# Package 'Biostrings'

October 31, 2025

Title Efficient manipulation of biological strings

**Description** Memory efficient string containers, string matching algorithms, and other utilities, for fast manipulation of large biological sequences or sets of sequences.

**biocViews** SequenceMatching, Alignment, Sequencing, Genetics, DataImport, DataRepresentation, Infrastructure

URL https://bioconductor.org/packages/Biostrings

BugReports https://github.com/Bioconductor/Biostrings/issues

Version 2.79.1

License Artistic-2.0

**Encoding UTF-8** 

**Depends** R (>= 4.1.0), BiocGenerics (>= 0.37.0), S4Vectors (>= 0.27.12), IRanges (>= 2.31.2), XVector (>= 0.37.1), Seqinfo

Imports methods, utils, grDevices, stats, crayon

LinkingTo S4Vectors, IRanges, XVector

**Suggests** graphics, pwalign, BSgenome (>= 1.13.14),

BSgenome.Celegans.UCSC.ce2 (>= 1.3.11),

BSgenome.Dmelanogaster.UCSC.dm3 (>= 1.3.11),

BSgenome. Hsapiens. UCSC. hg18, drosophila 2 probe, hgu 95 av 2 probe,

hgu133aprobe, GenomicFeatures (>= 1.3.14), hgu95av2cdf, affy

(>= 1.41.3), affydata (>= 1.11.5), RUnit, BiocStyle, knitr,

testthat (>= 3.0.0), covr

# VignetteBuilder knitr

Collate utils.R IUPAC\_CODE\_MAP.R AMINO\_ACID\_CODE.R GENETIC\_CODE.R

XStringCodec-class.R seqtype.R coloring.R XString-class.R XStringSet-class.R XStringSet-comparison.R XStringViews-class.R MaskedXString-class.R XStringSetList-class.R seqinfo-methods.R xscat.R XStringSet-io.R letter.R getSeq.R letterFrequency.R dinucleotideFrequencyTest.R chartr.R reverseComplement.R translate.R toComplex.R replaceAt.R replaceLetterAt.R injectHardMask.R padAndClip.R strsplit-methods.R misc.R

2 Contents

match-utils.R matchPattern.R maskMotif.R matchLRPatterns.R trimLRPatterns.R matchProbePair.R matchPWM.R findPalindromes.R PDict-class.R matchPDict.R XStringQuality-class.R QualityScaledXStringSet.R lcsuffix.R MultipleAlignment.R longestConsecutive.R predefined\_scoring\_matrices.R zzz.R git\_url https://git.bioconductor.org/packages/Biostrings git\_branch devel git\_last\_commit 778d5f9 git\_last\_commit\_date 2025-10-30 Repository Bioconductor 3.23 Date/Publication 2025-10-31 Author Hervé Pagès [aut, cre], Patrick Aboyoun [aut], Robert Gentleman [aut], Saikat DebRoy [aut], Vince Carey [ctb], Nicolas Delhomme [ctb], Felix Ernst [ctb], Wolfgang Huber [ctb] ('matchprobes' vignette), Beryl Kanali [ctb] (Converted 'MultipleAlignments' vignette from Sweave to RMarkdown), Haleema Khan [ctb] (Converted 'matchprobes' vignette from Sweave to RMarkdown), Aidan Lakshman [ctb], Kieran O'Neill [ctb], Valerie Obenchain [ctb], Marcel Ramos [ctb], Albert Vill [ctb], Jen Wokaty [ctb] (Converted 'matchprobes' vignette from Sweave to

SparseList-class.R MIndex-class.R lowlevel-matching.R

Maintainer Hervé Pagès <hpages.on.github@gmail.com>

# **Contents**

RMarkdown), Erik Wright [ctb]

AAString-clas	s									 								4
AMINO_ACI	D_CODE									 								5
Biostrings inte	rnals									 								6
chartr										 								6
coloring										 								8
detail										 								10
dinucleotideFi	equencyT	est								 								11
DNAString-cl	ass									 								12
findPalindrom																		

Contents 3

GENETIC_CODE	16
getSeq	19
gregexpr2	20
HNF4alpha	21
injectHardMask	22
IUPAC_CODE_MAP	23
lesuffix	25
letter	25
letterFrequency	26
longestConsecutive	32
lowlevel-matching	33
MaskedXString-class	38
maskMotif	40
match-utils	42
matchLRPatterns	44
matchPattern	46
matchPDict	<b>5</b> 0
matchPDict-inexact	<b>5</b> 9
matchProbePair	63
matchPWM	65
MIndex-class	67
misc	69
MultipleAlignment-class	70
nucleotideFrequency	74
padAndClip	<del>7</del> 9
PDict-class	81
predefined_scoring_matrices	85
QualityScaledXStringSet-class	86
replaceAt	88
replaceLetterAt	92
reverseComplement	95
RNAString-class	97
seqinfo-methods	99
toComplex	
translate	
trimLRPatterns	04
xscat	07
XString-class	08
XStringQuality-class	10
XStringSet-class	12
XStringSet-comparison	18
XStringSet-io	21
XStringSetList-class	28
XStringViews-class	30
yeastSEQCHR1	33

134

Index

4 AAString-class

AAString-class AAString objects

### **Description**

An AAString object allows efficient storage and manipulation of a long amino acid sequence.

#### Usage

```
AAString(x="", start=1, nchar=NA)

## Predefined constants:

AA_ALPHABET # full Amino Acid alphabet
AA_STANDARD # first 20 letters only
AA_PROTEINOGENIC # first 22 letters only
```

# **Arguments**

x A single string.

start, nchar Where to start reading from in x and how many letters to read.

#### **Details**

The AAString class is a direct XString subclass (with no additional slot). Therefore all functions and methods described in the XString man page also work with an AAString object (inheritance).

Unlike the BString container that allows storage of any single string (based on a single-byte character set) the AAString container can only store a string based on the Amino Acid alphabet (see below).

### The Amino Acid alphabet

This alphabet contains all letters from the Single-Letter Amino Acid Code (see ?AMINO\_ACID\_CODE) plus "\*" (the *stop* letter), "-" (the *gap* letter), "+" (the *hard masking* letter), and "." (the *not a letter* or *not available* letter). It is stored in the AA\_ALPHABET predefined constant (character vector).

The alphabet() function returns AA\_ALPHABET when applied to an AAString object.

#### Constructor-like functions and generics

In the code snippet below, x can be a single string (character vector of length 1) or a BString object.

AAString(x="", start=1, nchar=NA): Tries to convert x into an AAString object by reading nchar letters starting at position start in x.

# **Accessor methods**

In the code snippet below, x is an AAString object.

alphabet(x): If x is an AAString object, then return the Amino Acid alphabet (see above). See the corresponding man pages when x is a BString, DNAString or RNAString object.

AMINO\_ACID\_CODE

5

# **Display**

The letters in an AAString object are colored when displayed by the show() method. Set global option Biostrings.coloring to FALSE to turn off this coloring.

# Author(s)

H. Pagès

# See Also

```
AMINO_ACID_CODE, letter, XString-class, alphabetFrequency
```

# **Examples**

```
AA_ALPHABET
a <- AAString("MARKSLEMSIR*")
length(a)
alphabet(a)</pre>
```

AMINO\_ACID\_CODE

The Single-Letter Amino Acid Code

# Description

Named character vector mapping single-letter amino acid representations to 3-letter amino acid representations.

# See Also

```
AAString, GENETIC_CODE
```

```
## See all the 3-letter codes
AMINO_ACID_CODE

## Convert an AAString object to a vector of 3-letter amino acid codes
aa <- AAString("LANDEECQW")
AMINO_ACID_CODE[strsplit(as.character(aa), NULL)[[1]]]</pre>
```

6 chartr

Biostrings internals Biostrings internals

# Description

Biostrings objects, classes and methods that are not intended to be used directly.

### Author(s)

H. Pagès

chartr

Replace letters in a sequence or set of sequences

# **Description**

Replace letters in a sequence or set of sequences.

# Usage

```
## $4 method for signature 'ANY,ANY,XString'
chartr(old, new, x)
replaceAmbiguities(x, new="N")
```

# **Arguments**

old	A character string specifying the characters to be replaced.
new	A character string specifying the replacements. It must be a single letter for replaceAmbiguities.
Х	The sequence or set of sequences to translate. If x is an XString, XStringSet, XStringViews or MaskedXString object, then the appropriate chartr method is called, otherwise the standard chartr R function is called.

#### **Details**

See ?chartr for the details.

Note that, unlike the standard chartr R function, the methods for XString, XStringSet, XStringViews and MaskedXString objects do NOT support character ranges in the specifications.

replaceAmbiguities() is a simple wrapper around chartr() that replaces all IUPAC ambiguities with N for objects containing DNA or RNA sequence data.

# Value

An object of the same class and length as the original object.

chartr 7

#### See Also

- chartr in the base package.
- The replaceAt function for extracting or replacing arbitrary subsequences from/in a sequence or set of sequences.
- The replaceLetterAt function for a DNA-specific single-letter replacement functions useful for SNP injections.
- IUPAC\_CODE\_MAP for the mapping between IUPAC nucleotide ambiguity codes and their meaning.
- alphabetFrequency (and uniqueLetters) for tabulating letters in (and extracting the unique letters from) a sequence or set of sequences.
- The XString, XStringSet, XStringViews, and MaskedXString classes.

```
## A BASIC chartr() EXAMPLE
## -----
x <- BString("MiXeD cAsE 123")</pre>
chartr("iXs", "why", x)
## TRANSFORMING DNA WITH BISULFITE (AND SEARCHING IT...)
## -----
library(BSgenome.Celegans.UCSC.ce2)
chrII <- Celegans[["chrII"]]</pre>
alphabetFrequency(chrII)
pattern <- DNAString("TGGGTGTATTTA")</pre>
## Transforming and searching the + strand
plus_strand <- chartr("C", "T", chrII)</pre>
alphabetFrequency(plus_strand)
matchPattern(pattern, plus_strand)
matchPattern(pattern, chrII)
## Transforming and searching the - strand
minus_strand <- chartr("G", "A", chrII)</pre>
alphabetFrequency(minus_strand)
matchPattern(reverseComplement(pattern), minus_strand)
matchPattern(reverseComplement(pattern), chrII)
## -----
## replaceAmbiguities()
dna <- DNAStringSet(c("TTTKYTT-GR", "", "NAASACVT"))</pre>
dna
replaceAmbiguities(dna)
```

8 coloring

coloring

XString Display Colors

# **Description**

XString objects support custom coloring for display. Users can also set custom color palettes for XString objects using the update\_X\_palette functions.

# Usage

```
update_DNA_palette(colors=NULL)
update_RNA_palette(colors=NULL)
update_AA_palette(colors=NULL)
update_B_palette(colors=NULL)
```

#### **Arguments**

colors

A named list of colors to update, with entries fg and bg specifying the foreground and background colors, respectively. Colors can be specified in any way compatible with make\_style from the crayon package. Defaults to NULL, which resets the color palette to the default color scheme. See Details and Examples for more information.

#### **Details**

XString objects support the following default coloring for display.

- DNAString: A, C, G, and T are colored red, green, blue, and orange (respectively), N is colored light grey, other ambiguity codes are colored dark grey, and "-+." have no coloring.
- RNAString: All bases are colored identically to DNAString. U is colored yellow.
- AAString: Amino acids are colored according to JalView's Zappo color scheme, representing
  physicochemical properties. X is colored light grey, other ambiguity codes are colored dark
  grey, and "\*-+." are not colored.
- BStrings are not colored.

Users can change the default color scheme of Biostrings with the update\_X\_palette family functions. Each function expects a list with named entries corresponding to the values to update. Each entry can specify 'fg' and 'bg' values, corresponding to the foreground and background colors (respectively). If 'fg' is not specified, it defaults to rgb(1,1,1) (white). If 'bg' is not specified, it defaults to transparent.

These functions will only update the values passed, leaving the rest of the colors as-is. For example, calling update\_AA\_palette(list(A=list(fg="green"))) would update the coloring for A while leaving all other colors as the default schema.

To reset all colors to the default palette, call the function with no arguments (NULL).

To remove a coloring for a specific value, provide a named entry with value NULL. For example, update\_AA\_palette(list(A=NULL)) will remove the coloring for A.

coloring 9

update\_DNA\_palette and update\_RNA\_palette are identical internally, so either function can be used to update colorings for T,U.

See the Examples section for more examples of custom colorings.

# Value

For update\_X\_palette, Invisibly returns the new color mapping, consisting of a named character vector. Calling cat on the return value will print out all letters with their respective coloring.

#### Author(s)

Aidan Lakshman < AHL27@pitt.edu>

### See Also

XString-class

```
## display default colors
DNAString(paste(DNA_ALPHABET, collapse=''))
RNAString(paste(RNA_ALPHABET, collapse=''))
AAString(paste(AA_ALPHABET, collapse=''))
BString(paste(LETTERS, collapse=''))
## create new palettes
DNA_palette <- list(</pre>
  A=list(fg="blue",bg="black"),
  T=list(fg="red",bg='black'),
  G=list(fg='green',bg='black'),
  C=list(fg='yellow',bg='black')
update_DNA_palette(DNA_palette)
DNAString(paste(DNA_ALPHABET, collapse=''))
## reset to default palette
update_DNA_palette()
DNAString(paste(DNA_ALPHABET, collapse=''))
## colors can also be specified with `rgb()`
AA_palette <- list(
  A=list(fg="white", bg="purple"),
  B=list(fg=rgb(1,1,1), bg='orange')
update_AA_palette(AA_palette)
AAString(paste(AA_ALPHABET, collapse=''))
## remove all coloring for QEG
update_AA_palette(list(Q=NULL, E=NULL, G=NULL))
AAString(paste(AA_ALPHABET, collapse=''))
## reset to default
```

10 detail

```
update_AA_palette()
AAString(paste(AA_ALPHABET, collapse=''))
## We can also add colors to BStrings,
## which are normally not colored
## if 'fg' is not specified, defaults to rgb(1,1,1)
## if 'bg' is not specified, background is transparent
B_palette <- list(</pre>
  A=list(bg='green'),
  B=list(bg="red"),
  C=list(bg='blue'),
  D=list(fg="orange"),
  E=list(fg="yellow")
)
update_B_palette(B_palette)
BString(paste(LETTERS, collapse=''))
## can also directly view the changes with cat
cat(update_B_palette(B_palette), '\n')
## reset to default
update_B_palette()
BString(paste(LETTERS, collapse=''))
```

detail

Show (display) detailed object content

### **Description**

This is a variant of show, offering a more detailed display of object content.

# Usage

```
detail(x, ...)
```

#### **Arguments**

x An object. The default simply invokes show.

... Additional arguments. The default definition makes no use of these arguments.

### Value

None; the function is invoked for its side effect (detailed display of object content).

# Author(s)

Martin Morgan

### **Examples**

dinucleotideFrequencyTest

Pearson's chi-squared Test and G-tests for String Position Dependence

### **Description**

Performs Person's chi-squared test, G-test, or William's corrected G-test to determine dependence between two nucleotide positions.

# Usage

### **Arguments**

X	A DNAStringSet or RNAStringSet object.
i, j	Single integer values for positions to test for dependence.
test	One of "chisq" (Person's chi-squared test), "G" (G-test), or "adjG" (William's
	corrected G-test). See Details section.
simulate.p.val	ue
	a logical indicating whether to compute p-values by Monte Carlo simulation.
В	an integer specifying the number of replicates used in the Monte Carlo test.

#### **Details**

The null and alternative hypotheses for this function are:

```
H0: positions i and j are independentH1: otherwise
```

Let O and E be the observed and expected probabilities for base pair combinations at positions i and j respectively. Then the test statistics are calculated as:

```
test="chisq": stat = sum(abs(O - E)^2/E)  
test="G": stat = 2 * sum(O * log(O/E))  
test="adjG": stat = 2 * sum(O * log(O/E))/q, where q = 1 + ((df - 1)^2 - 1)/(6*length(x)*(df - 2))
```

Under the null hypothesis, these test statistics are approximately distributed chi-squared(df = ((distinct bases at i) - 1) \* ((distinct bases at j) - 1)).

DNAString-class

#### Value

An htest object. See help(chisq.test) for more details.

### Author(s)

P. Aboyoun

#### References

Ellrott, K., Yang, C., Sladek, F.M., Jiang, T. (2002) "Identifying transcription factor binding sites through Markov chain optimations", Bioinformatics, 18 (Suppl. 2), S100-S109.

Sokal, R.R., Rohlf, F.J. (2003) "Biometry: The Principle and Practice of Statistics in Biological Research", W.H. Freeman and Company, New York.

Tomovic, A., Oakeley, E. (2007) "Position dependencies in transcription factor binding sites", Bioinformatics, 23, 933-941.

Williams, D.A. (1976) "Improved Likelihood ratio tests for complete contingency tables", Biometrika, 63, 33-37.

#### See Also

nucleotideFrequencyAt, XStringSet-class, chisq.test

### **Examples**

```
data(HNF4alpha)
dinucleotideFrequencyTest(HNF4alpha, 1, 2)
dinucleotideFrequencyTest(HNF4alpha, 1, 2, test = "G")
dinucleotideFrequencyTest(HNF4alpha, 1, 2, test = "adjG")
```

DNAString-class

**DNAString** objects

# **Description**

A DNAString object allows efficient storage and manipulation of a long DNA sequence.

### **Details**

The DNAString class is a direct XString subclass (with no additional slot). Therefore all functions and methods described in the XString man page also work with a DNAString object (inheritance).

Unlike the BString container that allows storage of any single string (based on a single-byte character set) the DNAString container can only store a string based on the DNA alphabet (see below). In addition, the letters stored in a DNAString object are encoded in a way that optimizes fast search algorithms.

DNAString-class 13

### The DNA alphabet

This alphabet contains all letters from the IUPAC Extended Genetic Alphabet (see ?IUPAC\_CODE\_MAP) plus "-" (the *gap* letter), "+" (the *hard masking* letter), and "." (the *not a letter* or *not available* letter). It is stored in the DNA\_ALPHABET predefined constant (character vector).

The alphabet() function returns DNA\_ALPHABET when applied to a DNAString object.

#### Constructor-like functions and generics

In the code snippet below, x can be a single string (character vector of length 1), a BString object or an RNAString object.

DNAString(x="", start=1, nchar=NA): Tries to convert x into a DNAString object by reading nchar letters starting at position start in x.

#### Accessor methods

In the code snippet below, x is a DNAString object.

alphabet(x, baseOnly=FALSE): If x is a DNAString object, then return the DNA alphabet (see above). See the corresponding man pages when x is a BString, RNAString or AAString object.

### **Display**

The letters in a DNAString object are colored when displayed by the show() method. Set global option Biostrings.coloring to FALSE to turn off this coloring.

### Author(s)

H. Pagès

#### See Also

- The DNAStringSet class to represent a collection of DNAString objects.
- The XString and RNAString classes.
- reverseComplement
- alphabetFrequency
- IUPAC\_CODE\_MAP
- letter

```
DNA_BASES
DNA_ALPHABET
dna <- DNAString("TTGAAAA-CTC-N")
dna # 'options(Biostrings.coloring=FALSE)' to turn off coloring
length(dna)
alphabet(dna) # DNA_ALPHABET
alphabet(dna, baseOnly=TRUE) # DNA_BASES</pre>
```

14 findPalindromes

findPalindromes	Searching a sequence for palindromes		
-----------------	--------------------------------------	--	--

### **Description**

The findPalindromes function can be used to find palindromic regions in a sequence.

palindromeArmLength, palindromeLeftArm, and palindromeRightArm are utility functions for operating on palindromic sequences. They should typically be used on the output of findPalindromes.

### Usage

### **Arguments**

subject	An XString object containing the subject string, or an XStringViews object.
min.armlength	An integer giving the minimum length of the arms of the palindromes to search for.
max.looplength	An integer giving the maximum length of "the loop" (i.e the sequence separating the 2 arms) of the palindromes to search for. Note that by default (max.looplength=1), findPalindromes will search for strict palindromes only.
min.looplength	An integer giving the minimum length of "the loop" of the palindromes to search for.
max.mismatch	The maximum number of mismatching letters allowed between the 2 arms of the palindromes to search for.
allow.wobble	Logical indicating whether wobble base pairs (G/U or G/T base pairings) should be treated as mismatches (the default) or matches.
x	An XString object containing a 2-arm palindrome, or an XStringViews object containing a set of 2-arm palindromes.

#### **Details**

The findPalindromes function finds palindromic substrings in a subject string. The palindromes that can be searched for are either strict palindromes or 2-arm palindromes (the former being a particular case of the latter) i.e. palindromes where the 2 arms are separated by an arbitrary sequence called "the loop".

If the subject string is a nucleotide sequence (i.e. DNA or RNA), the 2 arms must contain sequences that are reverse complement from each other. Otherwise, they must contain sequences that are the same.

findPalindromes 15

#### Value

findPalindromes returns an XStringViews object containing all palindromes found in subject (one view per palindromic substring found).

palindromeArmLength returns the arm length (integer) of the 2-arm palindrome x. It will raise an error if x has no arms. Note that any sequence could be considered a 2-arm palindrome if we were OK with arms of length 0 but we are not: x must have arms of length greater or equal to 1 in order to be considered a 2-arm palindrome. When applied to an XStringViews object x, palindromeArmLength behaves in a vectorized fashion by returning an integer vector of the same length as x.

palindromeLeftArm returns an object of the same class as the original object x and containing the left arm of x.

palindromeRightArm does the same as palindromeLeftArm but on the right arm of x.

Like palindromeArmLength, both palindromeLeftArm and palindromeRightArm will raise an error if x has no arms. Also, when applied to an XStringViews object x, both behave in a vectorized fashion by returning an XStringViews object of the same length as x.

### Author(s)

H. Pagès, with contributions from Erik Wright and Thomas McCarthy

#### See Also

 $\verb|maskMotif|, \verb|matchPattern|, \verb|matchLRPattern|s|, \verb|matchProbePair|, XStringViews-class|, DNAString-class|$ 

```
x0 <- BString("abbbaabbcbbaccacabbbccbcaabbabacca")</pre>
pals0a <- findPalindromes(x0, min.armlength=3, max.looplength=5)</pre>
palindromeArmLength(pals0a)
palindromeLeftArm(pals0a)
palindromeRightArm(pals0a)
pals0b <- findPalindromes(x0, min.armlength=9, max.looplength=5,
                           max.mismatch=3)
pals0b
palindromeArmLength(pals0b, max.mismatch=3)
palindromeLeftArm(pals0b, max.mismatch=3)
palindromeRightArm(pals0b, max.mismatch=3)
## Whitespaces matter:
x1 <- BString("Delia saw I was aileD")</pre>
palindromeArmLength(x1)
palindromeLeftArm(x1)
palindromeRightArm(x1)
x2 <- BString("was it a car or a cat I saw")
```

16 GENETIC\_CODE

```
palindromeArmLength(x2)
palindromeLeftArm(x2)
palindromeRightArm(x2)
## On a DNA or RNA sequence:
x3 <- DNAString("CCGAAAACCATGATGGTTGCCAG")
findPalindromes(x3)
findPalindromes(RNAString(x3))
## Note that palindromes can be nested:
x4 <- DNAString("ACGTTNAACGTCCAAAATTTTCCACGTTNAACGT")</pre>
findPalindromes(x4, max.looplength=19)
## Treat wobble base pairings as matches:
x5 <- RNAString("AUGUCUNNNNAGGCGU")</pre>
findPalindromes(x5, max.looplength=4, min.looplength=4)
findPalindromes(x5, max.looplength=4, min.looplength=4, max.mismatch=2)
findPalindromes(x5, max.looplength=4, min.looplength=4, allow.wobble=TRUE)
## A real use case:
library(BSgenome.Dmelanogaster.UCSC.dm3)
chrX <- Dmelanogaster$chrX</pre>
chrX_pals0 <- findPalindromes(chrX, min.armlength=40, max.looplength=80)</pre>
chrX_pals0
palindromeArmLength(chrX_pals0) # 251 70 262
## Allowing up to 2 mismatches between the 2 arms:
chrX_pals2 <- findPalindromes(chrX, min.armlength=40, max.looplength=80,</pre>
                               max.mismatch=2)
chrX_pals2
palindromeArmLength(chrX_pals2, max.mismatch=2) # 254 77 44 48 40 264
```

GENETIC\_CODE

The Standard Genetic Code and its known variants

### Description

Two predefined objects (GENETIC\_CODE and RNA\_GENETIC\_CODE) that represent The Standard Genetic Code.

Other genetic codes are stored in predefined table GENETIC\_CODE\_TABLE from which they can conveniently be extracted with getGeneticCode.

# Usage

```
## The Standard Genetic Code:
GENETIC_CODE
RNA_GENETIC_CODE

## All the known genetic codes:
GENETIC_CODE_TABLE
getGeneticCode(id_or_name2="1", full.search=FALSE, as.data.frame=FALSE)
```

GENETIC\_CODE 17

#### **Arguments**

id\_or\_name2 A single string that uniquely identifies the genetic code to extract. Should be one of the values in the id or name2 columns of GENETIC\_CODE\_TABLE.

full.search By default, only the id and name2 columns of GENETIC\_CODE\_TABLE are searched

for an exact match with id\_or\_name2. If full.search is TRUE, then the search is extended to the name column of GENETIC\_CODE\_TABLE and id\_or\_name2 only needs to be a substring of one of the names in that column (also case is ignored).

as.data.frame Should the genetic code be returned as a data frame instead of a named character

vector?

#### **Details**

Formally, a *genetic code* is a mapping between the 64 tri-nucleotide sequences (called codons) and amino acids.

The Standard Genetic Code (a.k.a. The Canonical Genetic Code, or simply The Genetic Code) is the particular mapping that encodes the vast majority of genes in nature.

GENETIC\_CODE and RNA\_GENETIC\_CODE are predefined named character vectors that represent this mapping.

All the known genetic codes are summarized in GENETIC\_CODE\_TABLE, which is a predefined data frame with one row per known genetic code. Use getGeneticCode to extract one genetic code at a time from this object.

#### Value

GENETIC\_CODE and RNA\_GENETIC\_CODE are both named character vectors of length 64 (the number of all possible tri-nucleotide sequences) where each element is a single letter representing either an amino acid or the stop codon "\*" (aka termination codon).

The names of the GENETIC\_CODE vector are the DNA codons i.e. the tri-nucleotide sequences (directed 5' to 3') that are assumed to belong to the "coding DNA strand" (aka "sense DNA strand" or "non-template DNA strand") of the gene.

The names of the RNA\_GENETIC\_CODE are the RNA codons i.e. the tri-nucleotide sequences (directed 5' to 3') that are assumed to belong to the mRNA of the gene.

Note that the values in the GENETIC\_CODE and RNA\_GENETIC\_CODE vectors are the same, only their names are different. The names of the latter are those of the former where all occurrences of T (thymine) have been replaced by U (uracil).

Finally, both vectors have an alt\_init\_codons attribute on them, that lists the *alternative initiation codons*. Note that codons that always translate to M (Methionine) (e.g. ATG in GENETIC\_CODE or AUG in RNA\_GENETIC\_CODE) are omitted from the alt\_init\_codons attribute.

GENETIC\_CODE\_TABLE is a data frame that contains all the known genetic codes listed at <a href="ftp:ncbi.nih.gov/entrez/misc/data/gc.prt">ftp.ncbi.nih.gov/entrez/misc/data/gc.prt</a>. The data frame has one row per known genetic code and the 5 following columns:

- name: The long and very descriptive name of the genetic code.
- name2: The short name of the genetic code (not all genetic codes have one).
- id: The id of the genetic code.

18 GENETIC\_CODE

• AAs: A 64-character string representing the genetic code itself in a compact form (i.e. one letter per codon, the codons are assumed to be ordered like in GENETIC\_CODE).

• Starts: A 64-character string indicating the Initiation Codons.

By default (i.e. when as.data.frame is set to FALSE), getGeneticCode returns a named character vector of length 64 similar to GENETIC\_CODE i.e. it contains 1-letter strings from the Amino Acid alphabet (see ?AA\_ALPHABET) and its names are identical to names(GENETIC\_CODE). In addition it has an attribute on it, the alt\_init\_codons attribute, that lists the *alternative initiation codons*. Note that codons that always translate to M (Methionine) (e.g. ATG) are omitted from the alt\_init\_codons attribute.

When as . data . frame is set to TRUE, getGeneticCode returns a data frame with 64 rows (one per codon), rownames (3-letter strings representing the codons), and the 2 following columns:

- AA: A 1-letter string from the Amino Acid alphabet (see ?AA\_ALPHABET) representing the amino acid mapped to the codon ("\*" is used to mark the stop codon).
- Start: A 1-letter string indicating an alternative mapping for the codon i.e. what amino acid the codon is mapped to when it's the first tranlated codon.

The rownames of the data frame are identical to names (GENETIC\_CODE).

### Author(s)

H. Pagès

### References

All the known genetic codes are described here:

```
http://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi
```

The "official names" of the various codes ("Standard", "SGC0", "Vertebrate Mitochondrial", "SGC1", etc...) and their ids (1, 2, etc...) were taken from the print-form ASN.1 version of the above document (version 4.0 at the time of this writing):

ftp://ftp.ncbi.nih.gov/entrez/misc/data/gc.prt

### See Also

- AA\_ALPHABET and AMINO\_ACID\_CODE.
- The translate and trinucleotideFrequency functions.
- DNAString, RNAString, and AAString objects.

```
## -----
## THE STANDARD GENETIC CODE
## -----
GENETIC_CODE

## Codon ATG is *always* translated to M (Methionine)
GENETIC_CODE[["ATG"]]
```

getSeq 19

```
## Codons TTG and CTG are "normally" translated to L except when they are
## the first translated codon (a.k.a. start codon or initiation codon),
## in which case they are translated to M:
attr(GENETIC_CODE, "alt_init_codons")
GENETIC_CODE[["TTG"]]
GENETIC_CODE[["CTG"]]
sort(table(GENETIC_CODE)) # the same amino acid can be encoded by 1
                           # to 6 different codons
RNA_GENETIC_CODE
all(GENETIC_CODE == RNA_GENETIC_CODE) # TRUE
## ALL THE KNOWN GENETIC CODES
GENETIC_CODE_TABLE[1:3 , ]
getGeneticCode("SGC0") # The Standard Genetic Code, again
stopifnot(identical(getGeneticCode("SGC0"), GENETIC_CODE))
getGeneticCode("SGC1") # Vertebrate Mitochondrial
getGeneticCode("ascidian", full.search=TRUE) # Ascidian Mitochondrial
## EXAMINE THE DIFFERENCES BETWEEN THE STANDARD CODE AND A NON-STANDARD
idx <- which(GENETIC_CODE != getGeneticCode("SGC1"))</pre>
rbind(Standard=GENETIC_CODE[idx], SGC1=getGeneticCode("SGC1")[idx])
```

getSeq getSeq

# Description

A generic function for extracting a set of sequences (or subsequences) from a sequence container like a BSgenome object or other.

#### Usage

```
getSeq(x, ...)
```

20 gregexpr2

### **Arguments**

x A BSgenome object or any other supported object. Do showMethods("getSeq") to get the list of all supported types for x.

... Any additional arguments needed by the specialized methods.

### Value

An XString object or an XStringSet object or a character vector containing the extracted sequence(s).

See man pages of individual methods for the details e.g. with ?'getSeq,BSgenome-method' to access the man page of the method for BSgenome objects (make sure the BSgenome package is loaded first).

### See Also

```
getSeq,BSgenome-method, XString-class, XStringSet-class
```

### **Examples**

```
## Note that you need to load the package(s) defining the specialized
## methods to have showMethods() display them and to be able to access
## their man pages:
library(BSgenome)
showMethods("getSeq")
```

gregexpr2

A replacement for R standard gregexpr function

# **Description**

This is a replacement for the standard gregexpr function that does exact matching only. Standard gregexpr() misses matches when they are overlapping. The gregexpr2 function finds all matches but it only works in "fixed" mode i.e. for exact matching (regular expressions are not supported).

# Usage

```
gregexpr2(pattern, text)
```

# Arguments

pattern character string to be matched in the given character vector

text a character vector where matches are sought

HNF4alpha 21

### Value

A list of the same length as text each element of which is an integer vector as in gregexpr, except that the starting positions of all (even overlapping) matches are given. Note that, unlike gregexpr, gregexpr2 doesn't attach a "match.length" attribute to each element of the returned list because, since it only works in "fixed" mode, then all the matches have the length of the pattern. Another difference with gregexpr is that with gregexpr2, the pattern argument must be a single (non-NA, non-empty) string.

### Author(s)

H. Pagès

#### See Also

```
gregexpr, matchPattern
```

### **Examples**

```
gregexpr("aa", c("XaaaYaa", "a"), fixed=TRUE)
gregexpr2("aa", c("XaaaYaa", "a"))
```

HNF4alpha

Known HNF4alpha binding sequences

# **Description**

Seventy one known HNF4alpha binding sequences

### **Details**

A DNAStringSet containing 71 known binding sequences for HNF4alpha.

# Author(s)

P. Aboyoun

#### References

Ellrott, K., Yang, C., Sladek, F.M., Jiang, T. (2002) "Identifying transcription factor binding sites through Markov chain optimations", Bioinformatics, 18 (Suppl. 2), S100-S109.

```
data(HNF4alpha)
HNF4alpha
```

22 injectHardMask

injectHardMask

Injecting a hard mask in a sequence

#### **Description**

injectHardMask allows the user to "fill" the masked regions of a sequence with an arbitrary letter (typically the "+" letter).

### Usage

```
injectHardMask(x, letter="+")
```

### **Arguments**

x A MaskedXString or XStringViews object.

letter A single letter.

### **Details**

The name of the injectHardMask function was chosen because of the primary use that it is intended for: converting a pile of active "soft masks" into a "hard mask". Here the pile of active "soft masks" refers to the active masks that have been put on top of a sequence. In Biostrings, the original sequence and the masks defined on top of it are bundled together in one of the dedicated containers for this: the MaskedBString, MaskedDNAString, MaskedRNAString and MaskedAAString containers (this is the MaskedXString family of containers). The original sequence is always stored unmodified in a MaskedXString object so no information is lost. This allows the user to activate/deactivate masks without having to worry about losing the letters that are in the regions that are masked/unmasked. Also this allows better memory management since the original sequence never needs to be copied, even when the set of active/inactive masks changes.

However, there are situations where the user might want to *really* get rid of the letters that are in some particular regions by replacing them with a junk letter (e.g. "+") that is guaranteed to not interfer with the analysis that s/he is currently doing. For example, it's very likely that a set of motifs or short reads will not contain the "+" letter (this could easily be checked) so they will never hit the regions filled with "+". In a way, it's like the regions filled with "+" were masked but we call this kind of masking "hard masking".

Some important differences between "soft" and "hard" masking:

- injectHardMask creates a (modified) copy of the original sequence. Using "soft masking" does not.
- A function that is "mask aware" like alphabetFrequency or matchPattern will really skip
  the masked regions when "soft masking" is used i.e. they will not walk thru the regions that
  are under active masks. This might lead to some speed improvements when a high percentage
  of the original sequence is masked. With "hard masking", the entire sequence is walked thru.
- Matches cannot span over masked regions with "soft masking". With "hard masking" they
  can.

IUPAC\_CODE\_MAP 23

### Value

An XString object of the same length as the original object x if x is a MaskedXString object, or of the same length as subject(x) if it's an XStringViews object.

### Author(s)

H. Pagès

### See Also

maskMotif, MaskedXString-class, replaceLetterAt, chartr, XString, XStringViews-class

### **Examples**

```
## A. WITH AN XStringViews OBJECT
## -----
v2 <- Views("abCDefgHIJK", start=c(8, 3), end=c(14, 4))</pre>
injectHardMask(v2)
injectHardMask(v2, letter="=")
## -----
## B. WITH A MaskedXString OBJECT
## -----
mask0 < -Mask(mask.width=29, start=c(3, 10, 25), width=c(6, 8, 5))
x <- DNAString("ACACAACTAGATAGNACTNNGAGAGACGC")</pre>
masks(x) <- mask0</pre>
subject <- injectHardMask(x)</pre>
## Matches can span over masked regions with "hard masking":
matchPattern("ACggggggA", subject, max.mismatch=6)
## but not with "soft masking":
matchPattern("ACggggggA", x, max.mismatch=6)
```

IUPAC\_CODE\_MAP

The IUPAC Extended Genetic Alphabet

# **Description**

The IUPAC\_CODE\_MAP named character vector contains the mapping from the IUPAC nucleotide ambiguity codes to their meaning.

