	<i>Convert an AVASet object into a VCF object</i>
---------	---

---

**Description**

Converts all variants in a given AVASet object into a VCF object and writes it to a file in VCF format if filename is given.

**Usage**

```
ava2vcf(object, filename, annot)
```

**Arguments**

object	An object of class <a href="#">AVASet</a> .
filename	The name of the VCF file to write in, if omitted no file is written.
annot	An object of class AnnotatedVariants. Optional, if given variants are annotated with informations from dbSNP.

**Value**

An object of class [VCF-class](#)

**Author(s)**

Christian Ruckert

**See Also**

[AnnotatedVariants-class](#), [AVASet-class](#), [VCF-class](#), [writeVcf](#)

**Examples**

```
data("avaSetFiltered")
vcf <- ava2vcf(avaSetFiltered)
```

---

AVASet	<i>Creating an AVASet</i>
--------	---------------------------

---

**Description**

This function imports a project of Roche's Amplicon Variant Analyzer (AVA) Software. It stores all information into an extended version of the Biobase eSet.

**Usage**

```
AVASet(dirname, avaBin, file_sample, file_amp, file_reference, file_variant, file_variantHits)
```

## Arguments

dirname	The path of the AVA project. Without AVA-CLI (AVA version < 2.6): A directory that contains the files and subdirectories "Amplicons/ProjectDef/ampliconsProject.txt", "Amplicons/Results/Variants/currentV", "Amplicons/Results/Variants", "Amplicons/Results/Align". Using AVA-CLI (recommended): Path usually ends with directory "project-folder"
avaBin	The directory containing the AVA-CLI binary doAmplicon (usually "bin" in the AVA installation directory)
file_sample	Sample information exported with the AVA-CLI. File has to be in CSV format.
file_amp	Amplicons exported with the AVA-CLI. File has to be in CSV format.
file_reference	Reference sequences exported with the AVA-CLI. File has to be in CSV format.
file_variant	Variant information exported with the AVA-CLI. File has to be in CSV format.
file_variantHits	Report of variant hits exported with the AVA-CLI. File has to be in CSV format.

## Details

The five arguments for AVA command line interface (AVA-CLI) exports are optional and useful for exported projects, when no AVA software is installed. For exporting, start the AVA-CLI with the command "doAmplicon" and use the commands "open", then "list sample", "list amplicon", "list reference", "list variant" and "report variantHits". See [AVASet-class](#) for more details.

Giving only a project directory and the path to the AVA-CLI binary doAmplicon, AVASet will import all information by accessing the AVA-CLI from within R.

An AVASet object consists of three slots to store data about

1. variants

**variantForwCount/variantRevCount:** Data frames that contain the number of reads with the respective variant in forward/reverse direction.

**totalForwCount/totalRevCount:** Data frames that contain the total coverage for every variant location in forward/reverse direction.

**referenceSeq:** Gives the identifier of the reference sequence.

**variantBase/referenceBases:** The bases changed in each variant.

**start/end:** The position of the variant on the reference sequence.

**canonicalPattern/name:** Short identifiers of a variant including the position and the bases changed.

2. amplicons

**forwCount/revCount:** Data frames that contain the number of reads for every amplicon and each sample in forward/reverse direction.

**primer1,primer2:** The primer sequences for every amplicon.

**referenceSeqID:** The identifier of the reference sequence.

**targetStart/targetEnd:** The coordinates of the target region.

3. reference sequences

**If additional information has been loaded from Ensembl via alignShortReads, this slot knows about the chromos**



The structure of the variant and amplicon data is derived from the Biobase eSet and thus separated into assayData, phenoData and featureData. All information about the reference sequences is stored into an object of class AlignedRead.

The phenoData of the variants lists the sample-IDs and name, annotation and group of the read data for all samples. If available, the pico titer plate(s) (PTP) or MID(s) of each sample are shown as well (using the AVA-CLI, PTPs and MIDs cannot be imported at the moment).

### Value

An instance of the AVASet class.

### Note

It is recommended to use the import via AVA-CLI access. Although deprecated, the import for projects created with older version of the AVA software (< v2.6) is still possible.

### Author(s)

Christoph Bartenhagen

### See Also

[AVASet-class](#), [MapperSet-class](#), [alignShortReads](#)

### Examples

```
# Loading a project from AVA version < 2.6:
# Load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)
avaSetExample

# Loading exported data, that was exported via AVA-CLI
# Load an AVA dataset containing 6 samples, 4 amplicons and 222 variants
# by specifying each file exported from the AVA-CLI
projectDir = system.file("extdata", "AVASet_doAmplicon", package="R453Plus1Toolbox")
avaSetExample = AVASet(dirname=projectDir, file_sample="sample.csv", file_amp="amp.csv", file_reference="reference.csv")
avaSetExample

# In case AVA software is installed:
# Saying, for example, the AVA software was installed to the directory "/home/User/AVA",
# the easiest way to import a project via AVA-CLI would look like:
# avaSetExample = AVASet(dirname="myProjectDir", avaBin="/home/User/AVA/bin")
```

---

AVASet-class

*Class to contain Amplicon Variant Analyzer Output*

---

### Description

Container to store data imported from a project of Roche's Amplicon Variant Analyzer Software. It stores all information into an extended version of the Biobase ExpressionSet.

## Objects from the Class

Objects can be created by calls of the form `AVASet(dirname, avaBin)`. `dirname` is a character giving the project directory and `avaBin` is a character giving the path to the AVA software installation (i.e. the directory containing the `doAmplicon` binary). The constructor will start the AVA software command line and import all necessary data.

If the AVA software is not installed on the same machine that runs R, all data must be exported manually using the AVA Command Line Interface (AVA-CLI). After having exported all text files, the constructor `AVASet(dirname, avaBin, file_sample, file_amp, file_reference, file_variant, file_variantHits)` can be used to import them. See the example below.

Finally, old project folders generated by AVA software < 2.6 can be imported using `AVASet(dirname)`. Where `dirname` is the path to the project folder (i.e. a directory that contains the files and subdirectories "Amplicons/ProjectDef/ampliconsProject.txt", "Amplicons/Results/Variants/currentVariantDefs.txt", "Amplicons/Results/Variants", "Amplicons/Results/Align").

## Slots

**assayData:** Object of class `AssayData`. Contains the number of reads and the total read depth for every variant and each sample in forward and reverse direction. Its column number equals `nrow(phenoData)`.

**featureData:** Object of class `AnnotatedDataFrame`. Contains information about the type, position and reference of each variant.

**phenoData:** Object of class `AnnotatedDataFrame`. Contains the sample-IDs and name, annotation and group of the read data for all samples. If available, the lane, pico titer plate(s) (PTP) or MID(s) of each sample are shown as well.

**assayDataAmp:** Object of class `AssayData`. Contains the number of reads for every amplicon and each sample in forward/reverse direction. Its column number equals `nrow(featureDataAmp)`.

**featureDataAmp:** Object of class `AnnotatedDataFrame`. Contains the primer sequences, reference sequences and the coordinates of the target regions for every amplicon.

**referenceSequences:** Object of class `AlignedRead`. If additional alignment information were computed via `alignShortReads`, this slot knows about the chromosome, position and the strand of each reference sequence.

**variantFilterPerc:** Object of class `numeric`. Contains a threshold to display only those variants, whose coverage (in percent) in forward and reverse direction in at least one sample is higher than this filter value. See [setVariantFilter](#) for details about setting this value.

**variantFilter:** Object of class `character`. Contains a vector of variant names whose coverage (in percent) in forward and reverse direction in at least one sample is higher than the filter value in `variantFilterPerc`.

**dirs:** Object of class `character`. Based on a directory given at instantiation of the object, it contains a vector of several directories containing all relevant AVA-project files.

**experimentData:** Object of class `MIAME`. Contains details of the experiment.

**annotation:** Object of class `character`. Label associated with the annotation package used in the experiment.

**protocolData:** Object of class `annotatedDataFrame`. Contains additional information about the samples.

**.\_\_\_classVersion\_\_\_:** Object of class `Versions`. Remembers the R and R453Toolbox version numbers used to created the `AVASet` instance.

**Extends**

Class [eSet](#), directly. Class [VersionedBiobase](#), by class "eSet", distance 2. Class [Versioned](#), by class "eSet", distance 3.

**Methods**

**object[i,j :]** Allows subsetting an AVASet object by features (i) and samples (j).

**assayDataAmp(object), assayDataAmp(object)<-value:** Similar to assayData of the Biobase ExpressionSet, this function returns/replaces the amplicon assay data.

**fDataAmp(object):** Similar to fData of the Biobase ExpressionSet, this function returns the amplicon feature data as a data frame.

**featureDataAmp(object), featureDataAmp(object)<-value:** Similar to featureData of the Biobase ExpressionSet, this function returns/replaces the amplicon feature data and feature meta.

**referenceSequences(object), referenceSequences(object)<-value:** Returns/replaces the reference sequence slot.

**alignShortReads(object, bsGenome):** Retrieve the chromosomal positions of the amplicon sequences.

**setVariantFilter(object):** Sets the filter to display only those variants, whose coverage (in percent) in forward and reverse direction in at least one sample is higher than the given value.

**getVariantPercentages(object)** Computes the coverage for every variant over all reads (forward and/or reverse) and for each sample.

**annotateVariants(object):** Annotates given genomic variants. See [annotateVariants](#) for details.

**htmlReport(object):** Exports all (filtered) variant data into a html report. See [htmlReport](#) for details

**Author(s)**

Christoph Bartenhagen

**See Also**

[MapperSet-class](#), [annotateVariants](#), [alignShortReads](#), [htmlReport](#), [setVariantFilter](#), [getVariantPercentages](#)

**Examples**

```
# sum up class structure
showClass("AVASet")

# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)
avaSetExample

# show contents of assay, feature and pheno data
head(assayData(avaSetExample)$variantForwCount)
head(assayData(avaSetExample)$totalForwCount)
head(assayData(avaSetExample)$variantRevCount)
head(assayData(avaSetExample)$totalRevCount)
head(fData(avaSetExample))
pData(avaSetExample)
assayDataAmp(avaSetExample)
```

```
fDataAmp(avaSetExample)
referenceSequences(avaSetExample)

# Use these commands to export a project from within the AVA-CLI (doAmplicon):
# > list sample -outputFile sample.csv
# > list amplicon -outputFile amp.csv
# > list reference -outputFile reference.csv
# > list variant -outputFile variant.csv
# > report variantHits -outputFile variantHits.csv

# Load an AVA dataset containing 6 samples, 4 amplicons and 222 variants
# by specifying five files, that were exported with the AVA-CLI:
projectDir = system.file("extdata", "AVASet_doAmplicon", package="R453Plus1Toolbox")
avaSetExample = AVASet(dirname=projectDir, file_sample="sample.csv", file_amp="amp.csv", file_reference="re
```

---

avaSetExample

*Amplicon Variant Analyzer data import*


---

## Description

This is an example of an `link{AVASet-class}` object containing the output of Roche's Amplicon Variant Analyzer Software. It consists of 6 samples, 4 amplicons and 259 variants.

## Usage

```
data(avaSetExample)
```

## Format

Formal class 'AVASet'

## Source

'Next-generation sequencing technology reveals a characteristic pattern of molecular mutations in 72.8 leukemia by detecting frequent alterations in TET2, CBL, RAS, and RUNX1' (Kohlmann A et al., J Clin Oncol. 2010 Aug 20;28(24):3858-65. Epub 2010 Jul 19)

## Examples

```
data(avaSetExample)
avaSetExample
```

---

`avaSetFiltered`*Amplicon Variant Analyzer data import*

---

**Description**

This is an example of an `link{AVASet-class}` object containing the output of Roche's Amplicon Variant Analyzer Software. It consists of 6 samples, 4 amplicons and 4 variants. The variants were previously filtered according to the amplicon coverage (see [setVariantFilter](#) for details about filtering an AVASet object).

**Usage**

```
data(avaSetFiltered)
```

**Format**

Formal class 'AVASet'

**Source**

'Next-generation sequencing technology reveals a characteristic pattern of molecular mutations in 72.8 leukemia by detecting frequent alterations in TET2, CBL, RAS, and RUNX1' (Kohlmann A et al., J Clin Oncol. 2010 Aug 20;28(24):3858-65. Epub 2010 Jul 19)

**Examples**

```
data(avaSetFiltered)
avaSetFiltered
```

---

`avaSetFiltered_annot`*AVASet variant annotations*

---

**Description**

These are example annotations for 4 variants of an AVASet (try `data(avaSetFiltered)` to retrieve the corresponding `link{AVASet-class}` object). The annotations include affected genes, exons and codons as well as resulting amino acid changes and dbSNP identifiers (if the mutation is already known).

**Usage**

```
data(avaSetFiltered_annot)
```

**Format**

Formal class 'AnnotatedVariants'

**Source**

'Next-generation sequencing technology reveals a characteristic pattern of molecular mutations in 72.8 leukemia by detecting frequent alterations in TET2, CBL, RAS, and RUNX1' (Kohlmann A et al., J Clin Oncol. 2010 Aug 20;28(24):3858-65. Epub 2010 Jul 19)

**Examples**

```
data(avaSetFiltered_annot)
```

---

baseFrequency

*Absolute And Relative Frequency Of The Four Bases.*

---

**Description**

This function returns the absolute and the relative frequency of the four bases (A, C, G, T).

**Usage**

```
baseFrequency(object)
```

**Arguments**

object            An object of class [DNAStrngSet](#), [ShortRead](#) or [SFFContainer](#).

**Details**

This function makes use of the [alphabetFrequency](#) function from package Biostrings.

**Value**

A [data.frame](#) with two columns containing the absolute and relative frequencies respectively and six rows, one for each of the four bases (A, C, G, T), one for other symbols contained in the reads and one summarizing the five aforementioned rows.

**Author(s)**

Christian Ruckert

---

baseQualityHist	<i>Plot A Histogram Of The Base Qualities.</i>
-----------------	--

---

**Description**

Create a histogram based on the quality of every single base from all sequences.

