

配布するUSBメモリ中のhogeフォルダを  
デスクトップにコピーしておいてください。

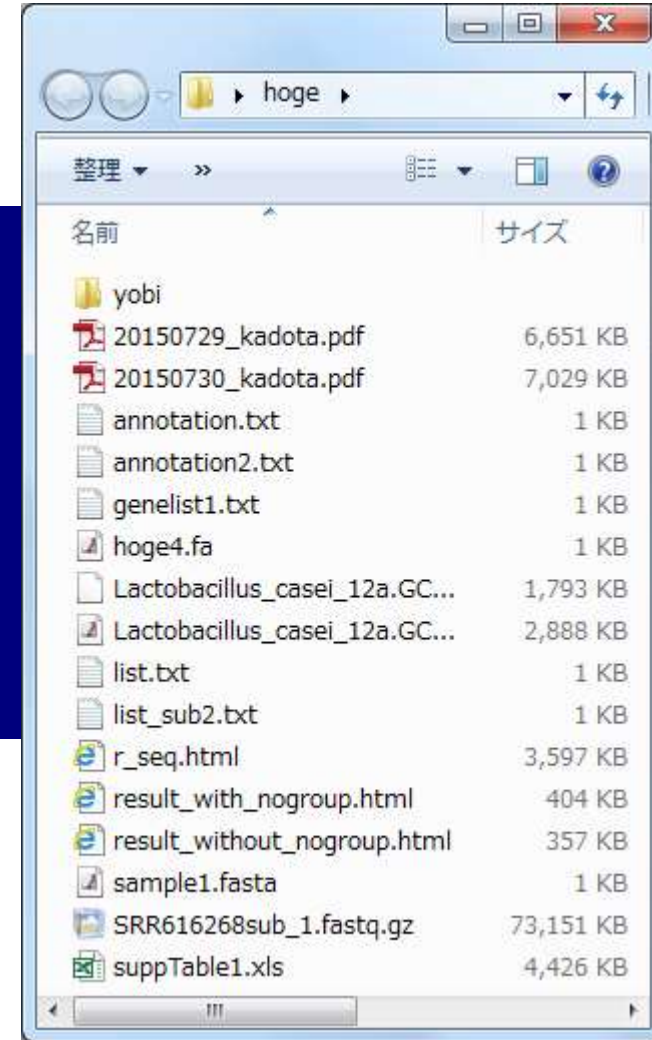
# NGSハンズオン講習会 Bioconductorの利用法

東京大学・大学院農学生命科学研究科  
アグリバイオインフォマティクス教育研究プログラム

門田幸二(かどた こうじ)

[kadota@iu.a.u-tokyo.ac.jp](mailto:kadota@iu.a.u-tokyo.ac.jp)

<http://www.iu.a.u-tokyo.ac.jp/~kadota/>



Rが導く未来  
は明るいのだ  
ろうか...



# Contents (全体)

- 7月22日(水): 84→83名。Bio-Linux 8とRのインストール状況確認。基本自習(門田・寺田先生)
- 7月23日(木): 92→90名。Linux基礎。LinuxコマンドなどUNIXの基礎の理解(門田)
- 7月24日(金): 85→83名。スクリプト言語。シェルスクリプト(アメリエフ株式会社 服部恵美先生)
- 7月27日(月): 93→91名。スクリプト言語。Perl(アメリエフ 服部先生)
- 7月28日(火): 91→90名。スクリプト言語。Python(アメリエフ 服部先生)
- 7月29日(水): 94→88名。データ解析環境R(門田)
  - R基礎1(初級): R言語の基礎(インストールから利用まで)
  - R基礎2(初級): ファイルの読み込み、行列演算の基本
  - R各種パッケージ(中級): パッケージのインストール法と代表的なパッケージの利用法
- 7月30日(木): 96→91名。データ解析環境R(門田)
  - Bioconductorの利用法1(中級): データの型やバージョンの違い
  - Bioconductorの利用法2(中級): FASTA/FASTQファイルの各種解析
- 8月3日(月): 89→84名。NGS解析。基礎(アメリエフ 山口昌雄先生)
- 8月4日(火): 85→80名。NGS解析。ゲノムReseq、変異解析(アメリエフ 山口先生)
- 8月5日(水): 86→81名。NGS解析。RNA-seq、統計解析(前半: 山口先生、後半: 門田)
- 8月6日(木): 104→98名。NGS解析。ChIP-seq(理研 森岡勝樹先生)

# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# パッケージ

R起動直後に「?関数名」と打ち込んでも、使用法を記したウェブページが開かずにエラーが出る場合があります。

```
R Console
R version 3.1.3 (2015-03-09) -- "Smooth Sidewalk"
Copyright (C) 2015 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R は、自由なソフトウェアであり、「完全に無保証」です。
一定の条件に従えば、自由にこれを再配布することができます。
配布条件の詳細に関しては、'license()' あるいは 'licence()' と入力し$

R は多くの貢献者による共同プロジェクトです。
詳しくは 'contributors()' と入力してください。
また、R や R のパッケージを出版物で引用する際の形式については
'citation()' と入力してください。

'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザによるヘルプがみられます。
'q()' と入力すれば R を終了します。

> ?subseq
No documentation for 'subseq' in specified packages and libraries:
you could try '??subseq'
> ?alphabetFrequency
No documentation for 'alphabetFrequency' in specified packages and$
you could try '??alphabetFrequency'
> |
```

# パッケージ

- ①「??alphabetFrequency」と打ち込むように勧められているので打ってみる。検索結果のウェブページが表示されるので、
- ②それっぽい関数名のところをクリック。

R Console

```
> ?subseq
No documentation for 'subseq' in specified packages and libraries:
you could try '??subseq'
> ?alphabetFrequency
No documentation for 'alphabetFrequency' in specified packages and$
you could try '??alphabetFrequency'
> ??alphabetFrequency
starting httpd help server ... done
> |
```



Search Results



The search string was "alphabetFrequency"

Help pages:

- |  |  |
|--|--|
| <a href="#">Biostrings::class:MultipleAlignment</a>    | MultipleAlignment objects  |
| <a href="#">Biostrings::letterFrequency</a>            | Calculate the frequency of letters in a biological sequence, or the consensus matrix of a set of sequences |
| <a href="#">GenomicAlignments::stackStringsFromBam</a> | Stack the read sequences stored in a BAM file on a region of interest                                      |
| <a href="#">ShortRead::QualityScore-class</a>          | Quality scores for short reads and their alignments  |



alphabetFrequency関数はBiostringsというパッケージから提供されているものと読み解く。「??関数名」は、関数名は既知だがどのパッケージから提供されているものかを知りたい場合などに利用する。

# パッケージ

letterFrequency {Biostrings} ←

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 tallies 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)
```

# パッケージ

multi-FASTAファイルを読み込んで様々な解析ができるのは、①Biostringsや②seqinrなどの塩基配列解析用パッケージのおかげです。

- インタロ | 一般 | [任意の位置の塩基を置換](#) (last modified 2013/09/12)
- インタロ | 一般 | [指定した範囲の配列を取得](#) (last modified 2015/04/06) **NEW**
- インタロ | 一般 | [指定したID\(染色体やdescription\)の配列を取得](#) (last modified 2014/03/10) **①**
- インタロ | 一般 | [翻訳配列\(translate\)を取得\(基礎\)](#) | [Biostrings](#) (last modified 2015/03/09)
- インタロ | 一般 | [翻訳配列\(translate\)を取得\(応用\)](#) | [seqinr\(Charif\\_2005\)](#) (last modified 2015/03/09) **②**
- インタロ | 一般 | [相補鎖\(complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)
- インタロ | 一般 | [逆相補鎖\(reverse complement\)を取得](#) (last modified 2013/06/14)

**インタロ | 一般 | 翻訳配列(translate)を取得(基礎) | Biostrings**

Biostringsパッケージを用いて塩基配列を読み込んでアミノ酸配列に翻訳するやり方を示します。翻訳のための [遺伝コード\(genetic code\)](#) は、Standard Genetic Codeだそうです。もちろん生物種!!によって多少違い(variants)があるようで、"Standard", "SGC0", "Vertebrate Mitochondrial", "SGC1"などいろいろ選べるようです。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

**1. FASTA形式ファイル(sample1.fasta)の場合:**

```
multi-FASTAではない single-FASTA形式ファイルです。

in_f <- "sample1.fasta"
out_f <- "hoge1.fasta"

#必要なパッケージをロード
library(Biostrings)

#入力ファイルの読み込み
fasta <- readDNAStringSet("sample1.fasta")
```

**インタロ | 一般 | 翻訳配列(translate)を取得(応用) | seqinr(Charif\_2005)**

seqinrパッケージを用いて塩基配列を読み込んでアミノ酸配列に翻訳するやり方を示します。本気で翻訳配列を取得する場合にはこちらの利用をお勧めします。翻訳できないコドン(不明なアミノ酸)に変換してくれたり、translate関数のオプションとしてambiguous=Tとすると、翻訳できるものは可能な限り翻訳してくれます(高橋 広夫 氏提供情報)。apply関数を用いるやり方(高橋 広夫 氏提供情報)とsapply関数を用いるやり方(甲斐 政親 氏提供情報)を示します。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

**1. FASTA形式ファイル(sample1.fasta)の場合:**

```
multi-FASTAではない single-FASTA形式ファイルです。

in_f <- "sample1.fasta"
out_f <- "hoge1.fasta"

#必要なパッケージをロード
library(seqinr)

#入力ファイルの読み込み
hoge <- read.fasta(in_f, seqtype="DNA")
hoge
```

#入力ファイル名を指定してin\_fに格納  
#出力ファイル名を指定してout\_fに格納  
#パッケージの読み込み  
#in\_fで指定したファイルの読み込み  
#確認してるだけです

# パッケージ

Biostringsというパッケージをlibrary関数を用いて読み込むことによって、alphabetFrequencyのようなBiostringsが提供する関数群を利用できるのです。ここでは、意図的に「library(Biostrings)」を2回実行して、2回目は何も表示されないということを思い出させています。実際には1回のみで大丈夫です。「?alp」まで打ってからTabキーを押すなどして「タブ補完」テクを有効利用。

```
R Console
> library(Biostrings)
要求されたパッケージ BiocGenerics をロード中です
要求されたパッケージ parallel をロード中です

次のパッケージを付け加えます: 'BiocGenerics'

The following objects are masked from 'package:parallel':
  clusterApply, clusterApplyLB, clusterCall, clusterExport, clusterMap, parApply, parCapply, parLapply, parLapplyLB, parRapply, parSapply, parSapplyLB

The following object is masked from 'package:stats':
  xtabs

The following objects are masked from 'package:base':
  anyDuplicated, append, as.data.frame, as.vector, colnames, do.call, duplicated, eval, evalq, Filter, get, intersect, is.unsorted, lapply, Map, mapply, mget, order, paste, pmax, pmax.int, pmin, pmin.int, Position, rank, rbind, Reduce, rep.int, rownames, setdiff, sort, table, tapply, union, unique, unlist, unsplit

要求されたパッケージ S4Vectors をロード中です
要求されたパッケージ stats4 をロード中です
要求されたパッケージ IRanges をロード中です
要求されたパッケージ XVector をロード中です
> library(Biostrings)
> ?alphabetFrequency
```

letterFrequency {Biostrings} R Documentation

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 tallies 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)
```



# R本体とパッケージの関係

「R本体」と「パッケージ」の関係は、「パソコン」と「ソフト」、「Microsoft EXCEL」と「アドイン」、「Cytoscape」と「プラグイン」のようなものという理解でよい。

- パソコンを購入しただけの状態では、できることが限られています。
  - 通常は、Officeやウイルス撃退ソフトなどをインストールして利用します。
- Linuxをインストールしただけの状態では、できることが限られています。
  - 通常は、マッピングなど各種プログラムをインストールして利用します。
- R本体をインストールしただけの状態では、できることが限られています。
  - NGS解析を行う各種パッケージ(またはライブラリ)をインストールして利用します。

# CRANとBioconductor

CRAN提供パッケージは生命科学を含む様々な分野で利用される。NGS解析は、主にBioconductor提供パッケージを利用。

## R上で利用可能なパッケージの2大リポジトリ(貯蔵庫)

- CRAN (The Comprehensive R Archive Network): 6,878パッケージ
- Bioconductor: 1,024パッケージ

- 作図 | ROC曲線 | 基礎編 | [7. 図の重ね書き\(new\)](#) (last modified 2015/02/15) NEW
- 作図 | ROC曲線 | 基礎編 | [8. 凡例を追加\(legend\)](#) (last modified 2015/02/15) NEW
- 作図 | ROC曲線 | 応用編 (last modified 2015/02/07) NEW
- 作図 | [SplicingGraphs](#) (last modified 2015/02/07) NEW
- [パイプライン](#) | [||](#)について (last modified 2015/02/07) NEW
- [パイプライン](#) | [ゲノム](#) | [発現変動](#) | 2群
- [パイプライン](#) | [ゲノム](#) | [機能解析](#) | 2群
- [パイプライン](#) | [ゲノム](#) | [機能解析](#) | 2群
- [パイプライン](#) | [ゲノム](#) | [small RNA](#) | [S](#)
- [リンク集](#) (last modified 2012/03/29)

## リンク集

- [R](#)
- [Bioconductor: Gentleman et al., Genome Biol., 2004](#)
- [CRAN](#)
- [RjpWiki](#)
- [R Tips](#)(竹澤様)
- [BioEdit](#)(フリーの配列編集ソフト)
- [BioMart: Smedley et al., BMC Genomics, 2009](#)
- [DDBJ Read Annotation Pipeline: Nagasaki et al., DNA Res., 2013](#)
- [EMBOSS explorer](#) (EMBOSSのウェブ版)
- [Biostar: Parnell et al., PLoS Comput Biol., 2011](#)
- [SEQanswers: Li et al., Bioinformatics, 2012](#)
- [NGS WikiBook: Li et al., Brief Bioinform., 2013](#)
- [HT Sequence Analysis with R and Bioconductor](#)

# 定期的にバージョンアップ

バグの修正や新たな機能がどんどん追加されている。最新版の利用をお勧め。毎年5月と11月ごろにバージョンアップするとよいだろう。

## ■ 近年のリリース頻度

### □ R本体 (<http://www.r-project.org/>)

- 2015-06-18にver. 3.2.1をリリース
- 2015-04-16にver. 3.2.0をリリース
- 2015-03-09にver. 3.1.3をリリース
- 2014-10-31にver. 3.1.2をリリース
- ...
- 2012-03-30にver. 2.15.0をリリース
- ...