The mergeIUPACLetters function provides the reverse mapping.

#### Usage

```
IUPAC_CODE_MAP
mergeIUPACLetters(x)
```

24 IUPAC\_CODE\_MAP

#### **Arguments**

Χ

A vector of non-empty character strings made of IUPAC letters.

### **Details**

IUPAC nucleotide ambiguity codes are used for representing sequences of nucleotides where the exact nucleotides that occur at some given positions are not known with certainty.

#### Value

IUPAC\_CODE\_MAP is a named character vector where the names are the IUPAC nucleotide ambiguity codes and the values are their corresponding meanings. The meaning of each code is described by a string that enumarates the base letters ("A", "C", "G" or "T") associated with the code.

The value returned by mergeIUPACLetters is an unnamed character vector of the same length as its argument x where each element is an IUPAC nucleotide ambiguity code.

#### Author(s)

H. Pagès

#### References

```
http://www.chick.manchester.ac.uk/SiteSeer/IUPAC\_codes.html
```

IUPAC-IUB SYMBOLS FOR NUCLEOTIDE NOMENCLATURE: Cornish-Bowden (1985) *Nucl. Acids Res.* 13: 3021-3030.

# See Also

```
DNAString, RNAString
```

```
IUPAC_CODE_MAP
some_iupac_codes <- c("R", "M", "G", "N", "V")
IUPAC_CODE_MAP[some_iupac_codes]
mergeIUPACLetters(IUPAC_CODE_MAP[some_iupac_codes])
mergeIUPACLetters(c("Ca", "Acc", "aA", "MAAMC", "gM", "AB", "bS", "mk"))</pre>
```

lcsuffix 25

lcsuffix

Find Longest Common Prefix/Suffix

# Description

Functions for searching the Longest Common Prefix/Suffix of two strings.

WARNING: These functions are experimental and might not work properly! Full documentation will come later.

Thanks for your comprehension!

# Usage

```
lcprefix(s1, s2)
lcsuffix(s1, s2)
```

# Arguments

- s1 1st string, a character string or an XString object.
- s2 2nd string, a character string or an XString object.

#### See Also

matchPattern, XStringViews-class, XString-class

letter

Subsetting a string

### **Description**

Extract a substring from a string by picking up individual letters by their position.

### Usage

```
letter(x, i)
```

# Arguments

- x A character vector, or an XString, XStringViews or MaskedXString object.
- i An integer vector with no NAs.

### **Details**

Unlike with the substr or substring functions, i must contain valid positions.

#### Value

A character vector of length 1 when x is an XString or MaskedXString object (the masks are ignored for the latter).

A character vector of the same length as x when x is a character vector or an XStringViews object. Note that, because i must contain valid positions, all non-NA elements in the result are guaranteed to have exactly length(i) characters.

#### See Also

subseq, XString-class, XStringViews-class, MaskedXString-class

#### **Examples**

```
x <- c("abcd", "ABC")
i <- c(3, 1, 1, 2, 1)

## With a character vector:
letter(x[1], 3:1)
letter(x, 3)
letter(x, i)
#letter(x, 4)  # Error!

## With a BString object:
letter(BString(x[1]), i) # returns a character vector
BString(x[1])[i] # returns a BString object

## With an XStringViews object:
x2 <- as(BStringSet(x), "Views")
letter(x2, i)</pre>
```

letterFrequency

Calculate the frequency of letters in a biological sequence, or the consensus matrix of a set of sequences

### **Description**

Given a biological sequence (or a set of biological sequences), the alphabetFrequency function computes the frequency of each letter of the relevant alphabet.

letterFrequency is similar, but more compact if one is only interested in certain letters. It can also tabulate letters "in common".

letterFrequencyInSlidingView is a more specialized version of letterFrequency for (non-masked) XString objects. It tallys the requested letter frequencies for a fixed-width view, or window, that is conceptually slid along the entire input sequence.

The consensusMatrix function computes the consensus matrix of a set of sequences, and the consensusString function creates the consensus sequence from the consensus matrix based upon specified criteria.

In this man page we call "DNA input" (or "RNA input") an XString, XStringSet, XStringViews or MaskedXString object of base type DNA (or RNA).

#### Usage

```
alphabetFrequency(x, as.prob=FALSE, ...)
hasOnlyBaseLetters(x)
uniqueLetters(x)
letterFrequency(x, letters, OR="|", as.prob=FALSE, ...)
letterFrequencyInSlidingView(x, view.width, letters, OR="|", as.prob=FALSE)
consensusMatrix(x, as.prob=FALSE, shift=0L, width=NULL, ...)
## S4 method for signature 'matrix'
consensusString(x, ambiguityMap="?", threshold=0.5)
## S4 method for signature 'DNAStringSet'
consensusString(x, ambiguityMap=IUPAC_CODE_MAP,
             threshold=0.25, shift=0L, width=NULL)
## S4 method for signature 'RNAStringSet'
consensusString(x,
             ambiguityMap=
            structure(as.character(RNAStringSet(DNAStringSet(IUPAC_CODE_MAP))),
                       names=
               as.character(RNAStringSet(DNAStringSet(names(IUPAC_CODE_MAP))))),
             threshold=0.25, shift=0L, width=NULL)
```

#### **Arguments**

Χ

An XString, XStringSet, XStringViews or MaskedXString object for alphabetFrequency, letterFrequency, or uniqueLetters.

DNA or RNA input for hasOnlyBaseLetters.

An XString object for letterFrequencyInSlidingView.

A character vector, or an XStringSet or XStringViews object for consensusMatrix.

A consensus matrix (as returned by consensusMatrix), or an XStringSet or

XStringViews object for consensusString.

as.prob

If TRUE then probabilities are reported, otherwise counts (the default).

view.width

For letterFrequencyInSlidingView, the constant (e.g. 35, 48, 1000) size of the "window" to slide along x. The specified letters are tabulated in each window of length view.width. The rows of the result (see value) correspond to the various windows.

letters

For letterFrequency or letterFrequencyInSlidingView, a character vector (e.g. "C", "CG", c("C", "G")) giving the letters to tabulate. When x is DNA or RNA input, letters must come from alphabet(x). Except with OR=0, multicharacter elements of letters ('nchar' > 1) are taken as groupings of letters into subsets, to be tabulated in common ("or"'d), as if their alphabetFrequency's were added (Arithmetic). The columns of the result (see value) correspond to the individual and sets of letters which are counted separately. Unrelated (and, with some post-processing, related) counts may of course be obtained in separate calls.

OR

For letterFrequency or letterFrequencyInSlidingView, the string (default |) to use as a separator in forming names for the "grouped" columns, e.g. "ClG". The otherwise exceptional value 0 (zero) disables or'ing and is provided for convenience, allowing a single multi-character string (or several strings) of letters that should be counted separately. If some but not all letters are to be counted separately, they must reside in separate elements of letters (with 'nchar' 1 unless they are to be grouped with other letters), and OR cannot be 0.

ambiguityMap

Either a single character to use when agreement is not reached or a named character vector where the names are the ambiguity characters and the values are the combinations of letters that comprise the ambiguity (e.g. link{IUPAC\_CODE\_MAP}). When ambiguityMap is a named character vector, occurrences of ambiguous letters in x are replaced with their base alphabet letters that have been equally weighted to sum to 1. (See Details for some examples.)

threshold

The minimum probability threshold for an agreement to be declared. When ambiguityMap is a single character, threshold is a single number in (0, 1]. When ambiguityMap is a named character vector (e.g. link{IUPAC\_CODE\_MAP}), threshold is a single number in (0, 1/sum(nchar(ambiguityMap) == 1)].

. . .

Further arguments to be passed to or from other methods.

For the XStringViews and XStringSet methods, the collapse argument is accepted.

Except for letterFrequency or letterFrequencyInSlidingView, and with DNA or RNA input, the baseOnly argument is accepted. If baseOnly is TRUE, the returned vector (or matrix) only contains the frequencies of the letters that belong to the "base" alphabet of x i.e. to the alphabet returned by alphabet(x, baseOnly=TRUE).

shift

An integer vector (recycled to the length of x) specifying how each sequence in x should be (horizontally) shifted with respect to the first column of the consensus matrix to be returned. By default (shift=0), each sequence in x has its first letter aligned with the first column of the matrix. A positive shift value means that the corresponding sequence must be shifted to the right, and a negative shift value that it must be shifted to the left. For example, a shift of 5 means that it must be shifted 5 positions to the right (i.e. the first letter in the sequence must be aligned with the 6th column of the matrix), and a shift of -3 means that it must be shifted 3 positions to the left (i.e. the 4th letter in the sequence must be aligned with the first column of the matrix).

width

The number of columns of the returned matrix for the consensusMatrix method for XStringSet objects. When width=NULL (the default), then this method returns a matrix that has just enough columns to have its last column aligned with the rightmost letter of all the sequences in x after those sequences have been shifted (see the shift argument above). This ensures that any wider consensus matrix would be a "padded with zeros" version of the matrix returned when width=NULL.

The length of the returned sequence for the consensusString method for XStringSet objects.

### **Details**

alphabetFrequency, letterFrequency, and letterFrequencyInSlidingView are generic func-

tions defined in the Biostrings package.

letterFrequency is similar to alphabetFrequency but specific to the letters of interest, hence more compact, especially with OR non-zero.

letterFrequencyInSlidingView yields the same result, on the sequence x, that letterFrequency would, if applied to the hypothetical (and possibly huge) XStringViews object consisting of all the intervals of length view.width on x. Taking advantage of the knowledge that successive "views" are nearly identical, for letter counting purposes, it is both lighter and faster.

For letterFrequencyInSlidingView, a masked (MaskedXString) object x is only supported through a cast to an (ordinary) XString such as unmasked (which includes its masked regions).

When consensusString is executed with a named character ambiguityMap argument, it weights each input string equally and assigns an equal probability to each of the base letters represented by an ambiguity letter. So for DNA and a threshold of 0.25, a "G" and an "R" would result in an "R" since 1/2 "G" + 1/2 "R" = 3/4 "G" + 1/4 "A" => "R"; two "G"'s and one "R" would result in a "G" since 2/3 "G" + 1/3 "R" = 5/6 "G" + 1/6 "A" => "G"; and one "A" and one "N" would result in an "N" since 1/2 "A" + 1/2 "N" = 5/8 "A" + 1/8 "C" + 1/8 "G" + 1/8 "T" => "N".

#### Value

alphabetFrequency returns an integer vector when x is an XString or MaskedXString object. When x is an XStringSet or XStringViews object, then it returns an integer matrix with length(x) rows where the i-th row contains the frequencies for x[[i]]. If x is a DNA, RNA, or AA input, then the returned vector is named with the letters in the alphabet. If the baseOnly argument is TRUE, then the returned vector has only 5 elements for DNA/RNA input (4 elements corresponding to the 4 nucleotides + the 'other' element) and 21 elements for AA input (20 elements corresponding to the 20 base amino acids + the 'other' element).

letterFrequency returns, similarly, an integer vector or matrix, but restricted and/or collated according to letters and OR.

letterFrequencyInSlidingView returns, for an XString object x of length (nchar) L, an integer matrix with L-view.width+1 rows, the i-th of which holding the letter frequencies of substring(x, i, i+view.width-1).

hasOnlyBaseLetters returns TRUE or FALSE indicating whether or not x contains only base letters (i.e. As, Cs, Gs and Ts for DNA input, As, Cs, Gs and Us for RNA input, or any of the 20 standard amino acids for AA input).

uniqueLetters returns a vector of 1-letter or empty strings. The empty string is used to represent the nul character if x happens to contain any. Note that this can only happen if the base class of x is BString.

An integer matrix with letters as row names for consensusMatrix.

A standard character string for consensusString.

#### Author(s)

H. Pagès and P. Aboyoun; H. Jaffee for letterFrequency and letterFrequencyInSlidingView

### See Also

alphabet, coverage, oligonucleotideFrequency, countPDict, XString-class, XStringSet-class, XStringViews-class, MaskedXString-class, strsplit

```
## -----
## alphabetFrequency()
## -----
data(yeastSEQCHR1)
yeast1 <- DNAString(yeastSEQCHR1)</pre>
alphabetFrequency(yeast1)
alphabetFrequency(yeast1, baseOnly=TRUE)
hasOnlyBaseLetters(yeast1)
uniqueLetters(yeast1)
## With input made of multiple sequences:
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
alphabetFrequency(probes[1:50], baseOnly=TRUE)
alphabetFrequency(probes, baseOnly=TRUE, collapse=TRUE)
## ------
## letterFrequency()
## -----
letterFrequency(probes[[1]], letters="ACGT", OR=0)
base_letters <- alphabet(probes, baseOnly=TRUE)</pre>
base_letters
letterFrequency(probes[[1]], letters=base_letters, OR=0)
base_letter_freqs <- letterFrequency(probes, letters=base_letters, OR=0)</pre>
head(base_letter_freqs)
GC_content <- letterFrequency(probes, letters="CG")
head(GC_content)
letterFrequency(probes, letters="CG", collapse=TRUE)
## letterFrequencyInSlidingView()
## -----
data(yeastSEQCHR1)
x <- DNAString(yeastSEQCHR1)</pre>
view.width <- 48
letters <- c("A", "CG")</pre>
two_columns <- letterFrequencyInSlidingView(x, view.width, letters)
head(two_columns)
tail(two_columns)
three_columns <- letterFrequencyInSlidingView(x, view.width, letters, OR=0)
head(three_columns)
tail(three_columns)
stopifnot(identical(two_columns[ , "C|G"],
                 three_columns[ , "C"] + three_columns[ , "G"]))
## Note that, alternatively, 'three_columns' can also be obtained by
## creating the views on 'x' (as a Views object) and by calling
## alphabetFrequency() on it. But, of course, that is be *much* less
## efficient (both, in terms of memory and speed) than using
```

```
## letterFrequencyInSlidingView():
v \leftarrow Views(x, start=seq\_len(length(x) - view.width + 1), width=view.width)
three_columns2 <- alphabetFrequency(v, baseOnly=TRUE)[ , c("A", "C", "G")]</pre>
stopifnot(identical(three_columns2, three_columns))
## Set the width of the view to length(x) to get the global frequencies:
letterFrequencyInSlidingView(x, letters="ACGTN", view.width=length(x), OR=0)
## consensus*()
## -----
## Read in ORF data:
file <- system.file("extdata", "someORF.fa", package="Biostrings")</pre>
orf <- readDNAStringSet(file)</pre>
## To illustrate, the following example assumes the ORF data
## to be aligned for the first 10 positions (patently false):
orf10 <- DNAStringSet(orf, end=10)</pre>
consensusMatrix(orf10, baseOnly=TRUE)
## The following example assumes the first 10 positions to be aligned
## after some incremental shifting to the right (patently false):
consensusMatrix(orf10, baseOnly=TRUE, shift=0:6)
consensusMatrix(orf10, baseOnly=TRUE, shift=0:6, width=10)
## For the character matrix containing the "exploded" representation
## of the strings, do:
as.matrix(orf10, use.names=FALSE)
## consensusMatrix() can be used to just compute the alphabet frequency
## for each position in the input sequences:
consensusMatrix(probes, baseOnly=TRUE)
## After sorting, the first 5 probes might look similar (at least on
## their first bases):
consensusString(sort(probes)[1:5])
consensusString(sort(probes)[1:5], ambiguityMap = "N", threshold = 0.5)
## Consensus involving ambiguity letters in the input strings
consensusString(DNAStringSet(c("NNNN", "ACTG")))
consensusString(DNAStringSet(c("AANN", "ACTG")))
consensusString(DNAStringSet(c("ACAG","ACAR")))
consensusString(DNAStringSet(c("ACAG", "ACAR", "ACAG")))
## -----
## C. RELATIONSHIP BETWEEN consensusMatrix() AND coverage()
## -----
## Applying colSums() on a consensus matrix gives the coverage that
## would be obtained by piling up (after shifting) the input sequences
## on top of an (imaginary) reference sequence:
cm <- consensusMatrix(orf10, shift=0:6, width=10)</pre>
colSums(cm)
```

32 longestConsecutive

```
## Note that this coverage can also be obtained with:
as.integer(coverage(IRanges(rep(1, length(orf)), width(orf)), shift=0:6, width=10))
```

longestConsecutive

Obtain the length of the longest substring containing only 'letter'

# **Description**

This function accepts a character vector and computes the length of the longest substring containing only letter for each element of x.

# Usage

```
longestConsecutive(seq, letter)
```

# **Arguments**

seq

Character vector.

letter

Character vector of length 1, containing one single character.

### **Details**

The elements of x can be in upper case, lower case or mixed. NAs are handled.

### Value

An integer vector of the same length as x.

### Author(s)

W. Huber

```
v <- c("AAACTGTGFG", "GGGAATT", "CCAAAAAAAAAATT")
longestConsecutive(v, "A")</pre>
```

lowlevel-matching

Low-level matching functions

#### **Description**

In this man page we define precisely and illustrate what a "match" of a pattern P in a subject S is in the context of the Biostrings package. This definition of a "match" is central to most pattern matching functions available in this package: unless specified otherwise, most of them will adhere to the definition provided here.

hasLetterAt checks whether a sequence or set of sequences has the specified letters at the specified positions.

neditAt, isMatchingAt and which.isMatchingAt are low-level matching functions that only look for matches at the specified positions in the subject.

### Usage

```
hasLetterAt(x, letter, at, fixed=TRUE)
## neditAt() and related utils:
neditAt(pattern, subject, at=1,
        with.indels=FALSE, fixed=TRUE)
neditStartingAt(pattern, subject, starting.at=1,
        with.indels=FALSE, fixed=TRUE)
neditEndingAt(pattern, subject, ending.at=1,
       with.indels=FALSE, fixed=TRUE)
## isMatchingAt() and related utils:
isMatchingAt(pattern, subject, at=1,
       max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE)
isMatchingStartingAt(pattern, subject, starting.at=1,
        max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE)
isMatchingEndingAt(pattern, subject, ending.at=1,
       max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE)
## which.isMatchingAt() and related utils:
which.isMatchingAt(pattern, subject, at=1,
       max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
        follow.index=FALSE, auto.reduce.pattern=FALSE)
which.isMatchingStartingAt(pattern, subject, starting.at=1,
       max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
        follow.index=FALSE, auto.reduce.pattern=FALSE)
which.isMatchingEndingAt(pattern, subject, ending.at=1,
        max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
        follow.index=FALSE, auto.reduce.pattern=FALSE)
```

#### **Arguments**

x A character vector, or an XString or XStringSet object.

letter A character string or an XString object containing the letters to check.

at, starting.at, ending.at

An integer vector specifying the starting (for starting.at and at) or ending (for ending.at) positions of the pattern relatively to the subject. With auto.reduce.pattern (below), either a single integer or a constant vector of length nchar(pattern) (below), to which the former is immediately converted. For the hasLetterAt function, letter and at must have the same length.

Tof the hastetter at function, fetter and at must have the same i

The pattern string (but see auto.reduce.pattern, below).

subject A character vector, or an XString or XStringSet object containing the subject

sequence(s).

max.mismatch, min.mismatch

Integer vectors of length >= 1 recycled to the length of the at (or starting.at, or ending.at) argument. More details below.

with.indels See details below.

fixed Only with a DNAString or RNAString-based subject can a fixed value other

than the default (TRUE) be used.

If TRUE (the default), an IUPAC ambiguity code in the pattern can only match the same code in the subject, and vice versa. If FALSE, an IUPAC ambiguity code in the pattern can match any letter in the subject that is associated with the code, and vice versa. See IUPAC\_CODE\_MAP for more information about the IUPAC Extended Genetic Alphabet.

fixed can also be a character vector, a subset of c("pattern", "subject"). fixed=c("pattern", "subject") is equivalent to fixed=TRUE (the default). An empty vector is equivalent to fixed=FALSE. With fixed="subject", ambiguities in the pattern only are interpreted as wildcards. With fixed="pattern", ambiguities in the subject only are interpreted as wildcards.

ambiguities in the subject only are interpreted as wildcards.

follow.index Whether the single integer returned by which.isMatchingAt (and related utils)

should be the first \*value\* in at for which a match occurred, or its \*index\* in

at (the default).

auto.reduce.pattern

Whether pattern should be effectively shortened by 1 letter, from its beginning for which.isMatchingStartingAt and from its end for which.isMatchingEndingAt, for each successive (at, max.mismatch) "pair".

### **Details**

A "match" of pattern P in subject S is a substring S' of S that is considered similar enough to P according to some distance (or metric) specified by the user. 2 distances are supported by most pattern matching functions in the Biostrings package. The first (and simplest) one is the "number of mismatching letters". It is defined only when the 2 strings to compare have the same length, so when this distance is used, only matches that have the same number of letters as P are considered. The second one is the "edit distance" (aka Levenshtein distance): it's the minimum number of operations needed to transform P into S', where an operation is an insertion, deletion, or substitution of a single letter. When this metric is used, matches can have a different number of letters than P.

The neditAt function implements these 2 distances. If with indels is FALSE (the default), then the first distance is used i.e. neditAt returns the "number of mismatching letters" between the pattern P and the substring S' of S starting at the positions specified in at (note that neditAt is vectorized so a long vector of integers can be passed thru the at argument). If with indels is TRUE, then the "edit distance" is used: for each position specified in at, P is compared to all the substrings S' of S starting at this position and the smallest distance is returned. Note that this distance is guaranteed to be reached for a substring of length < 2\*length(P) so, of course, in practice, P only needs to be compared to a small number of substrings for every starting position.

#### Value

hasLetterAt: A logical matrix with one row per element in x and one column per letter/position to check. When a specified position is invalid with respect to an element in x then the corresponding matrix element is set to NA.

neditAt: If subject is an XString object, then return an integer vector of the same length as at. If subject is an XStringSet object, then return the integer matrix with length(at) rows and length(subject) columns defined by:

```
sapply(unname(subject),
    function(x) neditAt(pattern, x, ...))
```

neditStartingAt is identical to neditAt except that the at argument is now called starting.at. neditEndingAt is similar to neditAt except that the at argument is now called ending.at and must contain the ending positions of the pattern relatively to the subject.

isMatchingAt: If subject is an XString object, then return the logical vector defined by:

```
min.mismatch <= neditAt(...) <= max.mismatch</pre>
```

If subject is an XStringSet object, then return the logical matrix with length(at) rows and length(subject) columns defined by:

isMatchingStartingAt is identical to isMatchingAt except that the at argument is now called starting.at. isMatchingEndingAt is similar to isMatchingAt except that the at argument is now called ending.at and must contain the ending positions of the pattern relatively to the subject.

which.isMatchingAt: The default behavior (follow.index=FALSE) is as follow. If subject is an XString object, then return the single integer defined by:

```
which(isMatchingAt(...))[1]
```

If subject is an XStringSet object, then return the integer vector defined by:

which.isMatchingStartingAt is identical to which.isMatchingAt except that the at argument is now called starting.at. which.isMatchingEndingAt is similar to which.isMatchingAt except that the at argument is now called ending.at and must contain the ending positions of the pattern relatively to the subject.

#### See Also

nucleotideFrequencyAt, matchPattern, matchPDict, matchLRPatterns, trimLRPatterns, IUPAC\_CODE\_MAP, XString-class, align-utils in the **pwalign** package

```
## -----
## hasLetterAt()
x <- DNAStringSet(c("AAACGT", "AACGT", "ACGT", "TAGGA"))</pre>
hasLetterAt(x, "AAAAAA", 1:6)
## hasLetterAt() can be used to answer questions like: "which elements
## in 'x' have an A at position 2 and a G at position 4?"
q1 \leftarrow hasLetterAt(x, "AG", c(2, 4))
which(rowSums(q1) == 2)
## or "how many probes in the drosophila2 chip have T, G, T, A at
## position 2, 4, 13 and 20, respectively?"
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
q2 <- hasLetterAt(probes, "TGTA", c(2, 4, 13, 20))</pre>
sum(rowSums(q2) == 4)
## or "what's the probability to have an A at position 25 if there is
## one at position 13?"
q3 <- hasLetterAt(probes, "AACGT", c(13, 25, 25, 25, 25))
sum(q3[, 1] & q3[, 2]) / sum(q3[, 1])
## Probabilities to have other bases at position 25 if there is an A
## at position 13:
sum(q3[, 1] & q3[, 3]) / sum(q3[, 1]) # C
sum(q3[ , 1] & q3[ , 4]) / sum(q3[ , 1]) # G
sum(q3[, 1] & q3[, 5]) / sum(q3[, 1]) # T
## See ?nucleotideFrequencyAt for another way to get those results.
## neditAt() / isMatchingAt() / which.isMatchingAt()
```

lowlevel-matching 37

```
subject <- DNAString("GTATA")</pre>
## Pattern "AT" matches subject "GTATA" at position 3 (exact match)
neditAt("AT", subject, at=3)
isMatchingAt("AT", subject, at=3)
## ... but not at position 1
neditAt("AT", subject)
isMatchingAt("AT", subject)
## ... unless we allow 1 mismatching letter (inexact match)
isMatchingAt("AT", subject, max.mismatch=1)
## Here we look at 6 different starting positions and find 3 matches if
## we allow 1 mismatching letter
isMatchingAt("AT", subject, at=0:5, max.mismatch=1)
## No match
neditAt("NT", subject, at=1:4)
isMatchingAt("NT", subject, at=1:4)
## 2 matches if N is interpreted as an ambiguity (fixed=FALSE)
neditAt("NT", subject, at=1:4, fixed=FALSE)
isMatchingAt("NT", subject, at=1:4, fixed=FALSE)
## max.mismatch != 0 and fixed=FALSE can be used together
neditAt("NCA", subject, at=0:5, fixed=FALSE)
isMatchingAt("NCA", subject, at=0:5, max.mismatch=1, fixed=FALSE)
some_starts <- c(10:-10, NA, 6)
subject <- DNAString("ACGTGCA")</pre>
is_matching <- isMatchingAt("CAT", subject, at=some_starts, max.mismatch=1)</pre>
some_starts[is_matching]
which.isMatchingAt("CAT", subject, at=some_starts, max.mismatch=1)
which.isMatchingAt("CAT", subject, at=some_starts, max.mismatch=1,
                  follow.index=TRUE)
## -----
## WITH INDELS
## -----
subject <- BString("ABCDEFxxxCDEFxxxABBCDE")</pre>
neditAt("ABCDEF", subject, at=9)
neditAt("ABCDEF", subject, at=9, with.indels=TRUE)
isMatchingAt("ABCDEF", subject, at=9, max.mismatch=1, with.indels=TRUE)
isMatchingAt("ABCDEF", subject, at=9, max.mismatch=2, with.indels=TRUE)
neditAt("ABCDEF", subject, at=17)
neditAt("ABCDEF", subject, at=17, with.indels=TRUE)
neditEndingAt("ABCDEF", subject, ending.at=22)
neditEndingAt("ABCDEF", subject, ending.at=22, with.indels=TRUE)
```

38 MaskedXString-class

MaskedXString-class MaskedXString objects

## **Description**

The MaskedBString, MaskedDNAString, MaskedRNAString and MaskedAAString classes are containers for storing masked sequences.

All those containers derive directly (and with no additional slots) from the MaskedXString virtual class.

#### **Details**

In Biostrings, a pile of masks can be put on top of a sequence. A pile of masks is represented by a MaskCollection object and the sequence by an XString object. A MaskedXString object is the result of bundling them together in a single object.

Note that, no matter what masks are put on top of it, the original sequence is always stored unmodified in a MaskedXString object. This allows the user to activate/deactivate masks without having to worry about losing the information stored in the masked/unmasked regions. Also this allows efficient memory management since the original sequence never needs to be copied (modifying it would require to make a copy of it first - sequences cannot and should never be modified in place in Biostrings), even when the set of active/inactive masks changes.

#### Accessor methods

In the code snippets below, x is a MaskedXString object. For masks(x) and masks(x) < -y, it can also be an XString object and y must be NULL or a MaskCollection object.

unmasked(x): Turns x into an XString object by dropping the masks.

masks(x): Turns x into a MaskCollection object by dropping the sequence.

masks(x) <- y: If x is an XString object and y is NULL, then this doesn't do anything.

If x is an XString object and y is a MaskCollection object, then this turns x into a MaskedXString object by putting the masks in y on top of it.

If x is a MaskedXString object and y is NULL, then this is equivalent to x < -unmasked(x).

If x is a MaskedXString object and y is a MaskCollection object, then this replaces the masks currently on top of x by the masks in y.

alphabet(x): Equivalent to alphabet(unmasked(x)). See ?alphabet for more information.

length(x): Equivalent to length(unmasked(x)). See ?`length,XString-method` for more
information.

#### "maskedwidth" and related methods

In the code snippets below, x is a MaskedXString object.

maskedwidth(x): Get the number of masked letters in x. A letter is considered masked iff it's masked by at least one active mask.

maskedratio(x): Equivalent to maskedwidth(x) / length(x).

nchar(x): Equivalent to length(x) - maskedwidth(x).

MaskedXString-class 39

## Coercion

In the code snippets below, x is a MaskedXString object.

as(x, "Views"): Turns x into a Views object where the views are the unmasked regions of the original sequence ("unmasked" means not masked by at least one active mask).

## Other methods

In the code snippets below, x is a MaskedXString object.

collapse(x): Collapses the set of masks in x into a single mask made of all active masks.

gaps(x): Reverses all the masks i.e. each mask is replaced by a mask where previously unmasked regions are now masked and previously masked regions are now unmasked.

#### Author(s)

H. Pagès

#### See Also

- maskMotif
- injectHardMask
- alphabetFrequency
- reverseComplement
- XString-class
- MaskCollection-class
- Views-class

40 maskMotif

maskMotif

*Masking by content (or by position)* 

# **Description**

Functions for masking a sequence by content (or by position).

# Usage

```
maskMotif(x, motif, min.block.width=1, ...)
mask(x, start=NA, end=NA, pattern)
```

## **Arguments**

The sequence to mask.

motif The motif to mask in the sequence.

min.block.width

The minimum width of the blocks to mask.

... Additional arguments for matchPattern.

An integer vector containing the starting positions of the regions to mask.

An integer vector containing the ending positions of the regions to mask.

pattern The motif to mask in the sequence.

maskMotif 41

#### Value

A MaskedXString object for maskMotif and an XStringViews object for mask.

#### Author(s)

H. Pagès

#### See Also

read. Mask, match Pattern, XString-class, Masked XString-class, XString Views-class, Mask Collection-class

```
## -----
## EXAMPLE 1
maskMotif(BString("AbcbbcbEEE"), "bcb")
maskMotif(BString("AbcbcbEEE"), "bcb")
## maskMotif() can be used in an incremental way to mask more than 1
## motif. Note that maskMotif() does not try to mask again what's
## already masked (i.e. the new mask will never overlaps with the
## previous masks) so the order in which the motifs are masked actually
## matters as it will affect the total set of masked positions.
x0 <- BString("AbcbEEEEbcbbEEEcbbcbc")</pre>
x1 <- maskMotif(x0, "E")</pre>
x1
x2 <- maskMotif(x1, "bcb")</pre>
x3 <- maskMotif(x2, "b")
## Note that inverting the order in which "b" and "bcb" are masked would
## lead to a different final set of masked positions.
## Also note that the order doesn't matter if the motifs to mask don't
## overlap (we assume that the motifs are unique) i.e. if the prefix of
## each motif is not the suffix of any other motif. This is of course
## the case when all the motifs have only 1 letter.
## EXAMPLE 2
x <- DNAString("ACACAACTAGATAGNACTNNGAGAGACGC")
## Mask the N-blocks
x1 <- maskMotif(x, "N")</pre>
as(x1, "Views")
gaps(x1)
as(gaps(x1), "Views")
```

42 match-utils

```
## Mask the AC-blocks
x2 <- maskMotif(x1, "AC")</pre>
x2
gaps(x2)
## Mask the GA-blocks
x3 <- maskMotif(x2, "GA", min.block.width=5)
x3 # masks 2 and 3 overlap
gaps(x3)
## -----
## EXAMPLE 3
library(BSgenome.Dmelanogaster.UCSC.dm3)
chrU <- Dmelanogaster$chrU</pre>
chrU
alphabetFrequency(chrU)
chrU <- maskMotif(chrU, "N")</pre>
chrU
alphabetFrequency(chrU)
as(chrU, "Views")
as(gaps(chrU), "Views")
mask2 <- Mask(mask.width=length(chrU),</pre>
            start=c(50000, 350000, 543900), width=25000)
names(mask2) <- "some ugly regions"</pre>
masks(chrU) <- append(masks(chrU), mask2)</pre>
chrU
as(chrU, "Views")
as(gaps(chrU), "Views")
## -----
## EXAMPLE 4
## -----
## Note that unlike maskMotif(), mask() returns an XStringViews object!
## masking "by position"
mask("AxyxyxBC", 2, 6)
## masking "by content"
mask("AxyxyxBC", "xyx")
noN_chrU <- mask(chrU, "N")</pre>
noN_chrU
alphabetFrequency(noN_chrU, collapse=TRUE)
```

match-utils

Utility functions operating on the matches returned by a high-level matching function

match-utils 43

#### **Description**

Miscellaneous utility functions operating on the matches returned by a high-level matching function like matchPattern, matchPDict, etc...

### Usage

```
mismatch(pattern, x, fixed=TRUE)
nmatch(pattern, x, fixed=TRUE)
nmismatch(pattern, x, fixed=TRUE)
## S4 method for signature 'MIndex'
coverage(x, shift=0L, width=NULL, weight=1L)
## S4 method for signature 'MaskedXString'
coverage(x, shift=0L, width=NULL, weight=1L)
```

### Arguments

The pattern string.

x An XStringViews object for mismatch (typically, one returned by matchPattern(pattern, subject)).

An MIndex object for coverage, or any object for which a coverage method is defined. See ?coverage.

fixed See ?`lowlevel-matching`.

shift, width See ?coverage.

weight An integer vector specifying how much each element in x counts.

#### **Details**

The mismatch function gives the positions of the mismatching letters of a given pattern relatively to its matches in a given subject.

The nmatch and nmismatch functions give the number of matching and mismatching letters produced by the mismatch function.

The coverage function computes the "coverage" of a subject by a given pattern or set of patterns.

#### Value

mismatch: a list of integer vectors.

nmismatch: an integer vector containing the length of the vectors produced by mismatch.

coverage: an Rle object indicating the coverage of x. See ?coverage for the details. If x is an MIndex object, the coverage of a given position in the underlying sequence (typically the subject used during the search that returned x) is the number of matches (or hits) it belongs to.

### See Also

lowlevel-matching, matchPattern, matchPDict, XString-class, XStringViews-class, MIndex-class, coverage, align-utils in the **pwalign** package

44 matchLRPatterns

## **Examples**

matchLRPatterns

Find paired matches in a sequence

# Description

The matchLRPatterns function finds paired matches in a sequence i.e. matches specified by a left pattern, a right pattern and a maximum distance between the left pattern and the right pattern.

## Usage

## **Arguments**

Lpattern The left part of the pattern. Rpattern The right part of the pattern. max.gaplength The max length of the gap in the middle i.e the max distance between the left and right parts of the pattern. subject An XString, XStringViews or MaskedXString object containing the target sequence. The maximum number of mismatching letters allowed in the left part of the patmax.Lmismatch tern. If non-zero, an inexact matching algorithm is used (see the matchPattern function for more information). max.Rmismatch Same as max. Lmismatch but for the right part of the pattern.

matchLRPatterns 45

with.Lindels 
If TRUE then indels are allowed in the left part of the pattern. In that case

max.Lmismatch is interpreted as the maximum "edit distance" allowed in the left part of the pattern.

See the with.indels argument of the matchPattern function for more information.

with.Rindels Same as with.Lindels but for the right part of the pattern.

Unly with a DNAString or RNAString subject can a Lfixed value other than the

default (TRUE) be used.

With Lfixed=FALSE, ambiguities (i.e. letters from the IUPAC Extended Genetic Alphabet (see IUPAC\_CODE\_MAP) that are not from the base alphabet) in the left pattern *and* in the subject are interpreted as wildcards i.e. they match any letter that they stand for.

Lfixed can also be a character vector, a subset of c("pattern", "subject"). Lfixed=c("pattern", "subject") is equivalent to Lfixed=TRUE (the default). An empty vector is equivalent to Lfixed=FALSE. With Lfixed="subject", ambiguities in the pattern only are interpreted as wildcards. With Lfixed="pattern",

ambiguities in the subject only are interpreted as wildcards.