**Usage**

```
baseQualityHist(object, xlab="Quality score", ylab="Number of bases", col="firebrick1", breaks=40, ...)
```

**Arguments**

object	An object of class <a href="#">QualityScaledDNStringSet</a> , <a href="#">ShortReadQ</a> or <a href="#">SFFContainer</a> .
xlab	The X axis label.
ylab	The Y axis label.
col	The plotting color.
breaks	The number of breaks in the histogram (see ‘hist’).
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

**Author(s)**

Christian Ruckert

---

baseQualityStats	<i>Statistics Of Base Quality</i>
------------------	-----------------------------------

---

**Description**

This function returns mean, minimum, maximum and standard deviation of the base quality scores over all sequences.

**Usage**

```
baseQualityStats(object)
```

**Arguments**

object	An object of class <a href="#">QualityScaledDNStringSet</a> , <a href="#">ShortReadQ</a> or <a href="#">SFFContainer</a> .
--------	--

**Value**

A numeric vector with four slots: mean, min, max, sd.

**Author(s)**

Christian Ruckert

---

breakpoints

*Putative breakpoints of chimeric reads*


---

### Description

This example holds two consensus (pathogenic and reciprocal) breakpoints of 12 chimeric reads indicating an inversion on chromosome 16. The Breakpoints object gives access to the breakpoint locations as well as alignment information for each of the 12 reads.

### Usage

```
data(breakpoints)
```

### Format

Formal class 'Breakpoints'

### Source

'Targeted next-generation sequencing detects point mutations, insertions, deletions, and balanced chromosomal rearrangements as well as identifies novel leukemia-specific fusion genes in a single procedure' (Leukemia, submitted)

### Examples

```
data(breakpoints)
```

---

Breakpoints-class

*Class "Breakpoints"*


---

### Description

Container to store chimeric reads that were clustered to putative breakpoints indicating structural variants. Related information like breakpoint position or alignment information about the chimeric reads is stored as well.

### Objects from the Class

Objects can be created by calls of the form `new("Breakpoints", ...)`. Usually, objects will be created by calling the `detectBreakpoints` method. It is not intended that users create objects of this class manually.

All slots of this class can be found twice. One slot name ends with "C1" and the other "C2". The slots labeled with "C2" are empty until `mergeBreakpoints` has been called and contain information about putatively associated breakpoints detected by `mergeBreakpoints`.



**Slots**

- seqsC1:** Object of class "list" with one data frame for each breakpoint. The data frame stores all chimeric reads covering the first breakpoint together with the alignment information.
- seqsC2:** Object of class "list" with one data frame for each breakpoint. The data frame stores all chimeric reads covering the second breakpoint together with the alignment information.
- commonBpsC1:** Object of class "list" with one data frame for each breakpoint. The data frame stores the consensus breakpoint sequence as well as the breakpoint coordinates of the first breakpoint.
- commonBpsC2:** Object of class "list" with one data frame for each breakpoint. The data frame stores the consensus breakpoint sequence as well as the breakpoint coordinates of the second breakpoint.
- commonAlignC1:** Object of class "list" with one object of class [PairwiseAlignedFixedSubject-class](#) for each breakpoint storing the alignments of the chimeric reads against the consensus breakpoint sequence for the first breakpoint.
- commonAlignC2:** Object of class "list" with one object of class [PairwiseAlignedFixedSubject-class](#) for each breakpoint storing the alignments of the chimeric reads against the consensus breakpoint sequence for the second breakpoint.
- alignedReadsC1:** Object of class "list" with one object of class [AlignedRead-class](#) storing all chimeric reads covering the first breakpoint and their alignments.
- alignedReadsC2:** Object of class "list" with one object of class [AlignedRead-class](#) storing all chimeric reads covering the second breakpoint and their alignments.

**Methods**

- alignedReadsC1<-** signature(object = "Breakpoints",value = "list"): Setter-method for the alignedReadsC1 slot.
- alignedReadsC1** signature(object = "Breakpoints"): Getter-method for the alignedReadsC1 slot.
- alignedReadsC2<-** signature(object = "Breakpoints",value = "list"): Setter-method for the alignedReadsC2 slot.
- alignedReadsC2** signature(object = "Breakpoints"): Getter-method for the alignedReadsC2 slot.
- commonAlignC1<-** signature(object = "Breakpoints",value = "list"): Setter-method for the commonAlignC1 slot.
- commonAlignC1** signature(object = "Breakpoints"): Getter-method for the commonAlignC1 slot.
- commonAlignC2<-** signature(object = "Breakpoints",value = "list"): Setter-method for the commonAlignC2 slot.
- commonAlignC2** signature(object = "Breakpoints"): Getter-method for the commonAlignC2 slot.
- commonBpsC1<-** signature(object = "Breakpoints",value = "list"): Setter-method for the commonBpsC1 slot.
- commonBpsC1** signature(object = "Breakpoints"): Getter-method for the commonBpsC1 slot.
- commonBpsC2<-** signature(object = "Breakpoints",value = "list"): Setter-method for the commonBpsC2 slot.
- commonBpsC2** signature(object = "Breakpoints"): Getter-method for the commonBpsC2 slot.

**seqsC1<-** signature(object = "Breakpoints", value = "list"): Setter-method for the seqsC1 slot.

**seqsC1** signature(object = "Breakpoints"): Getter-method for the seqsC1 slot.

**seqsC2<-** signature(object = "Breakpoints", value = "list"): Setter-method for the seqsC2 slot.

**seqsC2** signature(object = "Breakpoints"): Getter-method for the seqsC2 slot.

[ signature(x = "Breakpoints", i = "ANY", j = "ANY"): Subsetting a Breakpoints object.

**length** signature(x = "Breakpoints"): Returns the number of breakpoints stored.

**mergeBreakpoints** signature(breakpoints = "Breakpoints", maxDist = "missing", mergeBPs = "list"): Merge presumably related breakpoints.

**mergeBreakpoints** signature(breakpoints = "Breakpoints", maxDist = "missing", mergeBPs = "missing"): Merge presumably related breakpoints.

**mergeBreakpoints** signature(breakpoints = "Breakpoints", maxDist = "numeric", mergeBPs = "missing"): Merge presumably related breakpoints.

**names<-** signature(x = "Breakpoints", value = "ANY"): Set the names of the breakpoints.

**names** signature(x = "Breakpoints"): Get the names of the breakpoints.

**plotChimericReads** signature(brpData = "Breakpoints"): Plot the structural variant and the chimeric reads covering its breakpoints.

**summary** signature(object = "Breakpoints"): Create a data frame summarizing information about all breakpoints.

**table** signature(... = "Breakpoints"): Create a frequency table of cluster sizes.

**Author(s)**

Hans-Ulrich Klein, Christoph Bartenhagen

**See Also**

[filterChimericReads](#), [detectBreakpoints](#), [mergeBreakpoints](#), [plotChimericReads](#)

---

calculateTiTv

*Calculate transition transversion ratio*

---

**Description**

When many point mutations are detected, the ration of transitions to transversions can be used as quality measure to assess the number of false positives.

**Usage**

```
## S4 method for signature 'AVASet'
calculateTiTv(object)
## S4 method for signature 'MapperSet'
calculateTiTv(object)
```

**Arguments**

**object** An instance of AVASet or MapperSet storing the detected variants.

**Details**

For more information about the Ti/Tv ratio see [http://www.broadinstitute.org/gsa/wiki/index.php/QC\\_Methods](http://www.broadinstitute.org/gsa/wiki/index.php/QC_Methods)

**Value**

A list with two elements: A substitution matrix summarizing all observed substitutions and the transition/transversion ratio.

**Author(s)**

Hans-Ulrich Klein

**Examples**

```
data(avaSetExample)
ava = setVariantFilter(avaSetExample, c(0.03, 0.03))
calculateTiTv(ava)
```

---

captureArray

*Custom capture array design*

---

**Description**

Design of a custom Roche NimbleGen 385k capture array. The array captures short segments corresponding to all exon regions of 92 distinct target genes (genome build hg19). In addition, contiguous genomic regions were represented for three additional genes, i.e. CFBF, MLL, and RUNX1.

**Usage**

```
data(captureArray)
```

**Format**

Formal class 'CompressedIRangesList'

**Source**

'Targeted next-generation sequencing detects point mutations, insertions, deletions, and balanced chromosomal rearrangements as well as identifies novel leukemia-specific fusion genes in a single procedure' (Leukemia, submitted)

**Examples**

```
data(captureArray)
```

---

`complexity.dust`*Sequence Complexity Using The DUST Algorithm*

---

### Description

This function evaluates the sequence complexity using the DUST algorithm.

### Usage

```
complexity.dust(object, xlab="Complexity score (0=high, 100=low)", ylab="Number of sequences",  
xlim=c(0, 100), col="firebrick1", breaks=100, ...)
```

### Arguments

<code>object</code>	An object of class <a href="#">DNAStrngSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
<code>xlab</code>	The X axis label.
<code>ylab</code>	The Y axis label.
<code>xlim</code>	The limits of the X axis.
<code>col</code>	The plotting color.
<code>breaks</code>	The number of breaks in the histogram (see ‘hist’).
<code>...</code>	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

### Details

The complexity score is based on how often different trinucleotides occur and is scaled between 0 and 100. A sequence of homopolymer repeats (e.g. TTTTTTTTTT) has a score of 100, of dinucleotide repeats (e.g. TATATATATA) has a score around 49, and of trinucleotide repeats (e.g. TAGTAGTAG) has a score around 32. Scores above seven can be considered low-complexity.

### Value

A numeric vector containing the complexity score for each sequence.

### Author(s)

Christian Ruckert

### References

Schmieder R. (2011) Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, 2011 Mar 15;27(6):863-4.

---

complexity.entropy      *Sequence Complexity Using The Shannon-Wiener Algorithm*

---

## Description

This function evaluates the sequence complexity using the Shannon-Wiener Algorithm.

## Usage

```
complexity.entropy(object, xlab="Complexity score (0=low, 100=high)", ylab="Number of sequences",  
xlim=c(0, 100), col="firebrick1", breaks=100, ...)
```

## Arguments

object	An object of class <a href="#">DNAStrngSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
xlab	The X axis label.
ylab	The Y axis label.
xlim	The limits of the X axis.
col	The plotting color.
breaks	The number of breaks in the histogram (see ‘hist’).
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

## Details

The entropy approach evaluates the entropy of trinucleotides in a sequence. The entropy values are scaled from 0 to 100 and lower entropy values imply lower complexity. A sequence of homopolymer repeats (e.g. TTTTTTTTTT) has an entropy value of 0, of dinucleotide repeats (e.g. TATATATATA) has an entropy value around 16, and of trinucleotide repeats (e.g. TAGTAGTAG) has an entropy value around 26. Scores below 70 can be considered low-complexity.

## Value

A numeric vector containing the complexity score for each sequence.

## Author(s)

Christian Ruckert

## References

Schmieder R. (2011) Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, 2011 Mar 15;27(6):863-4.

---

convertCigar	<i>Basic functions for CIGAR strings</i>
--------------	--

---

### Description

These are temporary methods, that are likely to be replaced by methods from the Rsamtools package in near future.

### Usage

```
extendedCIGARToList(cigars)
listToExtendedCIGAR(cigarList)
```

### Arguments

cigars	A character vector with CIGAR strings.
cigarList	A list of converted CIGAR strings as produced by <a href="#">extendedCIGARToList</a>

---

coverageOnTarget	<i>Computes the coverage restricted to the target region.</i>
------------------	---

---

### Description

This method computes the approximate coverage of each base in a given region.

### Usage

```
coverageOnTarget(alnReads, targetRegion)
```

### Arguments

alnReads	A list as returned by scanBam storing aligned reads.
targetRegion	The target region as a RangesList. The chromosome names must fit to the chromosome names used in the alignment information of the given reads.

### Details

The detailed alignment information given by the CIGAR strings in .bam files are ignored by the function. Instead, it is assumed that the whole read alignes to the reference without indels. This is often not true for longer read (e.g. generated with Roche 454 Sequencing), but saves computation time.

### Value

A list of the same length as the alnReads argument. Each list element is an integer vector of the same length as the target region (in bases) and stores the coverage generated by the reads from the corresponding list element of alnReads.

**Author(s)**

Hans-Ulrich Klein

**See Also**[scanBam](#)**Examples**

```
library(Rsamtools)
bamFile = system.file("extdata", "StructuralVariantDetection", "bam", "N01.bam", package="R453Plus1Toolbox")
bam = scanBam(bamFile)
region = RangesList("11"=IRanges(start=118307205, end=118395936))
cov = coverageOnTarget(bam, region)
```

demultiplexReads

*Performs MID/Multiplex filtering***Description**

Roche's Genome Sequencer allows to load two or more samples on one region. To allocate sequences to samples, each sample has a unique multiplex sequence. The multiplex sequence should be the prefix of all sequences from that sample. This method demultiplexes a given set of sequences according to the given multiplex sequences (MIDs).

**Usage**

```
## S4 method for signature 'XStringSet,XStringSet,numeric,logical'
demultiplexReads(reads, mids, numMismatches, trim)
```

**Arguments**

reads	A DNASTringSet instance that contains reads starting with MIDs
mids	A DNASTringSet instance that contains the MIDs
numMismatches	The maximal number of mismatches allowed, default 2.
trim	Whether the MIDs should be cutted-out, default TRUE

**Details**

All given MIDs must have the same length. The algorithm computes the number of mismatches for each MID. The read is assigned to the MID with the lowest number of mismatches. If two or more MIDs have the same number of mismatches, or if the number of mismatches is greater than the given argument numMismatches, the read is not assigned to any MID. The default number of allowed mismatches is 2.