### □ Bioconductor (<http://bioconductor.org/>)は半年ごとにリリース

- 2015-04にver. 3.1をリリース (R ver. 3.2.1で動作確認)、提供パッケージ数: 1,024
- 2014-10にver. 3.0をリリース (R ver. 3.1.1で動作確認)、提供パッケージ数: 934
- 2014-04にver. 2.14をリリース (R ver. 3.1.0で動作確認)、提供パッケージ数: 824
- 2013-10にver. 2.13をリリース (R ver. 3.0で動作確認)、提供パッケージ数: 750
- 2013-04にver. 2.12をリリース (R ver. 3.0で動作確認)、提供パッケージ数: 672
- 2012-10にver. 2.11をリリース (R ver. 2.15.1で動作確認)、提供パッケージ数: 608
- 2012-04にver. 2.10をリリース (R ver. 2.15.0で動作確認)、提供パッケージ数: 553
- ...



# Bioconductor

Bioconductorに関する総説(Review)。ゲノム配列やアノテーションパッケージもBioconductorから提供されており、それらに関する言及もあり。

[Nat Methods](#). 2015 Feb;12(2):115-21. doi: 10.1038/nmeth.3252.

## Orchestrating high-throughput genomic analysis with Bioconductor.

[Huber W](#)<sup>1</sup>, [Carey VJ](#)<sup>2</sup>, [Gentleman R](#)<sup>3</sup>, [Anders S](#)<sup>1</sup>, [Carlson M](#)<sup>4</sup>, [Carvalho BS](#)<sup>5</sup>, [Bravo HC](#)<sup>6</sup>, [Davis S](#)<sup>7</sup>, [Gatto L](#)<sup>8</sup>, [Girke T](#)<sup>9</sup>, [Gottardo R](#)<sup>10</sup>, [Hahne F](#)<sup>11</sup>, [Hansen KD](#)<sup>12</sup>, [Irizarry RA](#)<sup>13</sup>, [Lawrence M](#)<sup>3</sup>, [Love MI](#)<sup>13</sup>, [MacDonald J](#)<sup>14</sup>, [Obenchain V](#)<sup>4</sup>, [Oleś AK](#)<sup>1</sup>, [Pagès H](#)<sup>4</sup>, [Reyes A](#)<sup>1</sup>, [Shannon P](#)<sup>4</sup>, [Smyth GK](#)<sup>15</sup>, [Tenenbaum D](#)<sup>4</sup>, [Waldron L](#)<sup>16</sup>, [Morgan M](#)<sup>4</sup>.

### + Author information

### Abstract

Bioconductor is an open-source, open-development software project for the analysis and comprehension of high-throughput data in genomics and molecular biology. The project aims to enable interdisciplinary research, collaboration and rapid development of scientific software. Based on the statistical programming language R, Bioconductor comprises 934 interoperable packages contributed by a large, diverse community of scientists. Packages cover a range of bioinformatic and statistical applications. They undergo formal initial review and continuous automated testing. We present an overview for prospective users and contributors.

# パッケージのインストール

「必要最小限プラスアルファ」の推奨インストール手順を行えば、「(Rで)塩基配列解析」で利用する多くのパッケージがインストールされます。

- [はじめに](#) (last modified 2015/03/31)
- [参考資料\(講義、講習会、本など\)](#) (last modified 2015/07/07) **NEW**
- [過去のお知らせ](#) (last modified 2015/07/08) **NEW**
- [インストール | について](#) (last modified 2015/04/04)
- [インストール | R本体 | 最新版 | Win用](#) (last modified 2015/03/22) 推奨
- [インストール | R本体 | 最新版 | Mac用](#) (last modified 2015/04/22) 推奨
- [インストール | R本体 | 過去版 | Win用](#) (last modified 2015/03/22)
- [インストール | R本体 | 過去版 | Mac用](#) (last modified 2015/03/22)
- [インストール | Rパッケージ | ほぼ全て\(20GB以上?!\)](#) (last modified 2015/05/25)
- [インストール | Rパッケージ | 必要最小限プラスアルファ\(数GB?!\)](#) (last modified 2015/06/02) 推奨
- [インストール | Rパッケージ | 必要最小限プラスアルファ\(アグリバイオ居室のみ\)](#) (last modified 2015/06/16)
- [インストール | Rパッケージ | 必要最小限\(数GB?!\)](#) (last modified 2015/05/25)
- [インストール | Rパッケージ | 個別](#) (last modified 2015/06/10)
- (削除予定)[Rのインストールと起動](#) (last modified 2015/04/02)
- (削除予定)[個別パッケージのインストール](#) (last modified 2015/02/20)
- [基本的な利用法](#) (last modified 2015/04/03)
- [サンプルデータ](#) (last modified 2015/06/15)
- [バイオインフォマティクス人材育成カリキュラム\(次世代シーケンサ\) | NGSハンズオン講習会2015](#) (last modified 2015/06/15)
- [バイオインフォマティクス人材育成カリキュラム\(次世代シーケンサ\) | 速習コース2014](#) (last modified 2015/06/15)
- [書籍 | トランスクリプトーム解析 | について](#) (last modified 2014/05/12)
- [書籍 | トランスクリプトーム解析 | 2.3.1 RNA-seqデータ\(FASTQファイル\)](#) (last modified 2014/04/15)
- [書籍 | トランスクリプトーム解析 | 2.3.2 RNA-seqデータ\(BAMファイル\)](#) (last modified 2014/04/16)

# パッケージのインストール

## インストール | Rパッケージ | 必要最小限プラスアルファ(数GB?!) **NEW**

(Rで)塩基配列解析、(Rで)マイクロアレイデータ解析中で利用するパッケージ、プラスアルファのパッケージをインストールするやり方です。Rパッケージの2大リポジトリであるCRANとBioconductorから提供されているパッケージ群のうち、一部のインストールに相当しますので、相当短時間でインストールが完了します。

### 1. R本体を起動

### 2. CRANから提供されているパッケージ群のインストール

以下を「R コンソール画面上」でコピー&ペースト。どこからダウンロードするか?と聞かれるので、その場合は自分のいる場所から近いサイトを指定しましょう。

#### #(Rで)塩基配列解析で主に利用

```
install.packages("limma")  
install.packages("samr")  
install.packages("seqinr")
```

```
#(Rで)マイクロアレイデータ解析でも利用  
#(Rで)マイクロアレイデータ解析でも利用  
#(Rで)マイクロアレイデータ解析でも利用
```

#### #(Rで)マイクロアレイデータ解析で利用

```
install.packages("cclust")  
install.packages("class")  
install.packages("e1071")  
install.packages("GeneCycle")  
install.packages("gptk")  
install.packages("GSA")  
install.packages("mixOmics")  
install.packages("pvclust")  
install.packages("RobLoxBioC")  
install.packages("som")  
install.packages("st")  
install.packages("varSelRF")
```

#### #アグリバイオの他の講義科目で利用予定

```
install.packages("ape")  
install.packages("cluster")  
install.packages("fields")
```

# パッケージのインストール

```
biocLite( DESeq , suppressUpdates=TRUE )
biocLite( "DESeq", suppressUpdates=TRUE )
biocLite( "DESeq2", suppressUpdates=TRUE )
biocLite( "DiffBind", suppressUpdates=TRUE )
biocLite( "doMC", suppressUpdates=TRUE )
biocLite( "EBSeq", suppressUpdates=TRUE )
biocLite( "EDASeq", suppressUpdates=TRUE )
biocLite( "edgeR", suppressUpdates=TRUE )
biocLite( "GenomicAlignments", suppressUpdates=TRUE )
```

①

## ①ゲノム情報のパッケージ群

(BSgenome...)はBioconductorから提供されています。ここでは計6パッケージをインストールしています。②例えば赤線部分は、マウスのmm10というバージョンのゲノム配列情報を含むパッケージの名前 (BSgenome.Mmusculus.UCSC.mm10) に相当します。biocLiteという関数を用いて該当パッケージをインストールしています。

## 4. Bioconductorから提供されているパッケージ群のインストール

ゲノム配列パッケージです。一つ一つの容量が尋常でないため、必要に応じてテキストエディタなどに予めコピーしておき、いらぬゲノムパッケージを削除してからお使いください。

```
source("http://bioconductor.org/biocLite.R")#おまじない
biocLite("BSgenome.Athaliana.TAIR.TAIR9", suppressUpdates=TRUE)#シロイヌナズナゲノム
biocLite("BSgenome.Celegans.UCSC.ce6", suppressUpdates=TRUE)#線虫ゲノム
biocLite("BSgenome.Drerio.UCSC.danRer7", suppressUpdates=TRUE)#ゼブラフィッシュゲノム
biocLite("BSgenome.Hsapiens.NCBI.GRCh38", suppressUpdates=TRUE)#ヒトゲノム(GRCh38)
biocLite("BSgenome.Hsapiens.UCSC.hg19", suppressUpdates=TRUE)#ヒトゲノム(hg19)
biocLite("BSgenome.Mmusculus.UCSC.mm10", suppressUpdates=TRUE)#マウスゲノム(mm10)
```

②

# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証



# Bioconductor概観

PubMed上で「R Bioconductor」でキーワード検索し原著論文があるパッケージのみ探すのも一つの戦略ですが、原著論文公開前のパッケージも見つかります。

- 作図 | ROC曲線 | 基礎編 | [7. 図の重ね書き\(new\)](#) (last modified 2015/02/15) NEW
- 作図 | ROC曲線 | 基礎編 | [8. 凡例を追加\(legend\)](#) (last modified 2015/02/15) NEW
- 作図 | ROC曲線 | 応用編 | (last modified 2015/02/07) NEW
- 作図 | [SplicingGraphs](#) (last modified 2015/02/07) NEW
- [パイプライン](#) | [||](#)について (last modified 2015/02/07) NEW
- [パイプライン](#) | [ゲノム](#) | [発](#)
- [パイプライン](#) | [ゲノム](#) | [機](#)
- [パイプライン](#) | [ゲノム](#) | [機](#)
- [パイプライン](#) | [ゲノム](#) | [sm](#)
- [リンク集](#) (last modified 2015/02/07) NEW

## リンク集

- [R](#)
- [Bioconductor: Gentleman et al., Genome Biol., 2004](#)
- [CRAN](#)
- [RjpWiki](#)
- [R Tips\(竹澤様\)](#)
- [BioEdit\(フリー\)](#)
- [BioMart: Smed](#)
- [DDBJ Read Ar](#)
- [EMBOSS expl](#)
- [Biostar: Parnell](#)
- [SEQanswers: L](#)
- [NGS WikiBook](#)
- [HT Sequence A](#)

# Bioconductor概観

ネットワーク環境にもよりますが、数秒～数分以内にこのような画面になります。

Home » BiocViews

## All Packages

**Bioconductor version 3.0 (Release)**

Autocomplete biocViews search:

Search:

Home Install Help Developers About

Packages found under Software:

Show  entries Search table:

Package	Maintainer	Title
<a href="#">a4</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Umbrella Package
<a href="#">a4Base</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Base Package
<a href="#">a4Classif</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Classification Package
<a href="#">a4Core</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Core Package
<a href="#">a4Preproc</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Preprocessing Package
<a href="#">a4Reporting</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Reporting Package
<a href="#">ABarray</a>	Yongming Andrew Sun	Microarray QA and statistical data analysis for Applied Biosystems Genome Survey Microarray (AB1700) gene expression data.
<a href="#">ABSSeq</a>	Wentao Yang	ABSSeq: a new RNA-Seq analysis method based on absolute expression differences and generalized Poisson

Software (936)

- AssayDomain (299)
- BiologicalQuestion (261)
- Infrastructure (187)
- ResearchField (193)
- StatisticalMethod (261)
- Technology (591)
- WorkflowStep (477)
- AnnotationData (895)
- ExperimentData (223)

# Bioconductor概観

基本的には左側のカテゴリ分けのところを眺めますが、Biostringsなど何を行うパッケージかがある程度分かっているものから逆引きして感覚をつかんでおくとよいでしょう。

「RNA-seq」など、フリーワード検索をやってもいいとは思いますが、経験上あまりうまく引かかってこないなので私はやりません。



Home

Install

Help

Developers

About

Home » BiocViews

## All Packages

Bioconductor version 3.0 (Release)

Packages found under Software:

Autocomplete biocViews search:

Show  entries

Search table:

- ▼ Software (936)
  - ▶ AssayDomain (299)
  - ▶ BiologicalQuestion (261)
  - ▶ Infrastructure (187)
  - ▶ ResearchField (193)
  - ▶ StatisticalMethod (261)
  - ▶ Technology (591)
  - ▶ WorkflowStep (477)
- ▶ AnnotationData (895)
- ▶ ExperimentData (223)

Package	Maintainer	Title
<a href="#">a4</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Umbrella Package
<a href="#">a4Base</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Base Package
<a href="#">a4Classif</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Classification Package
<a href="#">a4Core</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Core Package
<a href="#">a4Preproc</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Preprocessing Package
<a href="#">a4Reporting</a>	Tobias Verbeke, Willem Ligtenberg	Automated Affymetrix Array Analysis Reporting Package
<a href="#">ABarray</a>	Yongming Andrew Sun	Microarray QA and statistical data analysis for Applied Biosystems Genome Survey Microarray (AB1700) gene expression data.
<a href="#">ABSSeq</a>	Wentao Yang	ABSSeq: a new RNA-Seq analysis method based on absolute expression differences and generalized Poisson

# Bioconductor概観

「biostrings」と打つとすぐにリストアップされる。

The screenshot shows the Bioconductor website interface. At the top, there is a navigation bar with links for Home, Install, Help, Developers, and About. A search bar is located in the top right corner. Below the navigation bar, the page title is "All Packages". On the left side, there is a sidebar with a search input field and a list of package categories under "Bioconductor version 3.0 (Release)". The main content area displays "Packages found under Software:" with a search filter set to "All" and a search table containing "biostrings". A table lists two packages: "Biostrings" and "BSgenome". A red arrow points to the "Biostrings" package name in the table. Below the table, it says "Showing 1 to 2 of 2 entries (filtered from 939 total entries)".