Rfixed Same as Lfixed but for the right part of the pattern.

#### Value

An XStringViews object containing all the matches, even when they are overlapping (see the examples below), and where the matches are ordered from left to right (i.e. by ascending starting position).

# Author(s)

H. Pagès

#### See Also

matchPattern, matchProbePair, trimLRPatterns, findPalindromes, reverseComplement, XString-class, XStringViews-class, MaskedXString-class

```
library(BSgenome.Dmelanogaster.UCSC.dm3)
subject <- Dmelanogaster$chr3R
Lpattern <- "AGCTCCGAG"
Rpattern <- "TTGTTCACA"
matchLRPatterns(Lpattern, Rpattern, 500, subject) # 1 match
## Note that matchLRPatterns() will return all matches, even when they are
## overlapping:
subject <- DNAString("AAATTAACCCTT")
matchLRPatterns("AA", "TT", 0, subject) # 1 match
matchLRPatterns("AA", "TT", 1, subject) # 2 matches
matchLRPatterns("AA", "TT", 3, subject) # 3 matches
matchLRPatterns("AA", "TT", 7, subject) # 4 matches</pre>
```

matchPattern

String searching functions

## **Description**

A set of functions for finding all the occurrences (aka "matches" or "hits") of a given pattern (typically short) in a (typically long) reference sequence or set of reference sequences (aka the subject)

# Usage

```
matchPattern(pattern, subject,
             max.mismatch=0, min.mismatch=0,
             with.indels=FALSE, fixed=TRUE,
             algorithm="auto")
countPattern(pattern, subject,
             max.mismatch=0, min.mismatch=0,
             with.indels=FALSE, fixed=TRUE,
             algorithm="auto")
vmatchPattern(pattern, subject,
              max.mismatch=0, min.mismatch=0,
              with.indels=FALSE, fixed=TRUE,
              algorithm="auto", ...)
vcountPattern(pattern, subject,
              max.mismatch=0, min.mismatch=0,
              with.indels=FALSE, fixed=TRUE,
              algorithm="auto", ...)
```

# Arguments

pattern The pattern string.

subject An XString, XStringViews or MaskedXString object for matchPattern and

countPattern.

An XStringSet or XStringViews object for vmatchPattern and vcountPattern.

max.mismatch, min.mismatch

The maximum and minimum number of mismatching letters allowed (see ?`lowlevel-matching`

for the details). If non-zero, an algorithm that supports inexact matching is used.

with.indels

If TRUE then indels are allowed. In that case, min.mismatch must be  $\emptyset$  and max.mismatch is interpreted as the maximum "edit distance" allowed between the pattern and a match. Note that in order to avoid pollution by redundant matches, only the "best local matches" are returned. Roughly speaking, a "best local match" is a match that is locally both the closest (to the pattern P) and the shortest. More precisely, a substring S' of the subject S is a "best local match"

iff:

```
(a) nedit(P, S') <= max.mismatch
(b) for every substring S1 of S':
        nedit(P, S1) > nedit(P, S')
(c) for every substring S2 of S that contains S':
        nedit(P, S2) >= nedit(P, S')
```

One nice property of "best local matches" is that their first and last letters are

guaranteed to match the letters in P that they align with.

fixed If TRUE (the default), an IUPAC ambiguity code in the pattern can only match

the same code in the subject, and vice versa. If FALSE, an IUPAC ambiguity code in the pattern can match any letter in the subject that is associated with the

code, and vice versa. See ?`lowlevel-matching` for more information.

algorithm One of the following: "auto", "naive-exact", "naive-inexact", "boyer-moore",

"shift-or" or "indels".

... Additional arguments for methods.

#### **Details**

Available algorithms are: "naive exact", "naive inexact", "Boyer-Moore-like", "shift-or" and "indels". Not all of them can be used in all situations: restrictions apply depending on the "search criteria" i.e. on the values of the pattern, subject, max.mismatch, min.mismatch, with.indels and fixed arguments.

It is important to note that the algorithm argument is not part of the search criteria. This is because the supported algorithms are interchangeable, that is, if 2 different algorithms are compatible with a given search criteria, then choosing one or the other will not affect the result (but will most likely affect the performance). So there is no "wrong choice" of algorithm (strictly speaking).

Using algorithm="auto" (the default) is recommended because then the best suited algorithm will automatically be selected among the set of algorithms that are valid for the given search criteria.

#### Value

An XString Views object for matchPattern.

A single integer for countPattern.

An MIndex object for vmatchPattern.

An integer vector for vcountPattern, with each element in the vector corresponding to the number of matches in the corresponding element of subject.

# Note

Use matchPDict if you need to match a (big) set of patterns against a reference sequence.

Use pairwiseAlignment from the **pwalign** package if you need to solve a (Needleman-Wunsch) global alignment, a (Smith-Waterman) local alignment, or an (ends-free) overlap alignment problem.

## See Also

- lowlevel-matching
- matchPDict
- pairwiseAlignment in the **pwalign** package
- mismatch
- matchLRPatterns
- matchProbePair
- maskMotif
- alphabetFrequency
- XStringViews class
- MIndex class

```
## -----
## A. matchPattern()/countPattern()
## -----
## A simple inexact matching example with a short subject:
x <- DNAString("AAGCGCGATATG")</pre>
m1 <- matchPattern("GCNNNAT", x)</pre>
m2 <- matchPattern("GCNNNAT", x, fixed=FALSE)</pre>
as.matrix(m2)
## With DNA sequence of yeast chromosome number 1:
data(yeastSEQCHR1)
yeast1 <- DNAString(yeastSEQCHR1)</pre>
PpiI <- "GAACNNNNNCTC" # a restriction enzyme pattern</pre>
match1.PpiI <- matchPattern(PpiI, yeast1, fixed=FALSE)</pre>
match2.PpiI <- matchPattern(PpiI, yeast1, max.mismatch=1, fixed=FALSE)</pre>
## With a genome containing isolated Ns:
library(BSgenome.Celegans.UCSC.ce2)
chrII <- Celegans[["chrII"]]</pre>
alphabetFrequency(chrII)
matchPattern("N", chrII)
matchPattern("TGGGTGTCTTT", chrII) # no match
matchPattern("TGGGTGTCTTT", chrII, fixed=FALSE) # 1 match
## Using wildcards ("N") in the pattern on a genome containing N-blocks:
library(BSgenome.Dmelanogaster.UCSC.dm3)
chrX <- maskMotif(Dmelanogaster$chrX, "N")</pre>
as(chrX, "Views") # 4 non masked regions
matchPattern("TTTATGNTTGGTA", chrX, fixed=FALSE)
## Can also be achieved with no mask:
masks(chrX) <- NULL</pre>
```

```
matchPattern("TTTATGNTTGGTA", chrX, fixed="subject")
## B. vmatchPattern()/vcountPattern()
## -----
## Load Fly upstream sequences (i.e. the sequences 2000 bases upstream of
## annotated transcription starts):
dm3_upstream_filepath <- system.file("extdata",</pre>
                                   "dm3_upstream2000.fa.gz",
                                   package="Biostrings")
dm3_upstream <- readDNAStringSet(dm3_upstream_filepath)</pre>
dm3_upstream
Ebox <- DNAString("CANNTG")</pre>
subject <- dm3_upstream</pre>
mindex <- vmatchPattern(Ebox, subject, fixed="subject")</pre>
nmatch_per_seq <- elementNROWS(mindex) # Get the number of matches per</pre>
                                      # subject element.
sum(nmatch_per_seq) # Total number of matches.
table(nmatch_per_seq)
## Let's have a closer look at one of the upstream sequences with most
## matches:
i0 <- which.max(nmatch_per_seq)</pre>
subject0 <- subject[[i0]]</pre>
ir0 <- mindex[[i0]] # matches in 'subject0' as an IRanges object</pre>
Views(subject0, ir0) # matches in 'subject0' as a Views object
## -----
## C. WITH INDELS
## -----
library(BSgenome.Celegans.UCSC.ce2)
subject <- Celegans$chrI</pre>
pattern1 <- DNAString("ACGGACCTAATGTTATC")</pre>
pattern2 <- DNAString("ACGGACCTVATGTTRTC")</pre>
## Allowing up to 2 mismatching letters doesn't give any match:
m1a <- matchPattern(pattern1, subject, max.mismatch=2)</pre>
## But allowing up to 2 edit operations gives 3 matches:
system.time(m1b <- matchPattern(pattern1, subject, max.mismatch=2,</pre>
                              with.indels=TRUE))
m1b
## pwalign::pairwiseAlignment() returns the (first) best match only:
if (interactive()) {
 library(pwalign)
 mat <- nucleotideSubstitutionMatrix(match=1, mismatch=0, baseOnly=TRUE)</pre>
 ## Note that this call to pairwiseAlignment() will need to
 ## allocate 733.5 Mb of memory (i.e. length(pattern) * length(subject)
```

```
## * 3 bytes).
  system.time(pwa <- pairwiseAlignment(pattern1, subject, type="local",</pre>
                                      substitutionMatrix=mat,
                                      gapOpening=0, gapExtension=1))
 pwa
}
## With IUPAC ambiguities in the pattern:
m2a <- matchPattern(pattern2, subject, max.mismatch=2,</pre>
                   fixed="subject")
m2b <- matchPattern(pattern2, subject, max.mismatch=2,</pre>
                   with.indels=TRUE, fixed="subject")
## All the matches in 'm1b' and 'm2a' should also appear in 'm2b':
stopifnot(suppressWarnings(all(ranges(m1b) %in% ranges(m2b))))
stopifnot(suppressWarnings(all(ranges(m2a) %in% ranges(m2b))))
## D. WHEN 'with.indels=TRUE', ONLY "BEST LOCAL MATCHES" ARE REPORTED
## -----
## With deletions in the subject:
subject <- BString("ACDEFxxxCDEFxxxABCE")</pre>
matchPattern("ABCDEF", subject, max.mismatch=2, with.indels=TRUE)
matchPattern("ABCDEF", subject, max.mismatch=2)
## With insertions in the subject:
subject <- BString("AiBCDiEFxxxABCDiiFxxxAiBCDEFxxxABCiDEF")</pre>
matchPattern("ABCDEF", subject, max.mismatch=2, with.indels=TRUE)
matchPattern("ABCDEF", subject, max.mismatch=2)
## With substitutions (note that the "best local matches" can introduce
## indels and therefore be shorter than 6):
subject <- BString("AsCDEFxxxABDCEFxxxBACDEFxxxABCEDF")</pre>
matchPattern("ABCDEF", subject, max.mismatch=2, with.indels=TRUE)
matchPattern("ABCDEF", subject, max.mismatch=2)
```

matchPDict

Matching a dictionary of patterns against a reference

# Description

A set of functions for finding all the occurrences (aka "matches" or "hits") of a set of patterns (aka the dictionary) in a reference sequence or set of reference sequences (aka the subject)

The following functions differ in what they return: matchPDict returns the "where" information i.e. the positions in the subject of all the occurrences of every pattern; countPDict returns the "how many times" information i.e. the number of occurrences for each pattern; and whichPDict returns the "who" information i.e. which patterns in the input dictionary have at least one match.

vcountPDict and vwhichPDict are vectorized versions of countPDict and whichPDict, respectively, that is, they work on a set of reference sequences in a vectorized fashion.

This man page shows how to use these functions (aka the \*PDict functions) for exact matching of a constant width dictionary i.e. a dictionary where all the patterns have the same length (same number of nucleotides).

See ?'matchPDict-inexact' for how to use these functions for inexact matching or when the original dictionary has a variable width.

#### Usage

```
matchPDict(pdict, subject,
           max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
           algorithm="auto", verbose=FALSE)
countPDict(pdict, subject,
           max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
           algorithm="auto", verbose=FALSE)
whichPDict(pdict, subject,
           max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
           algorithm="auto", verbose=FALSE)
vcountPDict(pdict, subject,
            max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
            algorithm="auto", collapse=FALSE, weight=1L,
            verbose=FALSE, ...)
vwhichPDict(pdict, subject,
            max.mismatch=0, min.mismatch=0, with.indels=FALSE, fixed=TRUE,
            algorithm="auto", verbose=FALSE)
```

# **Arguments**

pdict

A PDict object containing the preprocessed dictionary.

All these functions also work with a dictionary that has not been preprocessed (in other words, the pdict argument can receive an XStringSet object). Of course, it won't be as fast as with a preprocessed dictionary, but it will generally be slightly faster than using matchPattern/countPattern or vmatchPattern/vcountPattern in a "lapply/sapply loop", because, here, looping is done at the C-level. However, by using a non-preprocessed dictionary, many of the restrictions that apply to preprocessed dictionaries don't apply anymore. For example, the dictionary doesn't need to be rectangular or to be a DNAStringSet object: it can be any type of XStringSet object and have a variable width.

subject

An XString or MaskedXString object containing the subject sequence for matchPDict, countPDict and whichPDict.

An XStringSet object containing the subject sequences for vcountPDict and vwhichPDict.

If pdict is a PDict object (i.e. a preprocessed dictionary), then subject must be of base class DNAString. Otherwise, subject must be of the same base class as pdict.

max.mismatch, min.mismatch

The maximum and minimum number of mismatching letters allowed (see ?isMatchingAt for the details). This man page focuses on exact matching of a constant width

> dictionary so max.mismatch=0 in the examples below. See ?`matchPDict-inexact` for inexact matching.

with.indels

Only supported by countPDict, whichPDict, vcountPDict and vwhichPDict at the moment, and only when the input dictionary is non-preprocessed (i.e. XStringSet).

If TRUE then indels are allowed. In that case, min.mismatch must be 0 and max.mismatch is interpreted as the maximum "edit distance" allowed between any pattern and any of its matches. See ?`matchPattern` for more information.

fixed

Whether IUPAC ambiguity codes should be interpreted literally or not (see ?isMatchingAt for more information). This man page focuses on exact matching of a constant width dictionary so fixed=TRUE in the examples below. See ?`matchPDict-inexact` for inexact matching.

algorithm

Ignored if pdict is a preprocessed dictionary (i.e. a PDict object). Otherwise, can be one of the following: "auto", "naive-exact", "naive-inexact", "boyer-moore" or "shift-or". See ?matchPattern for more information. Note that "indels" is not supported for now.

verbose

TRUE or FALSE.

collapse must be FALSE, 1, or 2.

If collapse=FALSE (the default), then weight is ignored and vcountPDict returns the full matrix of counts (M0). If collapse=1, then M0 is collapsed "horizontally" i.e. it is turned into a vector with length equal to length(pdict). If weight=1L (the default), then this vector is defined by rowSums(M0). If collapse=2, then M0 is collapsed "vertically" i.e. it is turned into a vector with length equal to length(subject). If weight=1L (the default), then this vector is defined by colSums (M0).

If collapse=1 or collapse=2, then the elements in subject (collapse=1) or in pdict (collapse=2) can be weighted thru the weight argument. In that case, the returned vector is defined by M0 %\*% rep(weight, length.out=length(subject)) and rep(weight, length.out=length(pdict)) %\*% M0, respectively.

Additional arguments for methods. . . .

#### **Details**

In this man page, we assume that you know how to preprocess a dictionary of DNA patterns that can then be used with any of the \*PDict functions described here. Please see ?PDict if you don't.

When using the \*PDict functions for exact matching of a constant width dictionary, the standard way to preprocess the original dictionary is by calling the PDict constructor on it with no extra arguments. This returns the preprocessed dictionary in a PDict object that can be used with any of the \*PDict functions.

### Value

If M denotes the number of patterns in the pdict argument (M <- length(pdict)), then matchPDict returns an MIndex object of length M, and countPDict an integer vector of length M.

whichPDict returns an integer vector made of the indices of the patterns in the pdict argument that have at least one match.

collapse, weight

If N denotes the number of sequences in the subject argument (N <- length(subject)), then vcountPDict returns an integer matrix with M rows and N columns, unless the collapse argument is used. In that case, depending on the type of weight, an integer or numeric vector is returned (see above for the details).

vwhichPDict returns a list of N integer vectors.

#### Author(s)

H. Pagès

#### References

Aho, Alfred V.; Margaret J. Corasick (June 1975). "Efficient string matching: An aid to bibliographic search". Communications of the ACM 18 (6): 333-340.

## See Also

PDict-class, MIndex-class, matchPDict-inexact, isMatchingAt, coverage, MIndex-method, matchPattern, alphabetFrequency, DNAStringSet-class, XStringViews-class, MaskedDNAString-class

```
## A. A SIMPLE EXAMPLE OF EXACT MATCHING
## -----
## Creating the pattern dictionary:
library(drosophila2probe)
dict0 <- DNAStringSet(drosophila2probe)</pre>
dict0
                                   # The original dictionary.
                                   # Hundreds of thousands of patterns.
length(dict0)
pdict0 <- PDict(dict0)</pre>
                                  # Store the original dictionary in
                                   # a PDict object (preprocessing).
## Using the pattern dictionary on chromosome 3R:
library(BSgenome.Dmelanogaster.UCSC.dm3)
chr3R <- Dmelanogaster$chr3R</pre>
                            # Load chromosome 3R
chr3R
mi0 <- matchPDict(pdict0, chr3R)
                                   # Search...
## Looking at the matches:
start_index <- startIndex(mi0)</pre>
                                # Get the start index.
length(start_index)
                                 # Same as the original dictionary.
                                 # Starts of the 8220th pattern.
start_index[[8220]]
end_index <- endIndex(mi0)  # Get the end index.
                                   # Ends of the 8220th pattern.
end_index[[8220]]
nmatch_per_pat <- elementNROWS(mi0) # Get the number of matches per pattern.</pre>
nmatch_per_pat[[8220]]
mi0[[8220]]
                                   # Get the matches for the 8220th pattern.
start(mi0[[8220]])
                                   # Equivalent to startIndex(mi0)[[8220]].
sum(nmatch_per_pat)
                                   # Total number of matches.
table(nmatch_per_pat)
```

```
i0 <- which(nmatch_per_pat == max(nmatch_per_pat))</pre>
                                  # The pattern with most occurrences.
pdict0[[i0]]
mi0[[i0]]
                                  # Its matches as an IRanges object.
Views(chr3R, mi0[[i0]])
                                  # And as an XStringViews object.
## Get the coverage of the original subject:
cov3R <- as.integer(coverage(mi0, width=length(chr3R)))</pre>
max(cov3R)
mean(cov3R)
sum(cov3R != 0) / length(cov3R) # Only 2.44% of chr3R is covered.
if (interactive()) {
 library(graphics)
 plotCoverage <- function(cx, start, end)</pre>
   graphics::plot.new()
   graphics::plot.window(c(start, end), c(0, 20))
   graphics::axis(1)
   graphics::axis(2)
   graphics::axis(4)
   graphics::lines(start:end, cx[start:end], type="1")
 plotCoverage(cov3R, 27600000, 27900000)
}
## -----
## B. NAMING THE PATTERNS
## -----
## The names of the original patterns, if any, are propagated to the
## PDict and MIndex objects:
names(dict0) <- mkAllStrings(letters, 4)[seq_len(length(dict0))]</pre>
dict0
dict0[["abcd"]]
pdict0n <- PDict(dict0)</pre>
names(pdict0n)[1:30]
pdict0n[["abcd"]]
mi0n <- matchPDict(pdict0n, chr3R)</pre>
names(mi0n)[1:30]
mi0n[["abcd"]]
## This is particularly useful when unlisting an MIndex object:
unlist(mi0)[1:10]
unlist(mi0n)[1:10] # keep track of where the matches are coming from
## C. PERFORMANCE
## -----
## If getting the number of matches is what matters only (without
## regarding their positions), then countPDict() will be faster,
## especially when there is a high number of matches:
nmatch_per_pat0 <- countPDict(pdict0, chr3R)</pre>
```

```
stopifnot(identical(nmatch_per_pat0, nmatch_per_pat))
if (interactive()) {
 ## What's the impact of the dictionary width on performance?
 ## Below is some code that can be used to figure out (will take a long
 ## time to run). For different widths of the original dictionary, we
      o pptime: preprocessing time (in sec.) i.e. time needed for
 ##
                  building the PDict object from the truncated input
 ##
                  sequences;
 ##
      o nnodes: nb of nodes in the resulting Aho-Corasick tree;
      o nupatt: nb of unique truncated input sequences;
      o matchtime: time (in sec.) needed to find all the matches;
      o totalcount: total number of matches.
 getPDictStats <- function(dict, subject)</pre>
    ans_width <- width(dict[1])</pre>
    ans_pptime <- system.time(pdict <- PDict(dict))[["elapsed"]]</pre>
    pptb <- pdict@threeparts@pptb</pre>
    ans_nnodes <- nnodes(pptb)</pre>
    ans_nupatt <- sum(!duplicated(pdict))</pre>
    ans_matchtime <- system.time(</pre>
                        mi0 <- matchPDict(pdict, subject)</pre>
                     )[["elapsed"]]
    ans_totalcount <- sum(elementNROWS(mi0))</pre>
    list(
      width=ans_width,
      pptime=ans_pptime,
      nnodes=ans_nnodes,
      nupatt=ans_nupatt,
      matchtime=ans_matchtime,
      totalcount=ans_totalcount
   )
 stats <- lapply(8:25,
                function(width)
                    getPDictStats(DNAStringSet(dict0, end=width), chr3R))
 stats <- data.frame(do.call(rbind, stats))</pre>
 stats
}
## D. USING A NON-PREPROCESSED DICTIONARY
dict3 <- DNAStringSet(mkAllStrings(DNA_BASES, 3)) # all trinucleotides</pre>
pdict3 <- PDict(dict3)</pre>
## The 3 following calls are equivalent (from faster to slower):
res3a <- countPDict(pdict3, chr3R)</pre>
res3b <- countPDict(dict3, chr3R)</pre>
res3c <- sapply(dict3,
```

```
function(pattern) countPattern(pattern, chr3R))
stopifnot(identical(res3a, res3b))
stopifnot(identical(res3a, res3c))
## One reason for using a non-preprocessed dictionary is to get rid of
## all the constraints associated with preprocessing, e.g., when
## preprocessing with PDict(), the input dictionary must be DNA and a
## Trusted Band must be defined (explicitly or implicitly).
## See '?PDict' for more information about these constraints.
## In particular, using a non-preprocessed dictionary can be
## useful for the kind of inexact matching that can't be achieved
## with a PDict object (if performance is not an issue).
## See '?`matchPDict-inexact`' for more information about inexact
## matching.
dictD <- xscat(dict3, "N", reverseComplement(dict3))</pre>
## The 2 following calls are equivalent (from faster to slower):
resDa <- matchPDict(dictD, chr3R, fixed=FALSE)</pre>
resDb <- sapply(dictD,</pre>
                function(pattern)
                 matchPattern(pattern, chr3R, fixed=FALSE))
stopifnot(all(sapply(seq_len(length(dictD)),
                     function(i)
                       identical(resDa[[i]], as(resDb[[i]], "IRanges")))))
## A non-preprocessed dictionary can be of any base class i.e. BString,
## RNAString, and AAString, in addition to DNAString:
matchPDict(AAStringSet(c("DARC", "EGH")), AAString("KMFPRNDEGHSTTWTEE"))
## E. vcountPDict()
## -----
## Load Fly upstream sequences (i.e. the sequences 2000 bases upstream of
## annotated transcription starts):
dm3_upstream_filepath <- system.file("extdata",</pre>
                                     "dm3_upstream2000.fa.gz",
                                     package="Biostrings")
dm3_upstream <- readDNAStringSet(dm3_upstream_filepath)</pre>
dm3_upstream
subject <- dm3_upstream[1:100]</pre>
mat1 <- vcountPDict(pdict0, subject)</pre>
dim(mat1) # length(pdict0) x length(subject)
nhit_per_probe <- rowSums(mat1)</pre>
table(nhit_per_probe)
## Without vcountPDict(), 'mat1' could have been computed with:
mat2 <- sapply(unname(subject), function(x) countPDict(pdict0, x))</pre>
stopifnot(identical(mat1, mat2))
## but using vcountPDict() is faster (10x or more, depending of the
## average length of the sequences in 'subject').
```

```
if (interactive()) {
 ## This will fail (with message "allocMatrix: too many elements
 ## specified") because, on most platforms, vectors and matrices in R
 ## are limited to 2^31 elements:
 subject <- dm3_upstream</pre>
 vcountPDict(pdict0, subject)
 length(pdict0) * length(dm3_upstream)
 1 * length(pdict0) * length(dm3_upstream) # > 2^31
 ## But this will work:
 nhit_per_seq <- vcountPDict(pdict0, subject, collapse=2)</pre>
 sum(nhit_per_seq >= 1) # nb of subject sequences with at least 1 hit
 table(nhit_per_seq) # max is 74
 which.max(nhit_per_seq) # 1133
 sum(countPDict(pdict0, subject[[1133]])) # 74
}
## -----
## F. RELATIONSHIP BETWEEN vcountPDict(), countPDict() AND
## vcountPattern()
## -----
subject <- dm3_upstream</pre>
## The 4 following calls are equivalent (from faster to slower):
mat3a <- vcountPDict(pdict3, subject)</pre>
mat3b <- vcountPDict(dict3, subject)</pre>
mat3c <- sapply(dict3,</pre>
              function(pattern) vcountPattern(pattern, subject))
mat3d <- sapply(unname(subject),</pre>
              function(x) countPDict(pdict3, x))
stopifnot(identical(mat3a, mat3b))
stopifnot(identical(mat3a, t(mat3c)))
stopifnot(identical(mat3a, mat3d))
## The 3 following calls are equivalent (from faster to slower):
nhitpp3a <- vcountPDict(pdict3, subject, collapse=1) # rowSums(mat3a)</pre>
nhitpp3b <- vcountPDict(dict3, subject, collapse=1)</pre>
nhitpp3c <- sapply(dict3,</pre>
                 function(pattern) sum(vcountPattern(pattern, subject)))
stopifnot(identical(nhitpp3a, nhitpp3b))
stopifnot(identical(nhitpp3a, nhitpp3c))
## The 3 following calls are equivalent (from faster to slower):
nhitps3a <- vcountPDict(pdict3, subject, collapse=2) # colSums(mat3a)</pre>
nhitps3b <- vcountPDict(dict3, subject, collapse=2)</pre>
nhitps3c <- sapply(unname(subject),</pre>
                 function(x) sum(countPDict(pdict3, x)))
stopifnot(identical(nhitps3a, nhitps3b))
stopifnot(identical(nhitps3a, nhitps3c))
## -----
## G. vwhichPDict()
## -----
```

```
subject <- dm3_upstream</pre>
## The 4 following calls are equivalent (from faster to slower):
vwp3a <- vwhichPDict(pdict3, subject)</pre>
vwp3b <- vwhichPDict(dict3, subject)</pre>
vwp3c <- lapply(seq_len(ncol(mat3a)), function(j) which(mat3a[ , j] != 0L))</pre>
vwp3d <- lapply(unname(subject), function(x) whichPDict(pdict3, x))</pre>
stopifnot(identical(vwp3a, vwp3b))
stopifnot(identical(vwp3a, vwp3c))
stopifnot(identical(vwp3a, vwp3d))
table(sapply(vwp3a, length))
which.min(sapply(vwp3a, length))
## Get the trinucleotides not represented in upstream sequence 21823:
dict3[-vwp3a[[21823]]] # 2 trinucleotides
## Sanity check:
tnf <- trinucleotideFrequency(subject[[21823]])</pre>
stopifnot(all(names(tnf)[tnf == 0] == dict3[-vwp3a[[21823]]]))
## -----
## H. MAPPING PROBE SET IDS BETWEEN CHIPS WITH vwhichPDict()
## -----
## Here we show a simple (and very naive) algorithm for mapping probe
## set IDs between the hgu95av2 and hgu133a chips (Affymetrix).
## 2 probe set IDs are considered mapped iff they share at least one
## WARNING: This example takes about 10 minutes to run.
if (interactive()) {
 library(hgu95av2probe)
 library(hgu133aprobe)
 probes1 <- DNAStringSet(hgu95av2probe)</pre>
 probes2 <- DNAStringSet(hgu133aprobe)</pre>
 pdict2 <- PDict(probes2)</pre>
 ## Get the mapping from probes1 to probes2 (based on exact matching):
 map1to2 <- vwhichPDict(pdict2, probes1)</pre>
 ## The following helper function uses the probe level mapping to induce
 ## the mapping at the probe set IDs level (from hgu95av2 to hgu133a).
 ## To keep things simple, 2 probe set IDs are considered mapped iff
 ## each of them contains at least one probe mapped to one probe of
 ## the other:
 mapProbeSetIDs1to2 <- function(psID)</pre>
   unique(hgu133aprobe$Probe.Set.Name[unlist(
     map1to2[hgu95av2probe$Probe.Set.Name == psID]
   )])
 ## Use the helper function to build the complete mapping:
 psIDs1 <- unique(hgu95av2probe$Probe.Set.Name)</pre>
 mapPSIDs1to2 <- lapply(psIDs1, mapProbeSetIDs1to2) # about 3 min.</pre>
 names(mapPSIDs1to2) <- psIDs1</pre>
```

```
## Do some basic stats:
 table(sapply(mapPSIDs1to2, length))
 ## [ADVANCED USERS ONLY]
 ## An alternative that is slightly faster is to put all the probes
 ## (hgu95av2 + hgu133a) in a single PDict object and then query its
 ## 'dups0' slot directly. This slot is a Dups object containing the
 ## mapping between duplicated patterns.
 ## Note that we can do this only because all the probes have the
 ## same length (25) and because we are doing exact matching:
 probes12 <- DNAStringSet(c(hgu95av2probe$sequence, hgu133aprobe$sequence))</pre>
 pdict12 <- PDict(probes12)</pre>
 dups0 <- pdict12@dups0</pre>
 mapProbeSetIDs1to2alt <- function(psID)</pre>
    ii1 <- unique(togroup(dups0, which(hgu95av2probe$Probe.Set.Name == psID)))</pre>
    ii2 <- members(dups0, ii1) - length(probes1)</pre>
    ii2 <- ii2[ii2 >= 1L]
    unique(hgu133aprobe$Probe.Set.Name[ii2])
 }
 mapPSIDs1to2alt <- lapply(psIDs1, mapProbeSetIDs1to2alt) # about 5 min.</pre>
 names(mapPSIDs1to2alt) <- psIDs1</pre>
 ## 'mapPSIDs1to2alt' and 'mapPSIDs1to2' contain the same mapping:
 stopifnot(identical(lapply(mapPSIDs1to2alt, sort),
                      lapply(mapPSIDs1to2, sort)))
}
```

matchPDict-inexact

Inexact matching with matchPDict()/countPDict()/whichPDict()

## **Description**

The matchPDict, countPDict and whichPDict functions efficiently find the occurrences in a text (the subject) of all patterns stored in a preprocessed dictionary.

This man page shows how to use these functions for inexact (or fuzzy) matching or when the original dictionary has a variable width.

See ?matchPDict for how to use these functions for exact matching of a constant width dictionary i.e. a dictionary where all the patterns have the same length (same number of nucleotides).

### **Details**

In this man page, we assume that you know how to preprocess a dictionary of DNA patterns that can then be used with matchPDict, countPDict or whichPDict. Please see ?PDict if you don't.

matchPDict and family support different kinds of inexact matching but with some restrictions. Inexact matching is controlled via the definition of a Trusted Band during the preprocessing step and/or via the max.mismatch, min.mismatch and fixed arguments. Defining a Trusted Band is also required when the original dictionary is not rectangular (variable width), even for exact matching. See ?PDict for how to define a Trusted Band.

Here is how matchPDict and family handle the Trusted Band defined on pdict:

- (1) Find all the exact matches of all the elements in the Trusted Band.
- (2) For each element in the Trusted Band that has at least one exact match, compare the head and the tail of this element with the flanking sequences of the matches found in (1).

Note that the number of exact matches found in (1) will decrease exponentially with the width of the Trusted Band. Here is a simple guideline in order to get reasonably good performance: if TBW is the width of the Trusted Band (TBW <- tb.width(pdict)) and L the number of letters in the subject (L <- nchar(subject)), then L / (4^TBW) should be kept as small as possible, typically < 10 or 20.

In addition, when a Trusted Band has been defined during preprocessing, then matchPDict and family can be called with fixed=FALSE. In this case, IUPAC ambiguity codes in the head or the tail of the PDict object are treated as ambiguities.

Finally, fixed="pattern" can be used to indicate that IUPAC ambiguity codes in the subject should be treated as ambiguities. It only works if the density of codes is not too high. It works whether or not a Trusted Band has been defined on pdict.

#### Author(s)

H. Pagès

#### References

Aho, Alfred V.; Margaret J. Corasick (June 1975). "Efficient string matching: An aid to bibliographic search". Communications of the ACM 18 (6): 333-340.

### See Also

PDict-class, MIndex-class, matchPDict

```
tail(pdict9)
sum(duplicated(pdict9))
table(patternFrequency(pdict9))
library(BSgenome.Dmelanogaster.UCSC.dm3)
chr3R <- Dmelanogaster$chr3R</pre>
chr3R
table(countPDict(pdict9, chr3R, max.mismatch=1))
table(countPDict(pdict9, chr3R, max.mismatch=3))
table(countPDict(pdict9, chr3R, max.mismatch=5))
## -----
## B. COMPARISON WITH EXACT MATCHING
## When the original dictionary is of constant width, exact matching
## (i.e. 'max.mismatch=0' and 'fixed=TRUE) will be more efficient with
## a full-width Trusted Band (i.e. a Trusted Band that covers the entire
## dictionary) than with a Trusted Band of width < width(dict0).
pdict0 <- PDict(dict0)</pre>
count0 <- countPDict(pdict0, chr3R)</pre>
count0b <- countPDict(pdict9, chr3R, max.mismatch=0)</pre>
identical(count0b, count0) # TRUE
## -----
## C. USING AN EXPLICIT TRUSTED BAND ON A VARIABLE WIDTH DICTIONARY
## -----
## Here is a small variable width dictionary that contains IUPAC
## ambiguities (pattern 1 and 3 contain an N):
dict0 <- DNAStringSet(c("TACCNG", "TAGT", "CGGNT", "AGTAG", "TAGT"))</pre>
## (Note that pattern 2 and 5 are identical.)
## If we only want to do exact matching, then it is recommended to use
## the widest possible Trusted Band i.e. to set its width to
## 'min(width(dict0))' because this is what will give the best
## performance. However, when 'dict0' contains IUPAC ambiguities (like
## in our case), it could be that one of them is falling into the
## Trusted Band so we get an error (only base letters can go in the
## Trusted Band for now):
 PDict(dict0, tb.end=min(width(dict0))) # Error!
## End(Not run)
## In our case, the Trusted Band cannot be wider than 3:
pdict <- PDict(dict0, tb.end=3)</pre>
tail(pdict)
subject <- DNAString("TAGTACCAGTTTCGGG")</pre>
m <- matchPDict(pdict, subject)</pre>
elementNROWS(m) # pattern 2 and 5 have 1 exact match
```

```
m[[2]]
## We can take advantage of the fact that our Trusted Band doesn't cover
## the entire dictionary to allow inexact matching on the uncovered parts
## (the tail in our case):
m <- matchPDict(pdict, subject, fixed=FALSE)</pre>
elementNROWS(m) # now pattern 1 has 1 match too
m[[1]]
m <- matchPDict(pdict, subject, max.mismatch=1)</pre>
elementNROWS(m) # now pattern 4 has 1 match too
m[[4]]
m <- matchPDict(pdict, subject, max.mismatch=1, fixed=FALSE)</pre>
elementNROWS(m) # now pattern 3 has 1 match too
m[[3]] # note that this match is "out of limit"
Views(subject, m[[3]])
m <- matchPDict(pdict, subject, max.mismatch=2)</pre>
elementNROWS(m) # pattern 4 gets 1 additional match
m[[4]]
## Unlist all matches:
unlist(m)
## D. WITH IUPAC AMBIGUITY CODES IN THE PATTERNS
## The Trusted Band cannot contain IUPAC ambiguity codes so patterns
## with ambiguity codes can only be preprocessed if we can define a
## Trusted Band with no ambiguity codes in it.
dict <- DNAStringSet(c("AAACAAKS", "GGGAAA", "TNCCGGG"))</pre>
pdict <- PDict(dict, tb.start=3, tb.width=4)</pre>
subject <- DNAString("AAACAATCCCGGGAAACAAGG")</pre>
matchPDict(pdict, subject)
matchPDict(pdict, subject, fixed="subject")
## Sanity checks:
res1 <- as.list(matchPDict(pdict, subject))</pre>
res2 <- as.list(matchPDict(dict, subject))</pre>
res3 <- lapply(dict,</pre>
  function(pattern)
    as(matchPattern(pattern, subject), "IRanges"))
stopifnot(identical(res1, res2))
stopifnot(identical(res1, res3))
res1 <- as.list(matchPDict(pdict, subject, fixed="subject"))</pre>
res2 <- as.list(matchPDict(dict, subject, fixed="subject"))</pre>
res3 <- lapply(dict,
  function(pattern)
```

matchProbePair 63

```
as(matchPattern(pattern, subject, fixed="subject"), "IRanges"))
stopifnot(identical(res1, res2))
stopifnot(identical(res1, res3))
## E. WITH IUPAC AMBIGUITY CODES IN THE SUBJECT
## -----
## 'fixed="pattern"' (or 'fixed=FALSE') can be used to indicate that
## IUPAC ambiguity codes in the subject should be treated as ambiguities.
pdict <- PDict(c("ACAC", "TCCG"))</pre>
matchPDict(pdict, DNAString("ACNCCGT"))
matchPDict(pdict, DNAString("ACNCCGT"), fixed="pattern")
matchPDict(pdict, DNAString("ACWCCGT"), fixed="pattern")
matchPDict(pdict, DNAString("ACRCCGT"), fixed="pattern")
matchPDict(pdict, DNAString("ACKCCGT"), fixed="pattern")
dict <- DNAStringSet(c("TTC", "CTT"))</pre>
pdict <- PDict(dict)</pre>
subject <- DNAString("CYTCACTTC")</pre>
mi1 <- matchPDict(pdict, subject, fixed="pattern")</pre>
mi2 <- matchPDict(dict, subject, fixed="pattern")</pre>
stopifnot(identical(as.list(mi1), as.list(mi2)))
```

matchProbePair

Find "theoretical amplicons" mapped to a probe pair

# **Description**

In the context of a computer-simulated PCR experiment, one wants to find the amplicons mapped to a given primer pair. The matchProbePair function can be used for this: given a forward and a reverse probe (i.e. the chromosome-specific sequences of the forward and reverse primers used for the experiment) and a target sequence (generally a chromosome sequence), the matchProbePair function will return all the "theoretical amplicons" mapped to this probe pair.