**Value**

demultiplexReads returns a list with one DNASTringSet instance for each MID.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[genomeSequencerMIDs](#), [DNAStrngSet](#)

**Examples**

```
library(Biostrings)
mids = genomeSequencerMIDs(c("MID1", "MID2", "MID3"))
reads = DNAStrngSet(c(
  paste(as.character(mids[["MID1"]]), "A", sep=""),
  paste(as.character(mids[["MID1"]]), "AA", sep=""),
  paste(as.character(mids[["MID2"]]), "AAA", sep="")))
demultiplexReads(reads, mids)
```

---

detectBreakpoints

*Clustering and consensus breakpoint detection for chimeric reads*

---

**Description**

Given a set of chimeric reads, this methods computes all putative breakpoints. First, chimeric reads are clustered such that all reads spanning the same breakpoint form a cluster. Then, a consensus breakpoint sequence and breakpoint position is computed for each cluster.

**Usage**

```
detectBreakpoints(chimericReads, bpDist=100, minClusterSize=4, removeSoftClips=TRUE, bsGenome)
```

**Arguments**

chimericReads	A list storing chimeric reads as returned by <a href="#">filterChimericReads</a> . The list must have the format as defined by the <a href="#">scanBam</a> method.
bpDist	The maximum distance in base pairs between the breakpoints of two chimeric reads at which the reads are merge to a cluster.
minClusterSize	Cluster whose size is below minClusterSize are be excluded from breakpoint detection.
removeSoftClips	If true, soft-clipped bases at the beginning or the end of a sequence are removed (see details below).
bsGenome	A bsGenome instance providing the reference sequences. If missing, the library BSgenome.Hsapiens.UCSC.hg19 is used by default.

**Details**

This method is usually invoked after calling [filterChimericReads](#) and before calling [mergeBreakpoints](#). It first forms clusters of chimeric reads (reads with exactly two local alignments) that span the same breakpoint and than computes a consensus breakpoint sequence for each cluster.

To carry out a hierarchical clustering, a measure for the distance between two chimeric reads must be defined. If reads span different chromosomes, their distance is set to infinity. The strand information of the local alignments may also indicate that two chimeric reads do not span the same breakpoint even if they span the same chromosomes. For example, the first reads has two local alignments on the positive strand whereas the second read has one local alignment on the positive strand and the other on the negative strand. In this case, the distance is set to infinity, too. Finally,



the distance measure distinguishes between the two breakpoints (sometimes called the pathogenic and the reciprocal breakpoint) that originate from the same structural variant. The distance between a read from the pathogenic and a read from the reciprocal breakpoint is infinity so that two different clusters will emerge. These two related breakpoints can be merged later using the [mergeBreakpoints](#) method. We observed that the breakpoints of these two cases often differ by a few ten or even a few hundred basepairs.

If the chromosome and strand information between two reads  $x$  and  $y$  are coherent, the Euclidean distance is used:

$$d(x, y) = (bp(x, ChrA) - bp(y, ChrA))^2 + (bp(x, ChrB) - bp(y, ChrB))^2$$

where  $bp$  gives the coordinates of the breakpoint for the given read and chromosome. Hierarchical clustering is applied with complete linkage and the dendrogram is cutted at a height of `bpDist` to obtain the final clusters. The `bpDist` argument does usually not influence the result, because we observed that reads spanning the same breakpoint have very little variation (only a few base pairs) in their local alignments due to sequencing errors or due to ambiguity caused by same/similar sequence of both chromosome near the breakpoint.

Although the given set of reads may belong to the same chimeric DNA, their individual breakpoints may differ in a few base pairs. Furthermore, a single read may have more than one possible breakpoint if a (small) part of the read was aligned to both parts.

The following step determines a consensus breakpoint for each cluster. It uses the supplied `bsGenome` to construct a chimeric reference sequence for all possible breakpoints over all reads within each cluster. After the reads were realigned to the chimeric reference sequences, the one that yields the highest alignment score is taken to represent best the chimeric DNA and its breakpoints.

As a preprocessing step, `detectBreakpoints` offers to remove soft clips occurring after the alignment:

Some reads may contain soft-clipped bases (e.g. linker sequences) at the beginning of the first part of the read or at the end of the second part. By default, `detectBreakpoints` removes these unaligned subsequences and adjusts the cigar string, the sequence, the sequence width (`qwidth`) and the local start/end coordinates.

## Value

`detectBreakpoints` returns an object of class `breakpoints`, which is a list of breakpoint clusters, which gives access to all alignments and consensus breakpoints:

<code>seqs</code>	This <a href="#">IRanges DataFrame</a> is mainly a rearranged version of the alignment input in <code>chimericReads</code> . In addition, it shows the corresponding breakpoints and local alignment coordinates.
<code>commonBps</code>	A dataframe listing the breakpoints for both parts of the chimeric reference, the associated chromosome, strand and the reference sequence itself, including positions " <code>localStart</code> "/" <code>localEnd</code> " indicating which part of the reference belongs to which breakpoint.
<code>commonAlign</code>	An object of class <a href="#">PairwiseAlignedFixedSubject</a> of the <code>Biostrings</code> package that contains the alignment to the (best) consensus reference sequence.
<code>alignedReads</code>	On the basis of <code>commonAlign</code> and <code>commonBps</code> , <code>alignedReads</code> is an instance of class <a href="#">AlignedRead</a> containing all aligned reads including their associated chromosomes, strands, and positions. Since the reference is a chimeric sequence each read has two chromosome and two strand entries.

**Author(s)**

Hans-Ulrich Klein, Christoph Bartenhagen

**See Also**

[filterChimericReads](#) [mergeBreakpoints](#) [plotChimericReads](#)

**Examples**

```
# Construct a small example with three chimeric reads
# (=6 local alignments) in bam format as given by
# aligners such as BWA-SW.
# The first two reads originate from the same case but
# from different strands. The third read originate from
# the reciprocal breakpoint.
library("BSgenome.Scerevisiae.UCSC.sacCer2")
bamReads = list()
bamReads[[1]] = list(
  qname=c("seq1", "seq1", "seq2", "seq2", "seq3", "seq3"),
  flag = as.integer(c(0, 0, 16, 16, 0, 0)),
  rname = factor(c("II", "III", "III", "II", "III", "II")),
  strand = factor(c("+", "+", "-", "-", "+", "+")),
  pos = as.integer(c(99951, 200000, 200000, 99951, 199950, 100001)),
  qwidth = as.integer(c(100, 100, 100, 100, 100, 100)),
  cigar = c("50M50S", "50S50M", "50S50M", "50M50S", "50M50S", "50S50M"),
  seq = DNASTringSet(c(
    paste(substr(Scerevisiae$chrII, start=99951, stop=100000),
          substr(Scerevisiae$chrIII, start=200000, stop=200049),
          sep=""),
    paste(substr(Scerevisiae$chrII, start=99951, stop=100000),
          substr(Scerevisiae$chrIII, start=200000, stop=200049),
          sep=""),
    paste(substr(Scerevisiae$chrIII, start=200000, stop=200049),
          substr(Scerevisiae$chrII, start=99951, stop=100000),
          sep=""),
    paste(substr(Scerevisiae$chrIII, start=200000, stop=200049),
          substr(Scerevisiae$chrII, start=99951, stop=100000),
          sep=""),
    paste(substr(Scerevisiae$chrIII, start=199950, stop=199999),
          substr(Scerevisiae$chrII, start=100001, stop=100050),
          sep=""),
    paste(substr(Scerevisiae$chrIII, start=199950, stop=199999),
          substr(Scerevisiae$chrII, start=100001, stop=100050),
          sep="")))
)

bps = detectBreakpoints(bamReads, minClusterSize=1, bsGenome=Scerevisiae)
summary(bps)
table(bps)

mergeBreakpoints(bps)
```

---

dinucleotideOddsRatio *Dinucleotide Odds Ratio*

---

### Description

This function calculates the dinucleotide odds ratio for each of the sixteen possible dinucleotides.

### Usage

```
dinucleotideOddsRatio(object, xlab="Under-/over-representation of dinucleotides",  
col="firebrick1", ...)
```

### Arguments

object	An object of class <a href="#">DNABStringSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
xlab	The X axis label.
col	The plotting color.
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

### Details

The dinucleotide odds ratio assigns a value between 0 and 2 to each of the sixteen possible dinucleotides (AA, AC, AG, AT, CA, CC, CG, CT, GA, GC, GG, GT, TA, TC, TG, TT). For values below 1 the dinucleotide is under-represented compared to the randomly expected frequency of this dinucleotide in a sequence of the given length and with the given frequencies of the four nucleotides (A, C, G, T). For values above 1 this dinucleotide is over-represented.

### Value

A matrix with sixteen columns, one for each dinucleotide, containing the dinucleotide odds ratio values for each sequence in a separate row.

### Author(s)

Christian Ruckert

### References

Schmieder R. (2011) Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, 2011 Mar 15;27(6):863-4.

---

fDataAmp	<i>Access the amplicon data of an AVASet.</i>
----------	---

---

**Description**

Similar to fData of the Biobase ExpressionSet, this function returns the feature data of the amplicon slot of an instance of the AVASet.

**Usage**

```
fDataAmp(object)
```

**Arguments**

object            An `link{AVASet-class}` object.

**Value**

The feature data of the amplicon slot contains the names, primers, start/end positions and reference sequences of all amplicons (see [AVASet-class](#) for details). It returns a data frame.

**Author(s)**

Christoph Bartenhagen

**See Also**

[featureDataAmp](#), [assayDataAmp](#), [AVASet-class](#)

**Examples**

```
# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)
avaSetExample

# show contents amplicon feature data
fDataAmp(avaSetExample)
```

---

featureDataAmp	<i>Access the amplicon data of an AVASet</i>
----------------	--

---

**Description**

Similar to featureData of the Biobase ExpressionSet, this function returns the feature data and feature meta of the amplicon slot of an instance of the AVASet.

**Usage**

```
featureDataAmp(object)
```

**Arguments**

object                    An link{AVASet-class} object.

**Value**

The feature data of the amplicon slot contains the names, primers, start/end positions and reference sequences of all amplicons (see [AVASet-class](#) for details). The returned object is of class [AnnotatedDataFrame](#).

**Author(s)**

Christoph Bartenhagen

**See Also**

[fDataAmp](#), [assayDataAmp](#), [AVASet-class](#),

**Examples**

```
# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)
avaSetExample

# show contents amplicon feature data
featureDataAmp(avaSetExample)
```

---

filterChimericReads	<i>Extract chimeric reads and apply filtering steps to remove artificial chimeric reads.</i>
---------------------	--

---

**Description**

Chimeric reads may be caused by sequencing a chromosomal aberration or by technical issues during sample preparation. This method implements several filter steps to remove false chimeric reads.

**Usage**

```
## S4 method for signature 'list,RangesList,DNAString,numeric,numeric'
filterChimericReads(alnReads, targetRegion, linkerSeq, minDist, dupReadDist)
```

**Arguments**

alnReads                    A list storing the aligned reads as produced by the function [scanBam](#).

targetRegion                A object of class [rangesList](#) containing the target region of e.g. a used capture array. The parameter may be omitted in case of a non targeted sequencing approach.

linkerSeq                    A linker sequence that was used during sample preparation. It may be omitted.

minDist                     The minimum distance between two local alignments (see details), default 1000

dupReadDist                 The maximum distance between the 5 prime start position of two duplicated reads (see details), default 1.

## Details

The following filter steps are performed:

1. All chimeric reads with exactly two local alignments are extracted. Reads with more than two local alignments are discarded.
2. If the targetRegion argument is given, chimeric reads must have one local alignment at least overlapping the target region. If both local alignments are outside the target region, the read is discarded.
3. If the linkerSeq argument is given, all chimeric reads that have the linker sequence between their local alignments are removed. When searching the linker sequence, 4 mismatches or indels are allowed and the linker sequence must not start or end within the first or last ten bases of the read. The function searches for the linkerSeq and for its reverse complement.
4. Two local alignment of a read must have minDist reads between the alignments (if both alignment are on the same chromosome). Otherwise, the read seems to span a deletion and not a chromosomal aberration and is discarded.
5. Duplicated reads are removed. Two reads are duplicated, if they lie on the same strand and have the same 5 prime start position. Due to sequencing and alignment errors, the start position may vary for a maximum of dupReadDist bases. In case of duplicated reads, only the longest read is kept.

Reads passing all filtering steps are returned in the list structure as given by the alnReads argument (as derived from the scanBam method). A data frame with information about the number of reads that passed each filter is added to the list.

## Value

A list containing only filtered chimeric reads. The list has the same structure like the given argument alnReads. Additionally, one element "log" with logging information of each filtering step is added.

## Author(s)

Hans-Ulrich Klein

## See Also

[detectBreakpoints](#), [mergeBreakpoints](#), [Breakpoints-class](#), [scanBam](#), [sequenceCaptureLinkers](#)

## Examples

```
library(Rsamtools)
bamFile = system.file("extdata", "StructuralVariantDetection", "bam", "N01.bam", package="R453Plus1Toolbox")
bam = scanBam(bamFile)
data(captureArray)
linker = sequenceCaptureLinkers("gSel3")[[1]]
filterReads = filterChimericReads(bam, targetRegion=captureArray, linkerSeq=linker)
```

---

flowgramBarplot	<i>Create A Barplot Of The Flow Intensities</i>
-----------------	---

---

**Description**

This function creates a barplot of the flow intensities with one bar for each nucleotide in the flow. With the height giving the measured intensity.

**Usage**

```
flowgramBarplot(x, range=c(0, length(flowgram(x))), xlab="Flow sequence",
  ylab="Flow intensity", col=c(A="black", C="red", G="blue", T="green"), ...)
```

**Arguments**

x	An object of class <a href="#">SFFRead</a> .
range	Two positions between which the flows should be plotted.
xlab	The X axis label.
ylab	The Y axis label.
col	The colors in which the four nucleotids should be plotted.
...	Arguments to be passed to methods, such as graphical parameters (see 'par').

**Author(s)**

Christian Ruckert

---

gcContent	<i>Calculate The Overall GC-Content</i>
-----------	---

---

**Description**

This function calculates the GC-content summarized over all sequences.

**Usage**

```
gcContent(object)
```

**Arguments**

object	An object of class <a href="#">DNAStrngSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
--------	--

**Details**

The GC-content is calculated as follows:

$$\frac{\#G + \#C}{\#G + \#C + \#A + \#T} * 100$$

Where #G is the number of base G in all sequences.