Bioconductor  
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

Home » BiocViews

## All Packages

Bioconductor version 3.0 (Release)

Autocomplete biocViews search:

Software (936)

- AssayDomain (299)
- BiologicalQuestion (261)
- Infrastructure (187)
- ResearchField (193)
- StatisticalMethod (261)
- Technology (591)
- WorkflowStep (477)
- AnnotationData (895)
- ExperimentData (223)

Packages found under Software:

Show All entries

Search table: biostrings

Package	Maintainer	Title
<a href="#">Biostrings</a>	H. Pages	String objects representing biological sequences, and matching algorithms
<a href="#">BSgenome</a>	H. Pages	Infrastructure for Biostrings-based genome data packages

Showing 1 to 2 of 2 entries (filtered from 939 total entries)

Previous Next

# Bioconductor概観

BioconductorのBiostringsパッケージのページに飛びます。

Home » [Bioconductor 3.0](#) » [Software Packages](#) » Biostrings

## Biostrings

String objects representing biological sequences, and matching algorithms

Bioconductor version: Release (3.0)

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

Author: H. Pages, P. Aboyoun, R. Gentleman, and S. DebRoy

Maintainer: H. Pages <hpages at fhcrc.org>

Citation (from within R, enter `citation("Biostrings")`):

Pages H, Aboyoun P, Gentleman R and DebRoy S. *Biostrings: String objects representing biological sequences, and matching algorithms*. R package version 2.34.1.

### Installation

To install this package, start R and enter:

```
source("http://bioconductor.org/biocLite.R")
biocLite("Biostrings")
```

### Documentation

To view documentation for the version of this package installed in your system, start R and enter:

### Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression](#) (parathyroideSE vignette)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)
- [RNA-Seq workflow: gene-level exploratory analysis and differential expression](#)
- [Changing genomic coordinate systems with rtracklayer::liftOver](#)
- [Mass spectrometry and proteomics data analysis](#)

### Mailing Lists »

Post questions about Bioconductor packages to our mailing lists. Read the [posting guide](#) before posting!

# Bioconductor概観

BioconductorのBiostringsパッケージのページで、ちょっと下のほうに移動。biocViewsのところで見えるキーワードっぽいのがさきほどのカテゴリ分けに相当。例えば、DataRepresentationをクリックするとカテゴリ分けの階層関係がわかる。

Documentation

To view documentation for the version of this package installed in your system, start R and enter:

```
browsevignettes("Biostrings")
```

PDF [R Script](#) A short presentation of the basic classes defined in Biostrings 2

PDF [R Script](#) Biostrings Quick Overview

PDF [R Script](#) Handling probe sequence information

PDF [R Script](#) Multiple Alignments

PDF [R Script](#) Pairwise Sequence Alignments

PDF [R Script](#) Reference Manual

Text [NEWS](#)

Details

biocViews [DataImport](#), [DataRepresentation](#), [Genetics](#), [Infrastructure](#), [SequenceMatching](#), [Sequencing](#), [Software](#)

Version 2.34.1

In Bioconductor since BioC 1.6 (R-2.1) or earlier

License Artistic-2.0

Depends R (>= 2.8.0), methods, [BiocGenerics](#)(>= 0.11.3), [S4Vectors](#)(>= 0.2.2), [IRanges](#)(>= 1.99.27), [XVector](#)(>= 0.5.8)

Imports graphics, methods, stats, utils, [BiocGenerics](#), [IRanges](#), [XVector](#), [zlibbioc](#)

Suggests [BSgenome](#)(>= 1.13.14), [BSgenome.Celegans.UCSC.ce2](#)(>= 1.3.11), [BSgenome.Dmelanogaster.UCSC.dm3](#)(>= 1.3.11), [BSgenome.Hsapiens.UCSC.hg18](#), [drosophila2probe](#), [hqu95av2probe](#), [hqu133aprobe](#), [GenomicFeatures](#)(>= 1.3.14), [hqu95av2cdf](#), [affv](#)(>= 1.41.3), [affydata](#)(>= 1.11.5), [RUnit](#)

System Requirements

URL [altcdfenvs](#), [Basic4Cseq](#), [BRAIN](#), [BSgenome](#), [ChIPpeakAnno](#), [ChIPsim](#), [cleaver](#), [CRISPRseek](#), [DASiR](#), [DECIPHER](#), [deepSNV](#), [Fdb.FANTOM4.promoters.hg19](#), [GeneRegionScan](#), [genomes](#), [GenomicAlignments](#), [GOTHIC](#), [harbChIP](#), [hiReadsProcessor](#), [IPAC](#), [JASPAR2014](#), [kebabs](#), [methVisual](#), [minfi](#), [MotifDb](#), [motifRG](#), [oliqo](#), [oneChannelGUI](#)

packages to our training lists. Read the [posting guide](#) before posting!

- [bioconductor](#)
- [bioc-devel](#)

# Bioconductor概観

DataRepresentationのカテゴリに含まれるのは34パッケージであることが分かる。赤枠部分がそのリスト。Biostringsは、①大分類はSoftware、②中分類はInfrastructureとなっており、その下の階層の③DataRepresentationに含まれていることがわかる。



Home

Install

Help

Developers

About

Home » BiocViews

## All Packages

Bioconductor version 3.0 (Release)

Packages found under DataRepresentation:

Autocomplete biocViews search:

Show  entries

Search table:

- ▼ Software (936)
  - ▶ AssayDomain (299)
  - ▶ BiologicalQuestion (261)
  - ▼ Infrastructure (187)
    - DataImport (82)
    - DataRepresentation (34)
    - GUI (20)
    - ThirdPartyClient (9)
  - ▶ ResearchField (193)
  - ▶ StatisticalMethod (261)
  - ▶ Technology (591)
  - ▶ WorkflowStep (477)
- ▶ AnnotationData (895)
- ▶ ExperimentData (223)

Package	Maintainer	Title
<a href="#">AtlasRDF</a>	James Malone	Gene Expression Atlas query and gene set enrichment package.
<a href="#">BaseSpaceR</a>	Adrian Alexa	R SDK for BaseSpace RESTful API
<a href="#">bigmemoryExtras</a>	Peter M. Haverty	An extension of the bigmemory package with added safety, convenience, and a factor class.
<a href="#">Biostrings</a>	H. Pages	String objects representing biological sequences, and matching algorithms
<a href="#">BSgenome</a>	H. Pages	Infrastructure for Biostrings-based genome data packages
<a href="#">cummeRbund</a>	Loyal A. Goff	Analysis, exploration, manipulation, and visualization of Cufflinks high-throughput sequencing data.
<a href="#">flowPlots</a>	N. Hawkins	flowPlots: analysis plots and data class for gated flow cytometry data
<a href="#">flowWorkspace</a>	Greg Finak, Mike Jiang	Import flowJo Workspaces into BioConductor and replicate flowJo gating with flowCore

# Bioconductor概観

①大分類のSoftware、②中分類のInfrastructureも存在する。Biostringsは塩基配列の切り出しや文字列検索系関数も提供しているので、③SequenceMatchingがあるのも妥当。

**Documentation**

To view documentation for the version of this package installed in your system, start R and enter:

```
browsevignettes("Biostrings")
```

[PDF](#) [R Script](#) A short presentation of the basic classes defined in Biostrings 2  
[PDF](#) [R Script](#) Biostrings Quick Overview  
[PDF](#) [R Script](#) Handling probe sequence information  
[PDF](#) [R Script](#) Multiple Alignments  
[PDF](#) [R Script](#) Pairwise Sequence Alignments  
[PDF](#) Reference Manual  
[Text](#) NEWS

**Details**

**biocViews** [DataImport](#), [DataRepresentation](#), [Genetics](#), [Infrastructure](#), [SequenceMatching](#), [Sequencing](#), [Software](#)

**Version** 2.34.1

**In** Bioconductor since BioC 1.6 (R-2.1) or earlier

**License** Artistic-2.0

**Depends** R (>= 2.8.0), methods, [BiocGenerics](#)(>= 0.11.3), [S4Vectors](#)(>= 0.2.2), [IRanges](#)(>= 1.99.27), [XVector](#)(>= 0.5.8)

**Imports** graphics, methods, stats, utils, [BiocGenerics](#), [IRanges](#), [XVector](#), [zlibbioc](#)

**Suggests** [BSgenome](#)(>= 1.13.14), [BSgenome.Celegans.UCSC.ce2](#)(>= 1.3.11), [BSgenome.Dmelanogaster.UCSC.dm3](#)(>= 1.3.11), [BSgenome.Hsapiens.UCSC.hg18](#), [drosophila2probe](#), [hgu95av2probe](#), [hgu133aprobe](#), [GenomicFeatures](#)(>= 1.3.14), [hgu95av2cdf](#), [affv](#)(>= 1.41.3), [affydata](#)(>= 1.11.5), [RUnit](#)

**System Requirements**

**URL** [altcdfenvs](#), [Basic4Cseq](#), [BRAIN](#), [BSgenome](#), [ChIPpeakAnno](#), [ChIPsim](#), [cleaver](#), [CRISPRseek](#), [DASiR](#), [DECIPHER](#), [deepSNV](#), [FDB.FANTOM4.promoters.hg19](#), [GeneRegionScan](#), [genomes](#), [GenomicAlignments](#), [GOTHIC](#), [harbChIP](#), [hiReadsProcessor](#), [IPAC](#), [JASPAR2014](#), [kebabs](#), [methVisual](#), [minfi](#), [MotifDb](#), [motifRG](#), [oligo](#), [oneChannelGUI](#)



# Bioconductor概観

①SequenceMatchingに含まれる17パッケージの一部しか表示されていないが、(Rで)塩基配列解析中にはない②CRISPR関連のパッケージなども存在することに気づく。また、③ゲノム配列パッケージBSgenomeもあることがわかる。



Home

Install

Help

Developers

About

Home » BiocViews

## All Packages

### Bioconductor version 3.0 (Release)

Autocomplete biocViews search:

- GenomeAnnotation (10)
- GenomicVariation (6)
- LinkageDisequilibrium (1)
- MotifAnnotation (3)
- MotifDiscovery (4)
- NetworkEnrichment (13)
- NetworkInference (17)
- SequenceMatching (17)** ①
- SomaticMutation (4)
- VariationDetection (1)
- ▶ Infrastructure (187)
- ▶ ResearchField (193)
- ▶ StatisticalMethod (261)
- ▶ Technology (591)
- ▶ WorkflowStep (477)
- ▶ AnnotationData (895)
- ▶ ExperimentData (223)

### Packages found under SequenceMatching:

Show  entries

Search table:

Package	Maintainer	Title
<a href="#">Biostrings</a>	H. Pages	String objects representing biological sequences, and matching algorithms
<a href="#">BSgenome</a> ③	H. Pages	Infrastructure for Biostrings-based genome data packages
<a href="#">cleanUpdTSeq</a>	Sarah Sheppard; Jianhong Ou; Lihua Julie Zhu	This package classifies putative polyadenylation sites as true or false/internally oligodT primed.
<a href="#">cobindR</a>	Manuela Benary	Finding Co-occurring motifs of transcription factor binding sites
<a href="#">CRISPRseek</a> ②	Lihua Julie Zhu	Design of target-specific guide RNAs in CRISPR-Cas9, genome-editing systems
<a href="#">dagLogo</a>	Jianhong Ou	dagLogo
<a href="#">FunciSNP</a>	Simon G. Coetzee	Integrating Functional Non-coding Datasets with Genetic Association Studies to Identify Candidate Regulatory SNPs
<a href="#">hapFabia</a>	Sepp Hochreiter	hapFabia: Identification of very short segments of identity by descent (IBD) characterized by rare variants in large sequencing data

# Bioconductor概観

ゲノム配列パッケージBSgenomeは、①内部的にBiostringsパッケージを利用していることがわかる。

The screenshot shows the Bioconductor website for the BSgenome package. The page title is "BSgenome" and the subtitle is "Infrastructure for Biostrings-based genome data packages". The page includes a search bar, navigation links (Home, Install, Help, Developers, About), and a list of workflows. A red arrow with the number 1 points to the package name "BSgenome".