# Usage

#### **Arguments**

Fprobe The forward probe.

Rprobe The reverse probe.

subject A DNAString object (or an XStringViews object with a DNAString subject)

containing the target sequence.

64 matchProbePair

algorithm One of the following: "auto", "naive-exact", "naive-inexact", "boyer-moore" or "shift-or". See matchPattern for more information.

logfile A file used for logging.

verbose TRUE or FALSE.

Additional arguments passed to matchPattern.

#### **Details**

The matchProbePair function does the following: (1) find all the "plus hits" i.e. the Fprobe and Rprobe matches on the "plus" strand, (2) find all the "minus hits" i.e. the Fprobe and Rprobe matches on the "minus" strand and (3) from the set of all (plus\_hit, minus\_hit) pairs, extract and return the subset of "reduced matches" i.e. the (plus\_hit, minus\_hit) pairs such that (a) plus\_hit <= minus\_hit and (b) there are no hits (plus or minus) between plus\_hit and minus\_hit. This set of "reduced matches" is the set of "theoretical amplicons".

Additional arguments can be passed to matchPattern via the . . . argument. This supports matching to ambiguity codes. See matchPattern for more information on supported arguments.

#### Value

An XStringViews object containing the set of "theoretical amplicons".

## Author(s)

H. Pagès

# See Also

matchPattern, matchLRPatterns, findPalindromes, reverseComplement, XStringViews-class

```
library(BSgenome.Dmelanogaster.UCSC.dm3)
subject <- Dmelanogaster$chr3R

## With 20-nucleotide forward and reverse probes:
Fprobe <- "AGCTCCGAGTTCCTGCAATA"
Rprobe <- "CGTTGTTCACAAATATGCGG"
matchProbePair(Fprobe, Rprobe, subject) # 1 "theoretical amplicon"

## With shorter forward and reverse probes, the risk of having multiple
## "theoretical amplicons" increases:
Fprobe <- "AGCTCCGAGTTCC"
Rprobe <- "CGTTGTTCACAA"
matchProbePair(Fprobe, Rprobe, subject) # 2 "theoretical amplicons"
Fprobe <- "AGCTCCGAGTT"
Rprobe <- "CGTTGTTCACA"
matchProbePair(Fprobe, Rprobe, subject) # 9 "theoretical amplicons"</pre>
```

matchPWM 65

matchPWM	PWM creating.	matching.	and related utilities
maccin mi	1 Will Cicalling,	maicing,	ana reiaica aiiiiics

Description

Position Weight Matrix (PWM) creating, matching, and related utilities for DNA data. (PWM for amino acid sequences are not supported.)

# Usage

```
PWM(x, type = c("log2probratio", "prob"),
    prior.params = c(A=0.25, C=0.25, G=0.25, T=0.25))

matchPWM(pwm, subject, min.score="80%", with.score=FALSE, ...)
countPWM(pwm, subject, min.score="80%", ...)
PWMscoreStartingAt(pwm, subject, starting.at=1)

## Utility functions for basic manipulation of the Position Weight Matrix
maxWeights(x)
minWeights(x)
maxScore(x)
minScore(x)
unitScale(x)
## S4 method for signature 'matrix'
reverseComplement(x, ...)
```

# **Arguments**

X	For PWM: a rectangular character vector or rectangular DNAStringSet object
	("rectangular" means that all elements have the same number of characters) with
	no IUPAC ambiguity letters, or a Position Frequency Matrix represented as an
	integer matrix with row names containing at least A, C, G and T (typically the

result of a call to consensusMatrix).

For maxWeights, minWeights, maxScore, minScore, unitScale and reverseComplement:

a Position Weight Matrix represented as a numeric matrix with row names A, C,

G and T.

type The type of Position Weight Matrix, either "log2probratio" or "prob". See De-

tails section for more information.

prior.params A positive numeric vector, which represents the parameters of the Dirichlet con-

jugate prior, with names A, C, G, and T. See Details section for more informa-

tion.

pwm A Position Weight Matrix represented as a numeric matrix with row names A,

C, G and T.

subject Typically a DNAString object. A Views object on a DNAString subject, a

MaskedDNAString object, or a single character string, are also supported.

66 matchPWM

	IUPAC ambiguity letters in subject are ignored (i.e. assigned weight 0) with a warning.
min.score	The minimum score for counting a match. Can be given as a character string containing a percentage (e.g. "85%") of the highest possible score or as a single number.
with.score	TRUE or FALSE. If TRUE, then the score of each hit is included in the returned object in a metadata column named score. Say the returned object is hits, this metadata column can then be accessed with mcols(hits)\$score.
starting.at	An integer vector specifying the starting positions of the Position Weight Matrix relatively to the subject.
	Additional arguments for methods.

## **Details**

The PWM function uses a multinomial model with a Dirichlet conjugate prior to calculate the estimated probability of base b at position i. As mentioned in the Arguments section, prior.params supplies the parameters for the DNA bases A, C, G, and T in the Dirichlet prior. These values result in a position independent initial estimate of the probabilities for the bases to be priorProbs = prior.params/sum(prior.params) and the posterior (data infused) estimate for the probabilities for the bases in each of the positions to be postProbs = (consensusMatrix(x) + prior.params)/(length(x) + sum(prior.params)). When type = "log2probratio", the PWM = unitScale(log2(postProbs/priorProbs)). When type = "prob", the PWM = unitScale(postProbs).

#### Value

A numeric matrix representing the Position Weight Matrix for PWM.

A numeric vector containing the Position Weight Matrix-based scores for PWMscoreStartingAt.

An XStringViews object for matchPWM.

A single integer for countPWM.

A vector containing the max weight for each position in pwm for maxWeights.

A vector containing the min weight for each position in pwm for minWeights.

The highest possible score for a given Position Weight Matrix for maxScore.

The lowest possible score for a given Position Weight Matrix for minScore.

The modified numeric matrix given by (x - minScore(x)/ncol(x))/(maxScore(x) - minScore(x)) for unitScale.

A PWM obtained by reverting the column order in PWM x and by reassigning each row to its complementary nucleotide for reverseComplement.

# Author(s)

H. Pagès and P. Aboyoun

# References

Wasserman, WW, Sandelin, A., (2004) Applied bioinformatics for the identification of regulatory elements, Nat Rev Genet., 5(4):276-87.

MIndex-class 67

## See Also

consensusMatrix, matchPattern, reverseComplement, DNAString-class, XStringViews-class

## **Examples**

```
## Data setup:
data(HNF4alpha)
library(BSgenome.Dmelanogaster.UCSC.dm3)
chr3R <- Dmelanogaster$chr3R</pre>
chr3R
## Create a PWM from a PFM or directly from a rectangular
## DNAStringSet object:
pfm <- consensusMatrix(HNF4alpha)</pre>
pwm <- PWM(pfm) # same as 'PWM(HNF4alpha)'</pre>
## Perform some general routines on the PWM:
round(pwm, 2)
maxWeights(pwm)
maxScore(pwm)
reverseComplement(pwm)
## Score the first 5 positions:
PWMscoreStartingAt(pwm, chr3R, starting.at=1:5)
## Match the plus strand:
hits <- matchPWM(pwm, chr3R)
nhit <- countPWM(pwm, chr3R) # same as 'length(hits)'</pre>
## Use 'with.score=TRUE' to get the scores of the hits:
hits <- matchPWM(pwm, chr3R, with.score=TRUE)</pre>
head(mcols(hits)$score)
min(mcols(hits)\$score / maxScore(pwm)) # should be >= 0.8
## The scores can also easily be post-calculated:
scores <- PWMscoreStartingAt(pwm, subject(hits), start(hits))</pre>
## Match the minus strand:
matchPWM(reverseComplement(pwm), chr3R)
```

MIndex-class

MIndex objects

# Description

The MIndex class is the basic container for storing the matches of a set of patterns in a subject sequence.

68 MIndex-class

#### **Details**

An MIndex object contains the matches (start/end locations) of a set of patterns found in an XString object called "the subject string" or "the subject sequence" or simply "the subject".

matchPDict function returns an MIndex object.

#### Accessor methods

In the code snippets below, x is an MIndex object.

length(x): The number of patterns that matches are stored for.

names(x): The names of the patterns that matches are stored for.

startIndex(x): A list containing the starting positions of the matches for each pattern.

endIndex(x): A list containing the ending positions of the matches for each pattern.

elementNROWS(x): An integer vector containing the number of matches for each pattern.

# Subsetting methods

In the code snippets below, x is an MIndex object.

x[[i]]: Extract the matches for the i-th pattern as an IRanges object.

#### Coercion

In the code snippets below, x is an MIndex object.

as(x, "CompressedIRangesList"): Turns x into an CompressedIRangesList object. This coercion changes x from one IntegerRangesList subtype to another with the underlying IntegerRanges values remaining unchanged.

# Other utility methods and functions

In the code snippets below, x and mindex are MIndex objects and subject is the XString object containing the sequence in which the matches were found.

unlist(x, recursive=TRUE, use.names=TRUE): Return all the matches in a single IRanges object. recursive and use.names are ignored.

extractAllMatches(subject, mindex): Return all the matches in a single XStringViews object.

## Author(s)

H. Pagès

# See Also

matchPDict, PDict-class, IRanges-class, XStringViews-class

## **Examples**

## See ?matchPDict and ?`matchPDict-inexact` for some examples.

misc 69

misc

Some miscellaneous stuff

# **Description**

Some miscellaneous stuff.

# Usage

```
N50(csizes)
```

# Arguments

csizes

A vector containing the contig sizes.

## Value

N50: The N50 value as an integer.

## The N50 contig size

**Definition** The N50 contig size of an assembly (aka the N50 value) is the size of the largest contig such that the contigs larger than that have at least 50% the bases of the assembly.

**How is it calculated?** It is calculated by adding the sizes of the biggest contigs until you reach half the total size of the contigs. The N50 value is then the size of the contig that was added last (i.e. the smallest of the big contigs covering 50% of the genome).

What for? The N50 value is a standard measure of the quality of a de novo assembly.

# Author(s)

Nicolas Delhomme <delhomme@embl.de>

#### See Also

XStringSet-class

```
my.size <- width(my.contig)
# Calculate the N50 value of this set of contigs:
my.contig.N50 <- N50(my.size)</pre>
```

MultipleAlignment-class

MultipleAlignment objects

# Description

The MultipleAlignment class is a container for storing multiple sequence alignments.

# Usage

```
## Constructors:
DNAMultipleAlignment(x=character(), start=NA, end=NA, width=NA,
    use.names=TRUE, rowmask=NULL, colmask=NULL)
RNAMultipleAlignment(x=character(), start=NA, end=NA, width=NA,
    use.names=TRUE, rowmask=NULL, colmask=NULL)

AAMultipleAlignment(x=character(), start=NA, end=NA, width=NA,
    use.names=TRUE, rowmask=NULL, colmask=NULL)

## Read functions:
readDNAMultipleAlignment(filepath, format)
readRNAMultipleAlignment(filepath, format)
readAAMultipleAlignment(filepath, format)
## Write funtions:
write.phylip(x, filepath)

## ... and more (see below)
```

#### **Arguments**

Х

Either a character vector (with no NAs), or an XString, XStringSet or XStringViews object containing strings with the same number of characters. If writing out a Phylip file, then x would be a MultipleAlignment object

start, end, width

Either NA, a single integer, or an integer vector of the same length as x specifying how x should be "narrowed" (see ?narrow in the **IRanges** package for the details).

use.names

TRUE or FALSE. Should names be preserved?

filepath

A character vector (of arbitrary length when reading, of length 1 when writing) containing the paths to the files to read or write. Note that special values like "" or "|cmd" (typically supported by other I/O functions in R) are not supported here. Also filepath cannot be a connection.

format Either "fasta" (the default), "stockholm", "phylip", or "clustal".

rowmask a NormalIRanges object that will set masking for rows colmask a NormalIRanges object that will set masking for columns

#### **Details**

The MultipleAlignment class is designed to hold and represent multiple sequence alignments. The rows and columns within an alignment can be masked for ad hoc analyses.

#### Accessor methods

In the code snippets below, x is a MultipleAlignment object.

unmasked(x): The underlying XStringSet object containing the multiple sequence alignment.

rownames(x): NULL or a character vector of the same length as x containing a short user-provided description or comment for each sequence in x.

rowmask(x), rowmask(x, append, invert) <- value: Gets and sets the NormalIRanges object representing the masked rows in x. The append argument takes union, replace or intersect to indicate how to combine the new value with rowmask(x). The invert argument takes a logical argument to indicate whether or not to invert the new mask. The value argument can be of any class that is coercible to a NormalIRanges via the as function.

colmask(x), colmask(x, append, invert) <- value: Gets and sets the NormalIRanges object representing the masked columns in x. The append argument takes union, replace or intersect to indicate how to combine the new value with colmask(x). The invert argument takes a logical argument to indicate whether or not to invert the new mask. The value argument can be of any class that is coercible to a NormalIRanges via the as function.

maskMotif(x, motif, min.block.width=1, ...): Returns a MultipleAlignment object with a modified column mask based upon motifs found in the consensus string where the consensus string keeps all the columns but drops the masked rows.

**motif** The motif to mask.

min.block.width The minimum width of the blocks to mask.

... Additional arguments for matchPattern.

maskGaps(x, min.fraction, min.block.width): Returns a MultipleAlignment object with a modified column mask based upon gaps in the columns. In particular, this mask is defined by min.block.width or more consecutive columns that have min.fraction or more of their non-masked rows containing gap codes.

**min.fraction** A value in [0, 1] that indicates the minimum fraction needed to call a gap in the consensus string (default is 0.5).

**min.block.width** A positive integer that indicates the minimum number of consecutive gaps to mask, as defined by min. fraction (default is 4).

nrow(x): Returns the number of sequences aligned in x.

ncol(x): Returns the number of characters for each alignment in x.

dim(x): Equivalent to c(nrow(x), ncol(x)).

maskednrow(x): Returns the number of masked aligned sequences in x.

```
maskedncol(x): Returns the number of masked aligned characters in x.
maskeddim(x): Equivalent to c(maskednrow(x), maskedncol(x)).
maskedratio(x): Equivalent to maskeddim(x) / dim(x).
nchar(x): Returns the number of unmasked aligned characters in x, i.e. ncol(x) - maskedncol(x).
alphabet(x): Equivalent to alphabet(unmasked(x)).
```

#### Coercion

In the code snippets below, x is a MultipleAlignment object.

- as(from, "DNAStringSet"), as(from, "RNAStringSet"), as(from, "AAStringSet"), as(from, "BStringSet"):
  Creates an instance of the specified XStringSet object subtype that contains the unmasked regions of the multiple sequence alignment in x.
- as.character(x, use.names): Convert x to a character vector containing the unmasked regions of the multiple sequence alignment. use.names controls whether or not rownames(x) should be used to set the names of the returned vector (default is TRUE).
- as.matrix(x, use.names): Returns a character matrix containing the "exploded" representation of the unmasked regions of the multiple sequence alignment. use.names controls whether or not rownames(x) should be used to set the row names of the returned matrix (default is TRUE).

#### **Utilities**

In the code snippets below, x is a MultipleAlignment object.

- consensusMatrix(x, as.prob, baseOnly): Creates an integer matrix containing the column frequencies of the underlying alphabet with masked columns being represented with NA values. If as.prob is TRUE, then probabilities are reported, otherwise counts are reported (the default). If baseOnly is TRUE, then the non-base letters are collapsed into an "other" category.
- consensusString(x, ...): Creates a consensus string for x with the symbol "#" representing a masked column. See consensusString for details on the arguments.
- consensusViews(x, ...): Similar to the consensusString method. It returns a XStringViews on the consensus string containing subsequence contigs of non-masked columns. Unlike the consensusString method, the masked columns in the underlying string contain a consensus value rather than the "#" symbol.
- alphabetFrequency(x, as.prob, collapse): Creates an integer matrix containing the row frequencies of the underlying alphabet. If as.prob is TRUE, then probabilities are reported, otherwise counts are reported (the default). If collapse is TRUE, then returns the overall frequency instead of the frequency by row.
- detail(x, invertColMask, hideMaskedCols): Allows for a full pager driven display of the object so that masked cols and rows can be removed and the entire sequence can be visually inspected. If hideMaskedCols is set to it's default value of TRUE then the output will hide all the the masked columns in the output. Otherwise, all columns will be displayed along with a row to indicate the masking status. If invertColMask is TRUE then any displayed mask will be flipped so as to represent things in a way consistent with Phylip style files instead of the mask that is actually stored in the MultipleAlignment object. Please notice that invertColMask will be ignored if hideMaskedCols is set to its default value of TRUE since in that case it will not make sense to show any masking information in the output. Masked rows are always hidden in the output.

## **Display**

The letters in a DNAMultipleAlignment or RNAMultipleAlignment object are colored when displayed by the show() method. Set global option Biostrings.coloring to FALSE to turn off this coloring.

### Author(s)

P. Aboyoun and M. Carlson

### See Also

XStringSet-class, MaskedXString-class

```
## create an object from file
origMAlign <-
 readDNAMultipleAlignment(filepath =
                            system.file("extdata",
                                         "msx2_mRNA.aln",
                                         package="Biostrings"),
                            format="clustal")
## list the names of the sequences in the alignment
rownames(origMAlign)
## rename the sequences to be the underlying species for MSX2
rownames(origMAlign) <- c("Human", "Chimp", "Cow", "Mouse", "Rat",</pre>
                           "Dog", "Chicken", "Salmon")
origMAlign
## See a detailed pager view
if (interactive()) {
detail(origMAlign)
## operations to mask rows
## For columns, just use colmask() and do the same kinds of operations
rowMasked <- origMAlign</pre>
rowmask(rowMasked) <- IRanges(start=1,end=3)</pre>
rowMasked
## remove rowumn masks
rowmask(rowMasked) <- NULL</pre>
rowMasked
## "select" rows of interest
rowmask(rowMasked, invert=TRUE) <- IRanges(start=4,end=7)</pre>
rowMasked
## or mask the rows that intersect with masked rows
rowmask(rowMasked, append="intersect") <- IRanges(start=1,end=5)</pre>
```

```
rowMasked
## TATA-masked
tataMasked <- maskMotif(origMAlign, "TATA")</pre>
colmask(tataMasked)
## automatically mask rows based on consecutive gaps
autoMasked <- maskGaps(origMAlign, min.fraction=0.5, min.block.width=4)</pre>
colmask(autoMasked)
autoMasked
## calculate frequencies
alphabetFrequency(autoMasked)
consensusMatrix(autoMasked, baseOnly=TRUE)[, 84:90]
## get consensus values
consensusString(autoMasked)
consensusViews(autoMasked)
## cluster the masked alignments
library(pwalign)
sdist <- pwalign::stringDist(as(autoMasked,"DNAStringSet"), method="hamming")</pre>
clust <- hclust(sdist, method = "single")</pre>
plot(clust)
fourgroups <- cutree(clust, 4)</pre>
fourgroups
## write out the alignement object (with current masks) to Phylip format
write.phylip(x = autoMasked, filepath = tempfile("foo.txt",tempdir()))
```

nucleotideFrequency

Calculate the frequency of oligonucleotides in a DNA or RNA sequence (and other related functions)

## **Description**

Given a DNA or RNA sequence (or a set of DNA or RNA sequences), the oligonucleotideFrequency function computes the frequency of all possible oligonucleotides of a given length (called the "width" in this particular context) in a sliding window that is shifted step nucleotides at a time.

The dinucleotideFrequency and trinucleotideFrequency functions are convenient wrappers for calling oligonucleotideFrequency with width=2 and width=3, respectively.

The nucleotideFrequencyAt function computes the frequency of the short sequences formed by extracting the nucleotides found at some fixed positions from each sequence of a set of DNA or RNA sequences.

In this man page we call "DNA input" (or "RNA input") an XString, XStringSet, XStringViews or MaskedXString object of base type DNA (or RNA).

nucleotideFrequency 75

## Usage

```
oligonucleotideFrequency(x, width, step=1,
                         as.prob=FALSE, as.array=FALSE,
                         fast.moving.side="right", with.labels=TRUE, ...)
## S4 method for signature 'XStringSet'
oligonucleotideFrequency(x, width, step=1,
                         as.prob=FALSE, as.array=FALSE,
                         fast.moving.side="right", with.labels=TRUE,
                         simplify.as="matrix")
dinucleotideFrequency(x, step=1,
                      as.prob=FALSE, as.matrix=FALSE,
                      fast.moving.side="right", with.labels=TRUE, ...)
trinucleotideFrequency(x, step=1,
                       as.prob=FALSE, as.array=FALSE,
                       fast.moving.side="right", with.labels=TRUE, ...)
nucleotideFrequencyAt(x, at,
                      as.prob=FALSE, as.array=TRUE,
                      fast.moving.side="right", with.labels=TRUE, ...)
## Some related functions:
oligonucleotideTransitions(x, left=1, right=1, as.prob=FALSE)
mkAllStrings(alphabet, width, fast.moving.side="right")
```

### **Arguments**

 $x \hspace{1cm} \textbf{Any DNA or RNA input for the *Frequency and oligonucleotideTransitions} \\$ 

functions.

An XStringSet or XStringViews object of base type DNA or RNA for nucleotideFrequencyAt.

width The number of nucleotides per oligonucleotide for oligonucleotideFrequency.

The number of letters per string for mkAllStrings.

step How many nucleotides should the window be shifted before counting the next

oligonucleotide (i.e. the sliding window step; default 1). If step is smaller than width, oligonucleotides will overlap; if the two arguments are equal, adjacent oligonucleotides will be counted (an efficient way to count codons in an ORF); and if step is larger than width, nucleotides will be sampled step nucleotides

apart.

at An integer vector containing the positions to look at in each element of x.

as.prob If TRUE then probabilities are reported, otherwise counts (the default).

as.array, as.matrix

Controls the "shape" of the returned object. If TRUE (the default for nucleotideFrequencyAt) then it's a numeric matrix (or array), otherwise it's just a "flat" numeric vector i.e. a vector with no dim attribute (the default for the \*Frequency functions).

fast.moving.side

Which side of the strings should move fastest? Note that, when as array is TRUE, then the supplied value is ignored and the effective value is "left".

with.labels If TRUE then the returned object is named.

... Further arguments to be passed to or from other methods.

simplify.as Together with the as.array and as.matrix arguments, controls the "shape"

of the returned object when the input x is an XStringSet or XStringViews object. Supported simplify.as values are "matrix" (the default), "list" and "collapsed". If simplify.as is "matrix", the returned object is a matrix with length(x) rows where the i-th row contains the frequencies for x[[i]]. If simplify.as is "list", the returned object is a list of the same length as length(x) where the i-th element contains the frequencies for x[[i]]. If simplify.as is "collapsed", then the the frequencies are computed for the entire object x as a whole (i.e. frequencies cumulated across all sequences in x).

left, right The number of nucleotides per oligonucleotide for the rows and columns respec-

tively in the transition matrix created by oligonucleotideTransitions.

alphabet The alphabet to use to make the strings.

### Value

If x is an XString or MaskedXString object, the \*Frequency functions return a numeric vector of length 4^width. If as.array (or as.matrix) is TRUE, then this vector is formatted as an array (or matrix). If x is an XStringSet or XStringViews object, the returned object has the shape specified by the simplify.as argument.

## Author(s)

H. Pagès and P. Aboyoun; K. Vlahovicek for the step argument

## See Also

alphabetFrequency, alphabet, hasLetterAt, XString-class, XStringSet-class, XStringViews-class, MaskedXString-class, GENETIC\_CODE, AMINO\_ACID\_CODE, reverseComplement, rev

```
## ------
## A. BASIC *Frequency() EXAMPLES
## -------
data(yeastSEQCHR1)
yeast1 <- DNAString(yeastSEQCHR1)

dinucleotideFrequency(yeast1)
trinucleotideFrequency(yeast1)
oligonucleotideFrequency(yeast1, 4)

## Get the counts of tetranucleotides overlapping by one nucleotide:
oligonucleotideFrequency(yeast1, 4, step=3)</pre>
```

nucleotideFrequency 77

```
## Get the counts of adjacent tetranucleotides, starting from the first
## nucleotide:
oligonucleotideFrequency(yeast1, 4, step=4)
## Subset the sequence to change the starting nucleotide (here we start
## counting from third nucleotide):
yeast2 <- subseq(yeast1, start=3)</pre>
oligonucleotideFrequency(yeast2, 4, step=4)
## Get the less and most represented 6-mers:
f6 <- oligonucleotideFrequency(yeast1, 6)</pre>
f6[f6 == min(f6)]
f6[f6 == max(f6)]
## Get the result as an array:
tri <- trinucleotideFrequency(yeast1, as.array=TRUE)</pre>
tri["A", "A", "C"] # == trinucleotideFrequency(yeast1)["AAC"]
tri["T", , ] # frequencies of trinucleotides starting with a "T"
## With input made of multiple sequences:
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
dfmat <- dinucleotideFrequency(probes) # a big matrix</pre>
dinucleotideFrequency(probes, simplify.as="collapsed")
dinucleotideFrequency(probes, simplify.as="collapsed", as.matrix=TRUE)
## B. OBSERVED DINUCLEOTIDE FREQUENCY VERSUS EXPECTED DINUCLEOTIDE
     FREQUENCY
## The expected frequency of dinucleotide "ab" based on the frequencies
## of its individual letters "a" and "b" is:
      exp_Fab = Fa * Fb / N if the 2 letters are different (e.g. CG)
      exp_Faa = Fa * (Fa-1) / N if the 2 letters are the same (e.g. TT)
## where Fa and Fb are the frequencies of "a" and "b" (respectively) and
## N the length of the sequence.
## Here is a simple function that implements the above formula for a
## DNAString object 'x'. The expected frequencies are returned in a 4x4
## matrix where the rownames and colnames correspond to the 1st and 2nd
## base in the dinucleotide:
expectedDinucleotideFrequency <- function(x)</pre>
{
    # Individual base frequencies.
    bf <- alphabetFrequency(x, baseOnly=TRUE)[DNA_BASES]</pre>
    (as.matrix(bf) %*% t(bf) - diag(bf)) / length(x)
}
## On Celegans chrI:
library(BSgenome.Celegans.UCSC.ce2)
chrI <- Celegans$chrI</pre>
obs_df <- dinucleotideFrequency(chrI, as.matrix=TRUE)</pre>
obs_df # CG has the lowest frequency
```

```
exp_df <- expectedDinucleotideFrequency(chrI)</pre>
## A sanity check:
stopifnot(as.integer(sum(exp_df)) == sum(obs_df))
## Ratio of observed frequency to expected frequency:
obs_df / exp_df # TA has the lowest ratio, not CG!
## C. nucleotideFrequencyAt()
nucleotideFrequencyAt(probes, 13)
nucleotideFrequencyAt(probes, c(13, 20))
nucleotideFrequencyAt(probes, c(13, 20), as.array=FALSE)
## nucleotideFrequencyAt() can be used to answer questions like: "how
## many probes in the drosophila2 chip have T, G, T, A at position
## 2, 4, 13 and 20, respectively?"
nucleotideFrequencyAt(probes, c(2, 4, 13, 20))["T", "G", "T", "A"]
## or "what's the probability to have an A at position 25 if there is
## one at position 13?"
nf <- nucleotideFrequencyAt(probes, c(13, 25))</pre>
sum(nf["A", "A"]) / sum(nf["A", ])
## Probabilities to have other bases at position 25 if there is an A
## at position 13:
sum(nf["A", "C"]) / sum(nf["A", ]) # C
sum(nf["A", "G"]) / sum(nf["A", ]) # G
sum(nf["A", "T"]) / sum(nf["A", ]) # T
## See ?hasLetterAt for another way to get those results.
## -----
## D. oligonucleotideTransitions()
## -----
## Get nucleotide transition matrices for yeast1
oligonucleotideTransitions(yeast1)
oligonucleotideTransitions(yeast1, 2, as.prob=TRUE)
## E. ADVANCED *Frequency() EXAMPLES
## -----
## Note that when dropping the dimensions of the 'tri' array, elements
## in the resulting vector are ordered as if they were obtained with
## 'fast.moving.side="left"':
triL <- trinucleotideFrequency(yeast1, fast.moving.side="left")</pre>
all(as.vector(tri) == triL) # TRUE
## Convert the trinucleotide frequency into the amino acid frequency
## based on translation:
tri1 <- trinucleotideFrequency(yeast1)</pre>
names(tri1) <- GENETIC_CODE[names(tri1)]</pre>
sapply(split(tri1, names(tri1)), sum) # 12512 occurrences of the stop codon
## When the returned vector is very long (e.g. width >= 10), using
```

padAndClip 79

```
## 'with.labels=FALSE' can improve performance significantly.
## Here for example, the observed speed up is between 25x and 500x:
f12 <- oligonucleotideFrequency(yeast1, 12, with.labels=FALSE) # very fast!</pre>
## With the use of 'step', trinucleotideFrequency() is a very fast way to
## calculate the codon usage table in an ORF (or a set of ORFs).
## Taking the same example as in '?codons':
file <- system.file("extdata", "someORF.fa", package="Biostrings")</pre>
my_ORFs <- readDNAStringSet(file)</pre>
## Strip flanking 1000 nucleotides around each ORF and remove first
## sequence as it contains an intron:
my_ORFs <- DNAStringSet(my_ORFs, start=1001, end=-1001)[-1]</pre>
## Codon usage for each ORF:
codon_usage <- trinucleotideFrequency(my_ORFs, step=3)</pre>
## Codon usage across all ORFs:
global_codon_usage <- trinucleotideFrequency(my_ORFs, step=3,</pre>
                                               simplify.as="collapsed")
stopifnot(all(colSums(codon_usage) == global_codon_usage)) # sanity check
## Some related functions:
dict1 <- mkAllStrings(LETTERS[1:3], 4)</pre>
dict2 <- mkAllStrings(LETTERS[1:3], 4, fast.moving.side="left")</pre>
stopifnot(identical(reverse(dict1), dict2))
```

padAndClip

Pad and clip strings

## Description

padAndClip first conceptually pads the supplied strings with an infinite number of padding letters on both sides, then clip them.

stackStrings is a convenience wrapper to padAndClip that turns a variable-width set of strings into a rectangular (i.e. constant-width) set, by padding and clipping the strings, after conceptually shifting them horizontally.

### Usage

## Arguments

x An XStringSet object containing the strings to pad and clip.

80 padAndClip

views

A IntegerRanges object (recycled to the length of x if necessary) defining the region to keep for each string. Because the strings are first conceptually padded with an infinite number of padding letters on both sides, regions can go beyond string limits.

Lpadding.letter, Rpadding.letter

A single letter to use for padding on the left, and another one to use for padding on the right. Note that the default letter ("") does not work if, for example, x is a DNAStringSet object, because the space is not a valid DNA letter (see ?DNA\_ALPHABET). So the Lpadding.letter and Rpadding.letter arguments *must* be supplied if x is not a BStringSet object. For example, if x is a DNAStringSet object, a typical choice is to use "+".

remove.out.of.view.strings

TRUE or FALSE. Whether or not to remove the strings that are out of view in the

returned object.

from, to Another way to specify the region to keep for each string, but with the restriction

that from and to must be single integers. So only 1 region can be specified, and

the same region is used for all the strings.

shift An integer vector (recycled to the length of x if necessary) specifying the amount of shifting (in number of letters) to apply to each string before doing pad and

clip. Positive values shift to the right and negative values to the left.

### Value

For padAndClip: An XStringSet object. If remove.out.of.view.strings is FALSE, it has the same length and names as x, and its "shape", which is described by the integer vector returned by width(), is the same as the shape of the views argument after recycling.

The class of the returned object is the direct concrete subclass of XStringSet that x belongs to or derives from. There are 4 direct concrete subclasses of the XStringSet virtual class: BStringSet, DNAStringSet, RNAStringSet, and AAStringSet. If x is an *instance* of one of those classes, then the returned object has the same class as x (i.e. in that case, padAndClip acts as an endomorphism). But if x *derives* from one of those 4 classes, then the returned object is downgraded to the class x derives from. In that case, padAndClip does not act as an endomorphism.

For stackStrings: Same as padAndClip. In addition it is guaranteed to have a rectangular shape i.e. to be a constant-width XStringSet object.

## Author(s)

H. Pagès

## See Also

- The stackStringsFromBam function in the **GenomicAlignments** package for stacking the read sequences (or their quality strings) stored in a BAM file on a region of interest.
- The XStringViews class to formally represent a set of views on a single string.
- The extractAt and replaceAt functions for extracting/replacing arbitrary substrings from/in
  a string or set of strings.
- The XStringSet class.
- The IntegerRanges class in the IRanges package.

## **Examples**

```
x <- BStringSet(c(seq1="ABCD", seq2="abcdefghijk", seq3="", seq4="XYZ"))
padAndClip(x, IRanges(3, 8:5), Lpadding.letter=">", Rpadding.letter="<")</pre>
padAndClip(x, IRanges(1:-2, 7), Lpadding.letter=">", Rpadding.letter="<")</pre>
stackStrings(x, 2, 8)
stackStrings(x, -2, 8, shift=c(0, -11, 6, 7),
             Lpadding.letter="#", Rpadding.letter=".")
stackStrings(x, -2, 8, shift=c(0, -14, 6, 7),
             Lpadding.letter="#", Rpadding.letter=".")
stackStrings(x, -2, 8, shift=c(0, -14, 6, 7),
             Lpadding.letter="#", Rpadding.letter=".",
             remove.out.of.view.strings=TRUE)
library(hgu95av2probe)
probes <- DNAStringSet(hgu95av2probe)</pre>
probes
stackStrings(probes, 0, 26,
             Lpadding.letter="+", Rpadding.letter="-")
options(showHeadLines=15)
stackStrings(probes, 3, 23, shift=6*c(1:5, -(1:5)),
             Lpadding.letter="+", Rpadding.letter="N",
             remove.out.of.view.strings=TRUE)
```

PDict-class

PDict objects

## **Description**

The PDict class is a container for storing a preprocessed dictionary of DNA patterns that can later be passed to the matchPDict function for fast matching against a reference sequence (the subject).