**Value**

A numeric vector of length one containing the overall GC-content.

**Author(s)**

Christian Ruckert

---

gcContentHist	<i>GC-Content Histogram</i>
---------------	-----------------------------

---

**Description**

This function creates a histogram of the relative GC-content per sequence.

**Usage**

```
gcContentHist(object, xlab="GC content", ylab="Number of sequences", col="firebrick1",
  breaks=50, ...)
```

**Arguments**

object	An object of class <a href="#">DNAStrngSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
xlab	The X axis label.
ylab	The Y axis label.
col	The plotting color.
breaks	The number of breaks in the histogram (see ‘hist’).
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

**Author(s)**

Christian Ruckert

---

gcPerPosition	<i>GC-Content Per Position</i>
---------------	--------------------------------

---

**Description**

This function plots the GC-content frequency per base position.

**Usage**

```
gcPerPosition(object, range=0.95, type="l", col="blue", xlab="Position", ylab="Frequency", ...)
```



**Arguments**

object	An object of class <a href="#">DNAStrngSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
range	An integer vector of length one or two. If length one only bases up to the percent quantile of read lengths defined by the given value are shown. A vector of length two gives the start and end base between which the GC-content is plotted.
type	The type of the plot (see 'plot').
col	The plotting color.
xlab	The X axis label.
ylab	The Y axis label.
...	Arguments to be passed to methods, such as graphical parameters (see 'par').

**Author(s)**

Christian Ruckert

---

genomeSequencerMIDs *Retrieve GS multiplex sequences*

---

**Description**

This method returns the standard multiplex sequences used by the Genome Sequence MID library kits.

**Usage**

```
## S4 method for signature 'missing'  
genomeSequencerMIDs()  
## S4 method for signature 'character'  
genomeSequencerMIDs(mid)
```

**Arguments**

mid                    Character vector with multiplex sequences' IDs (MIDs)

**Details**

If the argument mid is omitted, all 14 available multiplex sequences are returned.

**Value**

genomeSequencerMIDs returns a [DNAStrngSet](#) with the requested multiplex sequences.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[demultiplexReads](#)

**Examples**

```
genomeSequencerMIDs()
genomeSequencerMIDs(c("MID1", "MID3"))
```

---

```
getAlignedReads          Import reads from an Amplicon Variant Analyzer project
```

---

**Description**

For a given AVASet, this function imports all aligned reads belonging to all (or some selected) amplicons of all samples.

**Usage**

```
getAlignedReads(object, amplicons, dir)
```

**Arguments**

object	An instance of the link{AVASet-class}.
amplicons	An (optional) character vector of amplicon names as mentioned in the amplicon feature data (see <a href="#">fDataAmp</a> ).
dir	Usually, the method tries to retrieve the path to the AVA project from the given link{AVASet} object. However, if it fails to find the directory, dir can be used to set the root directory of the AVA project.

**Details**

This function reports all reads for all samples together. If you want to get the reads for some samples individually, try subsetting your AVASet as in the examples below.

**Value**

One [DNAStrngSet](#) that contains all aligned reads for all samples (eventually restricted to some given amplicons).

**Author(s)**

Christoph Bartenhagen

**See Also**

[AVASet-class](#), [fDataAmp](#)

**Examples**

```
# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)

# import all reads for amplicon "TET2_E11.04" of the first sample
avaProjectDir = system.file("extdata", "AVASet", package = "R453Plus1Toolbox")
alnReads = getAlignedReads(avaSetExample[, 1], dir=avaProjectDir, amplicons="TET2_E11.04")
show(alnReads)
```

---

getAminoAbbr	<i>Get amino acid abbreviations</i>
--------------	-------------------------------------

---

**Description**

This function returns a table with amino acid names as first column and the according abbreviations as row names.

**Usage**

```
getAminoAbbr()
```

---

getVariantPercentages	<i>Variant coverage</i>
-----------------------	-------------------------

---

**Description**

This function computes the coverage for each variant (in forward and/or reverse direction) for all samples. The coverage is defined as the percentual amount of reads that cover a variant.

**Usage**

```
getVariantPercentages(object, direction="both")
```

**Arguments**

object	An instance of class <a href="#">AVASet-class</a> or <a href="#">MapperSet-class</a> .
direction	A character indicating the direction ("forward", "reverse" or "both").

**Details**

If the direction was set to "both", the percentages are computed over the sum of both directions. Otherwise it is computed only over the occurrences in one direction (forward or reverse). The occurrences can be accessed via `assayData`.

**Value**

`getVariantPercentages` returns a data frame with all percentages/frequencies for all samples.

**Author(s)**

Christoph Bartenhagen

**See Also**

[setVariantFilter](#).

**Examples**

```

# load a (filtered) AVA dataset containing 6 samples, 4 amplicons and 4 variants
data(avaSetFiltered)
avaSetFiltered

# both directions
getVariantPercentages(avaSetFiltered, direction="both")
# this is equivalent to
(assayData(avaSetFiltered)[[1]] + assayData(avaSetFiltered)[[3]]) / (assayData(avaSetFiltered)[[2]] + assayData(avaSetFiltered)[[4]])

# forward direction only
getVariantPercentages(avaSetFiltered, direction="forward")
# this is equivalent to
assayData(avaSetFiltered)[[1]] / assayData(avaSetFiltered)[[2]]

# reverse direction only
getVariantPercentages(avaSetFiltered, direction="reverse")
# this is equivalent to
assayData(avaSetFiltered)[[3]] / assayData(avaSetFiltered)[[4]]

```

---

homopolymerHist

*Create A Histogram Of The Homopolymer Stretches*


---

**Description**

This function creates a histogram for the different lengths of the homopolymer stretches with one bar for each nucleotide in the flow. With the height giving the number of this homopolymer stretch in the flowgram.

**Usage**

```

homopolymerHist(x, range=c(0, length(flowgram(x))), xlab="Homopolymer length",
  ylab="Number of homopolymers", col=c(A="black", C="red", G="blue", T="green"), ...)

```

**Arguments**

x	An object of class <a href="#">SFFRead</a> .
range	Two positions between which the flows should be plotted.
xlab	The X axis label.
ylab	The Y axis label.
col	The colors in which the four nucleotids should be plotted.
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

**Author(s)**

Christian Ruckert

htmlReport

*HTML-Report Builder for the AVASet and MapperSet***Description**

This function creates a HTML variant and quality report for a given AVASet or MapperSet instance.

**Usage**

```
htmlReport(object, annot, blocks=c(), transcripts=c(), sampleCols, minMut=3, dir="HTMLReport", title="S
```

**Arguments**

object	An <a href="#">AVASet-class</a> or <a href="#">MapperSet-class</a> instance.
annot	An instance of class <a href="#">AnnotatedVariants</a> you get by calling <a href="#">annotateVariants</a> . If no such argument is supplied the data will be read from the Ensembl database automatically for all variants.
blocks	Character vector of block names for each variant. The variants will then be structured into several blocks. If no such list is supplied, the report consists of just one (big) table for all variants.
transcripts	Character vector containing Ensembl transcript-IDs that order the according entries on the transcript pages. Transcripts given in this argument will appear on top of the transcript page.
sampleCols	Character vector of column names of the sample data (phenoData) of the AVASet/MapperSet object to filter the sample output on the transcript pages. All columns will be listed if no such argument is given.
minMut	If the value of minMut is greater than zero, the report lists only variants, whose coverage for at least one sample is higher than minMut (percentage between 0 and 100).
dir	Character with the desired output directory. By defaultm the directory "HTML-Report" will be created in the current directory.
title	Heading for the first page with the variant information.

**Details**

The report is structured into two (MapperSet) or three (AVASet) parts containing variant and quality information:

- The main page sums up given variant information like the name, type, reference gene, position (see [fData](#), [annotateVariants](#)).  
Using the argument blocks, the main page can be individually structured by assigning a block name to each variant.  
The main page can be further structured by samples. For a given AVASet object, every sample links to another short quality report showing only the amplicon coverage for this sample.
- Every variant on the main page links to a page with further details about the affected genes and transcripts (e.g. Ensembl gene-IDs, transcript-IDs, codon sequences, changes of amino acids (if coding)).
- Only in case of AVASet object: A quality report shows the coverage of every amplicon in forward and/or reverse direction. Further plots display the coverage by MID and PTP (if this information is given in the pheno data of the object).

**Author(s)**

Christoph Bartenhagen, Hans-Ulrich Klein, Christian Ruckert

**See Also**

[annotateVariants](#).

**Examples**

```
# note: all examples save the report to the directory "htmlReportExample" in your current R working directory

# load a filtered AVA dataset containing 6 samples, 4 amplicons and 4 variants
# and its variant annotations
data("avaSetFiltered")
data("avaSetFiltered_annot")

# create a full report showing all (unfiltered) information
htmlReport(avaSetFiltered, avaSetFiltered_annot, dir="htmlReportExample", title="htmlReport Example", minMut=
# create a report that emphasizes on samples with variants covered by at least 50% of the reads
htmlReport(avaSetFiltered, avaSetFiltered_annot, dir="htmlReportExample", title="htmlReport Example", minMut=

# create a report that is structured by the reference genes
library("ShortRead")
refs = sapply(fData(avaSetFiltered)$referenceSeq, function(x)
  subset(pData(alignedData(referenceSequences(avaSetFiltered))), pData(alignedData(referenceSequences(avaSetFiltered)))
htmlReport(avaSetFiltered, avaSetFiltered_annot, dir="htmlReportExample", title="htmlReport Example", minMut=

# create a report whose sample information only lists the sample ids
pData(avaSetFiltered)
sampleCols = "SampleID"
htmlReport(avaSetFiltered, avaSetFiltered_annot, dir="htmlReportExample", title="htmlReport Example", minMut=
```

---

MapperSet

*Creating a MapperSet*

---

**Description**

This function imports a project of Roche's GS Reference Mapper Software. It stores all information into an instance of the Biobase ExpressionSet.

**Usage**

```
MapperSet(dirs, samplenames)
```

**Arguments**

dirs	A character vector containing all sample directories (i.e. directories that contain the files "mapping/454HCDiffs.txt" (required), "mapping/454ReadStatus.txt" (optional), "mapping/454NewblerMetrics.txt" (optional)).
samplenames	A character vector containing samplenames. The order and number of samplenames must be consistent with the filenames to ensure the correctness of the MapperSet. If no samplenames are given, the filenames are used for naming.

## Details

An instance of the MapperSet is derived from the Biobase eSet and thus structured into

1. assayData

**variantForwCount/variantRevCount:** Contain the number of reads with the respective difference in forward/reverse direction.

**totalForwCount/totalRevCount:** Contain the total coverage for every variant in forward/reverse direction.

2. featureData

**chromosome, start/end:** Give the location of each variant.

**referenceBases/variantBase:** Show the bases changed in each variant.

**regName:** The name of the region (gene) where the variant is located.

**knownSNP:** Lists Ensembl variant-ids for known SNPs (if any).

3. phenoData

**By default, the phenoData contains the accession number of every sample.**

## Value

An instance of the MapperSet.

## Author(s)

Christoph Bartenhagen

## See Also

[AVASet-class](#)

## Examples

```
# load a GS Mapper dataset containing 3 samples and 111 variants
data(mapperSetExample)
mapperSetExample
```

MapperSet-class

*Class to Contain GS Reference Mapper Output***Description**

Container to store data imported from a project of Roche's GS Reference Mapper Software. It stores all information into a Biobase ExpressionSet.

**Objects from the Class**

Objects can be created by calls of the form `MapperSet(filename)`. While `filename` is a vector containing all sample directories (i.e. directories that contain the files "mapping/454HCDiffs.txt" and "mapping/454NewblerMetrics.txt").

**Slots**

`assayData`: Object of class `AssayData`. Contains the number of reads with the respective difference and the total coverage for every variant in forward and reverse direction.

`featureData`: Object of class `AnnotatedDataFrame`. Contains information about the type, location and reference of each variant. If available, it shows further Ensembl variant-ids for known SNPs.

`phenoData`: Object of class `AnnotatedDataFrame`. By default, the `phenoData` contains the accession number of every sample.

`variantFilterPerc`: Object of class `numeric`. Contains a threshold to display only those variants, whose coverage (in percent) in forward and reverse direction in at least one sample is higher than this filter value. See [setVariantFilter](#) for details about setting this value.

`variantFilter`: Object of class `character`. Contains a vector of variant names whose coverage (in percent) in forward and reverse direction in at least one sample is higher than the filter value(s) in `variantFilterPerc`.

`dirs`: Object of class `character`. Based on a directory given at instantiation of the object, it contains a vector of several directories containing all relevant GS Mapper project files.

`experimentData`: Object of class `MIAME`. Contains details of the experiment.

`annotation`: Object of class `character` Label associated with the annotation package used in the experiment.

`protocolData`: Object of class `AnnotatedDataFrame`. Contains additional information about the samples.

`.___classVersion___`: Object of class `Versions`. Remembers the R and R453Toolbox version numbers used to created the `MapperSet` instance.

**Extends**

Class `eSet`, directly. Class `VersionedBiobase`, by class "eSet", distance 2. Class `Versioned`, by class "eSet", distance 3.



## Methods

**setVariantFilter(object):** Sets the filter to display only those variants, whose coverage (in percent) in forward and reverse direction in at least one sample is higher than the given value.

**getVariantPercentages(object)** Computes the coverage for every variant over all reads (forward and/or reverse) and for each sample.

**annotateVariants(object):** Annotates given genomic variants. See [annotateVariants](#) for details.

**htmlReport(object):** Exports all (filtered) variant data into a html report. See [htmlReport](#) for details

**getReadStatus(object):** Reads the file "454ReadStatus.txt" in the GSM project directory which contains information about the alignment of each read (chr, pos, strand, etc.)and returns its content in a dataframe.