Home » [Bioconductor 3.0](#) » [Software Packages](#) » BSgenome

## BSgenome

Infrastructure for Biostrings-based genome data packages

Bioconductor version: Release (3.0)

Infrastructure shared by all the Biostrings-based genome data packages

Author: Herve Pages

Maintainer: H. Pages <hpages at fhcrc.org>

Citation (from within R, enter `citation("BSgenome")`):

Pages H. *BSgenome: Infrastructure for Biostrings-based genome data packages*. R package version 1.34.1.

### Installation

To install this package, start R and enter:

```
source("http://bioconductor.org/biocLite.R")
biocLite("BSgenome")
```

### Documentation

To view documentation for the version of this package installed in your system, start R and enter:

### Workflows »

Common Bioconductor workflows include:

- [Oligonucleotide Arrays](#)
- [High-throughput Sequencing](#)
- [Counting Reads for Differential Expression](#) (parathyroideSE vignette)
- [Annotation](#)
- [Annotating Variants](#)
- [Annotating Ranges](#)
- [Flow Cytometry](#) and other assays
- [Candidate Binding Sites for Known Transcription Factors](#)
- [Cloud-enabled cis-eQTL search and annotation](#)
- [RNA-Seq workflow: gene-level exploratory analysis and differential expression](#)
- [Changing genomic coordinate systems with rtracklayer::liftOver](#)
- [Mass spectrometry and proteomics data analysis](#)

### Mailing Lists »

Post questions about Bioconductor packages to our mailing lists. Read the [posting guide](#) before posting!

# Bioconductor概観

ゲノム配列パッケージBSgenomeのちょっと下のほうに移動。Depends(依存)のところにBiostringsが存在することがわかる。BSgenomeパッケージを利用したい場合には、予めこれらの依存関係のあるパッケージ群のインストールが完了している必要がある。推奨手順通りにパッケージのインストールをしていけばBSgenomeを問題なく利用できるはずである。

Documentation

To view documentation for the version of this package installed in your system, start R and enter:

```
browseVignettes("BSgenome")
```

[PDF](#) [R Script](#) Efficient genome searching with Biostrings and the BSgenome data packages

[PDF](#) [R Script](#) How to forge a BSgenome data package

[PDF](#) Reference Manual

[Text](#) NEWS

Details

biocViews [Annotation](#), [DataRepresentation](#), [Genetics](#), [Infrastructure](#), [SNP](#), [SequenceMatching](#), [Software](#)

Version 1.34.1

In Bioconductor since BioC 1.9 (R-2.4)

License Artistic-2.0

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

Imports methods, stats, [BiocGenerics](#), [S4Vectors](#), [IRanges](#), [XVector](#), [GenomeInfoDb](#), [GenomicRanges](#), [Biostrings](#), [Rsamtools](#), [rtracklayer](#)

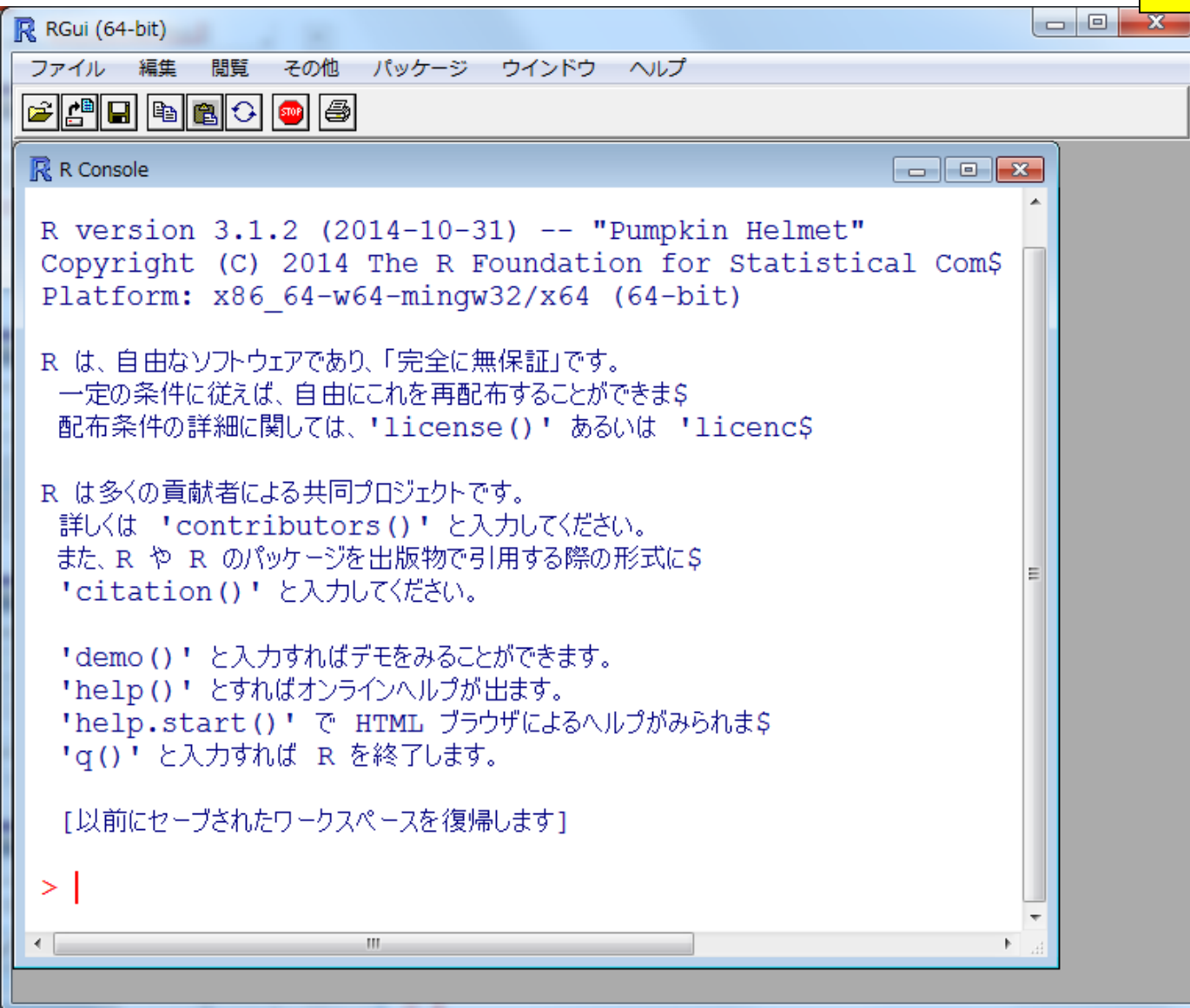
Suggests [BiocInstaller](#), [BSgenome.Celegans.UCSC.ce2](#)(>= 1.3.11), [BSgenome.Hsapiens.UCSC.hg19](#)(>= 1.3.11), [BSgenome.Hsapiens.UCSC.hg19.masked](#), [BSgenome.Rnorvegicus.UCSC.rm5](#), [SNPlocs.Hsapiens.dbSNP.20100427](#), [hqu95av2probe](#), [Biobase](#), [RUnit](#)

System Requirements

URL [BSgenome.Alvrata.JGI.v1](#), [BSgenome.Amellifera.BeeBase.assembly4](#), [BSgenome.Amellifera.UCSC.apiMel2](#), [BSgenome.Amellifera.UCSC.apiMel2.masked](#), [BSgenome.Athaliana.TAIR.04232008](#), [BSgenome.Athaliana.TAIR.TAIR9](#), [BSgenome.Btaurus.UCSC.bosTau3](#), [BSgenome.Btaurus.UCSC.bosTau3.masked](#), [BSgenome.Btaurus.UCSC.bosTau4](#), [BSgenome.Btaurus.UCSC.bosTau4.masked](#), [BSgenome.Btaurus.UCSC.bosTau6](#), [BSgenome.Btaurus.UCSC.bosTau6.masked](#)

# Bioconductor概観

BSgenomeを問題なく利用できるかは、library(BSgenome)が通るかどうかで判断。R Guiの新規画面を起動。



```
RGui (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ
R Console
R version 3.1.2 (2014-10-31) -- "Pumpkin Helmet"
Copyright (C) 2014 The R Foundation for Statistical Com$
Platform: x86_64-w64-mingw32/x64 (64-bit)

R は、自由なソフトウェアであり、「完全に無保証」です。
一定の条件に従えば、自由にこれを再配布することができます
配布条件の詳細に関しては、'license()' あるいは 'licenc$

R は多くの貢献者による共同プロジェクトです。
詳しくは 'contributors()' と入力してください。
また、R や R のパッケージを出版物で引用する際の形式に$
'citation()' と入力してください。

'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザによるヘルプがみられま$
'q()' と入力すれば R を終了します。

[以前にセーブされたワークスペースを復帰します]

> |
```

# Bioconductor概観

library(BSgenome)の実行結果中にエラーメッセージが出ていなければOK

```
RGui (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ
R Console

R version 3.1.2 (2014-10-31) -- "Pumpkin Helmet"
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R は、自由なソフトウェアであり、「完全に無保証」です。
一定の条件に従えば、自由にこれを再配布することができます。
配布条件の詳細に関しては、'license()' あるいは 'licenc$

R は多くの貢献者による共同プロジェクトです。
詳しくは 'contributors()' と入力してください。
また、R や R のパッケージを出版物で引用する際の形式に$
'citation()' と入力してください。

'demo()' と入力すればデモをみることができます。
'help()' とすればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザによるヘルプがみられま$
'q()' と入力すれば R を終了します。

[以前にセーブされたワークスペースを復帰します]

> library(BSgenome)
```

```
R Console

> library(BSgenome)
要求されたパッケージ BiocGenerics をロード中です
要求されたパッケージ parallel をロード中です

次のパッケージを付け加えます: 'BiocGenerics'

The following objects are masked from 'package:parallel$

clusterApply, clusterApplyLB, clusterCall,
clusterEvalQ, clusterExport, clusterMap,
parApply, parCapply, parLapply, parLapplyLB,
parRapply, parSapply, parSapplyLB

The following object is masked from 'package:stats':

xtabs

The following objects are masked from 'package:base':

anyDuplicated, append, as.data.frame,
as.vector, cbind, colnames, do.call,
duplicated, eval, evalq, Filter, Find, get,
intersect, is.unsorted, lapply, Map, mapply,
match, mget, order, paste, pmax, pmax.int,
pmin, pmin.int, Position, rank, rbind,
Reduce, rep.int, rownames, sapply, setdiff,
sort, table, tapply, union, unique, unlist

要求されたパッケージ IRanges をロード中です
要求されたパッケージ GenomicRanges をロード中です
要求されたパッケージ GenomeInfoDb をロード中です
要求されたパッケージ Biostrings をロード中です
要求されたパッケージ XVector をロード中です

> |
```

# Bioconductor概観

library(XXX)をやったときに、XXXパッケージ内部で利用するDependsやImportsにリストアップされているパッケージも同時にロードしている(読み込んでいる)。

**Documentation**

To view documentation for the version of this package installed in your system, start R and enter:

```
browseVignettes("BSgenome")
```

**Details**

biocViews [Annotation](#), [DataRepresentation](#), [Genetics](#), [Infrastructure](#), [SNP](#), [SequenceMatching](#), [Software](#)

Version 1.34.1

In Bioconductor BioC 1.9 (R-2.4) since

License Artistic-2.0

Depends R (>= 2.8.0), methods, [BiocGenerics](#) (>= 0.1.2), [S4Vectors](#) (>= 0.0.7), [IRanges](#) (>= 1.99.1), [GenomeInfoDb](#) (>= 1.1.4), [GenomicRanges](#) (>= 1.17.15), [Biostrings](#) (>= 2.33), [rtracklayer](#) (>= 1.25.8)

Imports methods, stats, [BiocGenerics](#), [S4Vectors](#), [IRanges](#), [XVector](#), [GenomeInfoDb](#), [GenomicRanges](#), [Biostrings](#), [Rsamtools](#), [rtracklayer](#)

Suggests [BiocInstaller](#), [BSgenome.Celegans.UCSC.ce2](#) (>= 1.3.11), [BSgenome.Hsapiens.UCSC.hg19](#) (>= 1.3.11), [BSgenome.Hsapiens.UCSC.hg19.masked](#), [BSgenome.Rnorvegicus.UCSC.rm5](#), [SNPlocs.Hsapiens.dbSNP.20100427](#), [hqu95av2pro](#), [Biobase](#), [RUnit](#)