PDict is the constructor function for creating new PDict objects.

## Usage

```
PDict(x, max.mismatch=NA, tb.start=NA, tb.end=NA, tb.width=NA, algorithm="ACtree2", skip.invalid.patterns=FALSE)
```

#### **Arguments**

x A character vector, a DNAStringSet object or an XStringViews object with a DNAString subject.

max.mismatch A single non-negative integer or NA. See the "Allowing a small number of mis-

matching letters" section below.

tb.start, tb.end, tb.width

A single integer or NA. See the "Trusted Band" section below.

algorithm "ACtree2" (the default) or "Twobit".

skip.invalid.patterns

This argument is not supported yet (and might in fact be replaced by the filter argument very soon).

#### **Details**

## THIS IS STILL WORK IN PROGRESS!

If the original dictionary x is a character vector or an XStringViews object with a DNAString subject, then the PDict constructor will first try to turn it into a DNAStringSet object.

By default (i.e. if PDict is called with max.mismatch=NA, tb.start=NA, tb.end=NA and tb.width=NA) the following limitations apply: (1) the original dictionary can only contain base letters (i.e. only As, Cs, Gs and Ts), therefore IUPAC ambiguity codes are not allowed; (2) all the patterns in the dictionary must have the same length ("constant width" dictionary); and (3) later matchPdict can only be used with max.mismatch=0.

A Trusted Band can be used in order to relax these limitations (see the "Trusted Band" section below).

If you are planning to use the resulting PDict object in order to do inexact matching where valid hits are allowed to have a small number of mismatching letters, then see the "Allowing a small number of mismatching letters" section below.

Two preprocessing algorithms are currently supported: algorithm="ACtree2" (the default) and algorithm="Twobit". With the "ACtree2" algorithm, all the oligonucleotides in the Trusted Band are stored in a 4-ary Aho-Corasick tree. With the "Twobit" algorithm, the 2-bit-per-letter signatures of all the oligonucleotides in the Trusted Band are computed and the mapping from these signatures to the 1-based position of the corresponding oligonucleotide in the Trusted Band is stored in a way that allows very fast lookup. Only PDict objects preprocessed with the "ACtree2" algo can then be used with matchPdict (and family) and with fixed="pattern" (instead of fixed=TRUE, the default), so that IUPAC ambiguity codes in the subject are treated as ambiguities. PDict objects obtained with the "Twobit" algo don't allow this. See ?\matchPDict-inexact\ for more information about support of IUPAC ambiguity codes in the subject.

## **Trusted Band**

What's a Trusted Band?

A Trusted Band is a region defined in the original dictionary where the limitations described above will apply.

Why use a Trusted Band?

Because the limitations described above will apply to the Trusted Band only! For example the Trusted Band cannot contain IUPAC ambiguity codes but the "head" and the "tail" can (see below for what those are). Also with a Trusted Band, if matchPdict is called with a non-null max.mismatch value then mismatching letters will be allowed in the head and the tail. Or, if matchPdict is called

with fixed="subject", then IUPAC ambiguity codes in the head and the tail will be treated as ambiguities.

How to specify a Trusted Band?

Use the tb.start, tb.end and tb.width arguments of the PDict constructor in order to specify a Trusted Band. This will divide each pattern in the original dictionary into three parts: a left part, a middle part and a right part. The middle part is defined by its starting and ending nucleotide positions given relatively to each pattern thru the tb.start, tb.end and tb.width arguments. It must have the same length for all patterns (this common length is called the width of the Trusted Band). The left and right parts are defined implicitely: they are the parts that remain before (prefix) and after (suffix) the middle part, respectively. Therefore three DNAStringSet objects result from this division: the first one is made of all the left parts and forms the head of the PDict object, the second one is made of all the middle parts and forms the Trusted Band of the PDict object, and the third one is made of all the right parts and forms the tail of the PDict object.

In other words you can think of the process of specifying a Trusted Band as drawing 2 vertical lines on the original dictionary (note that these 2 lines are not necessarily straight lines but the horizontal space between them must be constant). When doing this, you are dividing the dictionary into three regions (from left to right): the head, the Trusted Band and the tail. Each of them is a DNAStringSet object with the same number of elements than the original dictionary and the original dictionary could easily be reconstructed from those three regions.

The width of the Trusted Band must be >= 1 because Trusted Bands of width 0 are not supported.

Finally note that calling PDict with tb.start=NA, tb.end=NA and tb.width=NA (the default) is equivalent to calling it with tb.start=1, tb.end=-1 and tb.width=NA, which results in a full-width Trusted Band i.e. a Trusted Band that covers the entire dictionary (no head and no tail).

## Allowing a small number of mismatching letters

[TODO]

## Accessor methods

In the code snippets below, x is a PDict object.

length(x): The number of patterns in x.

width(x): A vector of non-negative integers containing the number of letters for each pattern in x.

names (x): The names of the patterns in x.

head(x): The head of x or NULL if x has no head.

tb(x): The Trusted Band defined on x.

tb.width(x): The width of the Trusted Band defined on x. Note that, unlike width(tb(x)), this is a single integer. And because the Trusted Band has a constant width, tb.width(x) is in fact equivalent to unique(width(tb(x))), or to width(tb(x))[1].

tail(x): The tail of x or NULL if x has no tail.

## Subsetting methods

In the code snippets below, x is a PDict object.

x[[i]]: Extract the i-th pattern from x as a DNAString object.

## Other methods

```
In the code snippet below, x is a PDict object.

duplicated(x): [TODO]

patternFrequency(x): [TODO]
```

### Author(s)

H. Pagès

## References

Aho, Alfred V.; Margaret J. Corasick (June 1975). "Efficient string matching: An aid to bibliographic search". Communications of the ACM 18 (6): 333-340.

#### See Also

matchPDict, DNA\_ALPHABET, IUPAC\_CODE\_MAP, DNAStringSet-class, XStringViews-class

```
## A. NO HEAD AND NO TAIL (THE DEFAULT)
## -----
library(drosophila2probe)
dict0 <- DNAStringSet(drosophila2probe)</pre>
dict0
                              # The original dictionary.
length(dict0)
                              # Hundreds of thousands of patterns.
unique(nchar(dict0))
                             # Patterns are 25-mers.
pdict0 <- PDict(dict0)</pre>
                             # Store the original dictionary in
                              # a PDict object (preprocessing).
pdict0
class(pdict0)
length(pdict0)
                              # Same as length(dict0).
tb.width(pdict0)
                              # The width of the (implicit)
                              # Trusted Band.
sum(duplicated(pdict0))
table(patternFrequency(pdict0)) # 9 patterns are repeated 3 times.
pdict0[[1]]
pdict0[[5]]
## -----
## B. NO HEAD AND A TAIL
## -----
dict1 <- c("ACNG", "GT", "CGT", "AC")</pre>
pdict1 <- PDict(dict1, tb.end=2)</pre>
pdict1
class(pdict1)
length(pdict1)
width(pdict1)
```

```
head(pdict1)
tb(pdict1)
tb.width(pdict1)
width(tb(pdict1))
tail(pdict1)
pdict1[[3]]
```

```
predefined_scoring_matrices
```

Predefined scoring matrices

## **Description**

Predefined scoring matrices for nucleotide and amino acid alignments.

WARNING: All the BLOSUM\* and PAM\* scoring matrices listed below are now located in the **pwalign** package and will soon be removed from the **Biostrings** package.

## Usage

```
data(BLOSUM45)
data(BLOSUM50)
data(BLOSUM62)
data(BLOSUM80)
data(BLOSUM100)
data(PAM30)
data(PAM40)
data(PAM70)
data(PAM120)
data(PAM250)
```

## **Format**

```
See ?pwalign::predefined_scoring_matrices in the pwalign package.
```

## **Details**

```
See ?pwalign::predefined_scoring_matrices in the pwalign package.
```

```
## See ?pwalign::predefined_scoring_matrices in the pwalign package.
```

```
QualityScaledXStringSet-class
```

QualityScaledBStringSet, QualityScaledDNAStringSet, QualityScaledRNAStringSet and QualityScaledAAStringSet objects

## **Description**

The QualityScaledBStringSet class is a container for storing a BStringSet object with an XStringQuality object.

Similarly, the QualityScaledDNAStringSet (or QualityScaledRNAStringSet, or QualityScaledAAStringSet) class is a container for storing a DNAStringSet (or RNAStringSet, or AAStringSet) objects with an XStringQuality object.

### Usage

### **Arguments**

Х

For the QualityScaled\*StringSet constructors: Either a character vector, or an XString, XStringSet or XStringViews object.

 $For write {\tt QualityScaledXStringSet:} \ A \ QualityScaledDNAStringSet \ object \ or \ other \ QualityScaledXStringSet \ derivative.$ 

quality An XStringQuality derivative.

filepath, nrec, skip, seek.first.rec, use.names, append, compress,
compression\_level

See ?`XStringSet-io`.

quality.scoring

Specify the quality scoring used in the FASTQ file. Must be one of "phred" (the default), "solexa", or "illumina". If set to "phred" (or "solexa" or "illumina"), the qualities will be stored in a PhredQuality (or SolexaQuality or IlluminaQuality, respectively) object.

### **Details**

The QualityScaledBStringSet, QualityScaledDNAStringSet, QualityScaledRNAStringSet and QualityScaledAAStringSet functions are constructors that can be used to "naturally" turn x into an QualityScaledXStringSet object of the desired base type.

### Accessor methods

The QualityScaledXStringSet class derives from the XStringSet class hence all the accessor methods defined for an XStringSet object can also be used on an QualityScaledXStringSet object. Common methods include (in the code snippets below, x is an QualityScaledXStringSet object):

```
length(x): The number of sequences in x.
```

width(x): A vector of non-negative integers containing the number of letters for each element in x.

```
nchar(x): The same as width(x).
```

names(x): NULL or a character vector of the same length as x containing a short user-provided description or comment for each element in x.

quality(x): The quality of the strings.

## Subsetting and appending

In the code snippets below, x and values are XStringSet objects, and i should be an index specifying the elements to extract.

x[i]: Return a new QualityScaledXStringSet object made of the selected elements.

## Author(s)

P. Aboyoun

## See Also

- BStringSet, DNAStringSet, RNAStringSet, and AAStringSet objects.
- XStringQuality objects.
- readDNAStringSet and writeXStringSet for reading/writing a DNAStringSet object (or other XStringSet derivative) from/to a FASTA or FASTQ file.

```
## ------
## QualityScaled*StringSet() CONSTRUCTORS
## ------
x1 <- DNAStringSet(c("TTGA", "CTCN"))
q1 <- PhredQuality(c("*+,-", "6789"))
qdna1 <- QualityScaledDNAStringSet(x1, q1)
qdna1</pre>
```

```
## READ/WRITE A QualityScaledDNAStringSet OBJECT FROM/TO A FASTQ FILE
filepath <- system.file("extdata", "s_1_sequence.txt",</pre>
                        package="Biostrings")
## By default, readQualityScaledDNAStringSet() assumes that the FASTQ
## file contains "Phred quality scores" (this is the standard Sanger
## variant to assess reliability of a base call):
qdna2 <- readQualityScaledDNAStringSet(filepath)</pre>
qdna2
outfile2a <- tempfile()</pre>
writeQualityScaledXStringSet(qdna2, outfile2a)
outfile2b <- tempfile()</pre>
writeQualityScaledXStringSet(qdna2, outfile2b, compress=TRUE)
## Use 'quality.scoring="solexa"' or 'quality.scoring="illumina"' if the
## quality scores are Solexa quality scores:
qdna3 <- readQualityScaledDNAStringSet(filepath, quality.scoring="solexa")</pre>
qdna3
outfile3a <- tempfile()</pre>
writeQualityScaledXStringSet(qdna3, outfile3a)
outfile3b <- tempfile()</pre>
writeQualityScaledXStringSet(qdna3, outfile3b, compress=TRUE)
## Sanity checks:
stopifnot(identical(readLines(outfile2a), readLines(filepath)))
stopifnot(identical(readLines(outfile2a), readLines(outfile2b)))
stopifnot(identical(readLines(outfile3a), readLines(filepath)))
stopifnot(identical(readLines(outfile3a), readLines(outfile3b)))
```

replaceAt

Extract/replace arbitrary substrings from/in a string or set of strings.

## **Description**

extractAt extracts multiple subsequences from XString object x, or from the individual sequences of XStringSet object x, at the ranges of positions specified thru at.

replaceAt performs multiple subsequence replacements (a.k.a. substitutions) in XString object x, or in the individual sequences of XStringSet object x, at the ranges of positions specified thru at.

#### Usage

```
extractAt(x, at)
replaceAt(x, at, value="")
```

### **Arguments**

at

An XString or XStringSet object. Х

> Typically a IntegerRanges object if x is an XString object, and an IntegerRanges-List object if x is an XStringSet object.

Alternatively, the ranges can be specified with only 1 number per range (its start position), in which case they are considered to be empty ranges (a.k.a. zerowidth ranges). So if at is a numeric vector, an IntegerList object, or a list of numeric vectors, each number in it is interpreted as the start position of a zerowidth range. This is useful when using replaceAt to perform insertions.

The following applies only if x is an XStringSet object:

at is recycled to the length of x if necessary. If at is a IntegerRanges object (or a numeric vector), it is first turned into a IntegerRangesList object of length 1 and then this IntegerRangesList object is recycled to the length of x. This is useful for specifying the same ranges across all sequences in x. The effective shape of at is described by its length together with the lengths of its list elements after recycling.

As a special case, extractAt accepts at and value to be both of length 0, in which case it just returns x unmodified (no-op).

The replacement sequences.

If x is an XString object, value is typically a character vector or an XStringSet object that is recycled to the length of at (if necessary).

If x is an XStringSet object, value is typically a list of character vectors or a CharacterList or XStringSetList object. If necessary, it is recycled "vertically" first and then "horizontally" to bring it into the *effective shape* of at (see above). "Vertical recycling" is the usual recycling whereas "horizontal recycling" recycles the individual list elements.

As a special case, extractAt accepts at and value to be both of length 0, in which case it just returns x unmodified (no-op).

### Value

For extractAt: An XStringSet object of the same length as at if x is an XString object. An XStringSetList object of the same length as x (and same effective shape as at) if x is an XStringSet object.

For replaceAt: An object of the same class as x. If x is an XStringSet object, its length and names and metadata columns are preserved.

### Note

Like subseq (defined and documented in the XVector package), extractAt does not copy the sequence data!

extractAt is equivalent to extractList (defined and documented in the IRanges package) when x is an XString object and at a IntegerRanges object.

## Author(s)

H. Pagès

value

### See Also

The subseq and subseq<- functions in the XVector package for simpler forms of subsequence extractions and replacements.</li>

- The extractList and unstrsplit functions defined and documented in the IRanges package.
- The replaceLetterAt function for a DNA-specific single-letter replacement functions useful for SNP injections.
- The padAndClip function for padding and clipping strings.
- The XString, XStringSet, and XStringSetList classes.
- The IntegerRanges, IntegerRangesList, IntegerList, and CharacterList classes defined and documented in the **IRanges** package.

```
## (A) ON AN XString OBJECT
## -----
x <- BString("abcdefghijklm")</pre>
at1 <- IRanges(5:1, width=3)
extractAt(x, at1)
names(at1) <- LETTERS[22:26]</pre>
extractAt(x, at1)
at2 <- IRanges(c(1, 5, 12), c(3, 4, 12), names=c("X", "Y", "Z"))
extractAt(x, at2)
extractAt(x, rev(at2))
value <- c("+", "-", "*")
replaceAt(x, at2, value=value)
replaceAt(x, rev(at2), value=rev(value))
at3 <- IRanges(c(14, 1, 1, 1, 1, 11), c(13, 0, 10, 0, 0, 10))
value <- 1:6
replaceAt(x, at3, value=value)
                                       # "24536klm1"
replaceAt(x, rev(at3), value=rev(value)) # "54236klm1"
## Deletions:
stopifnot(replaceAt(x, at2) == "defghijkm")
stopifnot(replaceAt(x, rev(at2)) == "defghijkm")
stopifnot(replaceAt(x, at3) == "klm")
stopifnot(replaceAt(x, rev(at3)) == "klm")
## Insertions:
at4 <- IRanges(c(6, 10, 2, 5), width=0)
stopifnot(replaceAt(x, at4, value="-") == "a-bcd-e-fghi-jklm")
stopifnot(replaceAt(x, start(at4), value="-") == "a-bcd-e-fghi-jklm")
at5 <- c(5, 1, 6, 5) # 2 insertions before position 5
replaceAt(x, at5, value=c("+", "-", "*", "/"))
```

```
## No-ops:
stopifnot(replaceAt(x, NULL, value=NULL) == x)
stopifnot(replaceAt(x, at2, value=extractAt(x, at2)) == x)
stopifnot(replaceAt(x, at3, value=extractAt(x, at3)) == x)
stopifnot(replaceAt(x, at4, value=extractAt(x, at4)) == x)
stopifnot(replaceAt(x, at5, value=extractAt(x, at5)) == x)
## The order of successive transformations matters:
    T1: insert "+" before position 1 and 4
    T2: insert "-" before position 3
## T1 followed by T2
x2a \leftarrow replaceAt(x, c(1, 4), value="+")
x3a <- replaceAt(x2a, 3, value="-")
## T2 followed by T1
x2b <- replaceAt(x, 3, value="-")</pre>
x3b \leftarrow replaceAt(x2b, c(1, 4), value="+")
## T1 and T2 simultaneously:
x3c \leftarrow replaceAt(x, c(1, 3, 4), value=c("+", "-", "+"))
## ==> 'x3a', 'x3b', and 'x3c' are all different!
## Append "**" to 'x3c':
replaceAt(x3c, length(x3c) + 1L, value="**")
## -----
## (B) ON AN XStringSet OBJECT
x <- BStringSet(c(seq1="ABCD", seq2="abcdefghijk", seq3="XYZ"))</pre>
at6 <- IRanges(c(1, 3), width=1)
extractAt(x, at=at6)
unstrsplit(extractAt(x, at=at6))
at7 <- IRangesList(IRanges(c(2, 1), c(3, 0)),
                   IRanges(c(7, 2, 12, 7), c(6, 5, 11, 8)),
                   IRanges(2, 2))
## Set inner names on 'at7'.
unlisted_at7 <- unlist(at7)</pre>
names(unlisted_at7) <-</pre>
    paste0("rg", sprintf("%02d", seq_along(unlisted_at7)))
at7 <- relist(unlisted_at7, at7)</pre>
extractAt(x, at7) # same as 'as(mapply(extractAt, x, at7), "List")'
extractAt(x, at7[3]) # same as 'as(mapply(extractAt, x, at7[3]), "List")'
replaceAt(x, at7, value=extractAt(x, at7)) # no-op
replaceAt(x, at7) # deletions
at8 <- IRangesList(IRanges(1:5, width=0),
                   IRanges(c(6, 8, 10, 7, 2, 5),
```

92 replaceLetterAt

```
width=c(0, 2, 0, 0, 0, 0),
                  IRanges(c(1, 2, 1), width=c(0, 1, 0)))
replaceAt(x, at8, value="-")
value8 <- relist(paste0("[", seq_along(unlist(at8)), "]"), at8)</pre>
replaceAt(x, at8, value=value8)
replaceAt(x, at8, value=as(c("+", "-", "*"), "List"))
## Append "**" to all sequences:
replaceAt(x, as(width(x) + 1L, "List"), value="**")
## (C) ADVANCED EXAMPLES
## -----
library(hgu95av2probe)
probes <- DNAStringSet(hgu95av2probe)</pre>
## Split the probes in 5-mer chunks:
at <- successiveIRanges(rep(5, 5))</pre>
extractAt(probes, at)
## Replace base 13 by its complement:
at <- IRanges(13, width=1)
base13 <- extractAt(probes, at)</pre>
base13comp <- relist(complement(unlist(base13)), base13)</pre>
replaceAt(probes, at, value=base13comp)
## See ?xscat for a more efficient way to do this.
## Replace all the occurences of a given pattern with another pattern:
midx <- vmatchPattern("VCGTT", probes, fixed=FALSE)</pre>
matches <- extractAt(probes, midx)</pre>
unlist(matches)
unique(unlist(matches))
probes2 <- replaceAt(probes, midx, value="-++-")</pre>
## See strings with 2 or more susbtitutions:
probes2[elementNROWS(midx) >= 2]
## 2 sanity checks:
stopifnot(all(replaceAt(probes, midx, value=matches) == probes))
probes2b <- gsub("[ACG]CGTT", "-++-", as.character(probes))</pre>
stopifnot(identical(as.character(probes2), probes2b))
```

replaceLetterAt

Replacing letters in a sequence (or set of sequences) at some specified locations

### **Description**

replaceLetterAt first makes a copy of a sequence (or set of sequences) and then replaces some of the original letters by new letters at the specified locations.

replaceLetterAt 93

.inplaceReplaceLetterAt is the IN PLACE version of replaceLetterAt: it will modify the original sequence in place i.e. without copying it first. Note that in place modification of a sequence is fundamentally dangerous because it alters all objects defined in your session that make reference to the modified sequence. NEVER use .inplaceReplaceLetterAt, unless you know what you are doing!

### Usage

```
replaceLetterAt(x, at, letter, if.not.extending="replace", verbose=FALSE)
## NEVER USE THIS FUNCTION!
.inplaceReplaceLetterAt(x, at, letter)
```

#### **Arguments**

x A DNAString or rectangular DNAStringSet object.

at The locations where the replacements must occur.

If x is a DNAString object, then at is typically an integer vector with no NAs but a logical vector or Rle object is valid too. Locations can be repeated and in this case the last replacement to occur at a given location prevails.

If x is a rectangular DNAStringSet object, then at must be a matrix of logicals with the same dimensions as x.

letter The new letters.

If x is a DNAString object, then letter must be a DNAString object or a character vector (with no NAs) with a total number of letters (sum(nchar(letter))) equal to the number of locations specified in at.

If x is a rectangular DNAStringSet object, then letter must be a DNAStringSet object or a character vector of the same length as x. In addition, the number of letters in each element of letter must match the number of locations specified in the corresponding row of at (all(width(letter) == rowSums(at))).

if.not.extending

What to do if the new letter is not "extending" the old letter? The new letter "extends" the old letter if both are IUPAC letters and the new letter is as specific or less specific than the old one (e.g. M extends A, Y extends Y, but Y doesn't extend S). Possible values are "replace" (the default) for replacing in all cases, "skip" for not replacing when the new letter does not extend the old letter, "merge" for merging the new IUPAC letter with the old one, and "error" for raising an error.

Note that the gap ("-") and hard masking ("+") letters are not extending or extended by any other letter.

Also note that "merge" is the only value for the if.not.extending argument that guarantees the final result to be independent on the order the replacement is performed (although this is only relevant when at contains duplicated locations, otherwise the result is of course always independent on the order, whatever the value of if.not.extending is).

verbose When TRUE, a warning will report the number of skipped or merged letters.

94 replaceLetterAt

### **Details**

.inplaceReplaceLetterAt semantic is equivalent to calling replaceLetterAt with if.not.extending="merge" and verbose=FALSE.

Never use .inplaceReplaceLetterAt! It is used by the injectSNPs function in the BSgenome package, as part of the "lazy sequence loading" mechanism, for altering the original sequences of a BSgenome object at "sequence-load time". This alteration consists in injecting the IUPAC ambiguity letters representing the SNPs into the just loaded sequence, which is the only time where in place modification of the external data of an XString object is safe.

### Value

A DNAString or DNAStringSet object of the same shape (i.e. length and width) as the original object x for replaceLetterAt.

### Author(s)

H. Pagès

### See Also

- The replaceAt function for extracting or replacing arbitrary subsequences from/in a sequence or set of sequences.
- IUPAC\_CODE\_MAP for the mapping between IUPAC nucleotide ambiguity codes and their meaning.
- The chartr and injectHardMask functions.
- The DNAString and DNAStringSet class.
- The injectSNPs function and the BSgenome class in the BSgenome package.

reverseComplement 95

reverseComplement	Sequence reversing and complementing	

## **Description**

Use these functions for reversing sequences and/or complementing DNA or RNA sequences.

## Usage

```
complement(x, ...)
reverseComplement(x, ...)
```

## **Arguments**

x A DNAString, RNAString, DNAStringSet, RNAStringSet, XStringViews (with DNAString or RNAString subject), MaskedDNAString or MaskedRNAString object for complement and reverseComplement.

. . . Additional arguments to be passed to or from methods.

### **Details**

See ?reverse for reversing an XString, XStringSet or XStringViews object.

If x is a DNAString or RNAString object, complement(x) returns an object where each base in x is "complemented" i.e. A, C, G, T in a DNAString object are replaced by T, G, C, A respectively and A, C, G, U in a RNAString object are replaced by U, G, C, A respectively.

Letters belonging to the IUPAC Extended Genetic Alphabet are also replaced by their complement (M <-> K, R <-> Y, S <-> S, V <-> B, W <-> W, H <-> D, N <-> N) and the gap ("-") and hard masking ("+") letters are unchanged.

reverseComplement(x) is equivalent to reverse(complement(x)) but is faster and more memory efficient.

## Value

An object of the same class and length as the original object.

### See Also

reverse, DNAString-class, RNAString-class, DNAStringSet-class, RNAStringSet-class, XStringViews-class, MaskedXString-class, chartr, findPalindromes, IUPAC\_CODE\_MAP

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

x <- DNAString("ACGT-YN-")
```

96 reverseComplement

```
reverseComplement(x)
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
probes
alphabetFrequency(probes, collapse=TRUE)
rcprobes <- reverseComplement(probes)</pre>
alphabetFrequency(rcprobes, collapse=TRUE)
## B. OBTAINING THE MISMATCH PROBES OF A CHIP
pm2mm <- function(probes)</pre>
{
   probes <- DNAStringSet(probes)</pre>
   subseq(probes, start=13, end=13) <- complement(subseq(probes, start=13, end=13))</pre>
}
mmprobes <- pm2mm(probes)</pre>
mmprobes
alphabetFrequency(mmprobes, collapse=TRUE)
## -----
## C. SEARCHING THE MINUS STRAND OF A CHROMOSOME
## -----
## Applying reverseComplement() to the pattern before calling
## matchPattern() is the recommended way of searching hits on the
## minus strand of a chromosome.
library(BSgenome.Dmelanogaster.UCSC.dm3)
chrX <- Dmelanogaster$chrX</pre>
pattern <- DNAString("ACCAACNNGGTTG")</pre>
matchPattern(pattern, chrX, fixed=FALSE) # 3 hits on strand +
rcpattern <- reverseComplement(pattern)</pre>
m0 <- matchPattern(rcpattern, chrX, fixed=FALSE)</pre>
m0 # 5 hits on strand -
## Applying reverseComplement() to the subject instead of the pattern is not
## a good idea for 2 reasons:
## (1) Chromosome sequences are generally big and sometimes very big
##
      so computing the reverse complement of the positive strand will
##
      take time and memory proportional to its length.
chrXminus <- reverseComplement(chrX) # needs to allocate 22M of memory!</pre>
## (2) Chromosome locations are generally given relatively to the positive
##
      strand, even for features located in the negative strand, so after
##
      doing this:
m1 <- matchPattern(pattern, chrXminus, fixed=FALSE)</pre>
      the start/end of the matches are now relative to the negative strand.
##
##
      You need to apply reverseComplement() again on the result if you want
```

RNAString-class 97

```
## them to be relative to the positive strand:
m2 <- reverseComplement(m1)  # allocates 22M of memory, again!
## and finally to apply rev() to sort the matches from left to right
## (5'3' direction) like in m0:
m3 <- rev(m2)  # same as m0, finally!

## WARNING: Before you try the example below on human chromosome 1, be aware
## that it will require the allocation of about 500Mb of memory!
if (interactive()) {
   library(BSgenome.Hsapiens.UCSC.hg18)
   chr1 <- Hsapiens$chr1
   matchPattern(pattern, reverseComplement(chr1))  # DON'T DO THIS!
   matchPattern(reverseComplement(pattern), chr1)  # DO THIS INSTEAD
}</pre>
```

RNAString-class

RNAString objects

## **Description**

An RNAString object allows efficient storage and manipulation of a long RNA sequence.

## **Details**

The RNAString class is a direct XString subclass (with no additional slot). Therefore all functions and methods described in the XString man page also work with an RNAString object (inheritance).

Unlike the BString container that allows storage of any single string (based on a single-byte character set) the RNAString container can only store a string based on the RNA alphabet (see below). In addition, the letters stored in an RNAString object are encoded in a way that optimizes fast search algorithms.

## The RNA alphabet

This alphabet is the same as the DNA alphabet, except that "T" is replaced by "U". See ?DNA\_ALPHABET for more information about the DNA alphabet. The RNA alphabet is stored in the RNA\_ALPHABET predefined constant (character vector).

The alphabet() function returns RNA\_ALPHABET when applied to an RNAString object.

## Constructor-like functions and generics

In the code snippet below, x can be a single string (character vector of length 1), a BString object or a DNAString object.

RNAString(x="", start=1, nchar=NA): Tries to convert x into an RNAString object by reading nchar letters starting at position start in x.

98 RNAString-class

## Accessor methods

In the code snippet below, x is an RNAString object.

alphabet(x, baseOnly=FALSE): If x is an RNAString object, then return the RNA alphabet (see above). See the corresponding man pages when x is a BString, DNAString or AAString object.

## **Display**

The letters in an RNAString object are colored when displayed by the show() method. Set global option Biostrings.coloring to FALSE to turn off this coloring.

## Author(s)

H. Pagès

#### See Also

- The RNAStringSet class to represent a collection of RNAString objects.
- The XString and DNAString classes.
- reverseComplement
- alphabetFrequency
- IUPAC\_CODE\_MAP
- letter

```
RNA_BASES
RNA_ALPHABET
dna <- DNAString("TTGAAAA-CTC-N")
rna <- RNAString(dna)
rna # 'options(Biostrings.coloring=FALSE)' to turn off coloring
alphabet(rna) # RNA_ALPHABET
alphabet(rna, baseOnly=TRUE) # RNA_BASES
## When comparing an RNAString object with a DNAString object,
## U and T are considered equals:
rna == dna # TRUE</pre>
```

seqinfo-methods 99

seqinfo-methods

seqinfo() method for DNAStringSet objects

## **Description**

seqinfo methods for extracting the sequence information stored in a DNAStringSet object.

## Usage

```
## S4 method for signature 'DNAStringSet'
seqinfo(x)
```

## **Arguments**

Х

A DNAStringSet object.

## Value

A Seqinfo object for the 'seqinfo' getter.

A DNAStringSet object containing sequence information for the 'seqinfo' setter.

# See Also

```
getSeq, DNAStringSet-class,
```

100 toComplex

toComplex

Turning a DNA sequence into a vector of complex numbers

# Description

The toComplex utility function turns a DNAString object into a complex vector.

# Usage

```
toComplex(x, baseValues)
```

# Arguments

```
x A DNAString object.
```

baseValues A named complex vector containing the values associated to each base e.g.

c(A=1+0i, G=0+1i, T=-1+0i, C=0-1i)

### Value

A complex vector of the same length as x.

# Author(s)

H. Pagès

## See Also

**DNAString** 

```
seq <- DNAString("accacctgaccattgtcct")
baseValues1 <- c(A=1+0i, G=0+1i, T=-1+0i, C=0-1i)
toComplex(seq, baseValues1)

## GC content:
baseValues2 <- c(A=0, C=1, G=1, T=0)
sum(as.integer(toComplex(seq, baseValues2)))
## Note that there are better ways to do this (see ?alphabetFrequency)</pre>
```

101 translate

translate

Translating DNA/RNA sequences

### **Description**

Functions for translating DNA or RNA sequences into amino acid sequences.

## Usage

```
## Translating DNA/RNA:
translate(x, genetic.code=GENETIC_CODE, no.init.codon=FALSE,
             if.fuzzy.codon="error")
## Extracting codons without translating them:
codons(x)
```

## **Arguments**

Х

A DNAStringSet, RNAStringSet, DNAString, RNAString, MaskedDNAString or MaskedRNAString object for translate.

A DNAString, RNAString, MaskedDNAString or MaskedRNAString object for codons.

genetic.code

The genetic code to use for the translation of codons into Amino Acid letters. It must be represented as a named character vector of length 64 similar to predefined constant GENETIC\_CODE. More precisely:

- it must contain 1-letter strings in the Amino Acid alphabet;
- its names must be identical to names (GENETIC\_CODE);
- it must have an alt\_init\_codons attribute on it, that lists the *alternative* initiation codons.

The default value for genetic.code is GENETIC\_CODE, which represents The Standard Genetic Code. See ?AA\_ALPHABET for the Amino Acid alphabet, and ?GENETIC\_CODE for The Standard Genetic Code and its known variants.

no.init.codon

By default, translate() assumes that the first codon in a DNA or RNA sequence is the initiation codon. This means that the alt\_init\_codons attribute on the supplied genetic.code will be used to translate the *alternative initiation* codons. This can be changed by setting no.init.codon to TRUE, in which case the alt\_init\_codons attribute will be ignored.

if.fuzzy.codon How fuzzy codons (i.e codon with IUPAC ambiguities) should be handled. Accepted values are:

- "error": An error will be raised on the first occurence of a fuzzy codon. This is the default.
- "solve": Fuzzy codons that can be translated non ambiguously to an amino acid or to \* (stop codon) will be translated. Ambiguous fuzzy codons will be translated to X.

102 translate

- "error.if.X": Fuzzy codons that can be translated non ambiguously to an amino acid or to \* (stop codon) will be translated. An error will be raised on the first occurence of an ambiguous fuzzy codon.
- "X": All fuzzy codons (ambiguous and non-ambiguous) will be translated to X.

Alternatively if.fuzzy.codon can be specified as a character vector of length 2 for more fine-grained control. The 1st string and 2nd strings specify how to handle non-ambiguous and ambiguous fuzzy codons, respectively. The accepted values for the 1st string are:

- "error": Any occurence of a non-ambiguous fuzzy codon will cause an error
- "solve": Non-ambiguous fuzzy codons will be translated to an amino acid or to \*.
- "X": Non-ambiguous fuzzy codons will be translated to X.

The accepted values for the 2nd string are:

- "error": Any occurence of an ambiguous fuzzy codon will cause an error.
- "X": Ambiguous fuzzy codons will be translated to X.

All the 6 possible combinations of 1st and 2nd strings are supported. Note that if.fuzzy.codon=c("error", "error") is equivalent to if.fuzzy.codon="error", if.fuzzy.codon=c("solve", "X") is equivalent to if.fuzzy.codon="solve", if.fuzzy.codon=c("solve", "error") is equivalent to if.fuzzy.codon="error.if.X", and if.fuzzy.codon=c("X", "X") is equivalent to if.fuzzy.codon="X".

## **Details**

translate reproduces the biological process of RNA translation that occurs in the cell. The input of the function can be either RNA or coding DNA. By default The Standard Genetic Code (see ?GENETIC\_CODE) is used to translate codons into amino acids but the user can supply a different genetic code via the genetic.code argument.

codons is a utility for extracting the codons involved in this translation without translating them.

### Value

For translate: An AAString object when x is a DNAString, RNAString, MaskedDNAString, or MaskedRNAString object. An AAStringSet object *parallel* to x (i.e. with 1 amino acid sequence per DNA or RNA sequence in x) when x is a DNAStringSet or RNAStringSet object. If x has names on it, they're propagated to the returned object.

For codons: An XStringViews object with 1 view per codon. When x is a MaskedDNAString or MaskedRNAString object, its masked parts are interpreted as introns and filled with the + letter in the returned object. Therefore codons that span across masked regions are represented by views that have a width > 3 and contain the + letter. Note that each view is guaranteed to contain exactly 3 base letters.

## See Also

• AA\_ALPHABET for the Amino Acid alphabet.

translate 103

- GENETIC\_CODE for The Standard Genetic Code and its known variants.
- The examples for extractTranscriptSeqs in the **GenomicFeatures** package for computing the full proteome of a given organism.
- The reverseComplement function.
- The DNAStringSet and AAStringSet classes.
- The XStringViews and MaskedXString classes.