## Author(s)

Christoph Bartenhagen

## See Also

[AVASet](#), [annotateVariants](#), [htmlReport](#), [setVariantFilter](#), [getVariantPercentages](#)

## Examples

```
# sum up class structure
showClass("MapperSet")

# load a GS Mapper dataset containing 3 samples and 111 variants
data(mapperSetExample)
mapperSetExample

# show contents of assay, feature and pheno data
assayData(mapperSetExample)
fData(mapperSetExample)
pData(mapperSetExample)
```

---

mapperSetExample

*GS Reference Mapper data import*

---

## Description

This is an example of an `link{MapperSet-class}` object containing the output of Roche's GS Reference Mapper Software. It consists of 3 samples and 111 variants.

## Usage

```
data(mapperSetExample)
```

## Format

Formal class 'MapperSet'

**Source**

'Targeted next-generation sequencing detects point mutations, insertions, deletions, and balanced chromosomal rearrangements as well as identifies novel leukemia-specific fusion genes in a single procedure' (Leukemia, submitted)

**Examples**

```
data(mapperSetExample)
mapperSetExample
```

---

mergeBreakpoints	<i>Identify and merge related breakpoints caused by the same variant.</i>
------------------	---

---

**Description**

Structural variation like transversions or inversion cause two breakpoints. In the context of fusion genes, these are called the pathogenic breakpoint and the reciproke breakpoint. The method [detectBreakpoints](#) processes each breakpoint individually and does explicitly not put reads from the pathogenic and reciproke breakpoint into the same cluster. Hence, it is usually sensible to call this methods afterward to search for related pairs of breakpoints to gain more confidence about the existence of a structural variation.

**Usage**

```
mergeBreakpoints(breakpoints, maxDist, mergeBPs)
```

**Arguments**

breakpoints	An object of class Breakpoints storing the breakpoints that will (potentially) be merged.
maxDist	The maximal distance in basepairs at which two breakpoints will be merged. Default value is 1000.
mergeBPs	An optional list of vectors of length two giving the breakpoints that should be merged. If this argument is given, the method will not search for related breakpoints.

**Details**

If the maxDist argument is given, the method compares each pair of breakpoints and checks, whether the two breakpoints may belong to the same structural variation. In addition to the spanned chromosomes, the orientation and the strand information of the reads spanning the breakpoints are also compared for this purpose. If chromosome, orientation and strand information of two breakpoints go well together, they will be merged, if the absolute distance of the breakpoints on chromosome A plus the absolute distance on chromosome B is smaller or equal to maxDist. If one breakpoint has more than one potential mate breakpoint for merging, it will be merged with the first candidate and a warning message is printed. The default value of maxDist is 1000.

If the mergeBPs argument is given, the method will not search for related breakpoints but simply merge the given breakpoints. mergeBPs must be a list with vectors of length two that either contain the names of the indices of the breakpoints that should be merged.

The arguments maxDist and mergeBPs cannot be given together. The given Breakpoints object must not contain breakpoints that have been merged before.

**Value**

An object of class Breakpoints storing merged and unmerged breakpoints.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[detectBreakpoints](#), [Breakpoints-class](#), [plotChimericReads](#)

**Examples**

```
# Load bam file and filter chimeric reads
library(Rsamtools)
bamFile = system.file("extdata", "StructuralVariantDetection", "bam", "N01.bam", package="R453Plus1Toolbox")
bam = scanBam(bamFile)
data(captureArray)
linker = sequenceCaptureLinkers("gSel3")[[1]]
filterReads = filterChimericReads(bam, targetRegion=captureArray, linkerSeq=linker)

# detect breakpoints of size >= 3
breakpoints = detectBreakpoints(filterReads, minClusterSize=3)
table(breakpoints)
summary(breakpoints)

# merge breakpoints
breakpoints = mergeBreakpoints(breakpoints)
summary(breakpoints)
```

---

mutationInfo

*Example data for [plotVariants](#)*

---

**Description**

This data.frame is part of the vignette example of the [plotVariants](#) function. It contains annotations for the different mutations occurring in the example data. It has columns "mutation", "legend" and "color".

**Usage**

```
data(plotVariantsExample)
```

**Format**

data.frame

**Examples**

```
data(plotVariantsExample)
mutationInfo
```

---

nucleotideCharts      *Nucleotide Charts*

---

### Description

This function plots the relative frequency of the four bases for each position in the sequences.

### Usage

```
nucleotideCharts(object, range=0.95, linetypes=c(A="l", C="l", G="l", T="l", N="l"),
  linecols=c(A="black", C="red", G="blue", T="green", N="grey70"), xlab="Position",
  ylab="Frequency", ...)
```

### Arguments

object	An object of class <a href="#">DNAStringSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
range	An integer vector of length one or two. If length one only bases up to the percent quantile of read lengths defined by the given value are shown. A vector of length two gives the start and end base between which the nucleotide frequencies are plotted.
linetypes	The line types used for the four nucleotids + N, see (see 'par')
linecols	The colors in which the four nucleotids + N should be plotted.
xlab	The X axis label.
ylab	The Y axis label.
...	Arguments to be passed to methods, such as graphical parameters (see 'par').

### Author(s)

Christian Ruckert

---

plotAmpliconCoverage      *Creates a plot visualizing the number of reads per amplicon*

---

### Description

A function for visualizing the number of reads per amplicon or per MID / pico titer plate.

### Usage

```
## S4 method for signature 'AVASet,character,logical'
plotAmpliconCoverage(avaSet, type="amplicon", bothDirections=TRUE, cex.names=0.8, cex.axis=0.8, las=3,
```

**Arguments**

avaSet	An instance of AVASet.
type	A character vector specifying the type of plot.
bothDirections	A logical value determining whether the plot sums forward and reverse reads or shows them separately.
cex.names	Font size of the amplicon name labels.
cex.axis	Font size of axes' labels.
las	Orientation of amplicon name labels.
col	Colors used in the plot.
...	Arguments to be passed to methods, such as graphical parameters (see 'par').

**Details**

If the argument type is "amplicon", the number of reads for each amplicon are visualized. In case of a AVASet with one sample, a barplot with one bar for each amplicon is created. In case of more than one sample, a boxplot with one box for each amplicon is plotted. If type is "mid", a boxplot with one box for each MID is created. If type is "ptp", a boxplot with one box for each pico titer plate is created.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[AVASet](#)

**Examples**

```
## Not run: data(avaSetExample)
plotAmpliconCoverage(avaSetExample)
plotAmpliconCoverage(avaSetExample[,1])
## End(Not run)
```

---

plotChimericReads      *Plots chimeric reads*

---

**Description**

This function plots a given set of aligned chimeric reads along a reference sequence. It plots the breakpoints of translations or inversions and marks deletions, insertions and mismatches. Optionally, it displays all base pairs in a given region around the breakpoint.

**Usage**

```
plotChimericReads(brpData, geneSymbols=FALSE, plotMut=TRUE, plotBasePairs=FALSE, maxBasePairs=
  col=c("red", "green", "black", "orange"))
```

**Arguments**

brpData	A Breakpoints object containing the consensus breakpoint of all reads and the consensus reference sequence as returned by the methods <a href="#">detectBreakpoints</a> and <a href="#">mergeBreakpoints</a> . Since only one plot is made, the function will only work for objects of class Breakpoints having length one.
geneSymbols	Boolean value whether to automatically load and plot the gene symbols from the Ensembl database. Additionally, geneSymbols can be a vector of two strings for an own annotation.
plotMut	Boolean value whether to mark deletions, insertions and mismatches.
plotBasePairs	Optionally, plotChimericReads displays all base pairs in a given region around the breakpoint (see maxBasePairs).
maxBasePairs	The maximum number of base pairs to be plotted. Only used in conjunction with plotBasePairs=TRUE.
legend	A logical value (TRUE/FALSE) whether to plot a legend that explains the colouration of the insertions, deletions, mismatches and breakpoints.
title	A title for the plot.
col	A vector of four colours to draw insertions, deletions, mismatches and breakpoints. In this order, the default colours are "red", "green", "black" and "orange" (use colours() to see a list of possible values).

**Details**

This method is intended to be run after the pipeline for structural variant detection. Therefore, see the methods [filterChimericReads](#), [detectBreakpoints](#) and [mergeBreakpoints](#) to correctly preprocess your alignment before running plotChimericReads.

**Note**

It is recommended to first create and resize the output device (e.g. the plotting window or a pdf file) before plotting. For example, on Unix systems you may try X11(width=w, height=h) or pdf(file="plotChimericReads.pdf", width=w, height=h) for some window width w (e.g. w=12) and window height h (e.g. h=6).

**Author(s)**

Christoph Bartenhagen

**See Also**

[Breakpoints-class](#), [detectBreakpoints](#), [mergeBreakpoints](#)

**Examples**

```
# load breakpoint data containing twelve chimeric reads describing an inversion in chromosome 16
data("breakpoints")
breakpoints

# standard plot
# (only arrangement of reads plotted; breakpoints in orange, deletions
# in red, insertions in green and mismatches in black by default)
plotChimericReads(breakpoints)
```

```
# plot base pairs in the breakpoint region (+/- 32bp)
## Not run: plotChimericReads(breakpoints, plotBasePairs=TRUE, maxBasePairs=32)

# use custom colours and display a legend:
# deletions="brown", insertions="blue", mismatches="yellow", breakpoints="gray"
plotChimericReads(breakpoints, col=c("brown", "blue", "yellow", "gray"), legend=TRUE)
```

---

plotVariants

*Plots variant positions*


---

## Description

This function illustrates the positions and types of mutations within a given gene and transcript. The plot shows only coding regions (thus, units are amino acids / codons). The coding region is further divided into exons labeled with their rank in the transcript.

## Usage

```
plotVariants(data, gene, transcript, regions, mutationInfo, groupBy, horiz=FALSE, cex=1, title="", legend=TRUE)
```

## Arguments

data	This can be either a data.frame or an instance of class <code>annotatedVariants</code> you get by calling <code>annotateVariants</code> . A data frame requires the columns "label", "pos" "mutation" and "color" specifying an annotation for each mutation, its position (a single numeric value), the mutation type and (optionally) the color of the mutation (e.g. depicting a certain group identity independent from the mutation type). The position needs to be given as amino acids / codons.
gene	A string containing the Ensembl id of the gene of interest (required).
transcript	A string containing an Ensembl transcript id (optional, but recommended). If no Ensembl transcript-id is passed, it is chosen automatically and the function will return an appropriate data frame with all annotated transcripts for the given gene.
regions	A data frame having columns "name", "start", "end" and "color". The plot will highlight these regions with the given colors and print their name in the legend.
mutationInfo	A data frame with annotations for the different mutations occurring in the data (optional but recommended). It requires the columns "mutation", "legend" and "color". The first column must list the exact mutation names that occur in the data column "mutation". The column "legend" allows for a more detailed name of the mutation that will appear in the legend. The color of each mutation type is optional and can also be assigned automatically.
groupBy	By default, the mutations will be grouped by their position (i.e. the column "pos"). If necessary, one might give the name of an other column in the <a href="#">data</a> here; for example the mutation labels. Please see the details below for more information.
horiz	In more comprehensive datasets, more than one mutation may be listed for a single position. If <code>horiz=FALSE</code> , these overlapping mutations will be aligned vertically. If <code>horiz=TRUE</code> they will be aligned horizontally in groups allowing a label for every single mutation.

cex	A numeric value $> 0$ giving the factor by which the labels are magnified relative to the default text size. If the width of the device is too small for all mutation labels, their size will be scaled automatically until it fits.
title	A title for the plot (optional).
legend	A logical value (TRUE/FALSE) whether to plot a legend for the mutations types (see mutationInfo) and the highlighted regions (if any).

### Details

The plot will show the coding part of the gene and its exons as x-axis at the bottom. Mutations will be marked and ordered above according to their genomic position. The axis units are amino acids / codons (hence, all given genomic positions should be divided by three if necessary).

Passing a data frame to this function allows a much more individual and detailed annotation of the mutations like labels, colors and user defined mutation types.

Passing an instance of class [annotatedVariants](#) is useful for integration into the R453Plus1Toolbox pipeline and for compatibility to older versions of the plot. Then, it will only distinguish deletions and missense, nonsense and silent mutations.

By default, the plot will group mutations (horizontally or vertically) by their position. It is possible to group by an other column in the data (see parameter groupBy), but in the current version this makes only sense if the mutations in one group are locally clustered, i.e. have the same or a similar position. The parameter groupBy is mainly useful to modify or even disable the automatic grouping of different mutations at the same position.

### Value

The function will return a data frame containing all Ensembl transcript information for the given gene ("ensembl\_transcript\_id", "rank", "cds\_start", "cds\_end" and "cds\_length"). This data frame may prove useful for retrieving a Ensembl transcript id for future plots.

### Note

Depending on the amount of mutations and the size of the gene, the plot may not fit into the device or the text may become too small. It is recommended to carefully select the right size of your device before starting this function to ensure a well scaled and beautiful plot.

This function requires the package TeachingDemos to work, which can be found at CRAN.

### Author(s)

Christoph Bartenhagen

### See Also

[annotateVariants](#).