System Requirements

URL [BSgenome.Alvrata.JGI.v1](#), [BSgenome.Amellifera.BeeBase.assembly4](#), [BSgenome.Amellifera.UCSC.apiMel2](#), [BSgenome.Amellifera.UCSC.apiMel2.masked](#), [BSgenome.Athaliana.TAIR.04232008](#), [BSgenome.Athaliana.TAIR.TAIR9](#), [BSgenome.Btaurus.UCSC.bosTau3](#), [BSgenome.Btaurus.UCSC.bosTau3.masked](#), [BSgenome.Btaurus.UCSC.bosTau4](#), [BSgenome.Btaurus.UCSC.bosTau4.masked](#), [BSgenome.Btaurus.UCSC.bosTau6](#), [BSgenome.Btaurus.UCSC.bosTau6.masked](#)

```
> library(BSgenome)
```

```
要求されたパッケージ BiocGenerics をロード中です ←
```

```
要求されたパッケージ parallel をロード中です ←
```

```
次のパッケージを付け加えます: 'BiocGenerics'
```

```
The following objects are masked from 'package:parallel':
```

```
clusterApply, clusterApplyLB, clusterCall,
clusterEvalQ, clusterExport, clusterMap,
parApply, parCapply, parLapply, parLapplyLB,
parRapply, parSapply, parSapplyLB
```

```
The following object is masked from 'package:stats':
```

```
xtabs
```

```
The following objects are masked from 'package:base':
```

```
anyDuplicated, append, as.data.frame,
as.vector, cbind, colnames, do.call,
duplicated, eval, evalq, Filter, Find, get,
intersect, is.unsorted, lapply, Map, mapply,
match, mget, order, paste, pmax, pmax.int,
pmin, pmin.int, Position, rank, rbind,
Reduce, rep.int, rownames, sapply, setdiff,
sort, table, tapply, union, unique, unlist
```

```
要求されたパッケージ IRanges をロード中です ←
```

```
要求されたパッケージ GenomicRanges をロード中です ←
```

```
要求されたパッケージ GenomeInfoDb をロード中です ←
```

```
要求されたパッケージ Biostrings をロード中です ←
```

```
要求されたパッケージ XVector をロード中です ←
```

```
> |
```

# Bioconductor概観

もう一度library(BSgenome)をやってもメッセージは出ません。エラーメッセージが出ていなければ特に気にする必要はありません。

```
R Console

The following object is masked from 'package:stats':

  xtabs

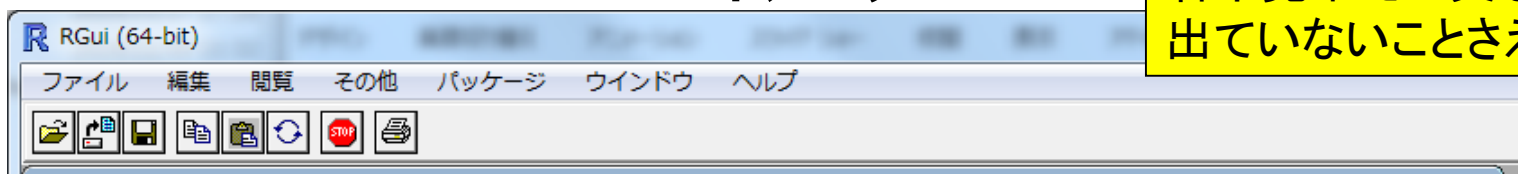
The following objects are masked from 'package:base':

  anyDuplicated, append, as.data.frame,
  as.vector, cbind, colnames, do.call,
  duplicated, eval, evalq, Filter, Find, get,
  intersect, is.unsorted, lapply, Map, mapply,
  match, mget, order, paste, pmax, pmax.int,
  pmin, pmin.int, Position, rank, rbind,
  Reduce, rep.int, rownames, sapply, setdiff,
  sort, table, tapply, union, unique, unlist

要求されたパッケージ IRanges をロード中です
要求されたパッケージ GenomicRanges をロード中です
要求されたパッケージ GenomeInfoDb をロード中です
要求されたパッケージ Biostrings をロード中です
要求されたパッケージ XVector をロード中です
> library(BSgenome)
> |
```

# Bioconductor概観

先にBSgenomeが内部的に利用している Biostringsパッケージのロードを行っておくと、若干見栄えが異なるが、エラーメッセージが出ていないことさえ確認できれば問題ない。



```
R version 3.1.2 (2014-10-31)
Copyright (C) 2014 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32

R は、自由なソフトウェアであり、「完全に無保証」である。ただし、特定の
条件下に於いて、特定の条件下に於いて、自由にこれを再配布することが
可能である。配布条件の詳細に関しては、'license()' または 'licence()'
を入力してください。

R は多くの貢献者による共同プロジェクトです。詳しくは 'contributors()'
を入力してください。また、R や R のパッケージを出版物で引用する場合は、
'citation()' と入力してください。

'demo()' と入力すればデモをみることもできます。
'help()' と入力すればオンラインヘルプが出ます。
'help.start()' で HTML ブラウザ版のヘルプを見ることができます。
'q()' と入力すれば R を終了します。

[以前にセーブされたワークスペースを復元しようとしています]

> library(Biostrings)|
```

```
R Console

clusterExport, clusterMap, parApply, parCapply, parLapply,
parLapplyLB, parRapply, parSapply, parSapplyLB

The following object is masked from 'package:stats':

xtabs

The following objects are masked from 'package:base':

anyDuplicated, append, as.data.frame, as.vector, cbind,
colnames, do.call, duplicated, eval, evalq, Filter, Find,
get, intersect, is.unsorted, lapply, Map, mapply, match,
mget, order, paste, pmax, pmax.int, pmin, pmin.int,
Position, rank, rbind, Reduce, rep.int, rownames, sapply,
setdiff, sort, table, tapply, union, unique, unlist

要求されたパッケージ IRanges をロード中です
要求されたパッケージ XVector をロード中です
> library(BSgenome)
要求されたパッケージ GenomicRanges をロード中です
要求されたパッケージ GenomeInfoDb をロード中です
> library(BSgenome)
> |
```



# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# BSgenome利用の意義

ゲノム配列情報はUCSCやEnsemblなどのウェブサイトから取得するのが一般的ではあるが、Rの生物種ごとに提供されているBSgenomeで取得、あるいは取り扱うことも可能。ChIP-seq用パッケージMEDIPSはBSgenomeを利用。

## 解析 | ChIP-seq | について

このあたりはほとんどノータッチです。SraTailor (Oki et al., 2014)は、[実験医学2014年12月号](#)の「Close Up 実験法」中で日本語による解説記事があります(沖真弥氏提供情報)。2015年2月に調査した結果をリストアップします。

### R用:

- ChIPsim: [Zhang et al., PLoS Comput. Biol., 2008](#)
- PeakSeq法: [Rozowsky et al., Nat Biotechnol., 2009](#)
- CSAR: [Kaufmann et al., PLoS Biol., 2009](#)
- rMAT: [Droit et al., Bioinformatics, 2010](#)
- ChIPpeakAnno: [Zhu et al., BMC Bioinformatics, 2010](#)
- PICS: [Zhang et al., Biometrics, 2011](#)
- ChIPseqR: [Humburg et al., BMC Bioinformatics, 2011](#)
- DiffBind: [Ross-Innes et al., Nature, 2012](#)
- MEDIPS: [Lienhard et al., Bioinformatics, 2014](#)
- DSS: [Feng et al., Nucleic Acids Res., 2014](#)
- methylSig: [Park et al., Bioinformatics, 2014](#)

### R以外:

- bwtool: [Pohl and Beato, Bioinformatics, 2014](#)
- SraTailor: [Oki et al., Genes Cells., 2014](#)

### Review、ガイドライン、パイプライン系:

- ガイドライン: [Bailey et al., PLoS Comput Biol., 2013](#)
- Review: [Robinson et al., Front Genet., 2014](#)

- 解析 | 機能解析 | パスウェイ(Pathway)解析 | [SeqGSEA\(Wang 2014\)](#) (last modified 2014/12/19)
- 解析 | 菌叢解析 | [phyloseq\(McMurdie 2013\)](#) (last modified 2014/05/29)
- 解析 | エクソーム解析 | [\(last modified 2014/12/19\)](#)
- 解析 | ChIP-seq | [\(last modified 2015/02/19\)](#)
- 解析 | ChIP-seq | [DiffBind\(Ross-Innes 2012\)](#) (last modified 2015/02/19)
- 解析 | ChIP-seq | [ChIPseqR\(Humburg 2011\)](#) (last modified 2015/02/19)
- 解析 | ChIP-seq | [chipseq](#) (last modified 2011/12/19)
- 解析 | ChIP-seq | [PICS\(Zhang 2011\)](#) (last modified 2011/12/19)

# BSgenome

- イントロ | 一般 | [任意の長さの連続塩基の出現頻度情報を取得](#) (last modified 2013/06/14)
- イントロ | 一般 | Tips | [任意の拡張子でファイルを保存](#) (last modified 2013/09/26)
- イントロ | 一般 | Tips | [拡張子は同じで任意の文字を追加して保存](#) (last modified 2013/09/26)
- イントロ | 一般 | 配列取得 | ゲノム配列 | [公共DBから](#) (last modified 2014/05/28)
- イントロ | 一般 | 配列取得 | ゲノム配列 | [BSgenome](#) ① (last modified 2015/02/18) **NEW**
- イントロ | 一般 | 配列取得 | プロモーター配列 | [公共DBから](#) (last modified 2014/04/02)
- イントロ | 一般 | 配列取得 | プロモーター配列 | [BSgenome](#) (last modified 2014/04/25)
- イントロ | 一般 | 配列取得 | プロモーター配列 | [GenomicFeatures\(Lawrence 2013\)](#) (last modified 2014/04/23)

## イントロ | 一般 | 配列取得 | ゲノム配列 | [BSgenome](#) **NEW**

[BSgenome](#)パッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ミヤマハタザオ (*A. lyrata*)、セイヨウミツバチ (*A. mellifera*)、シロイヌナズナ (*A. thaliana*)、ウシ (*B. taurus*)、線虫 (*C. elegans*)、犬 (*C. familiaris*)、キロショウジョウバエ (*D. melanogaster*)、ゼブラフィッシュ (*D. rerio*)、大腸菌 (*E. coli*)、イトヨ (*G. aculeatus*)、セキショクヤケイ (*G. gallus*)、ヒト (*H. sapiens*)、アカゲザル (*M. mulatta*)、マウス (*M. musculus*)、チンパンジー (*P. troglodytes*)、ラット (*R. norvegicus*)、出芽酵母 (*S. cerevisiae*)、トキソプラズマ (*T. gondii*)と実に様々な生物種が利用可能であることがわかります。`getSeq`関数はBSgenomeオブジェクト中の「single sequences」というあたりにリストアップされているchr...というものを全て抽出しています。したがって、例えばマウスゲノムは「chr1」以外に「chr1\_random」や「chrUn\_random」なども等価に取扱っている点に注意してください。「ファイル」-「ディレクトリの変更」でファイルを保存したいディレクトリに移動し以下をコピー。

### 1. 利用可能な生物種とRIにインストール済みの生物種をリストアップしたい場合:

```
#必要なパッケージをロード
library(BSgenome) #パッケージの読み込み

#本番 (利用可能なパッケージをリストアップ; インストール済みとは限らない)
available.genomes() #このパッケージ中で利用可能なゲノムをリストアップ

#本番 (インストール済みの生物種をリストアップ)
installed.genomes() #インストール済みの生物種をリストアップ

#後処理 (パッケージ名でだいたいわかるがproviderやversionを分割して表示したい場合)
installed.genomes(splitNameParts=TRUE) #インストール済みの生物種をリストアップ
```

# BSgenome

イントロ | 一般 | 配列取得 | ゲノム配列 | BSgenome NEW

BSgenomeパッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ミヤマバネ (A. mellifera)、シロイヌナズナ (A. thaliana)、ウシ (B. taurus)、線虫 (C. elegans)、ヒト (H. sapiens)、アカゲザル (M. mulatta)、マウス (M. musculus)、チンパンジー (P. troglodytes)、ラット (R. norvegicus)、出芽酵母 (S. cerevisiae)、トキソプラズマ (T. gondii) と実に様々な生物種が利用可能であることがわかります。getSeq関数はBSgenomeオブジェクト中の「single sequences」というあたりにリストアップされているchr...というものを全て抽出しています。したがって、例えばマウスゲノムは「chr1」以外に「chr1\_random」や「chrUn\_random」なども等価に取扱っている。 「ファイル」-「ディレクトリの変更」でファイルを保存したいディレクトリに

黒枠部分のコードをコピー。R ver. 3.1.3 (Bioconductor ver. 3.0)で利用可能な生物種のパッケージ名をリストアップ。71個あることが分かる。Rのバージョンが古いとパッケージ数は少なくなる。例えばR ver. 3.2.0 (Bioconductor ver. 3.1)では77個。

## 1. 利用可能な生物種とRにインストール済みの生物種をリストアップ

```
#必要なパッケージをロード
library(BSgenome) #パッケージの

#本番 (利用可能なパッケージをリストアップ; インストール済みの生物種をリストアップ)
available.genomes() #このパッケージ

#本番 (インストール済みの生物種をリストアップ)
installed.genomes() #インストール済みの生物種

#後処理 (パッケージ名でだいたいわかるがproviderやversionを抽出)
installed.genomes(splitNameParts=TRUE) #インストール済みの生物種
```