```
## -----
## 1. BASIC EXAMPLES
## -----
dna1 <- DNAString("TTGATATGGCCCTTATAA")</pre>
translate(dna1)
## TTG is an alternative initiation codon in the Standard Genetic Code:
translate(dna1, no.init.codon=TRUE)
SGC1 <- getGeneticCode("SGC1") # Vertebrate Mitochondrial code
translate(dna1, genetic.code=SGC1)
## TTG is NOT an alternative initiation codon in the Vertebrate
## Mitochondrial code:
translate(dna1, genetic.code=SGC1, no.init.codon=TRUE)
## All 6 codons except 4th (CCC) are fuzzy:
dna2 <- DNAString("HTGATHTGRCCCYTRTRA")</pre>
## Not run:
 translate(dna2) # error because of fuzzy codons
## End(Not run)
## Translate all fuzzy codons to X:
translate(dna2, if.fuzzy.codon="X")
## Or solve the non-ambiguous ones (3rd codon is ambiguous so cannot be
## solved):
translate(dna2, if.fuzzy.codon="solve")
## Fuzzy codons that are non-ambiguous with a given genetic code can
## become ambiguous with another genetic code, and vice versa:
translate(dna2, genetic.code=SGC1, if.fuzzy.codon="solve")
## -----
## 2. TRANSLATING AN OPEN READING FRAME
## -----
file <- system.file("extdata", "someORF.fa", package="Biostrings")</pre>
x <- readDNAStringSet(file)</pre>
```

104 trimLRPatterns

```
## The first and last 1000 nucleotides are not part of the ORFs:
x <- DNAStringSet(x, start=1001, end=-1001)</pre>
## Before calling translate() on an ORF, we need to mask the introns
## if any. We can get this information from the SGD database
## (http://www.yeastgenome.org/).
## According to SGD, the 1st ORF (YAL001C) has an intron at 71..160
## (see http://db.yeastgenome.org/cgi-bin/locus.pl?locus=YAL001C)
y1 <- x[[1]]
mask1 <- Mask(length(y1), start=71, end=160)</pre>
masks(y1) <- mask1
у1
translate(y1)
## Codons:
codons(y1)
which(width(codons(y1)) != 3)
codons(y1)[20:28]
## -----
## 3. AN ADVANCED EXAMPLE
## Translation on the '-' strand:
dna3 <- DNAStringSet(c("ATC", "GCTG", "CGACT"))</pre>
translate(reverseComplement(dna3))
## Translate sequences on both '+' and '-' strand across all
## possible reading frames (i.e., codon position 1, 2 or 3):
## First create a DNAStringSet of '+' and '-' strand sequences,
## removing the nucleotides prior to the reading frame start position.
dna3_subseqs <- lapply(1:3, function(pos)</pre>
    subseq(c(dna3, reverseComplement(dna3)), start=pos))
## Translation of 'dna3_subseqs' produces a list of length 3, each with
## 6 elements (3 '+' strand results followed by 3 '-' strand results).
lapply(dna3_subseqs, translate)
## Note that translate() throws a warning when the length of the sequence
## is not divisible by 3. To avoid this warning wrap the function in
## suppressWarnings().
```

trimLRPatterns

Trim Flanking Patterns from Sequences

## **Description**

The trimLRPatterns function trims left and/or right flanking patterns from sequences.

trimLRPatterns 105

## Usage

## **Arguments**

Lpattern The left pattern.

Rpattern The right pattern.

subject An XString object, XStringSet object, or character vector containing the target

sequence(s).

max.Lmismatch Either an integer vector of length nLp = nchar(Lpattern) representing an ab-

solute number of mismatches (or edit distance if with.Lindels is TRUE) or a single numeric value in the interval [0, 1) representing a mismatch rate when aligning terminal substrings (suffixes) of Lpattern with the beginning (prefix) of subject following the conventions set by nedit Starting At its Matching Starting

 $of \ subject following \ the \ conventions \ set \ by \ nedit Starting At, \ is \texttt{MatchingStartingAt},$ 

etc.

When max.Lmismatch is 0L or a numeric value in the interval [0, 1), it is taken as a "rate" and is converted to as.integer(1:nLp \* max.Lmismatch), analogous to agrep (which, however, employs ceiling).

Otherwise, max.Lmismatch is treated as an integer vector where negative numbers are used to prevent trimming at the i-th location. When an input integer vector is shorter than nLp, it is augmented with enough -1s at the beginning to bring its length up to nLp. Elements of max.Lmismatch beyond the first nLp are ignored.

Once the integer vector is constructed using the rules given above, when with.Lindels is FALSE, max.Lmismatch[i] is the number of acceptable mismatches (errors) between the suffix substring(Lpattern, nLp - i + 1, nLp) of Lpattern and the first i letters of subject. When with.Lindels is TRUE, max.Lmismatch[i] represents the allowed "edit distance" between that suffix of Lpattern and subject, starting at position 1 of subject (as in matchPattern and isMatchingStartingAt).

For a given element s of the subject, the initial segment (prefix) substring(s, 1, j) of s is trimmed if j is the largest i for which there is an acceptable match, if any.

max.Rmismatch Same as max.Lmismatch but with Rpattern, along with with.Rindels (below), and its initial segments (prefixes) substring(Rpattern, 1, i).

For a given element s of the subject, with nS = nchar(s), the terminal segment (suffix) substring(s, nS - j + 1, nS) of s is trimmed if j is the largest i for

which there is an acceptable match, if any.

with.Lindels If TRUE, indels are allowed in the alignments of the suffixes of Lpattern with the

subject, at its beginning. See the with.indels arguments of the matchPattern

and neditStartingAt functions for detailed information.

Same as with.Lindels but for alignments of the prefixes of Rpattern with the subject, at its end. See the with.indels arguments of the  ${\tt matchPattern}$  and

neditEndingAt functions for detailed information.

. . . . .

with.Rindels

106 trimLRPatterns

Lfixed, Rfixed Whether IUPAC extended letters in the left or right pattern should be interpreted

as ambiguities (see ?`lowlevel-matching` for the details).

ranges If TRUE, then return the ranges to use to trim subject. If FALSE, then returned

the trimmed subject.

### Value

A new XString object, XStringSet object, or character vector with the "longest" flanking matches removed, as described above.

## Author(s)

P. Aboyoun and H. Jaffee

### See Also

matchPattern, matchLRPatterns, lowlevel-matching, XString-class, XStringSet-class

```
Lpattern <- "TTCTGCTTG"</pre>
Rpattern <- "GATCGGAAG"</pre>
subject <- DNAString("TTCTGCTTGACGTGATCGGA")</pre>
subjectSet <- DNAStringSet(c("TGCTTGACGGCAGATCGG", "TTCTGCTTGGATCGGAAG"))</pre>
## Only allow for perfect matches on the flanks
trimLRPatterns(Lpattern = Lpattern, subject = subject)
trimLRPatterns(Rpattern = Rpattern, subject = subject)
trimLRPatterns(Lpattern = Lpattern, Rpattern = Rpattern, subject = subjectSet)
## Allow for perfect matches on the flanking overlaps
trimLRPatterns(Lpattern = Lpattern, Rpattern = Rpattern, subject = subjectSet,
               max.Lmismatch = 0, max.Rmismatch = 0)
## Allow for mismatches on the flanks
trimLRPatterns(Lpattern = Lpattern, Rpattern = Rpattern, subject = subject,
               max.Lmismatch = 0.2, max.Rmismatch = 0.2)
maxMismatches <- as.integer(0.2 * 1:9)</pre>
trimLRPatterns(Lpattern = Lpattern, Rpattern = Rpattern, subject = subjectSet,
               max.Lmismatch = maxMismatches, max.Rmismatch = maxMismatches)
## Produce ranges that can be an input into other functions
trimLRPatterns(Lpattern = Lpattern, Rpattern = Rpattern, subject = subjectSet,
               max.Lmismatch = 0, max.Rmismatch = 0, ranges = TRUE)
trimLRPatterns(Lpattern = Lpattern, Rpattern = Rpattern, subject = subject,
               max.Lmismatch = 0.2, max.Rmismatch = 0.2, ranges = TRUE)
```

xscat 107

xscat Concatenate sequences contained in XString, XStringSet and/or XStringViews objects

## **Description**

This function mimics the semantic of paste(..., sep="") but accepts XString, XStringSet or XStringViews arguments and returns an XString or XStringSet object.

## Usage

```
xscat(...)
```

## **Arguments**

One or more character vectors (with no NAs), XString, XStringSet or XStringViews objects.

## Value

An XString object if all the arguments are either XString objects or character strings. An XStringSet object otherwise.

### Author(s)

H. Pagès

#### See Also

XString-class, XStringSet-class, XStringViews-class, paste

```
## Return a BString object:
xscat(BString("abc"), BString("EF"))
xscat(BString("abc"), "EF")
xscat("abc", "EF")

## Return a BStringSet object:
xscat(BStringSet("abc"), "EF")

## Return a DNAStringSet object:
xscat(c("t", "a"), DNAString("N"))

## Arguments are recycled to the length of the longest argument:
res1a <- xscat("x", LETTERS, c("3", "44", "555"))
res1b <- paste0("x", LETTERS, c("3", "44", "555"))
stopifnot(identical(as.character(res1a), as.character(res1b)))</pre>
```

108 XString-class

```
## Concatenating big XStringSet objects:
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
mm <- complement(narrow(probes, start=13, end=13))</pre>
left <- narrow(probes, end=12)</pre>
right <- narrow(probes, start=14)</pre>
xscat(left, mm, right)
## Collapsing an XStringSet (or XStringViews) object with a small
## number of elements:
probes1000 <- as.list(probes[1:1000])</pre>
y1 <- do.call(xscat, probes1000)</pre>
y2 <- do.call(c, probes1000) # slightly faster than the above
y1 == y2 \# TRUE
## Note that this method won't be efficient when the number of
## elements to collapse is big (> 10000) so we need to provide a
## collapse() (or xscollapse()) function in Biostrings that will be
## efficient at doing this. Please request this on the Bioconductor
## mailing list (http://bioconductor.org/help/mailing-list/) if you
## need it.
```

XString-class

BString objects

## **Description**

The BString class is a general container for storing a big string (a long sequence of characters) and for making its manipulation easy and efficient.

The DNAString, RNAString and AAString classes are similar containers but with the more biologyoriented purpose of storing a DNA sequence (DNAString), an RNA sequence (RNAString), or a sequence of amino acids (AAString).

All those containers derive directly (and with no additional slots) from the XString virtual class.

### **Details**

The 2 main differences between an XString object and a standard character vector are: (1) the data stored in an XString object are not copied on object duplication and (2) an XString object can only store a single string (see the XStringSet container for an efficient way to store a big collection of strings in a single object).

Unlike the DNAString, RNAString and AAString containers that accept only a predefined set of letters (the alphabet), a BString object can be used for storing any single string based on a single-byte character set.

## Constructor-like functions and generics

In the code snippet below, x can be a single string (character vector of length 1) or an XString object.

BString(x="", start=1, nchar=NA): Tries to convert x into a BString object by reading nchar letters starting at position start in x.

XString-class 109

## **Accessor methods**

In the code snippets below, x is an XString object.

alphabet(x): NULL for a BString object. See the corresponding man pages when x is a DNAS-tring, RNAString or AAString object.

length(x): or nchar(x): Get the length of an XString object, i.e., its number of letters.

#### Coercion

In the code snippets below, x is an XString object.

```
as.character(x): Converts x to a character string. toString(x): Equivalent to as.character(x).
```

## Subsetting

In the code snippets below, x is an XString object.

x[i]: Return a new XString object made of the selected letters (subscript i must be an NA-free numeric vector specifying the positions of the letters to select). The returned object belongs to the same class as x.

Note that, unlike subseq, x[i] does copy the sequence data and therefore will be very inefficient for extracting a big number of letters (e.g. when i contains millions of positions).

# **Equality**

In the code snippets below, e1 and e2 are XString objects.

```
e1 == e2: TRUE if e1 is equal to e2. FALSE otherwise.
```

Comparison between two XString objects of different base types (e.g. a BString object and a DNAString object) is not supported with one exception: a DNAString object and an RNAString object can be compared (see RNAString-class for more details about this).

Comparison between a BString object and a character string is also supported (see examples below).

```
e1 != e2: Equivalent to ! (e1 == e2).
```

## Author(s)

H. Pagès

#### See Also

subseq, letter, DNAString-class, RNAString-class, AAString-class, XStringSet-class, XStringViews-class, reverseComplement, compact, XVector-class

110 XStringQuality-class

## **Examples**

```
b <- BString("I am a BString object")</pre>
b
length(b)
## Extracting a linear subsequence:
subseq(b)
subseq(b, start=3)
subseq(b, start=-3)
subseq(b, end=-3)
subseq(b, end=-3, width=5)
## Subsetting:
b2 <- b[length(b):1]
                        # better done with reverse(b)
as.character(b2)
b2 == b
                           # FALSE
b2 == as.character(b2)
                           # TRUE
## b[1:length(b)] is equal but not identical to b!
b == b[1:length(b)]
                     # TRUE
identical(b, 1:length(b)) # FALSE
## This is because subsetting an XString object with [ makes a copy
## of part or all its sequence data. Hence, for the resulting object,
## the internal slot containing the memory address of the sequence
## data differs from the original. This is enough for identical() to
## see the 2 objects as different.
## Compacting. As a particular type of XVector objects, XString
## objects can optionally be compacted. Compacting is done typically
## before serialization. See ?compact for more information.
```

XStringQuality-class PhredQuality, SolexaQuality and IlluminaQuality objects

## **Description**

Objects for storing string quality measures.

# Usage

```
## Constructors:
PhredQuality(x)
SolexaQuality(x)
IlluminaQuality(x)
## alphabet and encoding
## S4 method for signature 'XStringQuality'
```

XStringQuality-class 111

```
alphabet(x)
## S4 method for signature 'XStringQuality'
encoding(x)
```

## **Arguments**

Х

Either a character vector, BString, BStringSet, integer vector, or number vector of error probabilities.

#### **Details**

PhredQuality objects store characters that are interpreted as [0 - 99] quality measures by subtracting 33 from their ASCII decimal representation (e.g. ! = 0, " = 1, # = 2, ...). Quality measures q encode probabilities as  $-10 * \log 10(p)$ .

SolexaQuality objects store characters that are interpreted as [-5 - 99] quality measures by subtracting 64 from their ASCII decimal representation (e.g.; = -5, < = -4, = = -3, ...). Quality measures q encode probabilities as  $-10 * (\log 10(p) - \log 10(1 - p))$ .

IlluminaQuality objects store characters that are interpreted as [0 - 99] quality measures by subtracting 64 from their ASCII decimal representation (e.g. @ = 0, A = 1, B = 2, ...). Quality measures q encode probabilities as -10 \* log10(p)

# Alphabet and encoding

In the code snippets below, x is an XStringQuality object.

alphabet(x): Valid letters in this quality score; not all letters are encountered in actual sequencing runs.

encoding(x): Map between letters and their corresponding integer encoding. Use as.integer and as.numeric to coerce objects to their integer and probability representations.

## Author(s)

P. Aboyoun

#### See Also

pairwiseAlignment and PairwiseAlignments-class in the **pwalign** package, DNAString-class, BStringSet-class

# **Examples**

```
PhredQuality(0:40)
SolexaQuality(0:40)
IlluminaQuality(0:40)

pq <- PhredQuality(c("*+,-./", "0123456789:;"))
qs <- as(pq, "IntegerList") # quality scores
qs
as(qs, "PhredQuality")
p <- as(pq, "NumericList") # probabilities</pre>
```

```
as(p, "PhredQuality")
PhredQuality(seq(1e-4,0.5,length=10))
SolexaQuality(seq(1e-4,0.5,length=10))
IlluminaQuality(seq(1e-4,0.5,length=10))

x <- SolexaQuality(BStringSet(c(a="@ABC", b="abcd")))
as(x, "IntegerList")  # quality scores
as(x, "NumericList")  # probabilities
as.matrix(x)  # quality scores</pre>
```

XStringSet-class

XStringSet objects

## Description

The BStringSet class is a container for storing a set of BString objects and for making its manipulation easy and efficient.

Similarly, the DNAStringSet (or RNAStringSet, or AAStringSet) class is a container for storing a set of DNAString (or RNAString, or AAString) objects.

All those containers derive directly (and with no additional slots) from the XStringSet virtual class.

## Usage

```
## Constructors:
BStringSet(x=character(), start=NA, end=NA, width=NA, use.names=TRUE)
DNAStringSet(x=character(), start=NA, end=NA, width=NA, use.names=TRUE)
RNAStringSet(x=character(), start=NA, end=NA, width=NA, use.names=TRUE)
AAStringSet(x=character(), start=NA, end=NA, width=NA, use.names=TRUE)

## Accessor-like methods:
## S4 method for signature 'character'
width(x)
## S4 method for signature 'XStringSet'
nchar(x, type="chars", allowNA=FALSE)

## ... and more (see below)
```

# **Arguments**

```
x Either a character vector (with no NAs), or an XString, XStringSet or XStringViews object.

start, end, width

Either NA, a single integer, or an integer vector of the same length as x specifying how x should be "narrowed" (see ?narrow for the details).
```

use.names TRUE or FALSE. Should names be preserved?

type, allowNA Ignored.

#### **Details**

The BStringSet, DNAStringSet, RNAStringSet and AAStringSet functions are constructors that can be used to turn input x into an XStringSet object of the desired base type.

They also allow the user to "narrow" the sequences contained in x via proper use of the start, end and/or width arguments. In this context, "narrowing" means dropping a prefix or/and a suffix of each sequence in x. The "narrowing" capabilities of these constructors can be illustrated by the following property: if x is a character vector (with no NAs), or an XStringSet (or XStringViews) object, then the 3 following transformations are equivalent:

```
BStringSet(x, start=mystart, end=myend, width=mywidth):
subseq(BStringSet(x), start=mystart, end=myend, width=mywidth):
BStringSet(subseq(x, start=mystart, end=myend, width=mywidth)):
```

Note that, besides being more convenient, the first form is also more efficient on character vectors.

#### Accessor-like methods

In the code snippets below, x is an XStringSet object.

length(x): The number of sequences in x.

width(x): A vector of non-negative integers containing the number of letters for each element in x. Note that width(x) is also defined for a character vector with no NAs and is equivalent to nchar(x, type="bytes").

names(x): NULL or a character vector of the same length as x containing a short user-provided description or comment for each element in x. These are the only data in an XStringSet object that can safely be changed by the user. All the other data are immutable! As a general recommendation, the user should never try to modify an object by accessing its slots directly.

alphabet(x): Return NULL, DNA\_ALPHABET, RNA\_ALPHABET or AA\_ALPHABET depending on whether x is a BStringSet, DNAStringSet, RNAStringSet or AAStringSet object.

nchar(x): The same as width(x).

## Subsequence extraction and related transformations

In the code snippets below, x is a character vector (with no NAs), or an XStringSet (or XStringViews) object.

subseq(x, start=NA, end=NA, width=NA): Applies subseq on each element in x. See ?subseq for the details.

Note that this is similar to what substr does on a character vector. However there are some noticeable differences:

- (1) the arguments are start and stop for substr;
- (2) the SEW interface (start/end/width) interface of subseq is richer (e.g. support for negative start or end values); and (3) subseq checks that the specified start/end/width values are valid i.e., unlike substr, it throws an error if they define "out of limits" subsequences or subsequences with a negative width.

- narrow(x, start=NA, end=NA, width=NA, use.names=TRUE): Same as subseq. The only differences are: (1) narrow has a use.names argument; and (2) all the things narrow and subseq work on (IRanges, XStringSet or XStringViews objects for narrow, XVector or XStringSet objects for subseq). But they both work and do the same thing on an XStringSet object.
- threebands(x, start=NA, end=NA, width=NA): Like the method for IRanges objects, the threebands methods for character vectors and XStringSet objects extend the capability of narrow by returning the 3 set of subsequences (the left, middle and right subsequences) associated to the narrowing operation. See ?threebands in the IRanges package for the details.
- subseq(x, start=NA, end=NA, width=NA) <- value: A vectorized version of the subseq<- method for XVector objects. See ?`subseq<-` for the details.

# Subsetting and appending

In the code snippets below, x and values are XStringSet objects, and i should be an index specifying the elements to extract.

```
x[i]: Return a new XStringSet object made of the selected elements.
x[[i]]: Extract the i-th XString object from x.
append(x, values, after=length(x)): Add sequences in values to x.
```

## **Set operations**

In the code snippets below, x and y are XStringSet objects.

```
union(x, y): Union of x and y.
intersect(x, y): Intersection of x and y.
setdiff(x, y): Asymmetric set difference of x and y.
setequal(x, y): Set equality of x to y.
```

#### Other methods

In the code snippets below, x is an XStringSet object.

- unlist(x): Turns x into an XString object by combining the sequences in x together. Fast equivalent to do.call(c, as.list(x)).
- as.character(x, use.names=TRUE): Converts x to a character vector of the same length as x. The use.names argument controls whether or not names(x) should be propagated to the names of the returned vector.
- as.factor(x): Converts x to a factor, via as.character(x).
- as.matrix(x, use.names=TRUE): Returns a character matrix containing the "exploded" representation of the strings. Can only be used on an XStringSet object with equal-width strings. The use.names argument controls whether or not names(x) should be propagated to the row names of the returned matrix.
- toString(x): Equivalent to toString(as.character(x)).
- show(x): By default the show method displays 5 head and 5 tail lines. The number of lines can be altered by setting the global options showHeadLines and showTailLines. If the object length is less than the sum of the options, the full object is displayed. These options affect GRanges, GAlignments, IRanges, and XStringSet objects.

## **Display**

The letters in a DNAStringSet, RNAStringSet, or AAStringSet object are colored when displayed by the show() method. Set global option Biostrings.coloring to FALSE to turn off this coloring.

## Author(s)

H. Pagès

## See Also

- readDNAStringSet and writeXStringSet for reading/writing a DNAStringSet object (or other XStringSet derivative) from/to a FASTA or FASTQ file.
- XStringSet-comparison
- XString objects.
- XStringViews objects.
- XStringSetList objects.
- subseq, narrow, and substr.
- compact
- XVectorList objects.

## **Examples**

```
## A. USING THE XStringSet CONSTRUCTORS ON A CHARACTER VECTOR OR FACTOR
## Note that there is no XStringSet() constructor, but an XStringSet
## family of constructors: BStringSet(), DNAStringSet(), RNAStringSet(),
x0 <- c("#CTC-NACCAGTAT", "#TTGA", "TACCTAGAG")
width(x0)
x1 <- BStringSet(x0)</pre>
х1
## 3 equivalent ways to obtain the same BStringSet object:
BStringSet(x0, start=4, end=-3)
subseq(x1, start=4, end=-3)
BStringSet(subseq(x0, start=4, end=-3))
dna0 <- DNAStringSet(x0, start=4, end=-3)</pre>
dna0 # 'options(Biostrings.coloring=FALSE)' to turn off coloring
names(dna0)
names(dna0)[2] \leftarrow "seqB"
## When the input vector contains a lot of duplicates, turning it into
## a factor first before passing it to the constructor will produce an
## XStringSet object that is more compact in memory:
library(hgu95av2probe)
```

```
x2 <- sample(hgu95av2probe$sequence, 999000, replace=TRUE)
dna2a <- DNAStringSet(x2)</pre>
dna2b <- DNAStringSet(factor(x2)) # slower but result is more compact</pre>
object.size(dna2a)
object.size(dna2b)
## -----
## B. USING THE XStringSet CONSTRUCTORS ON A SINGLE SEQUENCE (XString
     OBJECT OR CHARACTER STRING)
x3 <- "abcdefghij"
BStringSet(x3, start=2, end=6:2) # behaves like 'substring(x3, 2, 6:2)'
BStringSet(x3, start=-(1:6))
x4 <- BString(x3)</pre>
BStringSet(x4, end=-(1:6), width=3)
## Randomly extract 1 million 40-mers from C. elegans chrI:
extractRandomReads <- function(subject, nread, readlength)</pre>
{
   if (!is.integer(readlength))
       readlength <- as.integer(readlength)</pre>
   start <- sample(length(subject) - readlength + 1L, nread,</pre>
                  replace=TRUE)
   DNAStringSet(subject, start=start, width=readlength)
}
library(BSgenome.Celegans.UCSC.ce2)
rndreads <- extractRandomReads(Celegans$chrI, 1000000, 40)</pre>
## Notes:
## - This takes only 2 or 3 seconds versus several hours for a solution
    using substring() on a standard character string.
## - The short sequences in 'rndreads' can be seen as the result of a
##
    simulated high-throughput sequencing experiment. A non-realistic
##
    one though because:
      (a) It assumes that the underlying technology is perfect (the
##
##
          generated reads have no technology induced errors).
##
      (b) It assumes that the sequenced genome is exactly the same as the
##
          reference genome.
      (c) The simulated reads can contain IUPAC ambiguity letters only
##
##
          because the reference genome contains them. In a real
          high-throughput sequencing experiment, the sequenced genome
##
##
          of course doesn't contain those letters, but the sequencer
##
          can introduce them in the generated reads to indicate ambiguous
##
          base-calling.
##
      (d) The simulated reads come from the plus strand only of a single
##
          chromosome.
## - See the getSeq() function in the BSgenome package for how to
    circumvent (d) i.e. how to generate reads that come from the whole
    genome (plus and minus strands of all chromosomes).
## -----
## C. USING THE XStringSet CONSTRUCTORS ON AN XStringSet OBJECT
## -----
library(drosophila2probe)
```

```
probes <- DNAStringSet(drosophila2probe)</pre>
probes
RNAStringSet(probes, start=2, end=-5) # does NOT copy the sequence data!
## -----
## D. USING THE XStringSet CONSTRUCTORS ON AN ORDINARY list OF XString
## -----
probes10 <- head(probes, n=10)</pre>
set.seed(33)
shuffled_nucleotides <- lapply(probes10, sample)</pre>
shuffled_nucleotides
DNAStringSet(shuffled_nucleotides) # does NOT copy the sequence data!
## Note that the same result can be obtained in a more compact way with
## just:
set.seed(33)
endoapply(probes10, sample)
## -----
## E. USING subseq() ON AN XStringSet OBJECT
## -----
subseq(probes, start=2, end=-5)
subseq(probes, start=13, end=13) <- "N"</pre>
probes
## Add/remove a prefix:
subseq(probes, start=1, end=0) <- "--"</pre>
probes
subseq(probes, end=2) <- ""</pre>
probes
## Do more complicated things:
subseq(probes, start=4:7, end=7) <- c("YYYY", "YYY", "YY", "Y")</pre>
subseq(probes, start=4, end=6) <- subseq(probes, start=-2:-5)</pre>
probes
## -----
## F. UNLISTING AN XStringSet OBJECT
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
unlist(probes)
## -----
## G. COMPACTING AN XStringSet OBJECT
## -----
## As a particular type of XVectorList objects, XStringSet objects can
## optionally be compacted. Compacting is done typically before
## serialization. See ?compact for more information.
```

```
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)

y <- subseq(probes[1:12], start=5)
probes@pool
y@pool
object.size(probes)
object.size(y)

y0 <- compact(y)
y0@pool
object.size(y0)</pre>
```

XStringSet-comparison Comparing and ordering the elements in one or more XStringSet objects

## **Description**

Methods for comparing and ordering the elements in one or more XStringSet objects.

## **Details**

Element-wise (aka "parallel") comparison of 2 XStringSet objects is based on the lexicographic order between 2 BString, DNAString, RNAString, or AAString objects.

For DNAStringSet and RNAStringSet objects, the letters in the respective alphabets (i.e. DNA\_ALPHABET and RNA\_ALPHABET) are ordered based on a predefined code assigned to each letter. The code assigned to each letter can be retrieved with:

```
dna_codes <- as.integer(DNAString(paste(DNA_ALPHABET, collapse="")))
names(dna_codes) <- DNA_ALPHABET

rna_codes <- as.integer(RNAString(paste(RNA_ALPHABET, collapse="")))
names(rna_codes) <- RNA_ALPHABET</pre>
```

Note that this order does NOT depend on the locale in use. Also note that comparing DNA sequences with RNA sequences is supported and in that case T and U are considered to be the same letter.

For BStringSet and AAStringSet objects, the alphabetical order is defined by the C collation. Note that, at the moment, AAStringSet objects are treated like BStringSet objects i.e. the alphabetical order is NOT defined by the order of the letters in AA\_ALPHABET. This might change at some point.

#### pcompare() and related methods

In the code snippets below, x and y are XStringSet objects.

pcompare(x, y): Performs element-wise (aka "parallel") comparison of x and y, that is, returns an integer vector where the i-th element is less than, equal to, or greater than zero if the i-th element in x is considered to be respectively less than, equal to, or greater than the i-th element in y. If x and y don't have the same length, then the shortest is recycled to the length of the longest (the standard recycling rules apply).

119

```
x == y, x != y, x <= y, x >= y, x < y, x > y: Equivalent to pcompare(x, y) == 0, pcompare(x, y) != 0, pcompare(x, y) <= 0, pcompare(x, y) >= 0, pcompare(x, y) < 0, and pcompare(x, y) > 0, respectively.
```

# order() and related methods

In the code snippets below, x is an XStringSet object.

is.unsorted(x, strictly=FALSE): Return a logical values specifying if x is unsorted. The strictly argument takes logical value indicating if the check should be for \_strictly\_ increasing values.

order(x, decreasing=FALSE): Return a permutation which rearranges x into ascending or descending order.

```
rank(x, ties.method=c("first", "min")): Rank x in ascending order. sort(x, decreasing=FALSE): Sort x into ascending or descending order.
```

## duplicated() and unique()

In the code snippets below, x is an XStringSet object.

duplicated(x): Return a logical vector whose elements denotes duplicates in x.

unique(x): Return the subset of x made of its unique elements.

#### match() and %in%

In the code snippets below, x and table are XStringSet objects.

match(x, table, nomatch=NA\_integer\_): Returns an integer vector containing the first positions of an identical match in table for the elements in x.

x %in% table: Returns a logical vector indicating which elements in x match identically with an element in table.

#### is.na() and related methods

In the code snippets below, x is an XStringSet object. An XStringSet object never contains missing values (these methods exist for compatibility).

```
is.na(x): Returns FALSE for every element. anyNA(x): Returns FALSE.
```

## Author(s)

H. Pagès

#### See Also

XStringSet-class, ==, is.unsorted, order, rank, sort, duplicated, unique, match, %in%

## **Examples**

```
## -----
## A. SIMPLE EXAMPLES
## -----
dna <- DNAStringSet(c("AAA", "TC", "", "TC", "AAA", "CAAC", "G"))</pre>
match(c("", "G", "AA", "TC"), dna)
library(drosophila2probe)
fly_probes <- DNAStringSet(drosophila2probe)</pre>
sum(duplicated(fly_probes)) # 481 duplicated probes
is.unsorted(fly_probes) # TRUE
fly_probes <- sort(fly_probes)</pre>
is.unsorted(fly_probes) # FALSE
is.unsorted(fly_probes, strictly=TRUE) # TRUE, because of duplicates
is.unsorted(unique(fly_probes), strictly=TRUE) # FALSE
## Nb of probes that are the reverse complement of another probe:
nb1 <- sum(reverseComplement(fly_probes) %in% fly_probes)</pre>
stopifnot(identical(nb1, 455L)) # 455 probes
## Probes shared between drosophila2probe and hgu95av2probe:
library(hgu95av2probe)
human_probes <- DNAStringSet(hgu95av2probe)</pre>
m <- match(fly_probes, human_probes)</pre>
stopifnot(identical(sum(!is.na(m)), 493L)) # 493 shared probes
## B. AN ADVANCED EXAMPLE
## -----
## We want to compare the first 5 bases with the 5 last bases of each
## probe in drosophila2probe. More precisely, we want to compute the
## percentage of probes for which the first 5 bases are the reverse
## complement of the 5 last bases.
library(drosophila2probe)
probes <- DNAStringSet(drosophila2probe)</pre>
first5 <- narrow(probes, end=5)</pre>
last5 <- narrow(probes, start=-5)</pre>
nb2 <- sum(first5 == reverseComplement(last5))</pre>
stopifnot(identical(nb2, 17L))
```

```
## Percentage:
100 * nb2 / length(probes) # 0.0064 %
## If the probes were random DNA sequences, a probe would have 1 chance
## out of 4^5 to have this property so the percentage would be:
100 / 4^5 # 0.098 %
## With randomly generated probes:
set.seed(33)
random_dna <- sample(DNAString(paste(DNA_BASES, collapse="")),</pre>
                      sum(width(probes)), replace=TRUE)
random_probes <- successiveViews(random_dna, width(probes))</pre>
random_probes
random_probes <- as(random_probes, "XStringSet")</pre>
random_probes
random_first5 <- narrow(random_probes, end=5)</pre>
random_last5 <- narrow(random_probes, start=-5)</pre>
nb3 <- sum(random_first5 == reverseComplement(random_last5))</pre>
100 * nb3 / length(random_probes) # 0.099 %
```

XStringSet-io

Read/write an XStringSet object from/to a file

# **Description**

Functions to read/write an XStringSet object from/to a file.

#### **Usage**

```
## Read FASTA (or FASTQ) files in an XStringSet object:
readBStringSet(filepath, format="fasta",
               nrec=-1L, skip=0L, seek.first.rec=FALSE,
               use.names=TRUE, with.qualities=FALSE)
readDNAStringSet(filepath, format="fasta",
               nrec=-1L, skip=0L, seek.first.rec=FALSE,
               use.names=TRUE, with.qualities=FALSE)
readRNAStringSet(filepath, format="fasta",
               nrec=-1L, skip=0L, seek.first.rec=FALSE,
               use.names=TRUE, with.qualities=FALSE)
readAAStringSet(filepath, format="fasta",
               nrec=-1L, skip=0L, seek.first.rec=FALSE,
               use.names=TRUE, with.qualities=FALSE)
## Extract basic information about FASTA (or FASTQ) files
## without actually loading the sequence data:
fasta.seqlengths(filepath,
               nrec=-1L, skip=0L, seek.first.rec=FALSE,
```