### Examples

```
# EXAMPLE 1: Working with instances of class annotatedVariants

# one missense, one nonsense point mutation and one deletion
variants = data.frame(
  row.names=c("missense", "deletion", "nonsense"),
  start=c(106157528, 106157635, 106193892),
  end=c(106157528, 106157635, 106193892),
```



```

chromosome=c("4", "4", "4"),
strand=c("+", "+", "+"),
seqRef=c("A", "G", "C"),
seqMut=c("G", "-", "T"),
seqSur=c("TACAGAA", "TAAGCAG", "CGGCGAA"),
stringsAsFactors=FALSE)

# annotate variants with affected genes, exons and codons (may take a minute to finish)
## Not run: varAnnot = annotateVariants(variants)

# plot variants for gene TET2 having the Ensembl id "ENSG00000168769"
# when passing no transcript, the largest transcript annotated in the Ensembl database for this gene will be selected automatically
## Not run: plotVariants(data=varAnnot, gene="ENSG00000168769", title="plotVariants Example", legend=TRUE)

# EXAMPLE 2: Working with a data frame

# two missense at one position, one nonsense point mutation and one deletion
# it is possible to assign a color to every single mutation independently from its type
variants = data.frame(
  label=c("A>G", "A>G(2)", "delG", "C>T"),
  pos=c(831,831,867,1437),
  mutation=c("M", "M", "D", "N"),
  color=c("black", "black", "green", "red"),
  stringsAsFactors=FALSE
)

# more detailed names for mutation abbreviations can be passed as mutationInfo
# this is useful for the legend, but can also be generated automatically
mutationInfo = data.frame(
  mutation=c("M", "D", "S", "N"),
  legend=c("Missense", "Nonsense", "Silent", "Deletion"),
  stringsAsFactors=FALSE
)

# regions of interest can be highlighted using the regions parameter
regions = data.frame(
  name = c("region1", "region2"),
  start = c(700, 1400),
  end = c(1000, 1900),
  color = c("red", "blue")
)

# using the horiz parameter, multiple mutations occurring at the same place can be either aligned ...
# ... vertically
## Not run: plotVariants(data=variants, gene="ENSG00000168769", transcript="ENST00000513237", regions=regions)

# ... or horizontally in groups
## Not run: plotVariants(data=variants, gene="ENSG00000168769", transcript="ENST00000513237", regions=regions, horiz=TRUE)

# group mutations by their label and not by their position (which is the default)
## Not run: plotVariants(data=variants, gene="ENSG00000168769", transcript="ENST00000513237", regions=regions, group="label")

```

## Description

This method creates a plot similar to the variation frequency plot in Roche's GS Amplicon Variant Analyzer. The plot shows the reference sequence along the x-axis and indicates variants as bars at the appropriate positions. The height of the bars corresponds to the percentage of reads carrying the variant. A second y-axis indicates the absolute number of reads covering the variant.

## Usage

```
plotVariationFrequency(object, plotRange, ...)
```

## Arguments

object	A character pointing to an Amplicon Variant Analyser Global Alignment export file.
plotRange	A two dimensional numeric vector giving the start and end base of the reference sequence that should be plotted.
...	Arguments passed to other plotting methods. Especially, argument col: Optional character vector of length 7 specifying the bars' colors indicating different base substitutions or deletions. See details. And argument sequenceCex: Optional numeric value specifying the size of the reference sequence's bases.

## Details

The text file used as input must have the format generated by the AVA export function. Such a file can be generated using the export button in the Global Alignment view of the AVA software. The col argument specifies the colours used for different bases and deletions. The following listing gives the meaning of the i-th position of the col vector (default values in braces):

1. A (green)
2. C (blue)
3. G (black)
4. T (red)
5. N (purple)
6. deletion (gray)

## Author(s)

Hans-Ulrich Klein

## Examples

```
## Not run:  
file = system.file("extdata", "AVAVarFreqExport", "AVAVarFreqExport.xls", package="R453Plus1Toolbox")  
plotVariationFrequency(file, plotRange=c(50, 150))  
## End(Not run)
```

---

positionQualityBoxplot    *Boxplot Of The Quality For Each Position*

---

**Description**

Creates a boxplot of the quality scores over all sequences at each position.

**Usage**

```
positionQualityBoxplot(object, range, binsize=10,
  xlab=paste("Read position in bp (Bin size: ", binsize, "bp)", sep=""), ylab="Quality score",
  col="firebrick1", ...)
```

**Arguments**

object	An object of class <a href="#">QualityScaledDNAStrngSet</a> , <a href="#">ShortReadQ</a> or <a href="#">SFFContainer</a> .
range	A numeric vector of length one or two. If length one only bases from the first until this position are plotted. If two all bases between these two positions are plotted.
binsize	Number of positions to summarize in one box in the plot.
xlab	The X axis label.
ylab	The Y axis label.
col	The plotting color.
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

**Author(s)**

Christian Ruckert

---

readLengthHist    *Histogram Of The Read Lengths*

---

**Description**

This function plots a histogram of the read lengths.

**Usage**

```
readLengthHist(object, cutoff=0.99, xlab="Read length", ylab="Number of sequences",
  col="firebrick1", breaks=100, ...)
```

**Arguments**

object	An object of class <a href="#">DNAStrngSet</a> , <a href="#">ShortRead</a> or <a href="#">SFFContainer</a> .
cutoff	Reads longer than the cutoff-percent quantile are omitted from the plot.
xlab	The X axis label.
ylab	The Y axis label.
col	The plotting color.
breaks	The number of breaks in the histogram (see ‘hist’).
...	Arguments to be passed to methods, such as graphical parameters (see ‘par’).

**Author(s)**

Christian Ruckert

---

readLengthStats	<i>Statistics For The Read Lengths</i>
-----------------	--

---

**Description**

This function returns the mean, median, minimum, maximum and standard deviation of the read lengths over a set of sequences.

**Usage**

```
readLengthStats(object)
```

**Arguments**

object            An object of class [DNAStrngSet](#), [ShortRead](#) or [SFFContainer](#).

**Value**

A [vector](#) with five entries: mean, median, min, max and sd.

**Author(s)**

Christian Ruckert

---

readSff	<i>Function To Read In Roche's .sff Files</i>
---------	---

---

**Description**

This function reads in files in Roche's Standard Flowgram Format (SFF) and store the contents in an [SFFContainer-class](#) object.

**Usage**

```
readSFF(files)
```

**Arguments**

files            The name of the .sff file to read in or a character vector of multiple file names or the name of a directory containing .sff files.

**Value**

An object or a list of objects of class [SFFContainer](#) storing all the information from the .sff file(s).

**Author(s)**

Christian Ruckert

**See Also**

[SFFContainer](#).

**Examples**

```
## Not run: sffContainer = readSFF("test.sff")
```

---

readsOnTarget	<i>Check for each read whether it aligns within the given region.</i>
---------------	---

---

**Description**

This methods checks (approximately) whether the given reads align within a given region.

**Usage**

```
readsOnTarget(alnReads, targetRegion)
```

**Arguments**

alnReads	A list as returned by scanBam storing aligned reads.
targetRegion	The target region as a RangesList. The chromosome names must fit to the chromosome names used in the alignment information of the given reads.

**Details**

The detailed alignment information given by the CIGAR strings in .bam files are ignored by the function. Instead, it is assumed that the whole read alignes to the reference without indels. This is often not true for longer read (e.g. generated with Roche 454 Sequencing), but saves computation time. Hence, this method is useful to approximate the number of reads that align in the target region of a targeted sequencing experiment.

**Value**

A list with one logical vector for each list entry in alnReads. The logical vector indicates for each read whether it overlaps with at least one base from any target region or not.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[scanBam](#)

**Examples**

```
library(Rsamtools)
bamFile = system.file("extdata", "StructuralVariantDetection", "bam", "N01.bam", package="R453Plus1Toolbox")
bam = scanBam(bamFile)
region = RangesList("11"=IRanges(start=118307205, end=118395936))
targetReads = readsOnTarget(bam, region)
sum(targetReads[[1]])
```

---

referenceSequences	<i>Access the reference sequences of an AVASet</i>
--------------------	--

---

**Description**

This function give access to a slot of an instance of the AVASet storing information about all reference sequences of the amplicons.

**Usage**

```
referenceSequences(object)
```

**Arguments**

object            An link{AVASet-class} object.

**Value**

The data is stored in an object of class AlignedRead and thus gives information about all reference sequences and their position on a chromosome (if [alignShortReads](#) has been called before).

**Author(s)**

Christoph Bartenhagen

**See Also**

[alignShortReads](#)

**Examples**

```
# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)

referenceSequences(avaSetExample)
```

---

regions	<i>Example data for <a href="#">plotVariants</a></i>
---------	--

---

**Description**

This data.frame is part of the vignette example of the [plotVariants](#) function. It has the columns "name", "start", "end" and "color". The plot will highlight these regions with the given colors and print their name in the legend.

**Usage**

```
data(plotVariantsExample)
```

**Format**

data.frame

**Examples**

```
data(plotVariantsExample)
regions
```

---

removeLinker	<i>Remove linker sequences located at the start of short reads</i>
--------------	--

---

**Description**

If linkers are attached during sample preparation, it may be useful to remove the linkers' sequences after sequencing. This method finds and removes linker sequences that are located at the start of the given reads.

**Usage**

```
## S4 method for signature 'XStringSet,DNAString,logical,numeric,numeric'
removeLinker(reads, linker, removeReadsWithoutLinker, minOverlap, penalty)
```

**Arguments**

reads	A DNAStringSet instance that contains reads possibly having linkers at their start site
linker	A DNAString instance with the linker's sequence
removeReadsWithoutLinker	Whether reads without linkers should be removed. Default is FALSE
minOverlap	The minimal score that must be achieved when aligning the linker. Default is length(linker)/2
penalty	The penalty for substitutions or indels. Default is -2

**Details**

The best alignment of the linker within the start (length of linker + 5) of each given sequence is computed. The following scoring schema is used: Each matching bases scores +1. Each substitution or indel scores the given penalty argument (default: penalty=-2). There are no penalties for gaps and the end of the linker (overlap). An alignment is considered as match, if the scores is larger or equal to minOverlap (default: minOverlap=round(length(linker)/2)). In cases of a successful match, the subsequence from position 1 until the end of the linker's alignment is removed.

**Value**

removeLinker returns a DNAStringSet with trimmed reads.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[sequenceCaptureLinkers](#), [DNAStrngSet](#), [pairwiseAlignment](#)

**Examples**

```
linker = sequenceCaptureLinkers()[[1]]
reads = DNAStrngSet(c(
  "CTCGAGAATTCTGGATCCTCAAA",
  "GAATTCTGGATCCTCAAA",
  "CTCGAGAAAAAAAAAATCCTCAAA"))
removeLinker(reads, linker)
```

---

sequenceCaptureLinkers *Retrieve NimbleGen's sequence capture linkers*

---

**Description**

This method returns the NimbleGen's linker sequences used with their sequence capture arrays. See pp.29-30 in the NimbleGen Arrays User's Guide.

**Usage**

```
## S4 method for signature 'character'
sequenceCaptureLinkers(name)
## S4 method for signature 'missing'
sequenceCaptureLinkers()
```

**Arguments**

name                    Character vector with linker sequences' names

**Details**

If the argument name is omitted, both linker sequences are returned.

**Value**

sequenceCaptureLinkers returns a DNAStrngSet with the requested linker sequences.

**Author(s)**

Hans-Ulrich Klein

**See Also**

[removeLinker](#)

**Examples**

```
sequenceCaptureLinkers()
```



---

sequenceQualityHist      *A Histogram Of The Sequence Qualities*

---

### Description

This function creates a histogram of the mean qualities of the sequences.

### Usage

```
sequenceQualityHist(object, xlab="Mean of quality scores per sequence",
  ylab="Number of sequences", col="firebrick1", ...)
```

### Arguments

object	An object of class <a href="#">QualityScaledDNAStrngSet</a> , <a href="#">ShortReadQ</a> or <a href="#">SFFContainer</a> .
xlab	The X axis label.
ylab	The Y axis label.
col	The plotting color.
...	Arguments to be passed to methods, such as graphical parameters (see 'par').

### Author(s)

Christian Ruckert

---

setVariantFilter      *Filters output of variant information*

---

### Description

This functions sets the filter to display only those variants, whose amplicon coverage (in percent) in forward and reverse direction in at least one sample is higher than a given value. The coverage is defined as the percentual amount of reads that cover a variant.

### Usage

```
setVariantFilter(object, filter=0)
```

### Arguments

object	An instance of an link{AVASet-class} or <a href="#">MapperSet-class</a> .
filter	A filter value between 0 and 1. If two values are given in a vector, the variants are filtered according to the forward (first value) and reverse direction (second value) separately. In this case, a variant has to meet both requirements.

### Details

Setting the filter affects the assayData and the featureData of the variant slot. See also [getVariantPercentages](#) for further details.

**Value**

setVariantFilter returns the given link{AVASet-class}/link{MapperSet-class} instance with an updated filter value.

**Author(s)**

Christoph Bartenhagen

**See Also**

link{AVASet-class}, link{MapperSet-class}, [getVariantPercentages](#).

**Examples**

```
# load an AVA dataset containing 6 samples, 4 amplicons and 259 variants
data(avaSetExample)
avaSetExample

# use only those variants that are covered by at least 10% of all reads in one sample in both directions together (259 -> 4 v
avaSetExample = setVariantFilter(avaSetExample, filter=0.1)
avaSetExample

# use only those variants that are covered by at least 0.1% of all reads in one sample in forward direction
# and by at least 0% in reverse direction (259 -> 6 variants)
avaSetExample = setVariantFilter(avaSetExample, filter=c(0.1, 0))
avaSetExample

# reset filter values to zero
avaSetExample = setVariantFilter(avaSetExample, filter=0)
# or simply
avaSetExample = setVariantFilter(avaSetExample)
```

---

sff2fastq

---

*Write A SFFContainer Object To A FASTQ File*


---

**Description**

This function takes a [SFFContainer](#) object and writes it to a file in FASTQ format.

**Usage**

```
sff2fastq(x, outdir, fname)
```

**Arguments**

x	An object of class <a href="#">SFFContainer</a> .
outdir	The directory where the file should be stored, defaults to the current working directory.
fname	The name of the file to write. Defaults to the filename slot of the <a href="#">SFFContainer</a> , with .sff substituted with .fastq.

**Author(s)**

Christian Ruckert

---

SFFContainer-class      *Class "SFFContainer"*


---

**Description**

This class is a container for data from files in Roche's Standard Flowgram Format (SFF).

**Objects from the Class**

Objects can be created by calls of the form `new("SFFContainer", ...)`. Usually, objects will be created by calling the `readSFF` method on a file in SFF format.