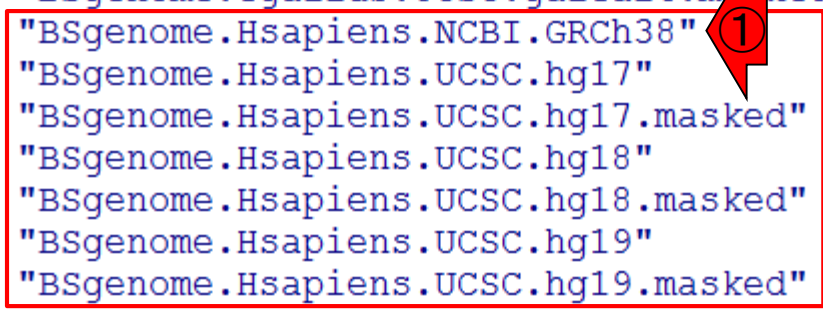
```
R Console

[56] "BSgenome.Ptroglydytes.UCSC.panTro2"
[57] "BSgenome.Ptroglydytes.UCSC.panTro2.masked"
[58] "BSgenome.Ptroglydytes.UCSC.panTro3"
[59] "BSgenome.Ptroglydytes.UCSC.panTro3.masked"
[60] "BSgenome.Rnorvegicus.UCSC.rn4"
[61] "BSgenome.Rnorvegicus.UCSC.rn4.masked"
[62] "BSgenome.Rnorvegicus.UCSC.rn5"
[63] "BSgenome.Rnorvegicus.UCSC.rn5.masked"
[64] "BSgenome.Scerevisiae.UCSC.sacCer1"
[65] "BSgenome.Scerevisiae.UCSC.sacCer2"
[66] "BSgenome.Scerevisiae.UCSC.sacCer3"
[67] "BSgenome.Sscrofa.UCSC.susScr3"
[68] "BSgenome.Sscrofa.UCSC.susScr3.masked"
[69] "BSgenome.Tgondii.ToxoDB.7.0"
[70] "BSgenome.Tguttata.UCSC.taeGut1"
[71] "BSgenome.Tguttata.UCSC.taeGut1.masked"
> |
```

```
> available.genomes()
[1] "BSgenome.Alyrata.JGI.v1"
[2] "BSgenome.Amellifera.BeeBase.asse
[3] "BSgenome.Amellifera.UCSC.apiMel2
[4] "BSgenome.Amellifera.UCSC.apiMel2
[5] "BSgenome.Athaliana.TAIR.04232008
[6] "BSgenome.Athaliana.TAIR.TAIR9"
[7] "BSgenome.Btaurus.UCSC.bosTau3"
[8] "BSgenome.Btaurus.UCSC.bosTau3.ma
[9] "BSgenome.Btaurus.UCSC.bosTau4"
[10] "BSgenome.Btaurus.UCSC.bosTau4.ma
[11] "BSgenome.Btaurus.UCSC.bosTau6"
[12] "BSgenome.Btaurus.UCSC.bosTau6.ma
[13] "BSgenome.Celegans.UCSC.ce10"
[14] "BSgenome.Celegans.UCSC.ce2"
[15] "BSgenome.Celegans.UCSC.ce6"
[16] "BSgenome.Cfamiliaris.UCSC.canFam
[17] "BSgenome.Cfamiliaris.UCSC.canFam
[18] "BSgenome.Cfamiliaris.UCSC.canFam
[19] "BSgenome.Cfamiliaris.UCSC.canFam
[20] "BSgenome.Dmelanogaster.UCSC.dm2"
[21] "BSgenome.Dmelanogaster.UCSC.dm2.
[22] "BSgenome.Dmelanogaster.UCSC.dm3"
[23] "BSgenome.Dmelanogaster.UCSC.dm3.
[24] "BSgenome.Drerio.UCSC.danRer5"
[25] "BSgenome.Drerio.UCSC.danRer5.mas
[26] "BSgenome.Drerio.UCSC.danRer6"
[27] "BSgenome.Drerio.UCSC.danRer6.mas
[28] "BSgenome.Drerio.UCSC.danRer7"
```

赤枠のヒトゲノムは4バージョン(hg17, 18, 19, GRCh38)を利用可能。2013年12月にリリースされた①最新版(GRCh38)のRパッケージも利用可能。

```
R Console
[29] "BSgenome.Drerio.UCSC.danRer7"
[30] "BSgenome.Ecoli.NCBI.20080805"
[31] "BSgenome.Gaculeatus.UCSC.gasAcu1"
[32] "BSgenome.Gaculeatus.UCSC.gasAcu1.masked"
[33] "BSgenome.Ggallus.UCSC.galGal3"
[34] "BSgenome.Ggallus.UCSC.galGal3.masked"
[35] "BSgenome.Ggallus.UCSC.galGal4"
[36] "BSgenome.Ggallus.UCSC.galGal4.masked"
[37] "BSgenome.Hsapiens.NCBI.GRCh38"
[38] "BSgenome.Hsapiens.UCSC.hg17"
[39] "BSgenome.Hsapiens.UCSC.hg17.masked"
[40] "BSgenome.Hsapiens.UCSC.hg18"
[41] "BSgenome.Hsapiens.UCSC.hg18.masked"
[42] "BSgenome.Hsapiens.UCSC.hg19"
[43] "BSgenome.Hsapiens.UCSC.hg19.masked"
[44] "BSgenome.Mfuro.UCSC.musFurl"
[45] "BSgenome.Mmulatta.UCSC.rheMac2"
[46] "BSgenome.Mmulatta.UCSC.rheMac2.masked"
[47] "BSgenome.Mmulatta.UCSC.rheMac3"
[48] "BSgenome.Mmulatta.UCSC.rheMac3.masked"
[49] "BSgenome.Mmusculus.UCSC.mm10"
[50] "BSgenome.Mmusculus.UCSC.mm10.masked"
[51] "BSgenome.Mmusculus.UCSC.mm8"
[52] "BSgenome.Mmusculus.UCSC.mm8.masked"
[53] "BSgenome.Mmusculus.UCSC.mm9"
[54] "BSgenome.Mmusculus.UCSC.mm9.masked"
[55] "BSgenome.Osativa.MSU.MSU7"
```



# BSgenome

イントロ | 一般 | 配列取得 | ゲノム配列 | BSgenome

黒枠部分のコードをコピー。10秒程度かかる。実際にインストール済みのものは(このPC環境では)7パッケージであることがわかる。講習会ではヒトゲノムの最新版(BSgenome.Hsapiens.NCBI.GRCh38)を使いますが、①これが入っていない場合は、②各自installed.genomes()で見られる他のパッケージに置き換えて実行してください。

BSgenomeパッケージを用いて様々な生物種のゲノム配列を取得するやり方を示す。lyrata)、セイヨウミツバチ (A. mellifera)、シロイヌナズナ (A. thaliana)、ウシ (B. taurus)、ヒト (C. familiaris)、キイロショウジョウバエ (D. melanogaster)、ゼブラフィッシュ (D. rerio)、アカゲザル (M. mulatta)、セキショクヤケイ (G. gallus)、ヒト (H. sapiens)、アカゲザル (M. mulatta)、パンジー (P. troglodytes)、ラット (R. norvegicus)、出芽酵母 (S. cerevisiae)、トキソプラズマ (Toxoplasma gondii) などの生物種が利用可能であることがわかります。getSeq関数はBSgenomeオブジェクト中の「single sequences」というあたりにリストアップされているchr...というものを全て抽出しています。したがって、例えばマウスゲノムは「chr1」以外に「chr1\_random」や「chrUn\_random」なども等価に取扱っている点に注意してください。「ファイル」-「ディレクトリの変更」でファイルを保存したいディレクトリに移動し以下をコピー。

## 1. 利用可能な生物種とRにインストール済みの生物種をリストアップしたい場合:

```
#必要なパッケージをロード
library(BSgenome) #パッケージの読み込み

#本番 (利用可能なパッケージをリストアップ; インストール済みと
available.genomes() #このパッケージ中

#本番 (インストール済みの生物種をリストアップ)
installed.genomes() #インストール済み

#後処理 (パッケージ名でだいたいわかるがproviderやversion
installed.genomes(splitNameParts=TRUE) #インストール済み
```

```
R Console

[68] "BSgenome.Sscrofa.UCSC.susScr3.masked"
[69] "BSgenome.Tgondii.ToxoDB.7.0"
[70] "BSgenome.Tguttata.UCSC.taeGut1"
[71] "BSgenome.Tguttata.UCSC.taeGut1.masked"
> installed.genomes() #イン$
[1] "BSgenome.Athaliana.TAIR.TAIR9"
[2] "BSgenome.Celegans.UCSC.ce2"
[3] "BSgenome.Celegans.UCSC.ce6"
[4] "BSgenome.Drerio.UCSC.danRer7"
[5] "BSgenome.Hsapiens.NCBI.GRCh38"
[6] "BSgenome.Hsapiens.UCSC.hg19"
[7] "BSgenome.Mmusculus.UCSC.mm10"
> |
```



# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# BSgenome

イントロ | 一般 | 配列取得 | ゲノム配列 | BSgenome **NEW**

BSgenomeパッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ライラタ (*A. lyrata*)、セイヨウミツバチ (*A. mellifera*)、シロイヌナズナ (*A. thaliana*)、ウシ (*B. taurus*)、ヒト (*C. familiaris*)、キイロショウジョウバエ (*D. melanogaster*)、ゼブラフィッシュ (*D. rerio*)、アカゲザル (*G. aculeatus*)、セキショクヤケイ (*G. gallus*)、ヒト (*H. sapiens*)、アカゲザル (*M. mulatta*)、マウス (*M. musculus*)、チンパンジー (*P. troglodytes*)、ラット (*R. norvegicus*)、出芽酵母 (*S. cerevisiae*)、トキソプラズマ (*T. gondii*) と実に様々な生物種が利用可能であることがわかります。getSeq関数はBSgenomeオブジェクト中の「single sequences」というあたりにリストアップされているchr\_ というものを全て抽出しています。したがって、例えばマウスゲノムには「chr1」以外に「chr1\_random」も含まれています。

①例題9。2013年12月にリリースされたヒトゲノム最新版(GRCh38)のRパッケージを入力、multi-FASTAファイルを出力として得る。作業ディレクトリはどこでもよいが基本はデスクトップ上のhoge。数分かかるが、約3.3GBのファイルが生成される。決してテキストエディタで開かないで!

## ① 9. インストール済みのヒト ("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム配列をmulti-FASTAファイルで保存したい場合:

2013年12月にリリースされたGenome Reference Consortium GRCh38です。R ver. 3.1.0とBioconductor ver. 2.14以上の環境で実行可能です。

### 1. 利用可能な生物種とR

```
#必要なパッケージを
library(BSgenome)

#本番 (利用可能なバ
available.genomes

#本番 (インストール
installed.genomes

#後処理 (パッケージ
installed.genomes
```

```
out_f <- "hoge9.fasta" #出力ファイル名を指定してout_fに格納
param <- "BSgenome.Hsapiens.NCBI.GRCh38"#パッケージ名を指定

#必要なパッケージをロード
library(param, character.only=T) #paramで指定したパッケージの読み込み

#前処理(paramで指定したパッケージ中のオブジェクト名をgenomeに統一)
#tmp <- unlist(strsplit(param, ".", fixed=TRUE))[2]#paramで指定した文字列からオブジェクト名を取得し
tmp <- ls(paste("package", param, sep=":"))#paramで指定したパッケージで利用可能なオブジェクト名を取
genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクトとしてgenomeに格納(パッケージ中には
genome #確認してるだけです

#本番
fasta <- getSeq(genome) #ゲノム塩基配列情報を抽出した結果をfastaに格納
names(fasta) <- seqnames(genome) #description情報を追加している

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fastaの中身を指定したファイル名で保存
```



# BSgenome

①出力ファイルの内容は、fastaオブジェクトに格納されている。慣れれば②fastaオブジェクトの中身を眺めるほうが全体像をつかみやすい。

## 9. インストール済みのヒト("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム配列をmulti-FASTAファイルで保存したい場合:

2013年12月にリリースされたGenome Reference Consortium GRCh38です。R ver. 3.1.0とBioconductor ver. 2.14以上の環境で実行可能です。

```
out_f <- "hoge9.fasta"
param <- "BSgenome.Hsapiens.NCBI.GRCh38"

#必要なパッケージをロード
library(param, character.only=T)

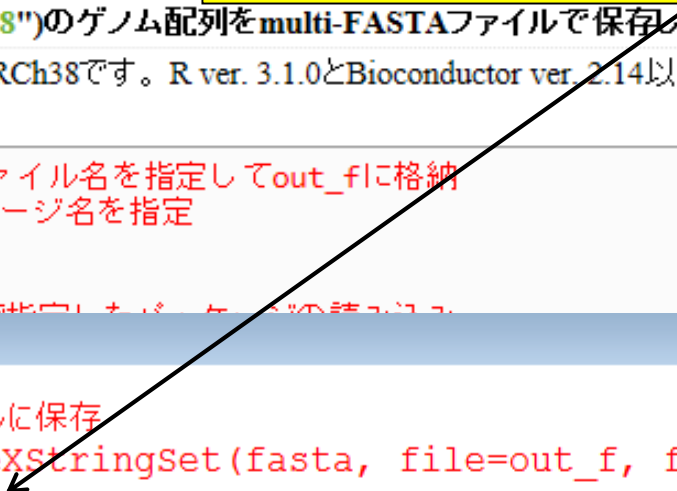
#前処理(paramで指定したパッケージ中のオ
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param, sep="."))
genome <- eval(parse(text=tmp))

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta")
```

#出力ファイル名を指定してout\_fに格納  
#パッケージ名を指定

```
R Console
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width$
> fasta
A DNASTringSet instance of length 455
      width seq          names
[1] 248956422 NNNNNNNNNN...NNNNNNNNNN 1
[2] 242193529 NNNNNNNNNN...NNNNNNNNNN 2
[3] 198295559 NNNNNNNNNN...NNNNNNNNNN 3
[4] 190214555 NNNNNNNNNN...NNNNNNNNNN 4
[5] 181538259 NNNNNNNNNN...NNNNNNNNNN 5
...
[451] 200773 TCTACTCTC...GGGGAATTC HSCHR19KIR_FH08_B...
[452] 170148 TTTCTTTCT...GGGGAATTC HSCHR19KIR_FH13_A...
[453] 215732 TGTGGTGAG...GGGGAATTC HSCHR19KIR_FH13_B...
[454] 170537 TCTACTCTC...GGGGAATTC HSCHR19KIR_FH15_A...
[455] 177381 GATCTATCT...GGGGAATTC HSCHR19KIR_RP5_B...
```