```
seqtype="B", use.names=TRUE)
fasta.index(filepath,
               nrec=-1L, skip=0L, seek.first.rec=FALSE,
               seqtype="B")
fastq.seqlengths(filepath,
               nrec=-1L, skip=0L, seek.first.rec=FALSE)
fastq.geometry(filepath,
               nrec=-1L, skip=0L, seek.first.rec=FALSE)
## Write an XStringSet object to a FASTA (or FASTQ) file:
writeXStringSet(x, filepath, append=FALSE,
                compress=FALSE, compression_level=NA, format="fasta", ...)
## Serialize an XStringSet object:
saveXStringSet(x, objname, dirpath=".", save.dups=FALSE, verbose=TRUE)
```

#### **Arguments**

filepath

A character vector (of arbitrary length when reading, of length 1 when writing) containing the path(s) to the file(s) to read or write. Reading files in gzip format (which usually have the '.gz' extension) is supported.

Note that special values like "" or "|cmd" (typically supported by other I/O functions in R) are not supported here.

Also filepath cannot be a standard connection. However filepath can be an object as returned by open\_input\_files. This object can be used to read files by chunks. See "READ FILES BY CHUNK" in the examples section for the details.

format

Either "fasta" (the default) or "fastq".

nrec

Single integer. The maximum of number of records to read in. Negative values are ignored.

skip

Single non-negative integer. The number of records of the data file(s) to skip before beginning to read in records.

seek.first.rec TRUE or FALSE (the default). If TRUE, then the reading function starts by setting the file position indicator at the beginning of the first line in the file that looks like the beginning of a FASTA (if format is "fasta") or FASTQ (if format is "fastq") record. More precisely this is the first line in the file that starts with a '>' (for FASTA) or a '@' (for FASTQ). An error is raised if no such line is found.

> Normal parsing then starts from there, and everything happens like if the file actually started there. In particular it will be an error if this first record is not a valid FASTA or FASTO record.

> Using seek.first.rec=TRUE is useful for example to parse GFF3 files with embedded FASTA data.

use names

TRUE (the default) or FALSE. If TRUE, then the returned vector is named. For FASTA the names are taken from the record description lines. For FASTQ they are taken from the record sequence ids. Dropping the names with use.names=FALSE

> can help reduce memory footprint e.g. for a FASTQ file containing millions of reads.

with qualities TRUE or FALSE (the default). This argument is only supported when reading a FASTQ file. If TRUE, then the quality strings are also read and returned in the qualities metadata column of the returned DNAStringSet object. Note that by default the quality strings are ignored. This helps reduce memory footprint if the FASTQ file contains millions of reads.

> Note that using readQualityScaledDNAStringSet() is the preferred way to load a set of DNA sequences and their qualities from a FASTQ file into Bioconductor. Its main advantage is that it will return a QualityScaledDNAStringSet object instead of a DNAStringSet object, which makes handling of the qualities more convenient and less error prone. See "READ A FASTQ FILE AS A QualityScaledDNAStringSet OBJECT" in the Examples section below for more information.

seqtype

A single string specifying the type of sequences contained in the FASTA file(s). Supported sequence types:

- "B" for anything i.e. any letter is a valid one-letter sequence code.
- "DNA" for DNA sequences i.e. only letters in DNA\_ALPHABET (case ignored) are valid one-letter sequence codes.
- "RNA" for RNA sequences i.e. only letters in RNA\_ALPHABET (case ignored) are valid one-letter sequence codes.
- "AA" for Amino Acid sequences i.e. only letters in AA\_ALPHABET (case ignored) are valid one-letter sequence codes.

Invalid one-letter sequence codes are ignored with a warning.

For writeXStringSet, the object to write to file. Х

For saveXStringSet, the object to serialize.

append

TRUE or FALSE. If TRUE output will be appended to file; otherwise, it will overwrite the contents of file. See ?cat for the details.

compress

Like for the save function in base R, must be TRUE or FALSE (the default), or a single string specifying whether writing to the file is to use compression. The only type of compression supported at the moment is "gzip".

Passing TRUE is equivalent to passing "gzip".

compression\_level

Not implemented yet.

Further format-specific arguments.

If format="fasta", the width argument can be used to specify the maximum number of letters per line of sequence. width must be a single integer.

If format="fastq", the qualities argument can be used to specify the quality strings, qualities must be a BStringSet object. If the argument is omitted, then the quality strings are taken from the qualities metadata column of x (i.e. from mcols(x)\$qualities). If x has no qualities metadata column and the qualities argument is omitted, then the fake quality ';' is assigned to each letter in x and written to the FASTQ file. If x is a QualityScaledXStringSet and qualities is not defined, the qualities contained in x are used automatically.

objname The name of the serialized object.

dirpath The path to the directory where to save the serialized object.

save.dups TRUE or FALSE. If TRUE then the Dups object describing how duplicated elements

in x are related to each other is saved too. For advanced users only.

verbose TRUE or FALSE.

## **Details**

gzip compression is supported by reading and writing functions on all platforms.

readDNAStringSet and family (i.e. readBStringSet, readDNAStringSet, readRNAStringSet and readAAStringSet) load sequences from an input file (or multiple input files) into an XStringSet object. When multiple input files are specified, all must have the same format (i.e. FASTA or FASTQ) and files with different compression types can be mixed with non-compressed files. The files are read in the order they were specified and the sequences are stored in the returned object in the order they were read.

Only FASTA and FASTQ files are supported for now.

The fasta. seqlengths utility returns an integer vector with one element per FASTA record in the input files. Each element is the length of the sequence found in the corresponding record, that is, the number of valid one-letter sequence codes in the record. See description of the seqtype argument above for how to control the set of valid one-letter sequence codes.

The fasta.index utility returns a data frame with 1 row per FASTA record in the input files and the following columns:

- recno: The rank of the record in the (virtually) concatenated input files.
- fileno: The rank of the file where the record is located.
- offset: The offset of the record relative to the start of the file where it's located. Measured in bytes.
- desc: The description line (a.k.a. header) of the record.
- seqlength: The length of the sequence in the record (not counting invalid letters).
- filepath: The path to the file where the record is located. Always a local file, so if the user specified a remote file, this column will contain the path to the downloaded file.

A subset of this data frame can be passed to readDNAStringSet and family for direct access to an arbitrary subset of sequences. More precisely, if fai is a FASTA index that was obtained with fasta.index(filepath, ..., seqtype="DNA"), then readDNAStringSet(fai[i, ]) is equivalent to readDNAStringSet(filepath, ...)[i] for any valid subscript i, except that the former only loads the requested sequences in memory and thus will be more memory efficient if only a small subset of sequences is requested.

The fastq.seqlengths utility returns the read lengths in an integer vector with one element per FASTQ record in the input files.

The fastq.geometry utility is a convenience wrapper around fastq.seqlengths that returns an integer vector of length 2 describing the *geometry* of the FASTQ files. The first integer gives the total number of FASTQ records in the files and the second element the common length of the reads (this common length is set to NA in case of variable length reads or if no FASTQ record was found).

This compact representation of the geometry can be useful if the FASTQ files are known to contain fixed length reads.

writeXStringSet writes an XStringSet object to a file. Like with readDNAStringSet and family, only FASTA and FASTQ files are supported for now. WARNING: Please be aware that using writeXStringSet on a BStringSet object that contains the '\n' (LF) or '\r' (CR) characters or the FASTA markup characters '>' or ';' is almost guaranteed to produce a broken FASTA file!

Serializing an XStringSet object with saveXStringSet is equivalent to using the standard save mechanism. But it will try to reduce the size of x in memory first before calling save. Most of the times this leads to a much reduced size on disk.

#### References

http://en.wikipedia.org/wiki/FASTA\_format

#### See Also

- BStringSet, DNAStringSet, RNAStringSet, and AAStringSet objects.
- readQualityScaledDNAStringSet and writeQualityScaledXStringSet for reading/writing a QualityScaledDNAStringSet object (or other QualityScaledXStringSet derivative) from/to a FASTQ file.

## **Examples**

```
## A. READ/WRITE FASTA FILES
## -----
## Read a non-compressed FASTA files:
filepath1 <- system.file("extdata", "someORF.fa", package="Biostrings")</pre>
fasta.seqlengths(filepath1, seqtype="DNA")
x1 <- readDNAStringSet(filepath1)</pre>
х1
## Read a gzip-compressed FASTA file:
filepath2 <- system.file("extdata", "someORF.fa.gz", package="Biostrings")</pre>
fasta.seqlengths(filepath2, seqtype="DNA")
x2 <- readDNAStringSet(filepath2)</pre>
x2
## Sanity check:
stopifnot(identical(as.character(x1), as.character(x2)))
## Read 2 FASTA files at once:
filepath3 <- system.file("extdata", "fastaEx.fa", package="Biostrings")</pre>
fasta.seqlengths(c(filepath2, filepath3), seqtype="DNA")
x23 <- readDNAStringSet(c(filepath2, filepath3))</pre>
x23
## Sanity check:
x3 <- readDNAStringSet(filepath3)</pre>
stopifnot(identical(as.character(x23), as.character(c(x2, x3))))
```

```
## Use a FASTA index to load only an arbitrary subset of sequences:
filepath4 <- system.file("extdata", "dm3_upstream2000.fa.gz",</pre>
                          package="Biostrings")
fai <- fasta.index(filepath4, seqtype="DNA")</pre>
head(fai)
head(fai$desc)
i <- sample(nrow(fai), 10) # randomly pick up 10 sequences</pre>
x4 <- readDNAStringSet(fai[i, ])</pre>
## Sanity check:
stopifnot(identical(as.character(readDNAStringSet(filepath4)[i]),
                     as.character(x4)))
## Write FASTA files:
out23a <- tempfile()</pre>
writeXStringSet(x23, out23a)
out23b <- tempfile()</pre>
writeXStringSet(x23, out23b, compress=TRUE)
file.info(c(out23a, out23b))$size
## Sanity checks:
stopifnot(identical(as.character(readDNAStringSet(out23a)),
                     as.character(x23)))
stopifnot(identical(readLines(out23a), readLines(out23b)))
## B. READ/WRITE FASTQ FILES
filepath5 <- system.file("extdata", "s_1_sequence.txt",</pre>
                          package="Biostrings")
fastq.geometry(filepath5)
## The quality strings are ignored by default:
reads <- readDNAStringSet(filepath5, format="fastq")</pre>
reads
mcols(reads)
## Use 'with.qualities=TRUE' to load them:
reads <- readDNAStringSet(filepath5, format="fastq", with.qualities=TRUE)</pre>
reads
mcols(reads)
mcols(reads)$qualities
## Each quality string contains one letter per nucleotide in the
## corresponding read:
stopifnot(identical(width(mcols(reads)$qualities), width(reads)))
## Write the reads to a FASTQ file:
outfile <- tempfile()</pre>
writeXStringSet(reads, outfile, format="fastq")
```

```
outfile2 <- tempfile()</pre>
writeXStringSet(reads, outfile2, compress=TRUE, format="fastq")
## Sanity checks:
stopifnot(identical(readLines(outfile), readLines(filepath5)))
stopifnot(identical(readLines(outfile), readLines(outfile2)))
## C. READ FILES BY CHUNK
## readDNAStringSet() supports reading an arbitrary number of FASTA or
## FASTQ records at a time in a loop. This can be useful to process
## big FASTA or FASTQ files by chunk and thus avoids loading the entire
## file in memory. To achieve this the files to read from need to be
## opened with open_input_files() first. Note that open_input_files()
## accepts a vector of file paths and/or URLs.
## With FASTA files:
files <- open_input_files(filepath4)</pre>
i <- 0
while (TRUE) {
    i < -i + 1
   ## Load 4000 records at a time. Each new call to readDNAStringSet()
   ## picks up where the previous call left.
   dna <- readDNAStringSet(files, nrec=4000)</pre>
    if (length(dna) == 0L)
       break
   cat("processing chunk", i, "...\n")
    ## do something with 'dna' ...
}
## With FASTQ files:
files <- open_input_files(filepath5)</pre>
i <- 0
while (TRUE) {
   i < -i + 1
   \#\# Load 75 records at a time.
   reads <- readDNAStringSet(files, format="fastq", nrec=75)</pre>
    if (length(reads) == 0L)
       break
    cat("processing chunk", i, "...\n")
    ## do something with 'reads' ...
}
## IMPORTANT NOTE: Like connections, the object returned by
## open_input_files() can NOT be shared across workers in the
## context of parallelization!
## D. READ A FASTQ FILE AS A QualityScaledDNAStringSet OBJECT
## -----
## Use readQualityScaledDNAStringSet() if you want the object to be
```

128 XStringSetList-class

```
## returned as a QualityScaledDNAStringSet instead of a DNAStringSet
## object. See ?readQualityScaledDNAStringSet for more information.
## Note that readQualityScaledDNAStringSet() is a simple wrapper around
## readDNAStringSet() that does the following if the file contains
## "Phred quality scores" (which is the standard Sanger variant to
## assess reliability of a base call):
reads <- readDNAStringSet(filepath5, format="fastq", with.qualities=TRUE)</pre>
quals <- PhredQuality(mcols(reads)$qualities)</pre>
QualityScaledDNAStringSet(reads, quals)
## The call to PhredQuality() is replaced with a call to SolexaQuality()
## or IlluminaQuality() if the quality scores are Solexa quality scores.
## E. GENERATE FAKE READS AND WRITE THEM TO A FASTQ FILE
library(BSgenome.Celegans.UCSC.ce2)
## Create a "sliding window" on chr I:
sw_start <- seq.int(1, length(Celegans$chrI)-50, by=50)</pre>
sw <- Views(Celegans$chrI, start=sw_start, width=10)</pre>
my_fake_reads <- as(sw, "XStringSet")</pre>
my_fake_ids <- sprintf("ID%06d", seq_len(length(my_fake_reads)))</pre>
names(my_fake_reads) <- my_fake_ids</pre>
my_fake_reads
## Fake quality ';' will be assigned to each base in 'my_fake_reads':
out2 <- tempfile()</pre>
writeXStringSet(my_fake_reads, out2, format="fastq")
## Passing qualities thru the 'qualities' argument:
my_fake_quals <- rep.int(BStringSet("DCBA@?>=<;"), length(my_fake_reads))</pre>
my_fake_quals
out3 <- tempfile()
writeXStringSet(my_fake_reads, out3, format="fastq",
               qualities=my_fake_quals)
## -----
## F. SERIALIZATION
## -----
saveXStringSet(my_fake_reads, "my_fake_reads", dirpath=tempdir())
```

# Description

The XStringSetList class is a virtual container for storing a list of XStringSet objects.

XStringSetList-class 129

## Usage

```
## Constructors:
BStringSetList(..., use.names=TRUE)
DNAStringSetList(..., use.names=TRUE)
RNAStringSetList(..., use.names=TRUE)
AAStringSetList(..., use.names=TRUE)
```

# **Arguments**

... Character vector(s) (with no NAs), or XStringSet object(s), or XStringViews

object(s) to be concatenated into a XStringSetList.

use.names TRUE or FALSE. Should names be preserved?

#### **Details**

Concrete flavors of the XStringSetList container are the BStringSetList, DNAStringSetList, RNAStringSetList and AAStringSetList containers for storing a list of BStringSet, DNAStringSet, RNAStringSet and AAStringSet objects, respectively. These four containers are direct subclasses of XStringSetList with no additional slots. The XStringSetList class itself is virtual and has no constructor.

#### Methods

The XStringSetList class extends the List class defined in the IRanges package. Using a less technical jargon, this just means that an XStringSetList object is a list-like object that can be manipulated like an ordinary list. Or, said otherwise, most of the operations that work on an ordinary list (e.g. length, names, [, [[, c, unlist, etc...) should work on an XStringSetList object. In addition, Bioconductor specific list operations like elementNROWS and PartitioningByEnd (defined in the IRanges package) are supported too.

## Author(s)

H. Pagès

## See Also

XStringSet-class, List-class

## **Examples**

```
## ------
## A. THE XStringSetList CONSTRUCTORS
## ------
dna1 <- c("AAA", "AC", "", "T", "GGATA")
dna2 <- c("G", "TT", "C")

x <- DNAStringSetList(dna1, dna2)
x</pre>
```

130 XString Views-class

```
DNAStringSetList(DNAStringSet(dna1), DNAStringSet(dna2))
DNAStringSetList(dna1, DNAStringSet(dna2))
DNAStringSetList(DNAStringSet(dna1), dna2)
DNAStringSetList(dna1, RNAStringSet(DNAStringSet(dna2)))
DNAStringSetList(list(dna1, dna2))
DNAStringSetList(CharacterList(dna1, dna2))
## Empty object (i.e. zero-length):
DNAStringSetList()
## Not empty (length is 1):
DNAStringSetList(character(0))
## B. UNLISTING AN XStringSetList OBJECT
length(x)
elementNROWS(x)
unlist(x)
x[[1]]
x[[2]]
as.list(x)
names(x) <- LETTERS[1:length(x)]</pre>
x[["A"]]
x[["B"]]
as.list(x) # named list
```

XStringViews-class

The XStringViews class

## **Description**

The XStringViews class is the basic container for storing a set of views (start/end locations) on the same sequence (an XString object).

## **Details**

An XStringViews object contains a set of views (start/end locations) on the same XString object called "the subject string" or "the subject sequence" or simply "the subject". Each view is defined by its start and end locations: both are integers such that start <= end. An XStringViews object is in fact a particular case of an Views object (the XStringViews class contains the Views class) so it can be manipulated in a similar manner: see ?Views for more information. Note that two views can overlap and that a view can be "out of limits" i.e. it can start before the first letter of the subject or/and end after its last letter.

XStringViews-class 131

#### Constructor

Views(subject, start=NULL, end=NULL, width=NULL, names=NULL): See ?Views in the IRanges package for the details.

#### Accessor-like methods

All the accessor-like methods defined for Views objects work on XStringViews objects. In addition, the following accessors are defined for XStringViews objects:

nchar(x): A vector of non-negative integers containing the number of letters in each view. Values in nchar(x) coincide with values in width(x) except for "out of limits" views where they are lower.

#### Other methods

In the code snippets below, x, object, e1 and e2 are XStringViews objects, and i can be a numeric or logical vector.

e1 == e2: A vector of logicals indicating the result of the view by view comparison. The views in the shorter of the two XStringViews object being compared are recycled as necessary.

Like for comparison between XString objects, comparison between two XStringViews objects with subjects of different classes is not supported with one exception: when the subjects are DNAString and RNAString instances.

Also, like with XString objects, comparison between an XStringViews object with a BString subject and a character vector is supported (see examples below).

- e1 != e2: Equivalent to ! (e1 == e2).
- as.character(x, use.names=TRUE, check.limits=TRUE): Converts x to a character vector of the same length as x. The use.names argument controls whether or not names(x) should be propagated to the names of the returned vector. The check.limits argument controls whether or not an error should be raised if x has "out of limit" views. If check.limits is FALSE then "out of limit" views are trimmed with a warning.
- as.data.frame(x, row.names = NULL, optional = FALSE, ...) Equivalent of as.data.frame(as.character(x)).
- as.matrix(x, use.names=TRUE): Returns a character matrix containing the "exploded" representation of the views. Can only be used on an XStringViews object with equal-width views. The use.names argument controls whether or not names(x) should be propagated to the row names of the returned matrix.

toString(x): Equivalent to toString(as.character(x)).

# Author(s)

H. Pagès

#### See Also

Views-class, gaps, XString-class, XStringSet-class, letter, MIndex-class

132 XString Views-class

## **Examples**

```
## One standard way to create an XStringViews object is to use
## the Views() constructor.
## Views on a DNAString object:
s <- DNAString("-CTC-N")</pre>
v4 <- Views(s, start=3:0, end=5:8)
٧4
subject(v4)
length(v4)
start(v4)
end(v4)
width(v4)
## Attach a comment to views #3 and #4:
names(v4)[3:4] <- "out of limits"</pre>
names(v4)
## A more programatical way to "tag" the "out of limits" views:
names(v4)[start(v4) < 1 \mid nchar(subject(v4)) < end(v4)] <- "out of limits"
names(v4)[nchar(v4) < width(v4)] <- "out of limits"</pre>
## Two equivalent ways to extract a view as an XString object:
s2a <- v4[[2]]
s2b <- subseq(subject(v4), start=start(v4)[2], end=end(v4)[2])</pre>
identical(s2a, s2b) # TRUE
## It is an error to try to extract an "out of limits" view:
#v4[[3]] # Error!
v12 <- Views(DNAString("TAATAATG"), start=-2:9, end=0:11)</pre>
v12 == DNAString("TAA")
v12[v12 == v12[4]]
v12[v12 == v12[1]]
v12[3] == Views(RNAString("AU"), start=0, end=2)
## Here the first view doesn't even overlap with the subject:
Views(BString("aaa--b"), start=-3:4, end=-3:4 + c(3:6, 6:3))
## 'start' and 'end' are recycled:
subject <- "abcdefghij"</pre>
Views(subject, start=2:1, end=4)
Views(subject, start=5:7, end=nchar(subject))
Views(subject, start=1, end=5:7)
## Applying gaps() to an XStringViews object:
v2 <- Views("abCDefgHIJK", start=c(8, 3), end=c(14, 4))</pre>
gaps(v2)
## Coercion:
as(v12, "XStringSet") # same as 'as(v12, "DNAStringSet")'
```

yeastSEQCHR1 133

```
rna <- as(v12, "RNAStringSet")
as(rna, "Views")</pre>
```

yeastSEQCHR1

An annotation data file for CHR1 in the yeastSEQ package

# **Description**

This is a single character string containing DNA sequence of yeast chromosome number 1. The data were obtained from the Saccharomyces Genome Database (ftp://genome-ftp.stanford.edu/pub/yeast/data\\_download/sequence/genomic\\_sequence/chromosomes/fasta/).

## **Details**

Annotation based on data provided by Yeast Genome project.

Source data built: Yeast Genome data are built at various time intervals. Sources used were downloaded Fri Nov 21 14:00:47 2003 Package built: Fri Nov 21 14:00:47 2003

# References

http://www.yeastgenome.org/DownloadContents.shtml

# **Examples**