**Slots**

- name:** Object of class "character" containing the name of the file this SFFContainer was created from.
- flowgramFormat:** Object of class "numeric" representing the format used to encode each of the flowgram values for each read. Currently, only one flowgram format has been adopted and is coded by the value 1.
- flowChars:** Object of class "character" containing the array of nucleotide bases ('A', 'C', 'G', 'T') that correspond to the nucleotides used for each flow of each read.
- keySequence:** Object of class "character" representing the nucleotide bases of the key sequence used for these reads.
- clipQualityLeft:** Object of class "numeric" representing the position of the first base after the clipping point for an attached quality sequence for each read. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityLeft. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.
- clipQualityRight:** Object of class "numeric" representing the position of the last base before the clipping point for an attached quality sequence for each read. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityRight. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.
- clipAdapterLeft:** Object of class "numeric" representing the position of the first base after the clipping point for an attached adapter sequence for each read. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityLeft. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.
- clipAdapterRight:** Object of class "numeric" representing the position of the last base before the clipping point for an attached adapter sequence for each read. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityRight. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.
- flowgrams:** Object of class "list" containing the homopolymer stretch estimates for each flow using one list item for each read.
- flowIndexes:** Object of class "list" containing the flow positions for each base in the called sequence, i.e. for each base, the position in the flowgram whose estimate resulted in that base being called. Each read has its own list item.
- reads:** Object of class "QualityScaledDNASet" containing the basecalled nucleotide sequences of each read together with the quality scores for each of the bases in the sequence using the standard  $-\log_{10}$  probability scale.

**Methods**

- addRead** signature(object = "SFFContainer", read = "SFFRead"): Adds an object of class [SFFRead](#) to the [SFFContainer](#)
- getRead** signature(object = "SFFContainer", readname = "character"): Returns the read with the given name as an object of class [SFFRead](#).
- clipAdapterLeft<-** signature(object = "SFFContainer", value = "numeric"): Setter-method for the clipAdapterLeft slot.
- clipAdapterLeft** signature(object = "SFFContainer"): Getter-method for the clipAdapterLeft slot.
- clipAdapterRight<-** signature(object = "SFFContainer", value = "numeric"): Setter-method for the clipAdapterRight slot.
- clipAdapterRight** signature(object = "SFFContainer"): Getter-method for the clipAdapterRight slot.
- clipQualityLeft<-** signature(object = "SFFContainer", value = "numeric"): Setter-method for the clipQualityLeft slot.
- clipQualityLeft** signature(object = "SFFContainer"): Getter-method for the clipQualityLeft slot.
- clipQualityRight<-** signature(object = "SFFContainer", value = "numeric"): Setter-method for the clipQualityRight slot.
- clipQualityRight** signature(object = "SFFContainer"): Getter-method for the clipQualityRight slot.
- name<-** signature(object = "SFFContainer", value = "character"): Setter-method for the name slot.
- name** signature(object = "SFFContainer"): Getter-method for the name slot.
- flowChars<-** signature(object = "SFFContainer", value = "character"): Setter-method for the flowChars slot.
- flowChars** signature(object = "SFFContainer"): Getter-method for the flowChars slot.
- flowgramFormat<-** signature(object = "SFFContainer", value = "numeric"): Setter-method for the flowgramFormat slot.
- flowgramFormat** signature(object = "SFFContainer"): Getter-method for the flowgramFormat slot.
- flowgrams<-** signature(object = "SFFContainer", value = "list"): Setter-method for the flowgrams slot.
- flowgrams** signature(object = "SFFContainer"): Getter-method for the flowgrams slot.
- flowIndexes<-** signature(object = "SFFContainer", value = "list"): Setter-method for the flowIndexes slot.
- flowIndexes** signature(object = "SFFContainer"): Getter-method for the flowIndexes slot.
- keySequence<-** signature(object = "SFFContainer", value = "character"): Setter-method for the keySequence slot.
- keySequence** signature(object = "SFFContainer"): Getter-method for the keySequence slot.
- reads<-** signature(object = "SFFContainer", value = "QualityScaledDNAStringSet"): Setter-method for the reads slot.
- reads** signature(object = "SFFContainer"): Getter-method for the reads slot.
- [ signature(x = "SFFContainer", i = "ANY", j = "ANY"): Subsetting a SFFContainer object.

**Author(s)**

Christian Ruckert

**See Also**[readSFF](#), [SFFRead](#)**Examples**

```
showClass("SFFContainer")
```

---

*SFFRead-class**Class "SFFRead"*

---

**Description**

This class is a container for a single read from files in Roche's Standard Flowgram Format (SFF).

**Objects from the Class**

Objects can be created by calls of the form `new("SFFRead", ...)`. Usually, objects will be created by calling the [getRead](#) method on an object of class [SFFContainer](#).

**Slots**

**name:** Object of class "character" representing the name of the read.

**read:** Object of class "DNASTring" containing the basecalled nucleotide sequence of the read.

**flowgramFormat:** Object of class "numeric" representing the format used to encode each of the flowgram values for each read. Currently, only one flowgram format has been adopted and is coded by the value 1.

**flowChars:** Object of class "character" containing the array of nucleotide bases ('A', 'C', 'G', 'T') that correspond to the nucleotides used for each flow of each read.

**keySequence:** Object of class "character" representing the nucleotide bases of the key sequence used for these reads.

**clipQualityLeft:** Object of class "numeric" representing the position of the first base after the clipping point for an attached quality sequence. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityLeft. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.

**clipQualityRight:** Object of class "numeric" representing the position of the last base before the clipping point for an attached quality sequence. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityRight. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.

**clipAdapterLeft:** Object of class "numeric" representing the position of the first base after the clipping point for an attached adapter sequence. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityLeft. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.

**clipAdapterRight:** Object of class "numeric" representing the position of the last base before the clipping point for an attached adapter sequence. If only a combined (quality+adapter) clipping position is computed it should be stored in clipQualityRight. If no clipping value is computed the field is set to 0. The position values use 1-based indexing.

**flowgram**: Object of class "numeric" containing the homopolymer stretch estimates for each flow.

**flowIndexes**: Object of class "numeric" containing the flow positions for each base in the called sequence, i.e. for each base, the position in the flowgram whose estimate resulted in that base being called.

**quality**: Object of class "BString" containing the quality scores for each of the bases in the sequence, where the values use the standard  $-\log_{10}$  probability scale.

## Methods

**read<-** signature(object = "SFFRead", value = "DNAStrng"): Setter-method for the read slot.

**read** signature(object = "SFFRead"): Getter-method for the read slot.

**flowChars<-** signature(object = "SFFContainer", value = "character"): Setter-method for the flowChars slot.

**flowChars** signature(object = "SFFContainer"): Getter-method for the flowChars slot.

**flowgramFormat<-** signature(object = "SFFContainer", value = "numeric"): Setter-method for the flowgramFormat slot.

**flowgramFormat** signature(object = "SFFContainer"): Getter-method for the flowgramFormat slot.

**keySequence<-** signature(object = "SFFContainer", value = "character"): Setter-method for the keySequence slot.

**keySequence** signature(object = "SFFContainer"): Getter-method for the keySequence slot.

**clipAdapterLeft<-** signature(object = "SFFRead", value = "numeric"): Setter-method for the clipAdapterLeft slot.

**clipAdapterLeft** signature(object = "SFFRead"): Getter-method for the clipAdapterLeft slot.

**clipAdapterRight<-** signature(object = "SFFRead", value = "numeric"): Setter-method for the clipAdapterRight slot.

**clipAdapterRight** signature(object = "SFFRead"): Getter-method for the clipAdapterRight slot.

**clipQualityLeft<-** signature(object = "SFFRead", value = "numeric"): Setter-method for the clipQualityLeft slot.

**clipQualityLeft** signature(object = "SFFRead"): Getter-method for the clipQualityLeft slot.

**clipQualityRight<-** signature(object = "SFFRead", value = "numeric"): Setter-method for the clipQualityRight slot.

**clipQualityRight** signature(object = "SFFRead"): Getter-method for the clipQualityRight slot.

**flowgram<-** signature(object = "SFFRead", value = "numeric"): Setter-method for the flowgram slot.

**flowgram** signature(object = "SFFRead"): Getter-method for the flowgram slot.

**flowIndexes<-** signature(object = "SFFRead", value = "numeric"): Setter-method for the flowIndexes slot.

**flowIndexes** signature(object = "SFFRead"): Getter-method for the flowIndexes slot.

**name<-** signature(object = "SFFRead", value = "character"): Setter-method for the name slot.

**name** signature(object = "SFFRead"): Getter-method for the name slot.

**quality<-** signature(object = "SFFRead", value = "BString"): Setter-method for the quality slot.

**quality** signature(object = "SFFRead"): Getter-method for the quality slot.

**Author(s)**

Christian Ruckert

**See Also**[readSFF](#), [SFFContainer](#)**Examples**

```
showClass("SFFRead")
```

---

`variants`*Example data for [plotVariants](#)*

---

**Description**

This data.frame is part of the vignette example of the [plotVariants](#) function. It has the columns "label", "pos" "mutation" and "color" specifying an annotation for each mutation, its position, the mutation type and an individual color. The position is given as amino acids / codons.

**Usage**

```
data(plotVariantsExample)
```

**Format**

data.frame

**Examples**

```
data(plotVariantsExample)
variants
```

# Index

- \*Topic **TiTv, Transition, Transversion**
  - calculateTiTv, 18
- \*Topic **Variation Frequency Coverage**
  - plotVariationFrequency, 49
- \*Topic **alignShortReads**
  - alignShortReads, 3
- \*Topic **amplicon coverage**
  - plotAmpliconCoverage, 44
- \*Topic **annotateVariants**
  - annotateVariants, 5
- \*Topic **classes**
  - AnnotatedVariants-class, 4
  - AVASet-class, 9
  - Breakpoints-class, 16
  - MapperSet-class, 40
  - SFFContainer-class, 59
  - SFFRead-class, 61
- \*Topic **datasets**
  - avaSetExample, 12
  - avaSetFiltered, 13
  - avaSetFiltered\_annot, 13
  - breakpoints, 16
  - captureArray, 19
  - mapperSetExample, 41
  - mutationInfo, 43
  - regions, 54
  - variants, 63
- \*Topic **demultiplex**
  - demultiplexReads, 23
  - genomeSequencerMIDs, 33
- \*Topic **filterChimericReads**
  - filterChimericReads, 29
- \*Topic **linker**
  - sequenceCaptureLinkers, 56
- \*Topic **methods**
  - ava2vcf, 7
  - baseFrequency, 14
  - baseQualityHist, 15
  - baseQualityStats, 15
  - complexity.dust, 20
  - complexity.entropy, 21
  - dinucleotideOddsRatio, 27
  - flowgramBarplot, 31
  - gcContent, 31
  - gcContentHist, 32
  - gcPerPosition, 32
  - homopolymerHist, 36
  - nucleotideCharts, 44
  - positionQualityBoxplot, 51
  - readLengthHist, 51
  - readLengthStats, 52
  - sequenceQualityHist, 57
  - sff2fastq, 58
- \*Topic **plotChimericReads, detectBreakpoints, mergeBreakpoints, Breakpoints**
  - plotChimericReads, 45
- \*Topic **readsOnTarget**
  - readsOnTarget, 53
- \*Topic **removeLinker**
  - removeLinker, 55
- [,AVASet,ANY,ANY-method (AVASet-class), 9
- [,Breakpoints,ANY,ANY-method (Breakpoints-class), 16
- [,SFFContainer,ANY,ANY-method (SFFContainer-class), 59
- addRead (SFFContainer-class), 59
- addRead,SFFContainer,SFFRead-method (SFFContainer-class), 59
- AlignedRead, 3, 25
- alignedReadsC1 (Breakpoints-class), 16
- alignedReadsC1,Breakpoints-method (Breakpoints-class), 16
- alignedReadsC1<- (Breakpoints-class), 16
- alignedReadsC1<- ,Breakpoints,list-method (Breakpoints-class), 16
- alignedReadsC2 (Breakpoints-class), 16
- alignedReadsC2,Breakpoints-method (Breakpoints-class), 16
- alignedReadsC2<- (Breakpoints-class), 16
- alignedReadsC2<- ,Breakpoints,list-method (Breakpoints-class), 16
- alignShortReads, 3, 9, 11, 54