# BSgenome

①1~22番染色体のみ取扱いたい場合。  
染色体番号の数が大きくなるほど配列長  
が短くなっている傾向が一目瞭然ですね。

## 9. インストール済みのヒト("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム配列をmulti-FASTAファイルで保存したい場合:

2013年12月にリリースされたGenome Reference Consortium GRCh38です。R ver. 3.1.0とBioconductor ver. 2.14以上の環境で実行可能です。

```
out_f <- "hoge9.fasta" #出力ファイル名を指定してout_fに格納
param <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定

#必要なパッケージをロード
library(param, character.only=T)

#前処理(paramで指定したパッケージ中のオブジェクトをリスト化)
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param, sep="."))
genome <- eval(parse(text=tmp))

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta")
```

```
R Console
[454] 170537 TCTACTCTC...GGGGAATTC HSCHR19KIR_FH15_A...
[455] 177381 GATCTATCT...GGGGAATTC HSCHR19KIR_RP5_B...
> fasta[1:22]
A DNASTringSet instance of length 22
      width seq names
[1] 248956422 NNNNNNNNNN...NNNNNNNNNN 1
[2] 242193529 NNNNNNNNNN...NNNNNNNNNN 2
[3] 198295559 NNNNNNNNNN...NNNNNNNNNN 3
[4] 190214555 NNNNNNNNNN...NNNNNNNNNN 4
[5] 181538259 NNNNNNNNNN...NNNNNNNNNN 5
...
[18] 80373285 NNNNNNNNNN...NNNNNNNNNN 18
[19] 58617616 NNNNNNNNNN...NNNNNNNNNN 19
[20] 64444167 NNNNNNNNNN...NNNNNNNNNN 20
[21] 46709983 NNNNNNNNNN...NNNNNNNNNN 21
[22] 50818468 NNNNNNNNNN...NNNNNNNNNN 22
> |
```

①X, Y, およびミトコンドリア配列も含めたい場合。②配列の並びの確認は試行錯誤。最初からわかっていたわけではありません。R画面上で眺めるほうが、全体像を把握しやすい。

# BSgenome

## 9. インストール済みのヒト("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム配列をmulti-FASTAファイルで保存したい場合:

2013年12月にリリースされたGenome Reference Consortium GRCh38です。R ver. 3.1.0とBioconductor ver. 2.14以上の環境で実行可能です。

```

out_f <- "hoge9.fasta"
param <- "BSgenome.Hsapiens.NCBI.GRCh38"

#必要なパッケージをロード
library(param, character.only=T)

#前処理(paramで指定したパッケージ中のオ
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param, sep="."))
genome <- eval(parse(text=tmp))

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, fo

```

#出力ファイル名を指定してout\_fに格納  
#パッケージ名を指定

```

R Console
[21] 46709983 NNNNNNNNNN...NNNNNNNNN 21
[22] 50818468 NNNNNNNNNN...NNNNNNNNN 22
> fasta[1:25]
A DNASTringSet instance of length 25
width seq
[1] 248956422 NNNNNNNNNN...NNNNNNNNN 1
[2] 242193529 NNNNNNNNNN...NNNNNNNNN 2
[3] 198295559 NNNNNNNNNN...NNNNNNNNN 3
[4] 190214555 NNNNNNNNNN...NNNNNNNNN 4
[5] 181538259 NNNNNNNNNN...NNNNNNNNN 5
...
[21] 46709983 NNNNNNNNNN...NNNNNNNNN 21
[22] 50818468 NNNNNNNNNN...NNNNNNNNN 22
[23] 156040895 NNNNNNNNNN...NNNNNNNNN X
[24] 57227415 NNNNNNNNNN...NNNNNNNNN Y
[25] 16569 GATCACAGGT...ATCACGATG MT
> |

```

# BSgenome

X, Y, およびミトコンドリア配列までのサブセットを hoge10.fasta で保存したい場合。①上矢印キーを何回か押して、ファイルに保存するためのコマンドを出し、②と③の水色下線部分を変更すればよい。

## 9. インストール済みのヒト("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム

2013年12月にリリースされた Genome Reference Consortium GRCh38 です。R ver. 3.1.0 と Bioconductor ver. 2.14 以上の環境で実行可能です。

```

out_f <- "hoge9.fasta"
param <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定

#必要なパッケージをロード
library(param, character.only=T)

#前処理(paramで指定したパッケージ中のオブジェクト)
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param, sep="."))
genome <- eval(parse(text=tmp))

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=80)

```

```

R Console
[21] 46709983 NNNNNNNNNN...NNNNNNNNN 21
[22] 50818468 NNNNNNNNNN...NNNNNNNNN 22
> fasta[1:25]
A DNASTringSet instance of length 25
      width seq names
[1] 248956422 NNNNNNNNNN...NNNNNNNNN 1
[2] 242193529 NNNNNNNNNN...NNNNNNNNN 2
[3] 198295559 NNNNNNNNNN...NNNNNNNNN 3
[4] 190214555 NNNNNNNNNN...NNNNNNNNN 4
[5] 181538259 NNNNNNNNNN...NNNNNNNNN 5
...
[21] 46709983 NNNNNNNNNN...NNNNNNNNN 21
[22] 50818468 NNNNNNNNNN...NNNNNNNNN 22
[23] 156040895 NNNNNNNNNN...NNNNNNNNN X
[24] 57227415 NNNNNNNNNN...NNNNNNNNN Y
[25] 16569 GATCACAGGT...ATCACGATG MT
> writeXStringSet(fasta, file=out_f, format="fasta", width=80)

```



# BSgenome

②と③をこんな感じで変更。実行後にhoge9.fastaよりも若干ファイルサイズの小さいhoge10.fastaが生成されていることが確認できるはず。決してテキストエディタで開かないで!

## 9. インストール済みのヒト("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム配列をmulti-FASTAファイルで保存したい場合:

2013年12月にリリースされたGenome Reference Consortium GRCh38です。R ver. 3.1.0とBioconductor ver. 2.14以上の環境で実行可能です。

```
out_f <- "hoge9.fasta" #出力ファイル名を指定してout_fに格納
param <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定

#必要なパッケージをロード
library(param, character.only=T)

#前処理(paramで指定したパッケージ中のオ
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param, sep="."))
genome <- eval(parse(text=tmp))

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta")
```

```
R Console
[22] 50818468 NNNNNNNNNN...NNNNNNNNN 22
> fasta[1:25]
A DNAStringSet instance of length 25
      width seq          names
[1] 248956422 NNNNNNNNNN...NNNNNNNNN 1
[2] 242193529 NNNNNNNNNN...NNNNNNNNN 2
[3] 198295559 NNNNNNNNNN...NNNNNNNNN 3
[4] 190214555 NNNNNNNNNN...NNNNNNNNN 4
[5] 181538259 NNNNNNNNNN...NNNNNNNNN 5
...
[21] 46709983 NNNNNNNNNN...NNNNNNNNN 21
[22] 50818468 NNNNNNNNNN...NNNNNNNNN 22
[23] 156040895 NNNNNNNNNN...NNNNNNNNN X
[24] 57227415 NNNNNNNNNN...NNNNNNNNN Y
[25] 16569 GATCACAGGT...ATCACGATG MT
> writeXStringSet(fasta[1:25], file="hoge10.fasta", format="fasta")
> |
```



①例題10。②様々な記述形式があります。やらなくていいです。決してテキストエディタで開かないで!

# BSgenome

イントロ | 一般 | 配列取得 | ゲノム配列 | BSgenome NEW

BSgenomeパッケージを用いて様々な生物種のゲノム配列を取得するやり方を示します。ミヤマハタザオ (A. mellifera)、シロイヌナズナ (A. thaliana)、ウシ (B. taurus)、線虫 (C. elegans)、犬

①

10. インストール済みのヒト ("BSgenome.Hsapiens.NCBI.GRCh38")のゲノム配列のmulti-FASTAファイルで保存したい場合:

一部を抽出して保存するやり方です。このパッケージ中の染色体の並びが既知(chr1, 2, ..., chr22, chrX, chrY, and MT)であるという前提です。

```

out_f <- "hoge10.fasta" #出力ファイル名を指定してout_fに格納
param <- "BSgenome.Hsapiens.NCBI.GRCh38" #パッケージ名を指定
param_range <- 1:25 #抽出したい範囲を指定

```

②

```

#必要なパッケージロード
library(param, character.only=T)

#前処理(paramで指定したパッケージ中)
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param))
genome <- eval(parse(text=tmp))
genome

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#後処理(フィルタリング)
obj <- param_range
fasta <- fasta[obj]
fasta

```

```

R Console
width seq names
[1] 248956422 NNNNNNNNNNNN...NNNNNNNNNNN 1
[2] 242193529 NNNNNNNNNNNN...NNNNNNNNNNN 2
[3] 198295559 NNNNNNNNNNNN...NNNNNNNNNNN 3
[4] 190214555 NNNNNNNNNNNN...NNNNNNNNNNN 4
[5] 181538259 NNNNNNNNNNNN...NNNNNNNNNNN 5
...
[21] 46709983 NNNNNNNNNNNN...NNNNNNNNNNN 21
[22] 50818468 NNNNNNNNNNNN...NNNNNNNNNNN 22
[23] 156040895 NNNNNNNNNNNN...NNNNNNNNNNN X
[24] 57227415 NNNNNNNNNNNN...NNNNNNNNNNN Y
[25] 16569 GATCACAGGTC...CATCACGATG MT
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=$
> |

```

# BSgenome

例題9。26番目以降の配列は、ヒトゲノムの一部ではあるものの、まだ割り当てられる染色体が定まっていないものたちです。メタゲノム解析などでヒトゲノムにマップされないリードのみ取扱いたい場合には、利用可能な全配列をマッピング時のリファレンスとして用いるのが自然だと思います。

## 9. インストール済みのヒト ("BSgenome.Hsapiens.NCBI.GRCh38")

2013年12月にリリースされた Genome Reference Consortium (GRCh38) 実行可能です。

```
out_f <- "hoge9.fasta"
param <- "BSgenome.Hsapiens.NCBI.GRCh38"

#必要なパッケージをロード
library(param, character.only=T)

#前処理(paramで指定したパッケージ中のオブジェクト)
#tmp <- unlist(strsplit(param, "."))
tmp <- ls(paste("package", param, sep="."))
genome <- eval(parse(text=tmp))

#本番
fasta <- getSeq(genome)
names(fasta) <- seqnames(genome)

#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta")
```

#出力ファイル名を指定してout\_fに格納  
#パッケージ名を指定

```
R Console
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width$width)
> fasta
A DNAStringSet instance of length 455
      width seq          names
[1] 248956422 NNNNNNNNNN...NNNNNNNNNN 1
[2] 242193529 NNNNNNNNNN...NNNNNNNNNN 2
[3] 198295559 NNNNNNNNNN...NNNNNNNNNN 3
[4] 190214555 NNNNNNNNNN...NNNNNNNNNN 4
[5] 181538259 NNNNNNNNNN...NNNNNNNNNN 5
...
[451] 200773 TCTACTCTC...GGGGAATTC HSCHR19KIR_FH08_B...
[452] 170148 TTTCTTTCT...GGGGAATTC HSCHR19KIR_FH13_A...
[453] 215732 TGTGGTGAG...GGGGAATTC HSCHR19KIR_FH13_B...
[454] 170537 TCTACTCTC...GGGGAATTC HSCHR19KIR_FH15_A...
[455] 177381 GATCTATCT...GGGGAATTC HSCHR19KIR_RP5_B...
```

# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証



①ゲノム配列と②アノテーション情報があれば、任意の領域の配列を取り出すことができます。

# プロモーター配列取得

- ・ [イントロ](#) | [一般](#) | [Tips](#) | [拡張子は同じで任意の文字を追加して保存](#) (last modified 2013/09/26)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [公共DBから](#) (last modified 2014/05/28)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [BSgenome](#) (last modified 2015/02/19) **NEW**
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [公共DBから](#) (last modified 2014/04/02)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [BSgenomeとTxDbから](#) (last modified 2015/02/20) **NEW**
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [GenomicFeatures\(Lawrence 2013\)](#) (last modified 2015/02/20)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [トランスクリプトーム配列](#) | [公共DBから](#) (last modified 2014/04/02)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [トランスクリプトーム配列](#) | [biomaRt\(Durink 2009\)](#) (last modified 2015/02/20) **NEW**

## [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [BSgenomeとTxDbから](#)

**NEW**

①

②

ゲノム配列([BSgenome](#))パッケージとアノテーション情報([TxDb](#))パッケージを用いて様々な生物種のプロモーター配列(転写開始点近傍配列; 上流配列)を取得するやり方を示します。2014年4月リリースの [Bioconductor 2.14](#)以降の推奨手順では、ゲノムのパッケージ(例:[BSgenome.Hsapiens.UCSC.hg19](#))と対応するアノテーションパッケージ(例:[TxDb.Hsapiens.UCSC.hg19.knownGene](#))の両方を読み込ませる必要がありますので、2015年2月に記述内容を大幅に変更しました。ヒトなどの主要なパッケージ以外はおそらくデフォルトではインストールされていないので、「パッケージがインストールされていない」的なエラーが出た場合は、[個別パッケージのインストール](#)を参考にして予め利用したいパッケージのインストールを行ってから再挑戦してください。出力はmulti-FASTAファイルです。現状では、ゼブラフィッシュ([danRer7](#))はゲノムパッケージ([BSgenome.Drerio.UCSC.danRer7](#))は存在しますが、対応するTxDbパッケージが存在しないので、どこかからGFFファイルを取得して[makeTranscriptDbFromGFF](#)関数などを利用してTranscriptDbオブジェクトを得るなどする必要があります。