```
data(yeastSEQCHR1)
nchar(yeastSEQCHR1)
```

# **Index**

!=,BString,character-method	* distribution
(XString-class), 108	dinucleotideFrequencyTest, 11
!=,XString,XString-method	* htest
(XString-class), 108	dinucleotideFrequencyTest, 11
!=,XString,XStringViews-method	* internal
(XStringViews-class), 130	Biostrings internals, 6
!=,XStringViews,XString-method	* manip
(XStringViews-class), 130	chartr, 6
!=,XStringViews,XStringViews-method	detail, 10
(XStringViews-class), 130	getSeq, 19
!=,XStringViews,character-method	gregexpr2, 20
(XStringViews-class), 130	injectHardMask, 22
!=,character,BString-method	letterFrequency, 26
(XString-class), 108	longestConsecutive, 32
!=,character,XStringViews-method	maskMotif, 40
(XStringViews-class), 130	matchPWM, 65
* classes	misc, 69
AAString-class,4	nucleotideFrequency, 74
DNAString-class, 12	padAndClip, 79
MaskedXString-class, 38	replaceAt, 88
MIndex-class, 67	replaceLetterAt, 92
MultipleAlignment-class, 70	reverseComplement, 95
PDict-class, 81	translate, 101
QualityScaledXStringSet-class, 86	xscat, 107
RNAString-class, 97	XStringSet-io, 121
XString-class, 108	* methods
XStringQuality-class, 110	AAString-class, 4
XStringSet-class, 112	chartr, 6
XStringSetList-class, 128	DNAString-class, 12
XStringViews-class, 130	findPalindromes, 14
* datasets	lcsuffix, 25
HNF4alpha, 21	letter, 25
<pre>predefined_scoring_matrices, 85</pre>	letterFrequency, 26
yeastSEQCHR1, 133	lowlevel-matching, 33
* data	MaskedXString-class, 38
AMINO_ACID_CODE, 5	maskMotif, 40
GENETIC_CODE, 16	match-utils, 42
IUPAC_CODE_MAP, 23	matchLRPatterns, 44
<pre>predefined_scoring_matrices, 85</pre>	matchPattern, 46
	,

matchPDict, 50	(XString-class), 108
matchPDict-inexact, 59	==,character,XStringViews-method
matchProbePair, 63	(XStringViews-class), 130
matchPWM, 65	<pre>[[,ByPos_MIndex-method(MIndex-class),</pre>
MIndex-class, 67	67
misc, 69	[[,PDict-method(PDict-class), 81
MultipleAlignment-class, 70	[[,SparseList-method(Biostrings
nucleotideFrequency, 74	internals), 6
padAndClip, 79	%in%, <i>120</i>
PDict-class, 81	
QualityScaledXStringSet-class, 86	AA_ALPHABET, 18, 101, 102, 113, 118, 123
replaceAt, 88	AA_ALPHABET (AAString-class),4
reverseComplement, 95	AA_PROTEINOGENIC(AAString-class),4
RNAString-class, 97	AA_STANDARD (AAString-class), 4
toComplex, 100	AAMultipleAlignment
translate, 101	(MultipleAlignment-class), 70
trimLRPatterns, 104	AAMultipleAlignment-class
xscat, 107	(MultipleAlignment-class), 70
	AAString, 4, 5, 13, 18, 98, 102, 108, 109, 112,
XString-class, 108	118
XStringQuality-class, 110	AAString (AAString-class), 4
XStringSet-class, 112	AAString-class, 4, 109
XStringSet-comparison, 118	AAStringSet, 80, 86, 87, 102, 103, 118, 125,
XStringSetList-class, 128	129
XStringViews-class, 130	AAStringSet(XStringSet-class), 112
utilities	AAStringSet-class(XStringSet-class),
AMINO_ACID_CODE, 5	112
GENETIC_CODE, 16	AAStringSetList(XStringSetList-class),
injectHardMask, 22	128
IUPAC_CODE_MAP, 23	AAStringSetList-class
matchPWM, 65	(XStringSetList-class), 128
replaceLetterAt, 92	ACtree2 (PDict-class), 81
XStringSet-io, 121	ACtree2-class (PDict-class), 81
inplaceReplaceLetterAt	agrep, <i>105</i>
(replaceLetterAt), 92	align-utils, <i>36</i> , <i>43</i>
==, 120	alphabet, 26, 27, 29, 38, 76
==,BString,character-method	alphabet(XString-class), 108
(XString-class), 108	alphabet, ANY-method (XString-class), 108
==,XString,XString-method	alphabet,XStringQuality-method
(XString-class), 108	(XStringQuality-class), 110
==,XString,XStringViews-method	alphabetFrequency, 5, 7, 13, 39, 48, 53, 76,
(XStringViews-class), 130	98
==,XStringViews,XString-method	alphabetFrequency (letterFrequency), $26$
(XStringViews-class), 130	${\tt alphabetFrequency,AAString-method}$
==,XStringViews,XStringViews-method	(letterFrequency), 26
(XStringViews-class), 130	alphabetFrequency,AAStringSet-method
==,XStringViews,character-method	(letterFrequency), 26
(XStringViews-class), 130	alphabetFrequency,DNAString-method
==,character,BString-method	(letterFrequency), 26

alphabetFrequency,DNAStringSet-method	as.vector,XString-method
(letterFrequency), 26	(XString-class), 108
alphabetFrequency,MaskedXString-method	as.vector,XStringQuality-method
(letterFrequency), 26	(XStringQuality-class), 110
alphabetFrequency, MultipleAlignment-method	as.vector,XStringSet-method
(MultipleAlignment-class), 70	(XStringSet-class), 112
alphabetFrequency,RNAString-method	
(letterFrequency), 26	Biostrings internals, 6
alphabetFrequency,RNAStringSet-method	BLOSUM100
(letterFrequency), 26	<pre>(predefined_scoring_matrices),</pre>
alphabetFrequency,XString-method	85
(letterFrequency), 26	BLOSUM45 (predefined_scoring_matrices),
alphabetFrequency,XStringSet-method	85
(letterFrequency), 26	BLOSUM50 (predefined_scoring_matrices),
alphabetFrequency,XStringViews-method	85
(letterFrequency), 26	<pre>BLOSUM62 (predefined_scoring_matrices),</pre>
AMINO_ACID_CODE, 4, 5, 5, 18, 76	85
anyNA,XStringSet-method	BLOSUM80 (predefined_scoring_matrices),
(XStringSet-comparison), 118	85
Arithmetic, 27	BSgenome, 19, 20, 94
as.character,MaskedXString-method	BString, 4, 12, 13, 29, 97, 98, 111, 112, 118
(MaskedXString-class), 38	BString (XString-class), 108
as.character,MultipleAlignment-method	BString-class (XString-class), 108
(MultipleAlignment-class), 70	BStringSet, 80, 86, 87, 111, 118, 123, 125,
as.character,XString-method	129
(XString-class), 108	BStringSet (XStringSet-class), 112
as.character,XStringSet-method	BStringSet-class, 111
(XStringSet-class), 112	BStringSet-class (XStringSet-class), 112
as.character,XStringViews-method	BStringSetList (XStringSetList-class),
(XStringViews-class), 130	128
as.data.frame,XStringSet-method	BStringSetList-class
(XStringSet-class), 112	(XStringSetList-class), 128
as.data.frame,XStringViews-method	ByPos_MIndex-class (MIndex-class), 67
(XStringViews-class), 130	
as.factor,XStringSet-method	c, 27
(XStringSet-class), 112	cat, <i>123</i>
as.list,MTB_PDict-method(PDict-class),	ceiling, 105
81	CharacterList, 89, 90
as.list,SparseList-method(Biostrings	chartr, 6, 6, 7, 23, 94, 95
internals), 6	chartr, ANY, ANY, MaskedXString-method
as.matrix,MultipleAlignment-method	(chartr), 6
(MultipleAlignment-class), 70	chartr, ANY, ANY, XString-method (chartr),
as.matrix,XStringQuality-method	6
(XStringQuality-class), 110	chartr, ANY, ANY, XStringSet-method
as.matrix,XStringSet-method	(chartr), 6
(XStringSet-class), 112	chartr, ANY, ANY, XStringViews-method
as.matrix,XStringViews-method	(chartr), 6
(XStringViews-class), 130	chisq.test, 12

class:AAMultipleAlignment	86
(MultipleAlignment-class), 70	class:QualityScaledDNAStringSet
<pre>class:AAString (AAString-class), 4</pre>	(QualityScaledXStringSet-class)
<pre>class:AAStringSet (XStringSet-class),</pre>	86
112	class:QualityScaledRNAStringSet
class:AAStringSetList	(QualityScaledXStringSet-class)
(XStringSetList-class), 128	86
class: ACtree2 (PDict-class), 81	class:QualityScaledXStringSet
class:BString (XString-class), 108	(QualityScaledXStringSet-class)
<pre>class:BStringSet (XStringSet-class), 112</pre>	86
class:BStringSetList	class:RNAMultipleAlignment
(XStringSetList-class), 128	(MultipleAlignment-class), 70
<pre>class:ByPos_MIndex (MIndex-class), 67</pre>	<pre>class:RNAString(RNAString-class), 97</pre>
class:DNAMultipleAlignment	<pre>class:RNAStringSet(XStringSet-class),</pre>
(MultipleAlignment-class), 70	112
<pre>class:DNAString (DNAString-class), 12</pre>	class:RNAStringSetList
<pre>class:DNAStringSet (XStringSet-class),</pre>	(XStringSetList-class), 128
112	class:SolexaQuality
class:DNAStringSetList	(XStringQuality-class), 110
(XStringSetList-class), 128	<pre>class:SparseList(Biostrings</pre>
<pre>class:Expanded_TB_PDict (PDict-class),</pre>	internals), 6
81	<pre>class:TB_PDict(PDict-class), 81</pre>
class:IlluminaQuality	<pre>class:Twobit (PDict-class), 81</pre>
(XStringQuality-class), 110	class:XString(XString-class), 108
class:MaskedAAString	<pre>class:XStringCodec (Biostrings</pre>
(MaskedXString-class), 38	internals), 6
class:MaskedBString	class:XStringQuality
(MaskedXString-class), 38	(XStringQuality-class), 110
class:MaskedDNAString	<pre>class:XStringSet (XStringSet-class), 112</pre>
(MaskedXString-class), 38	class:XStringSetList
class:MaskedRNAString	(XStringSetList-class), 128
(MaskedXString-class), 38	class:XStringViews
class:MaskedXString	(XStringViews-class), 130
(MaskedXString-class), 38	codons (translate), 101
class:MIndex (MIndex-class), 67	codons, DNAString-method (translate), 101
<pre>class:MTB_PDict (PDict-class), 81</pre>	codons,MaskedDNAString-method
class:MultipleAlignment	(translate), 101
(MultipleAlignment-class), 70	codons,MaskedRNAString-method
<pre>class:PDict (PDict-class), 81</pre>	(translate), 101
class:PDict3Parts(PDict-class), 81	codons, RNAString-method(translate), 101
class:PhredQuality	coerce,AAString,MaskedAAString-method
(XStringQuality-class), 110	(MaskedXString-class), 38
<pre>class:PreprocessedTB (PDict-class), 81</pre>	coerce, ANY, AAStringSet-method
class:QualityScaledAAStringSet	(XStringSet-class), 112
(QualityScaledXStringSet-class),	coerce, ANY, BStringSet-method
86	(XStringSet-class), 112
class:QualityScaledBStringSet	coerce, ANY, DNAStringSet-method
(QualityScaledXStringSet-class),	(XStringSet-class), 112

coerce, ANY, RNAStringSet-method	<pre>coerce,IntegerList,IlluminaQuality-method</pre>
(XStringSet-class), 112	(XStringQuality-class), 110
coerce, ANY, XStringSet-method	<pre>coerce,IntegerList,PhredQuality-method</pre>
(XStringSet-class), 112	(XStringQuality-class), 110
coerce, BString, IlluminaQuality-method	<pre>coerce,IntegerList,SolexaQuality-method</pre>
(XStringQuality-class), 110	(XStringQuality-class), 110
coerce, BString, MaskedBString-method	coerce, List, AAStringSetList-method
(MaskedXString-class), 38	(XStringSetList-class), 128
coerce, BString, PhredQuality-method	coerce, list, AAStringSetList-method
(XStringQuality-class), 110	(XStringSetList-class), 128
coerce, BString, SolexaQuality-method	coerce,List,BStringSetList-method
(XStringQuality-class), 110	(XStringSetList-class), 128
coerce, BStringSet, IlluminaQuality-method	coerce, list, BStringSetList-method
(XStringQuality-class), 110	(XStringSetList-class), 128
coerce, BStringSet, PhredQuality-method	coerce, List, DNAStringSetList-method
(XStringQuality-class), 110	(XStringSetList-class), 128
coerce, BStringSet, SolexaQuality-method	coerce, list, DNAStringSetList-method
(XStringQuality-class), 110	(XStringSetList-class), 128
coerce, character, AAMultipleAlignment-method	coerce,List,RNAStringSetList-method
(MultipleAlignment-class), 70	(XStringSetList-class), 128
coerce, character, AAString-method	coerce, list, RNAStringSetList-method
(XString-class), 108	(XStringSetList-class), 128
coerce, character, BString-method	coerce, List, XStringSetList-method
(XString-class), 108	(XStringSetList-class), 128
coerce, character, DNAMultipleAlignment-method	
(MultipleAlignment-class), 70	(XStringSetList-class), 128
coerce, character, DNAString-method	coerce, MaskedAAString, AAString-method
(XString-class), 108	(MaskedXString-class), 38
coerce, character, IlluminaQuality-method	coerce, MaskedBString, BString-method
(XStringQuality-class), 110	(MaskedXString-class), 38
coerce, character, PhredQuality-method	coerce, MaskedDNAString, DNAString-method
(XStringQuality-class), 110	(MaskedXString-class), 38
coerce, character, RNAMultipleAlignment-method	
(MultipleAlignment-class), 70	(MaskedXString-class), 38
coerce, character, RNAString-method	coerce, MaskedXString, MaskCollection-method
(XString-class), 108	(MaskedXString-class), 38
coerce, character, SolexaQuality-method	coerce, MaskedXString, MaskedAAString-method
(XStringQuality-class), 110	(MaskedXString-class), 38
coerce, character, XString-method	coerce, MaskedXString, MaskedBString-method
(XString-class), 108	(MaskedXString-class), 38
coerce, DNAString, MaskedDNAString-method (MaskedXString-class), 38	coerce, MaskedXString, MaskedDNAString-method
	(MaskedXString-class), 38
coerce, integer, IlluminaQuality-method	coerce, MaskedXString, MaskedRNAString-method
(XStringQuality-class), 110	(MaskedXString-class), 38
coerce, integer, PhredQuality-method	coerce, MaskedXString, NormalIRanges-method
(XStringQuality-class), 110	(MaskedXString-class), 38
coerce, integer, SolexaQuality-method	coerce, MaskedXString, Views-method
(XStringQuality-class), 110	(MaskedXString-class), 38

coerce, Masked XString, XString Views-method	coerce,XStringViews,AAStringSet-method
(MaskedXString-class), 38	(XStringViews-class), 130
<pre>coerce,MIndex,CompressedIRangesList-method</pre>	<pre>coerce,XStringViews,BStringSet-method</pre>
(MIndex-class), 67	(XStringViews-class), 130
<pre>coerce,MultipleAlignment,AAStringSet-method</pre>	<pre>coerce,XStringViews,DNAStringSet-method</pre>
(MultipleAlignment-class), 70	(XStringViews-class), 130
<pre>coerce,MultipleAlignment,BStringSet-method</pre>	<pre>coerce,XStringViews,RNAStringSet-method</pre>
(MultipleAlignment-class), 70	(XStringViews-class), 130
<pre>coerce,MultipleAlignment,DNAStringSet-method</pre>	<pre>coerce,XStringViews,XStringSet-method</pre>
(MultipleAlignment-class), 70	(XStringViews-class), 130
<pre>coerce,MultipleAlignment,RNAStringSet-method</pre>	collapse, MaskedXString-method
(MultipleAlignment-class), 70	(MaskedXString-class), 38
coerce, numeric, IlluminaQuality-method	<pre>colmask (MultipleAlignment-class), 70</pre>
(XStringQuality-class), 110	colmask, MultipleAlignment-method
coerce, numeric, PhredQuality-method	(MultipleAlignment-class), 70
(XStringQuality-class), 110	<pre>colmask&lt;- (MultipleAlignment-class), 70</pre>
coerce, numeric, SolexaQuality-method	<pre>colmask&lt;-,MultipleAlignment,ANY-method</pre>
(XStringQuality-class), 110	(MultipleAlignment-class), 70
coerce, NumericList, IlluminaQuality-method	<pre>colmask&lt;-,MultipleAlignment,NULL-method</pre>
(XStringQuality-class), 110	(MultipleAlignment-class), 70
coerce, NumericList, PhredQuality-method	coloring, 8
(XStringQuality-class), 110	compact, 109, 115
coerce, NumericList, SolexaQuality-method	complement (reverseComplement), 95
(XStringQuality-class), 110	complement, DNAString-method
coerce, RNAString, MaskedRNAString-method	(reverseComplement), 95
(MaskedXString-class), 38	complement, DNAStringSet-method
coerce, XString, AAString-method	(reverseComplement), 95
(XString-class), 108	complement, MaskedDNAString-method
coerce, XString, BString-method	(reverseComplement), 95
(XString-class), 108	complement, MaskedRNAString-method
coerce, XString, DNAString-method	(reverseComplement), 95
(XString-class), 108	complement, RNAString-method
coerce, XString, RNAString-method	(reverseComplement), 95
(XString-class), 108	complement, RNAStringSet-method
coerce, XStringQuality, CompressedIntegerList-	
(XStringQuality-class), 110	complement, XStringViews-method
coerce, XStringQuality, CompressedNumericList-	
(XStringQuality-class), 110	CompressedIRangesList, 68
coerce, XStringQuality, IntegerList-method	computeAllFlinks(PDict-class), 81
(XStringQuality-class), 110	computeAllFlinks,ACtree2-method
coerce, XStringQuality, matrix-method	(PDict-class), 81
(XStringQuality-class), 110	connection, 122
coerce, XStringQuality, NumericList-method	consensusMatrix, 65, 67
(XStringQuality-class), 110	consensusMatrix (letterFrequency), 26
coerce, XStringSet, Views-method	consensusMatrix, character-method
(XStringViews-class), 130	(letterFrequency), 26
coerce, XStringSet, XStringViews-method	consensusMatrix, matrix-method
(XStringViews-class), 130	(letterFrequency), 26

consensusMatrix, MultipleAlignment-method	countPattern, XStringViews-method
(MultipleAlignment-class), 70	(matchPattern), 46
consensusMatrix,XStringSet-method	countPDict, 29
(letterFrequency), 26	countPDict (matchPDict), 50
consensusMatrix,XStringViews-method	countPDict, MaskedXString-method
(letterFrequency), 26	(matchPDict), 50
consensusString, 72	countPDict, XString-method (matchPDict),
consensusString (letterFrequency), 26	50
consensusString,AAMultipleAlignment-method	countPDict, XStringSet-method
(MultipleAlignment-class), 70	(matchPDict), 50
consensusString,ANY-method	countPDict, XStringViews-method
(1etterFrequency), 26	(matchPDict), 50
consensusString,BStringSet-method	countPWM (matchPWM), 65
(letterFrequency), 26	countPWM, character-method (matchPWM), 65
<pre>consensusString,DNAMultipleAlignment-method</pre>	countPWM, DNAString-method (matchPWM), 65
(MultipleAlignment-class), 70	countPWM,MaskedDNAString-method
consensusString,DNAStringSet-method	(matchPWM), 65
(letterFrequency), 26	countPWM,XStringViews-method
consensusString,matrix-method	(matchPWM), 65
(letterFrequency), 26	coverage, 29, 43
consensusString,MultipleAlignment-method	coverage,MaskedXString-method
(MultipleAlignment-class), 70	(match-utils), 42
$consensus {\tt String,RNAMultipleAlignment-method}$	coverage, MIndex-method (match-utils), 42
(MultipleAlignment-class), 70	
consensusString,RNAStringSet-method	detail, 10
(letterFrequency), 26	detail,MultipleAlignment-method
consensusString,XStringViews-method	(MultipleAlignment-class), 70
(letterFrequency), 26	dim,MultipleAlignment-method
consensusViews	(MultipleAlignment-class), 70
(MultipleAlignment-class), 70	dinucleotideFrequency
consensusViews, AAMultipleAlignment-method	(nucleotideFrequency), 74
(MultipleAlignment-class), 70	dinucleotideFrequencyTest, 11
consensusViews,DNAMultipleAlignment-method	dinucleotideFrequencyTest,DNAStringSet-method
(MultipleAlignment-class), 70	<pre>(dinucleotideFrequencyTest), 11</pre>
consensusViews,MultipleAlignment-method	dinucleotideFrequencyTest,RNAStringSet-method
(MultipleAlignment-class), 70	<pre>(dinucleotideFrequencyTest), 11</pre>
consensusViews,RNAMultipleAlignment-method	DNA_ALPHABET, 80, 84, 97, 113, 118, 123
(MultipleAlignment-class), 70	DNA_ALPHABET (DNAString-class), 12
countPattern, 51	DNA_BASES (DNAString-class), 12
countPattern (matchPattern), 46	DNAMultipleAlignment
countPattern,character-method	(MultipleAlignment-class), 70
(matchPattern), 46	DNAMultipleAlignment-class
countPattern,MaskedXString-method	(MultipleAlignment-class), 70
(matchPattern), 46	DNAString, 4, 18, 24, 34, 45, 51, 63, 65,
countPattern,XString-method	81–83, 93–95, 97, 98, 100–102, 108,
(matchPattern), 46	109, 112, 118, 131
countPattern,XStringSet-method	DNAString (DNAString-class), 12
(matchPattern), 46	DNAString-class, 12, 15, 67, 95, 109, 111

DNAStringSet, 11, 13, 51, 80-83, 86, 87,	fastq.geometry(XStringSet-io), 121
93–95, 99, 101–103, 115, 118, 123,	fastq.seqlengths(XStringSet-io), 121
125, 129	findPalindromes, 14, 45, 64, 95
DNAStringSet (XStringSet-class), 112	findPalindromes, DNAString-method
DNAStringSet-class, 53, 84, 95, 99	(findPalindromes), 14
DNAStringSet-class (XStringSet-class),	findPalindromes,MaskedXString-method (findPalindromes), 14
DNAStringSetList	findPalindromes,RNAString-method
(XStringSetList-class), 128	(findPalindromes), 14
DNAStringSetList-class	findPalindromes, XString-method
(XStringSetList-class), 128	(findPalindromes), 14
duplicated, 120	findPalindromes, XStringViews-method
duplicated, 720 duplicated, PDict-method (PDict-class),	(findPalindromes), 14
81	
duplicated, PreprocessedTB-method	gaps, <i>131</i>
(PDict-class), 81	gaps,MaskedXString-method
Dups, <i>124</i>	(MaskedXString-class), 38
	GENETIC_CODE, 5, 16, 76, 101–103
elementNROWS, 129	GENETIC_CODE_TABLE (GENETIC_CODE), 16
elementNROWS,MIndex-method	<pre>get_seqtype_conversion_lookup</pre>
(MIndex-class), 67	(Biostrings internals), 6
encoding (XStringQuality-class), 110	<pre>getGeneticCode (GENETIC_CODE), 16</pre>
encoding,XStringQuality-method	getSeq, 19, 99
(XStringQuality-class), 110	getSeq,BSgenome-method,20
endIndex (MIndex-class), 67	gregexpr, 21
<pre>endIndex,ByPos_MIndex-method</pre>	gregexpr2, 20
(MIndex-class), 67	
<pre>Expanded_TB_PDict (PDict-class), 81</pre>	hasAllFlinks (PDict-class), 81
<pre>Expanded_TB_PDict-class (PDict-class),</pre>	hasAllFlinks,ACtree2-method
81	(PDict-class), 81
<pre>extract_character_from_XString_by_positions</pre>	hasLetterAt, 76
(Biostrings internals), 6	hasLetterAt (lowlevel-matching), 33
<pre>extract_character_from_XString_by_positions,</pre>	xsasqnbysasabetters (letterFrequency), 26
(Biostrings internals), 6	hasOnlyBaseLetters,DNAString-method
extract_character_from_XString_by_ranges	(letterFrequency), 26
(Biostrings internals), 6	hasOnlyBaseLetters,DNAStringSet-method
<pre>extract_character_from_XString_by_ranges,XSt</pre>	ring-methodetterFrequency), 26
(Biostrings internals), 6	hasOnlyBaseLetters,MaskedXString-method
extractAllMatches (matchPDict), 50	(letterFrequency), 26
extractAt, 80	hasOnlyBaseLetters,RNAString-method
extractAt (replaceAt), 88	(letterFrequency), 26
extractAt, XString-method (replaceAt), 88	hasOnlyBaseLetters,RNAStringSet-method
extractAt, XStringSet-method	(letterFrequency), 26
(replaceAt), 88	hasOnlyBaseLetters,XStringViews-method
extractList, 89, 90	(letterFrequency), 26
extractTranscriptSeqs, 103	head, PDict3Parts-method (PDict-class),
	81
<pre>fasta.index (XStringSet-io), 121</pre>	head, TB_PDict-method (PDict-class), 81
fasta.seqlengths(XStringSet-io), 121	HNF4alpha, 21

IlluminaQuality, 86	lcprefix(lcsuffix), 25
IlluminaQuality (XStringQuality-class),	lcprefix,character,character-method
110	(lcsuffix), 25
IlluminaQuality-class	lcprefix,character,XString-method
(XStringQuality-class), 110	(lcsuffix), 25
initialize, ACtree2-method	lcprefix,XString,character-method
(PDict-class), 81	(lcsuffix), 25
initialize, PreprocessedTB-method	lcprefix,XString,XString-method
(PDict-class), 81	(lcsuffix), 25
<pre>initialize, Twobit-method (PDict-class),</pre>	lcsuffix, 25
81	lcsuffix,character,character-method
initialize, XStringCodec-method	(lcsuffix), 25
(Biostrings internals), 6	lcsuffix,character,XString-method
injectHardMask, 22, 39, 94	(lcsuffix), 25
injectHardMask,MaskedXString-method	lcsuffix,XString,character-method
(injectHardMask), 22	(lcsuffix), 25
injectHardMask,XStringViews-method	lcsuffix,XString,XString-method
(injectHardMask), 22	(lcsuffix), 25
injectSNPs, 94	length, MaskedXString-method
IntegerList, 89, 90	(MaskedXString-class), 38
IntegerRanges, 68, 80, 89, 90	length, MIndex-method (MIndex-class), 67
IntegerRangesList, 68, 89, 90	length, PDict-method (PDict-class), 81
IRanges, 68, 114	length, PDict3Parts-method
IRanges-class, 68	(PDict-class), 81
is.na,XStringSet-method	length, PreprocessedTB-method
(XStringSet-comparison), 118	(PDict-class), 81
is.unsorted, 120	length, SparseList-method (Biostrings
isMatchingAt, $51-53$	internals), 6
isMatchingAt (lowlevel-matching), 33	letter, 5, 13, 25, 98, 109, 131
isMatchingEndingAt (lowlevel-matching), 33	letter, character-method (letter), 25
33	letter, MaskedXString-method (letter), 25
isMatchingEndingAt,character-method	letter, XString-method (letter), 25
(lowlevel-matching), 33	letter, XString Views-method (letter), 25
isMatchingEndingAt,XString-method	letterFrequency, 26
(lowlevel-matching), 33	letterFrequency, MaskedXString-method
isMatchingEndingAt,XStringSet-method	(letterFrequency), 26
(lowlevel-matching), 33	letterFrequency, XString-method
isMatchingStartingAt, 105	(letterFrequency), 26
isMatchingStartingAt	letterFrequency, XStringSet-method
(lowlevel-matching), 33	(letterFrequency), 26
	letterFrequency, XStringViews-method
<pre>isMatchingStartingAt,character-method     (lowlevel-matching), 33</pre>	(letterFrequency), 26
isMatchingStartingAt,XString-method	letterFrequencyInSlidingView
(lowlevel-matching), 33	(letterFrequency), 26
isMatchingStartingAt,XStringSet-method	<pre>letterFrequencyInSlidingView,XString-method</pre>
(lowlevel-matching), 33	List, 129
IUPAC_CODE_MAP, 7, 13, 23, 34, 36, 45, 84, 94, 95, 98	List, 129 List-class, 129
75, 70	LI31 CIG33, 147

longestConsecutive, 32	MaskedXString, 6, 7, 22, 23, 25–27, 29, 41,
lowlevel-matching, 33, 43, 48, 106	44, 46, 51, 74, 76, 103
	MaskedXString (MaskedXString-class), 38
make_style, 8	MaskedXString-class, 23, 26, 29, 38, 41, 45,
<pre>make_XString_from_string (Biostrings</pre>	73, 76, 95
internals), 6	maskGaps (MultipleAlignment-class), 70
<pre>make_XString_from_string,XString-method</pre>	maskGaps,MultipleAlignment-method
(Biostrings internals), 6	(MultipleAlignment-class), 70
<pre>make_XStringSet_from_strings</pre>	maskMotif, <i>15</i> , <i>23</i> , <i>39</i> , 40, 48
(Biostrings internals), 6	maskMotif, MaskedXString, character-method
${\tt make\_XStringSet\_from\_strings}, {\tt XStringSet\_methology} \\$	od (maskMotif), 40
(Biostrings internals), 6	maskMotif, MaskedXString, XString-method
mask (maskMotif), 40	(maskMotif), 40
MaskCollection, 38	maskMotif,MultipleAlignment,ANY-method
MaskCollection-class, 39, 41	(MultipleAlignment-class), 70
MaskedAAString, 22	maskMotif,XString,ANY-method
MaskedAAString (MaskedXString-class), 38	(maskMotif), 40
MaskedAAString-class	masks (MaskedXString-class), 38
(MaskedXString-class), 38	masks, MaskedXString-method
MaskedBString, 22	(MaskedXString-class), 38
MaskedBString (MaskedXString-class), 38	masks, XString-method
MaskedBString-class	(MaskedXString-class), 38
(MaskedXString-class), 38	masks<- (MaskedXString-class), 38
maskeddim (MultipleAlignment-class), 70	masks<-, MaskedXString, MaskCollection-method
maskeddim, MultipleAlignment-method	(MaskedXString-class), 38
(MultipleAlignment-class), 70	
MaskedDNAString, 22, 65, 95, 101, 102	masks<-, MaskedXString, NULL-method
MaskedDNAString (MaskedXString-class),	(MaskedXString-class), 38
38	masks<-, XString, ANY-method
MaskedDNAString-class, 53	(MaskedXString-class), 38
MaskedDNAString-class	masks<-,XString,NULL-method
(MaskedXString-class), 38	(MaskedXString-class), 38
maskedncol (MultipleAlignment-class), 70	match, 120
maskedncol, MultipleAlignment-method	match, ANY, XStringSet-method
(MultipleAlignment-class), 70	(XStringSet-comparison), 118
maskednrow (MultipleAlignment-class), 70	match, XStringSet, ANY-method
maskednrow, MultipleAlignment-method	(XStringSet-comparison), 118
(MultipleAlignment-class), 70	match, XStringSet, XStringSet-method
maskedratio, MaskedXString-method	(XStringSet-comparison), 118
(MaskedXString-class), 38	match-utils, 42
maskedratio, MultipleAlignment-method	matchLRPatterns, 15, 36, 44, 48, 64, 106
(MultipleAlignment-class), 70	matchLRPatterns, MaskedXString-method
MaskedRNAString, 22, 95, 101, 102	(matchLRPatterns), 44
<pre>MaskedRNAString (MaskedXString-class),</pre>	matchLRPatterns,XString-method
38	(matchLRPatterns), 44
MaskedRNAString-class	matchLRPatterns,XStringViews-method
(MaskedXString-class), 38	(matchLRPatterns), 44
maskedwidth, MaskedXString-method	matchPattern, 15, 21, 25, 36, 41, 43–45, 46,
(MaskedXString-class), 38	51–53, 64, 67, 105, 106

matchPattern,character-method	misc, 69
(matchPattern), 46	mismatch, 48
matchPattern,MaskedXString-method	mismatch (match-utils), 42
(matchPattern), 46	mismatch, ANY, XStringViews-method
matchPattern,XString-method	(match-utils), 42
(matchPattern), 46	mkAllStrings(nucleotideFrequency),74
matchPattern,XStringSet-method	MTB_PDict(PDict-class), 81
(matchPattern), 46	MTB_PDict-class (PDict-class), 81
matchPattern,XStringViews-method	MultipleAlignment, $70$
(matchPattern), 46	MultipleAlignment
matchPDict, 36, 43, 47, 48, 50, 59, 60, 68, 81,	(MultipleAlignment-class), 70
84	MultipleAlignment-class, 70
matchPDict,MaskedXString-method	
(matchPDict), 50	N50 (misc), 69
<pre>matchPDict,XString-method(matchPDict),</pre>	names, MIndex-method (MIndex-class), 67
50	names, PDict-method (PDict-class), 81
matchPDict,XStringSet-method	<pre>names&lt;-,MIndex-method(MIndex-class), 67</pre>
(matchPDict), 50	<pre>names&lt;-,PDict-method(PDict-class), 81</pre>
matchPDict, XStringViews-method	narrow, <i>70</i> , <i>112</i> , <i>115</i>
(matchPDict), 50	narrow,character-method
matchPDict-exact (matchPDict), 50	(XStringSet-class), 112
matchPDict-inexact, 53, 59	narrow,QualityScaledXStringSet-method
matchProbePair, 15, 45, 48, 63	(QualityScaledXStringSet-class)
matchProbePair, DNAString-method	86
(matchProbePair), 63	nchar, 29
matchProbePair, MaskedDNAString-method	nchar,MaskedXString-method
(matchProbePair), 63	(MaskedXString-class), 38
matchProbePair,XStringViews-method	nchar, MultipleAlignment-method
(matchProbePair), 63	(MultipleAlignment-class), 70
matchPWM, 65	nchar, XString-method (XString-class),
matchPWM, character-method (matchPWM), 65	108
matchPWM, DNAString-method (matchPWM), 65	nchar, XStringSet-method
matchPWM, MaskedDNAString-method	(XStringSet-class), 112
(matchPWM), 65	nchar,XStringSetList-method
matchPWM,XStringViews-method	(XStringSetList-class), 128
(matchPWM), 65	nchar,XStringViews-method
maxScore (matchPWM), 65	(XStringViews-class), 130
maxScore, ANY-method (matchPWM), 65	ncol, MultipleAlignment-method
maxWeights (matchPWM), 65	(MultipleAlignment-class), 70
maxWeights, matrix-method (matchPWM), 65	neditAt (lowlevel-matching), 33
mergeIUPACLetters (IUPAC_CODE_MAP), 23	neditEndingAt, 105
MIndex, 43, 47, 48, 52	neditEndingAt (lowlevel-matching), 33
MIndex (MIndex-class), 67	neditEndingAt, character-method
	(lowlevel-matching), 33
MIndex-class, 43, 53, 60, 67, 131	neditEndingAt,XString-method
minScore (matchPWM), 65	(lowlevel-matching), 33
minScore, ANY-method (matchPWM), 65	neditEndingAt,XStringSet-method
minWeights (matchPWM), 65	(lowlevel-matching), 33
minWeights, matrix-method (matchPWM), 65	neditStartingAt, <i>105</i>

neditStartingAt (lowlevel-matching), 33	palindromeArmLength,RNAString-method
neditStartingAt,character-method	(findPalindromes), 14
(lowlevel-matching), 33	palindromeArmLength,XString-method
neditStartingAt,XString-method	(findPalindromes), 14
(lowlevel-matching), 33	palindromeArmLength,XStringSet-method
neditStartingAt,XStringSet-method	(findPalindromes), 14
(lowlevel-matching), 33	palindromeArmLength,XStringViews-method
nmatch (match-utils), 42	(findPalindromes), 14
nmatch, ANY, XStringViews-method	palindromeLeftArm (findPalindromes), 14
(match-utils), 42	palindromeLeftArm, XString-method
nmismatch (match-utils), 42	(findPalindromes), 14
nmismatch, ANY, XStringViews-method	palindromeLeftArm,XStringViews-method
(match-utils), 42	(findPalindromes), 14
nnodes (PDict-class), 81	palindromeRightArm (findPalindromes), 14
nnodes, ACtree2-method (PDict-class), 81	palindromeRightArm, XString-method
NormalIRanges, 71	(findPalindromes), 14
nrow,MultipleAlignment-method	palindromeRightArm, XStringViews-method
(MultipleAlignment-class), 70	(findPalindromes), 14
nucleotideFrequency, 74	PAM120 (predefined_scoring_matrices), 85
nucleotideFrequencyAt, 12, 36	PAM250 (predefined_scoring_matrices), 85
nucleotideFrequencyAt	PAM30 (predefined_scoring_matrices), 85
(nucleotideFrequency), 74	PAM40 (predefined_scoring_matrices), 85
nucleotideFrequencyAt,XStringSet-method	PAM70 (predefined_scoring_matrices), 85
(nucleotideFrequency), 74	parallel_slot_names,QualityScaledXStringSet-method
nucleotideFrequencyAt,XStringViews-method	(QualityScaledXStringSet-class),
(nucleotideFrequency), 74	86
	PartitioningByEnd, 129
oligonucleotideFrequency, 29	paste, <i>107</i>
oligonucleotideFrequency	patternFrequency (PDict-class), 81
(nucleotideFrequency), 74	
oligonucleotideFrequency, MaskedXString-metho	(PDict-class), 81
(nucleotideFrequency), 74	pcompare, ANY, XStringSet-method
oligonucleotideFrequency,XString-method	(XStringSet-comparison), 118
(nucleotideFrequency), 74	pcompare, XStringSet, ANY-method
oligonucleotideFrequency,XStringSet-method	(XStringSet-comparison), 118
(nucleotideFrequency), 74	pcompare, XStringSet, XStringSet-method
oligonucleotideFrequency,XStringViews-method	(XStringSet-comparison), 118
(nucleotideFrequency), 74	PDict, 51, 52, 59, 60
oligonucleotideTransitions	PDict (PDict-class), 81
(nucleotideFrequency), 74	
order, <i>120</i>	PDict, AsIs-method (PDict-class), 81
	PDict, character-method (PDict-class), 81
padAndClip, 79, 90	PDict, DNAStringSet-method
pairwiseAlignment, 47, 48, 111	(PDict-class), 81
PairwiseAlignments-class, 111	PDict, probetable-method (PDict-class),
palindromeArmLength(findPalindromes),	81
14	PDict, XStringViews-method
palindromeArmLength, DNAString-method	(PDict-class), 81
(findPalindromes), 14	PDict-class, 53, 60, 68, 81

PDict3Parts (PDict-class), 81	QualityScaledXStringSet, 123, 125
PDict3Parts-class (PDict-class), 81	QualityScaledXStringSet
PhredQuality, 86	(QualityScaledXStringSet-class),
PhredQuality (XStringQuality-class), 110	86
PhredQuality-class	QualityScaledXStringSet-class, 86
(XStringQuality-class), 110	
<pre>predefined_scoring_matrices, 85, 85</pre>	rank, <i>120</i>
PreprocessedTB (PDict-class), 81	read.Mask, <i>41</i>
PreprocessedTB-class (PDict-class), 81	readAAMultipleAlignment
<pre>print.moved_to_pwalign_pkg</pre>	(MultipleAlignment-class), 70
<pre>(predefined_scoring_matrices),</pre>	readAAStringSet (XStringSet-io), 121
85	readBStringSet (XStringSet-io), 121
PWM (matchPWM), 65	readDNAMultipleAlignment
PWM, character-method (matchPWM), 65	(MultipleAlignment-class), 70
PWM, DNAStringSet-method (matchPWM), 65	readDNAStringSet, 87, 115
PWM, matrix-method (matchPWM), 65	_
PWMscoreStartingAt (matchPWM), 65	readDNAStringSet (XStringSet-io), 121
	readQualityScaledDNAStringSet, 123, 125
quality	readQualityScaledDNAStringSet
(QualityScaledXStringSet-class),	(QualityScaledXStringSet-class),
86	86
quality,QualityScaledXStringSet-method	readRNAMultipleAlignment
(QualityScaledXStringSet-class),	(MultipleAlignment-class), 70
86	readRNAStringSet (XStringSet-io), 121
QualityScaledAAStringSet	relistToClass,XStringSet-method
(QualityScaledXStringSet-class),	(XStringSetList-class), 128
86	replaceAmbiguities (chartr), 6
QualityScaledAAStringSet-class	replaceAt, 7, 80, 88, 94
(QualityScaledXStringSet-class),	replaceAt, XString-method (replaceAt), 88
86	replaceAt,XStringSet-method
QualityScaledBStringSet	(replaceAt), 88
(QualityScaledXStringSet-class),	replaceLetterAt, 7, 23, 90, 92
86	replaceLetterAt,DNAString-method
QualityScaledBStringSet-class	(replaceLetterAt), 92
(QualityScaledXStringSet-class),	replaceLetterAt,DNAStringSet-method
86	(replaceLetterAt), 92
QualityScaledDNAStringSet, 123, 125	rev, 76
QualityScaledDNAStringSet	reverse, 95
(QualityScaledXStringSet-class),	reverse,MaskedXString-method
86	(reverseComplement), 95
QualityScaledDNAStringSet-class	reverse,QualityScaledXStringSet-method
(QualityScaledXStringSet-class),	$({\tt QualityScaledXStringSet-class}),$
86	86
QualityScaledRNAStringSet	reverseComplement, 13, 39, 45, 64, 67, 76,
(QualityScaledXStringSet-class),	95, 98, 103, 109
86	${\tt reverseComplement,DNAString-method}$
QualityScaledRNAStringSet-class	(reverseComplement), 95
(QualityScaledXStringSet-class),	$reverse {\tt Complement}, {\tt DNAStringSet-method}$
86	(reverseComplement), 95

(MultipleAlignment-class), 70
rownames, MultipleAlignment-method
(MultipleAlignment-class), 70
rownames<-,MultipleAlignment-method
(MultipleAlignment-class), 70
meare XStringSet (XStringSet-io), 121
seqinfo, 99
seqinfo(seqinfo-methods), 99
meaninfo, DNAStringSet-method
(Seq11110-lilethous), 99
seqinfo-methods, 99
seqinfo<-,DNAStringSet-method
(seqinfo-methods), 99
seqtype (Biostrings internals), 6
seqtype,AAString-method
(XString-class), 108
seqtype, BString-method (XString-class)
108
seqtype,DNAString-method
(XString-class), 108
seqtype,MaskedXString-method
(MaskedXString-class), 38
seqtype,MultipleAlignment-method
(MultipleAlignment-class), 70
seqtype,RNAString-method
(XString-class), 108
seqtype,XStringSet-method
(XStringSet-class), 112
seqtype,XStringSetList-method
(XStringSetList-class), 128
seqtype,XStringViews-method
(XStringViews-class), 130
seqtype<- (Biostrings internals), 6
seqtype<-,MaskedXString-method
(MaskedXString-class), 38
seqtype<-,XString-method
(XString-class), 108
seqtype<-,XStringSet-method
(XStringSet-class), 112
seqtype<-,XStringSetList-method
(XStringSetList-class), 128
seqtype<-,XStringViews-method
(XStringViews-class), 130
show, <i>10</i>
show, ACtree2-method (PDict-class), 81
show,MaskedXString-method
(MaskedXString-class), 38
show, MIndex-method (MIndex-class), 67

show, MTB_PDict-method (PDict-class), 81	substring, 29
show, MultipleAlignment-method	substring,XString-method
(MultipleAlignment-class), 70	(XString-class), 108
show,QualityScaledXStringSet-method	, , , , , , , , , , , , , , , , , , , ,
$({\tt QualityScaledXStringSet-class}),$	tail, PDict3Parts-method (PDict-class), 81
86	tail, TB_PDict-method (PDict-class), 81
show, TB_PDict-method (PDict-class), 81	tb (PDict-class), 81
show, Twobit-method (PDict-class), 81	tb,PDict3Parts-method(PDict-class), 81
show, XString-method (XString-class), 108	
show,XStringSet-method	tb, PreprocessedTB-method (PDict-class), 81
(XStringSet-class), 112	tb,TB_PDict-method(PDict-class),81
show,XStringSetList-method	· · · · · · · · · · · · · · · · · · ·
(XStringSetList-class), 128	tb.width (PDict-class), 81
show,XStringViews-method	tb.width,PDict3Parts-method
(XStringViews-class), 130	(PDict-class), 81
showAsCell,XString-method	tb.width,PreprocessedTB-method
(XString-class), 108	(PDict-class), 81
showAsCell,XStringSet-method	<pre>tb.width,TB_PDict-method(PDict-class),</pre>
(XStringSet-class), 112	81
showAsCell,XStringSetList-method	TB_PDict (PDict-class), 81
(XStringSetList-class), 128	TB_PDict-class (PDict-class), 81
SolexaQuality, 86	threebands, 114
SolexaQuality, 80 SolexaQuality (XStringQuality-class),	threebands,character-method
110	(XStringSet-class), 112
	toComplex, 100
SolexaQuality-class	toComplex, DNAString-method (toComplex),
(XStringQuality-class), 110	100
sort, <i>120</i>	toString,MaskedXString-method
SparseList (Biostrings internals), 6	(MaskedXString-class), 38
SparseList-class (Biostrings	toString,XString-method
internals), 6	(XString-class), 108
stackStrings(padAndClip), 79	toString,XStringSet-method
stackStringsFromBam, $80$	(XStringSet-class), 112
startIndex (MIndex-class), 67	toString,XStringViews-method
startIndex,ByPos_MIndex-method	(XStringViews-class), 130
(MIndex-class), 67	translate, 18, 101
strsplit, 29	translate, DNAString-method (translate),
subseq, 26, 89, 90, 109, 113, 115	101
subseq, character-method	translate,DNAStringSet-method
(XStringSet-class), 112	(translate), 101
subseq, MaskedXString-method	translate, MaskedDNAString-method
(MaskedXString-class), 38	(translate), 101
subseq<-,character-method	
(XStringSet-class), 112	translate, MaskedRNAString-method
	(translate), 101
subseq<-,XStringSet-method	translate, RNAString-method (translate),
(XStringSet-class), 112	101
substr, 113, 115	translate,RNAStringSet-method
substr,XString-method(XString-class),	(translate), 101
108	trimLRPatterns. 36, 45, 104

trimLRPatterns, character-method	vcountPattern, 51
(trimLRPatterns), 104	vcountPattern (matchPattern), 46
trimLRPatterns, XString-method	vcountPattern, character-method
(trimLRPatterns), 104	(matchPattern), 46
trimLRPatterns, XStringSet-method	vcountPattern, MaskedXString-method
(trimLRPatterns), 104	(matchPattern), 46
trinucleotideFrequency, 18	vcountPattern, XString-method
trinucleotideFrequency	(matchPattern), 46
(nucleotideFrequency), 74	vcountPattern, XStringSet-method
Twobit (PDict-class), 81	(matchPattern), 46
Twobit-class (PDict-class), 81	vcountPattern, XStringViews-method
//	(matchPattern), 46
unique, <i>120</i>	vcountPDict (matchPDict), 50
uniqueLetters, $7$	vcountPDict, MaskedXString-method
uniqueLetters (letterFrequency), 26	(matchPDict), 50
uniqueLetters,MaskedXString-method	vcountPDict, XString-method
(letterFrequency), 26	(matchPDict), 50
uniqueLetters,XString-method	vcountPDict, XStringSet-method
(1etterFrequency), 26	(matchPDict), 50
uniqueLetters,XStringSet-method	vcountPDict, XStringViews-method
(1etterFrequency), 26	(matchPDict), 50
uniqueLetters,XStringViews-method	Views, 39, 65, 130, 131
(1etterFrequency), 26	Views, character-method
unitScale (matchPWM), 65	(XStringViews-class), 130
unlist, MIndex-method (MIndex-class), 67	Views, MaskedXString-method
unlist,XStringSet-method	(MaskedXString-class), 38
(XStringSet-class), 112	Views, XString-method
unmasked, 29	(XStringViews-class), 130
unmasked (MaskedXString-class), 38	Views-class, 39, 131
unmasked, MaskedXString-method	vmatchPattern, 51
(MaskedXString-class), 38	vmatchPattern (matchPattern), 46
unmasked, MultipleAlignment-method	vmatchPattern,character-method
(MultipleAlignment-class), 70	(matchPattern), 46
unmasked, XString-method	vmatchPattern,MaskedXString-method
(MaskedXString-class), 38	(matchPattern), 46
unstrsplit, 90	vmatchPattern, XString-method
update_AA_palette (coloring), 8	(matchPattern), 46
update_B_palette (coloring), 8	vmatchPattern, XStringSet-method
update_DNA_palette (coloring), 8	(matchPattern), 46
update_RNA_palette (coloring), 8	vmatchPattern, XStringViews-method
update_X_palette (coloring), 8	(matchPattern), 46
updateObject, AAString-method	vmatchPDict (matchPDict), 50
(XString-class), 108	vmatchPDict, ANY-method (matchPDict), 50
updateObject, AAStringSet-method	
(XStringSet-class), 112	vmatchPDict, MaskedXString-method
updateObject,XString-method	<pre>(matchPDict), 50 vmatchPDict, XString-method</pre>
(XString-class), 108	(matchPDict, AString-method
updateObject, XStringSet-method	
(XStringSet-class), 112	vwhichPDict (matchPDict), 50

<pre>vwhichPDict,MaskedXString-method</pre>	<pre>windows,QualityScaledXStringSet-method      (QualityScaledXStringSet-class),</pre>
vwhichPDict,XString-method	86
(matchPDict), 50	write.phylip (MultipleAlignment-class),
vwhichPDict,XStringSet-method	70
<pre>(matchPDict), 50 vwhichPDict,XStringViews-method</pre>	writeQualityScaledXStringSet, 125 writeQualityScaledXStringSet
(matchPDict), 50	(QualityScaledXStringSet-class),
which.isMatchingAt(lowlevel-matching),	writeXStringSet, 87, 115
33	<pre>writeXStringSet(XStringSet-io), 121</pre>
which.isMatchingEndingAt	
(lowlevel-matching), 33	xscat, 107
which.isMatchingEndingAt,character-method	xscodes (Biostrings internals), 6
(lowlevel-matching), 33	xscodes, ANY-method (Biostrings
which.isMatchingEndingAt,XString-method	internals), 6
(lowlevel-matching), 33	XString, 4, 6–8, 12–14, 20, 23, 25–27, 29, 34,
which.isMatchingEndingAt,XStringSet-method	35, 38, 44, 46, 51, 68, 70, 74, 76, 86,
(lowlevel-matching), 33	88–90, 94, 95, 97, 98, 105–107, 112,
which.isMatchingStartingAt	114, 115, 130, 131
(lowlevel-matching), 33	XString (XString-class), 108
which.isMatchingStartingAt,character-method	XString-class, 5, 9, 20, 25, 26, 29, 36, 39,
(lowlevel-matching), 33	41, 43, 45, 76, 106, 107, 108, 131
which.isMatchingStartingAt,XString-method	XStringCodec (Biostrings internals), 6
(lowlevel-matching), 33	XStringCodec-class (Biostrings
which.isMatchingStartingAt,XStringSet-method	internals), 6
(lowlevel-matching), 33	XStringQuality, 86, 87
whichPDict, 59	XStringQuality (XStringQuality-class),
whichPDict(matchPDict), 50	110
whichPDict,MaskedXString-method	XStringQuality-class, 110
(matchPDict), 50	XStringSet, 6, 7, 20, 26–29, 34, 35, 46, 51,
whichPDict,XString-method(matchPDict),	70–72, 74–76, 79, 80, 86–90, 95,
50	105–108, 115, 118, 119, 121, 124,
whichPDict,XStringSet-method	125, 128, 129
(matchPDict), 50	XStringSet (XStringSet-class), 112
whichPDict,XStringViews-method	XStringSet-class, 12, 20, 29, 69, 73, 76,
(matchPDict), 50	106, 107, 109, 112, 120, 129, 131
width, character-method	XStringSet-comparison, 115, 118
(XStringSet-class), 112	XStringSet-io, 121
width, PDict-method (PDict-class), 81	XStringSetList, 89, 90, 115, 129
width, PDict3Parts-method (PDict-class),	XStringSetList (XStringSetList-class),
81	128
width, PreprocessedTB-method	XStringSetList-class, 128
(PDict-class), 81	XStringViews, 6, 7, 14, 15, 22, 23, 25–29, 41,
width0 (MIndex-class), 67	43–48, 63, 64, 66, 68, 70, 72, 74–76,
width0, MIndex-method (MIndex-class), 67	80–82, 86, 95, 102, 103, 107,
windows, character-method	112–115, 129
(XStringSet-class), 112	XStringViews (XStringViews-class), 130

```
XStringViews-class, 15, 23, 25, 26, 29, 41,
43, 45, 53, 64, 67, 68, 76, 84, 95,
107, 109, 130

XVector, 114

XVector-class, 109

XVectorList, 115

yeastSEQCHR1, 133
```