- clipAdapterRight,SFFContainer-method  
(SFFContainer-class), 59
- clipAdapterRight,SFFRead-method  
(SFFRead-class), 61
- clipAdapterRight<- (SFFContainer-class),  
59
- clipAdapterRight<-,SFFContainer,numeric-method  
(SFFContainer-class), 59
- clipAdapterRight<-,SFFRead,numeric-method  
(SFFRead-class), 61
- clipQualityLeft (SFFContainer-class), 59
- clipQualityLeft,SFFContainer-method  
(SFFContainer-class), 59
- clipQualityLeft,SFFRead-method  
(SFFRead-class), 61
- clipQualityLeft<- (SFFContainer-class), 59
- clipQualityLeft<-,SFFContainer,numeric-method  
(SFFContainer-class), 59
- clipQualityLeft<-,SFFRead,numeric-method  
(SFFRead-class), 61
- clipQualityRight (SFFContainer-class), 59
- clipQualityRight,SFFContainer-method  
(SFFContainer-class), 59
- clipQualityRight,SFFRead-method  
(SFFRead-class), 61
- clipQualityRight<- (SFFContainer-class),  
59
- clipQualityRight<-,SFFContainer,numeric-method  
(SFFContainer-class), 59
- clipQualityRight<-,SFFRead,numeric-method  
(SFFRead-class), 61
- commonAlignC1 (Breakpoints-class), 16
- commonAlignC1,Breakpoints-method  
(Breakpoints-class), 16
- commonAlignC1<- (Breakpoints-class), 16
- commonAlignC1<-,Breakpoints,list-method  
(Breakpoints-class), 16
- commonAlignC2 (Breakpoints-class), 16
- commonAlignC2,Breakpoints-method  
(Breakpoints-class), 16
- commonAlignC2<- (Breakpoints-class), 16
- commonAlignC2<-,Breakpoints,list-method  
(Breakpoints-class), 16
- commonBpsC1 (Breakpoints-class), 16
- commonBpsC1,Breakpoints-method  
(Breakpoints-class), 16
- commonBpsC1<- (Breakpoints-class), 16
- commonBpsC1<-,Breakpoints,list-method  
(Breakpoints-class), 16
- commonBpsC2 (Breakpoints-class), 16
- commonBpsC2,Breakpoints-method  
(Breakpoints-class), 16
- commonBpsC2<- (Breakpoints-class), 16
- commonBpsC2<-,Breakpoints,list-method  
(Breakpoints-class), 16
- complexity.dust, 20
- complexity.dust,DNAStringSet-method  
(complexity.dust), 20
- complexity.dust,SFFContainer-method  
(complexity.dust), 20
- complexity.dust,ShortRead-method  
(complexity.dust), 20
- complexity.entropy, 21
- complexity.entropy,DNAStringSet-method  
(complexity.entropy), 21
- complexity.entropy,SFFContainer-method  
(complexity.entropy), 21
- complexity.entropy,ShortRead-method  
(complexity.entropy), 21
- convertCigar, 22
- coverageOnTarget, 22
- coverageOnTarget,list,RangesList-method  
(coverageOnTarget), 22
- data, 47
- data.frame, 14
- DataFrame, 25
- demultiplexReads, 23, 33
- demultiplexReads,XStringSet,XStringSet,missing,logical-method  
(demultiplexReads), 23
- demultiplexReads,XStringSet,XStringSet,missing,missing-method  
(demultiplexReads), 23
- demultiplexReads,XStringSet,XStringSet,numeric,logical-method  
(demultiplexReads), 23
- demultiplexReads,XStringSet,XStringSet,numeric,missing-method  
(demultiplexReads), 23
- detectBreakpoints, 16, 18, 24, 30, 42, 43, 46
- detectBreakpoints,list-method  
(detectBreakpoints), 24
- dinucleotideOddsRatio, 27
- dinucleotideOddsRatio,DNAStringSet-method  
(dinucleotideOddsRatio), 27
- dinucleotideOddsRatio,SFFContainer-method  
(dinucleotideOddsRatio), 27
- dinucleotideOddsRatio,ShortRead-method  
(dinucleotideOddsRatio), 27
- DNAStringSet, 3, 14, 20, 21, 24, 27, 31–34,  
44, 51, 52, 56
- eSet, 11, 40
- extendedCIGARToList, 22
- extendedCIGARToList (convertCigar), 22
- fData, 37
- fDataAmp, 6, 28, 29, 34

- fDataAmp, AVASet-method (AVASet-class), 9
- featureDataAmp, 6, 28, 28
- featureDataAmp, AVASet-method (AVASet-class), 9
- featureDataAmp<- (AVASet-class), 9
- featureDataAmp<-, AVASet, AnnotatedDataFrame-method (AVASet-class), 9
- filterChimericReads, 18, 24, 26, 29, 46
- filterChimericReads, list, missing, DNASString, missing-method (filterChimericReads), 29
- filterChimericReads, list, missing, DNASString, numeric, numeric-method (filterChimericReads), 29
- filterChimericReads, list, missing, missing, missing, missing-method (filterChimericReads), 29
- filterChimericReads, list, missing, missing, numeric, numeric-method (filterChimericReads), 29
- filterChimericReads, list, RangesList, DNASString, missing, missing-method (filterChimericReads), 29
- filterChimericReads, list, RangesList, DNASString, numeric, numeric-method (filterChimericReads), 29
- filterChimericReads, list, RangesList, missing, missing, missing, missing-method (filterChimericReads), 29
- filterChimericReads, list, RangesList, missing, numeric, numeric-method (filterChimericReads), 29
- flowChars (SFFContainer-class), 59
- flowChars, SFFContainer-method (SFFContainer-class), 59
- flowChars, SFFRead-method (SFFRead-class), 61
- flowChars<- (SFFContainer-class), 59
- flowChars<-, SFFContainer, character-method (SFFContainer-class), 59
- flowChars<-, SFFRead, character-method (SFFRead-class), 61
- flowgram (SFFRead-class), 61
- flowgram, SFFRead-method (SFFRead-class), 61
- flowgram<- (SFFRead-class), 61
- flowgram<-, SFFRead, numeric-method (SFFRead-class), 61
- flowgramBarplot, 31
- flowgramBarplot, SFFRead-method (flowgramBarplot), 31
- flowgramFormat (SFFContainer-class), 59
- flowgramFormat, SFFContainer-method (SFFContainer-class), 59
- flowgramFormat, SFFRead-method (SFFRead-class), 61
- flowgramFormat<- (SFFContainer-class), 59
- flowgramFormat<-, SFFContainer, numeric-method (SFFContainer-class), 59
- flowgramFormat<-, SFFRead, numeric-method (SFFRead-class), 61
- flowgrams (SFFContainer-class), 59
- flowgrams, SFFContainer-method (SFFContainer-class), 59
- flowgrams<- (SFFContainer-class), 59
- flowgrams<-, SFFContainer, list-method (SFFContainer-class), 59
- flowIndexes (SFFContainer-class), 59
- flowIndexes, SFFContainer-method (SFFContainer-class), 59
- flowIndexes, SFFRead-method (SFFRead-class), 61
- flowIndexes<- (SFFContainer-class), 59
- flowIndexes<-, SFFContainer, list-method (SFFContainer-class), 59
- flowIndexes<-, SFFRead, numeric-method (SFFRead-class), 61
- gcContent, 31
- gcContent, DNASStringSet-method (gcContent), 31
- gcContent, SFFContainer-method (gcContent), 31
- gcContent, ShortRead-method (gcContent), 31
- gcContentHist, 32
- gcContentHist, DNASStringSet-method (gcContentHist), 32
- gcContentHist, SFFContainer-method (gcContentHist), 32
- gcContentHist, ShortRead-method (gcContentHist), 32
- gcPerPosition, 32
- gcPerPosition, DNASStringSet-method (gcPerPosition), 32
- gcPerPosition, SFFContainer-method (gcPerPosition), 32
- gcPerPosition, ShortRead-method (gcPerPosition), 32
- genomeSequencerMIDs, 24, 33
- genomeSequencerMIDs, character-method (genomeSequencerMIDs), 33
- genomeSequencerMIDs, missing-method (genomeSequencerMIDs), 33
- getAlignedReads, 34
- getAlignedReads, AVASet-method (getAlignedReads), 34
- getAminoAbbr, 35
- getRead, 61
- getRead (SFFContainer-class), 59

- getRead,SFFContainer,character-method (SFFContainer-class), 59
- getReadStatus (MapperSet-class), 40
- getReadStatus,MapperSet-method (MapperSet-class), 40
- getVariantPercentages, 11, 35, 41, 57, 58
- getVariantPercentages,AVASet-method (AVASet-class), 9
- getVariantPercentages,MapperSet-method (MapperSet-class), 40
  
- homopolymerHist, 36
- homopolymerHist,SFFRead-method (homopolymerHist), 36
- htmlReport, 4, 6, 11, 37, 41
- htmlReport,AVASet-method (AVASet-class), 9
- htmlReport,MapperSet-method (MapperSet-class), 40
  
- IRanges, 25
  
- keySequence (SFFContainer-class), 59
- keySequence,SFFContainer-method (SFFContainer-class), 59
- keySequence,SFFRead-method (SFFRead-class), 61
- keySequence<- (SFFContainer-class), 59
- keySequence<-,SFFContainer,character-method (SFFContainer-class), 59
- keySequence<-,SFFRead,character-method (SFFRead-class), 61
  
- length,Breakpoints-method (Breakpoints-class), 16
- listToExtendedCIGAR (convertCigar), 22
  
- MapperSet, 5, 38
- MapperSet,character-method (MapperSet), 38
- MapperSet-class, 40
- mapperSetExample, 41
- matchPDict, 3
- mergeBreakpoints, 16, 18, 24–26, 30, 42, 46
- mergeBreakpoints,Breakpoints,missing,list-method (Breakpoints-class), 16
- mergeBreakpoints,Breakpoints,missing,missing-method (Breakpoints-class), 16
- mergeBreakpoints,Breakpoints,numeric,missing-method (Breakpoints-class), 16
- mutationInfo, 43
- name (SFFContainer-class), 59
- name,SFFContainer-method (SFFContainer-class), 59
- name,SFFRead-method (SFFRead-class), 61
- name<- (SFFContainer-class), 59
- name<-,SFFContainer,character-method (SFFContainer-class), 59
- name<-,SFFRead,character-method (SFFRead-class), 61
- names,AnnotatedVariants-method (AnnotatedVariants-class), 4
- names,Breakpoints-method (Breakpoints-class), 16
- names<-,AnnotatedVariants,character-method (AnnotatedVariants-class), 4
- names<-,Breakpoints,ANY-method (Breakpoints-class), 16
- nucleotideCharts, 44
- nucleotideCharts,DNAStringSet-method (nucleotideCharts), 44
- nucleotideCharts,SFFContainer-method (nucleotideCharts), 44
- nucleotideCharts,ShortRead-method (nucleotideCharts), 44
  
- PairwiseAlignedFixedSubject, 25
- pairwiseAlignment, 56
- plotAmpliconCoverage, 44
- plotAmpliconCoverage,AVASet,character,logical-method (plotAmpliconCoverage), 44
- plotAmpliconCoverage,AVASet,character,missing-method (plotAmpliconCoverage), 44
- plotAmpliconCoverage,AVASet,missing,logical-method (plotAmpliconCoverage), 44
- plotAmpliconCoverage,AVASet,missing,missing-method (plotAmpliconCoverage), 44
- plotChimericReads, 18, 26, 43, 45
- plotChimericReads,Breakpoints-method (Breakpoints-class), 16
- plotVariants, 43, 47, 54, 63
- plotVariants,AnnotatedVariants,character-method (plotVariants), 47
- plotVariants,data.frame,character-method (plotVariants), 47
- plotVariationFrequency, 49
- plotVariationFrequency,character,numeric-method (plotVariationFrequency), 49
- positionQualityBoxplot, 51
- positionQualityBoxplot,QualityScaledDNAStringSet-method (positionQualityBoxplot), 51
- positionQualityBoxplot,SFFContainer-method (positionQualityBoxplot), 51

- positionQualityBoxplot,ShortReadQ-method  
(positionQualityBoxplot), 51
- quality,SFFRead-method (SFFRead-class),  
61
- quality<- (SFFRead-class), 61
- quality<-,SFFRead,BString-method  
(SFFRead-class), 61
- QualityScaledDNAStrngSet, 15, 51, 57
- read (SFFRead-class), 61
- read,SFFRead-method (SFFRead-class), 61
- read<- (SFFRead-class), 61
- read<-,SFFRead,QualityScaledDNAStrngSet-method  
(SFFRead-class), 61
- readLengthHist, 51
- readLengthHist,DNAStrngSet-method  
(readLengthHist), 51
- readLengthHist,SFFContainer-method  
(readLengthHist), 51
- readLengthHist,ShortRead-method  
(readLengthHist), 51
- readLengthStats, 52
- readLengthStats,DNAStrngSet-method  
(readLengthStats), 52
- readLengthStats,SFFContainer-method  
(readLengthStats), 52
- readLengthStats,ShortRead-method  
(readLengthStats), 52
- reads (SFFContainer-class), 59
- reads,SFFContainer-method  
(SFFContainer-class), 59
- reads<- (SFFContainer-class), 59
- reads<-,SFFContainer,QualityScaledDNAStrngSet-method  
(SFFContainer-class), 59
- readSFF, 59, 61, 63
- readSFF (readSff), 52
- readSff, 52
- readsOnTarget, 53
- readsOnTarget,list,RangesList-method  
(readsOnTarget), 53
- referenceSequences, 54
- referenceSequences,AVASet-method  
(AVASet-class), 9
- referenceSequences<- (AVASet-class), 9
- referenceSequences<-,AVASet,AlignedRead-method  
(AVASet-class), 9
- regions, 54
- removeLinker, 55, 56
- removeLinker,XStringSet,DNAStrng,logical,numeric,numeric-method  
(removeLinker), 55
- removeLinker,XStringSet,DNAStrng,missing,missing,missing-method  
(removeLinker), 55
- scanBam, 23, 24, 29, 30, 53
- seqsC1 (Breakpoints-class), 16
- seqsC1,Breakpoints-method  
(Breakpoints-class), 16
- seqsC1<- (Breakpoints-class), 16
- seqsC1<-,Breakpoints,list-method  
(Breakpoints-class), 16
- seqsC2 (Breakpoints-class), 16
- seqsC2,Breakpoints-method  
(Breakpoints-class), 16
- seqsC2<- (Breakpoints-class), 16
- seqsC2<-,Breakpoints,list-method  
(Breakpoints-class), 16
- sequenceCaptureLinkers, 30, 56, 56
- sequenceCaptureLinkers,character-method  
(sequenceCaptureLinkers), 56
- sequenceCaptureLinkers,missing-method  
(sequenceCaptureLinkers), 56
- sequenceQualityHist, 57
- sequenceQualityHist,QualityScaledDNAStrngSet-method  
(sequenceQualityHist), 57
- sequenceQualityHist,SFFContainer-method  
(sequenceQualityHist), 57
- sequenceQualityHist,ShortReadQ-method  
(sequenceQualityHist), 57
- setVariantFilter, 10, 11, 13, 35, 40, 41, 57
- setVariantFilter,AVASet-method  
(AVASet-class), 9
- setVariantFilter,MapperSet-method  
(MapperSet-class), 40
- sff2fastq, 58
- sff2fastq,SFFContainer-method (sff2fastq),  
58
- SFFContainer, 14, 15, 20, 21, 27, 31–33, 44,  
51–53, 57, 58, 60, 61, 63
- SFFContainer (SFFContainer-class), 59
- SFFContainer-class, 59
- SFFRead, 31, 36, 60, 61
- SFFRead (SFFRead-class), 61
- SFFRead-class, 61
- ShortRead, 14, 20, 21, 27, 31–33, 44, 51, 52
- ShortReadQ, 15, 51, 57
- variants, 63
- vector, 52
- Versioned, 11, 40
- VersionedBiobase, 11, 40
- writeVcf, 7