# プロモーター配列取得

例題1。作業ディレクトリはどこでもよいが基本はデスクトップ上のhoge。転写開始点上流1000塩基を取得したい場合。width列は配列長に相当。出力ファイルを見なくてもうまくいっているだろうと思える。

## 1. ヒト (hg19) の場合:

ゲノムパッケージ ([BSgenome.Hsapiens.UCSC.hg19](#)) と対応するアノテーションパッケージ ([TxDb.Hsapiens.UCSC.hg19.knownGene](#)) を読み込んで、転写開始点上流1000塩基分を取得する

```
out_f <- "hoge1.fasta" #出力ファイル名を指定してout_fに格納
param_bsgenome <- "BSgenome.Hsapiens.UCSC.hg19" #パッケージ名を指定(BSgenome系のゲノムパ
param_txdb <- "TxDb.Hsapiens.UCSC.hg19.knownGene" #パッケージ名を指定(TxDB系のアノテーシ
param_upstream <- 1000 #転写開始点上流の塩基配列数を指定
```

```
#前処理(指定したパッケージ中のオブジェクト)
library(param_bsgenome, character.only=T)
tmp <- ls(paste("package", param_bsgenome))
genome <- eval(parse(text=tmp))
```

```
library(param_txdb, character.only=T)
tmp <- ls(paste("package", param_txdb))
txdb <- eval(parse(text=tmp))
```

```
#本番
gn <- sort(genes(txdb))
hoge <- flank(gn, width=param_upstream)
fasta <- getSeq(genome, hoge)
fasta
```

```
#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=param_upstream)
```

```
R Console
> fasta <- getSeq(genome, hoge) #指定した範囲の$
> fasta #確認してるだけ$
A DNAStringSet instance of length 23056
  width seq names
[1] 1000 ACACATGCTA...TTTCGTTAA 100287102
[2] 1000 GCTATTATCA...AATAACTCT 79501
[3] 1000 GAATTAGGCT...AAGGCGGGG 643837
[4] 1000 CGGGGAGCCC...CTTGGGCCA 148398
[5] 1000 CGGCGGGGCT...AGCGGCGGG 339451
...
[23052] 1000 GGTGAGCCAA...ATCTCAGCC 283788
[23053] 1000 AGCCCTCCAC...CTCTCCAAC 100507412
[23054] 1000 CGGGGCCAG...CCTGGCTGC 728410
[23055] 1000 CGGGGCCAG...CCTGGCTGC 100653046
[23056] 1000 CAGGCTGAGC...CTCACC GCG 100288687
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=param_upstream)
> |
```

# プロモーター配列取得

例題3。線虫はBSgenomeとTxDbの両方が提供されています。予め2つのパッケージを個別にインストールしておけば、上流500～下流20塩基の範囲の配列を取得できる。

## 3. 線虫(ce6)の場合:

ゲノムパッケージ([BSgenome.Celegans.UCSC.ce6](#))と対応するアノテーションパッケージ([TxDb.Celegans.UCSC.ce6.ensGene](#))を読み込んで、転写開始点上流500塩基から下流20塩基までの範囲を取得するやり方です。

```
out_f <- "hoge3.fasta" #出力ファイル名を指定してout_fに格納
param_bsgenome <- "BSgenome.Celegans.UCSC.ce6" #パッケージ名を指定(BSgenome系のゲノムパッ
param_txdb <- "TxDb.Celegans.UCSC.ce6.ensGene" #パッケージ名を指定(TxDb系のアノテーション
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20
```

```
#前処理(指定したパッケージ中のオブジェクト)
library(param_bsgenome, character.only=T)
tmp <- ls(paste("package", param_bsgenome))
genome <- eval(parse(text=tmp))
```

```
library(param_txdb, character.only=T)
tmp <- ls(paste("package", param_txdb))
txdb <- eval(parse(text=tmp))
```

```
#本番
gn <- sort(genes(txdb))
hoge <- promoters(gn, upstream=param_upstream,
                 downstream=param_downstream)
fasta <- getSeq(genome, hoge)
fasta
```

```
R Console
> fasta <- getSeq(genome, hoge) #指定した範囲の$
> fasta #確認してるだけ$
A DNAStringSet instance of length 27928
  width seq names
[1] 520 TTGCGCGTAA...CCGCGTCAC Y74C9A.2
[2] 520 GCAGATAATT...ACGATGGAA Y74C9A.1
[3] 520 AGAGGAATTT...CGAGACAAA Y48G1C.12
[4] 520 TTACCTCCAG...CCCGATCCC Y48G1C.4
[5] 520 CTTCAATCCC...GCAGCAGCA Y48G1C.2
... ..
[27924] 520 GGGTGTTACA...ATAAAAAAG MTCE.29
[27925] 520 ATATCTGCAG...ATAATTCAG MTCE.30
[27926] 520 GTAGTATTTT...TTATATTAA MTCE.32
[27927] 520 TTTATATTAT...TGTTTAAAT MTCE.33
[27928] 520 GATTAATATT...TTAATAAAA MTCE.36
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width$
> |
```

# Tips: sessionInfo

自分の解析環境を知る。エラーが出てどうにもならない場合は、`sessionInfo()`実行結果をメール添付か本文中に記載して助けを求める。見る側は、②R本体のバージョン(この場合、3.2.0)や、パッケージが読み込まれているか、あるいはパッケージのバージョンをチェックする。

```
R Console
> sessionInfo()
R version 3.2.0 (2015-04-16)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows 7 x64 (build 7601) Service Pack 1

locale:
[1] LC_COLLATE=Japanese_Japan.932  LC_CTYPE=Japanese_Japan.932
[3] LC_MONETARY=Japanese_Japan.932 LC_NUMERIC=C
[5] LC_TIME=Japanese_Japan.932

attached base packages:
[1] stats4      parallel    stats       graphics    grDevices   utils       datasets    methods
[9] base

other attached packages:
 [1] TxDb.Celegans.UCSC.ce6.ensGene_3.1.2      BSgenome.Celegans.UCSC.ce6_1.4.0
 [3] BSgenome.Hsapiens.UCSC.hg19_1.4.0        BSgenome_1.36.0
 [5] rtracklayer_1.28.3                       Biostrings_2.36.1
 [7] XVector_0.8.0                             RCurl_1.95-4.6
 [9] bitops_1.0-6                              TxDb.Hsapiens.UCSC.hg19.knownGene_3.1.2
[11] GenomicFeatures_1.20.1                   AnnotationDbi_1.30.1
[13] Biobase_2.28.0                           GenomicRanges_1.20.3
[15] GenomeInfoDb_1.4.0                       IRanges_2.2.1
[17] S4Vectors_0.6.0                          BiocGenerics_0.14.0

loaded via a namespace (and not attached):
```

# Tips: sessionInfo

これは、うまくプロモーター配列を取得できたときのR環境。①例題1のヒト(H. sapiens)プロモーター配列を取得したときと、②例題3の線虫(C. elegans)プロモーター配列を取得したときに読み込んだパッケージが表示されているのがわかる。

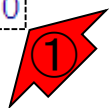
```
R Console
> sessionInfo()
R version 3.2.0 (2015-04-16)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows 7 x64 (build 7601) Service Pack 1

locale:
[1] LC_COLLATE=Japanese_Japan.932  LC_CTYPE=Japanese_Japan.932
[3] LC_MONETARY=Japanese_Japan.932 LC_NUMERIC=C
[5] LC_TIME=Japanese_Japan.932

attached base packages:
[1] stats4      parallel    stats       graphics    grDevices   utils       datasets    methods
[9] base

other attached packages:
[1] TxDb.Celegans.UCSC.ce6.ensGene_3.1.2   BSgenome.Celegans.UCSC.ce6_1.4.0
[3] BSgenome.Hsapiens.UCSC.hg19_1.4.0     BSgenome_1.36.0
[5] rtracklayer_1.28.3                    Biostrings_2.36.1
[7] XVector_0.8.0                          RCurl_1.95-4.6
[9] bitops_1.0-6                           TxDb.Hsapiens.UCSC.hg19.knownGene_3.1.2
[11] GenomicFeatures_1.20.1                 AnnotationDbi_1.30.1
[13] Biobase_2.28.0                         GenomicRanges_1.20.3
[15] GenomeInfoDb_1.4.0                     IRanges_2.2.1
[17] S4Vectors_0.6.0                       BiocGenerics_0.14.0

loaded via a namespace (and not attached):
```



# パッケージのインストール

「必要最小限プラスアルファ」の推奨インストール手順通りにインストールしたヒトの多くはうまくいっている。

- インストール | R本体 | 過去版 | [Mac用](#) (last modified 2015/03/22)
- インストール | Rパッケージ | [ほぼ全て\(20GB以上?!\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [必要最小限プラスアルファ\(数GB?!\)](#) (last modified 2015/07/24) 推奨
- インストール | Rパッケージ | [必要最小限プラスアルファ\(アグリバイオ/居室のみ\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [必要最小限\(数GB?!\)](#) (last modified 2015/05/25)

## インストール | Rパッケージ | 必要最小限プラスアルファ(数GB?!) NEW

- (削除予定) [基本的な利用](#)

(Rで)塩基配列解析、(Rで)マイクロアレイデータ解析 中で利用するパッケージ、プラスアルファのパッケージをインストールするやり方です。Rパッケージの2大リポジトリであるCRANとBioconductor から提供されているパッケージ群のうち、一部のインストールに相当しますので、相当短時間でインストールが完了します。

SAFEではなくsafeパッケージの間違い  
( )m limmaの取得先が間違っってCRAN追加しました。

### 1. R本体を起動

### 2. CRANから提供されているパッケージ

以下を「R コンソール画面」上でコピー  
合は自分のいる場所から近いサイトを

```
#(Rで)塩基配列解析で主に利用
install.packages("PoissonSeq")
install.packages("samr")
install.packages("seqinr")
```

```
#(Rで)マイクロアレイデータ解析
```

```
biocLite("SKADD", suppressUpdates=TRUE)
biocLite("TCC", suppressUpdates=TRUE)
biocLite("TxDb.Celegans.UCSC.ce6.ensGene", suppressUpdates=TRUE)
biocLite("TxDb.Hsapiens.UCSC.hg19.knownGene", suppressUpdates=TRUE)
biocLite("TxDb.Hsapiens.UCSC.hg38.knownGene", suppressUpdates=TRUE)
biocLite("TxDb.Mmusculus.UCSC.mm10.knownGene", suppressUpdates=TRUE)
biocLite("TxDb.Rnorvegicus.UCSC.rn5.refGene", suppressUpdates=TRUE)
#biocLite("yeastRNASeq", suppressUpdates=TRUE)
```

```
#(Rで)マイクロアレイデータ解析で利用
biocLite("affy", suppressUpdates=TRUE)
biocLite("agilo", suppressUpdates=TRUE)
```

### 4. Bioconductorから提供されているパッケージ群のインストール

ゲノム配列パッケージです。一つ一つの容量が尋常でないため、必要に応じてテキストエディタなどに予めコピーしておき、いらぬゲノムパッケージを削除してからお使いください。

```
source("http://bioconductor.org/biocLite.R")#おまじない
biocLite("BSgenome.Athaliana.TAIR.TAIR2", suppressUpdates=TRUE)#シロイヌナズナゲノム
biocLite("BSgenome.Celegans.UCSC.ce6", suppressUpdates=TRUE)#線虫ゲノム
biocLite("BSgenome.Drerio.UCSC.danRer7", suppressUpdates=TRUE)#ゼブラフィッシュゲノム
biocLite("BSgenome.Hsapiens.NCBI.GRCh38", suppressUpdates=TRUE)#ヒトゲノム(GRCh38)
biocLite("BSgenome.Hsapiens.UCSC.hg19", suppressUpdates=TRUE)#ヒトゲノム(hg19)
biocLite("BSgenome.Mmusculus.UCSC.mm10", suppressUpdates=TRUE)#マウスゲノム(mm10)
```

# 全体像

リポジトリ(パッケージ提供元)	パッケージ名	R本体	
		R ver. 3.1.3	R ver. 3.2.0
Bioconductor	<a href="#">TxDb.Celegans.UCSC.ce6.ensGene</a>		3.1.2
Bioconductor	<a href="#">BSgenome.Celegans.UCSC.ce6</a>		1.4.0
CRAN	XML		
CRAN	RCurl		1.95-4.6
CRAN	checkmate		-

# R ver. 3.1.3でこけた

①例題3を門田のR ver. 3.1.3でやった結果。②XMLおよび③RCurlパッケージがないためにエラー祭りとなった(爆)

```
R Console
> out_f <- "hoge3.fasta" #出力ファイル名を指定してo$
> param_bsgenome <- "BSgenome.Celegans.UCSC.ce6" #パッケージ名を指定$
> param_txdb <- "TxDb.Celegans.UCSC.ce6.ensGene" #パッケージ名を指定$
> param_upstream <- 500 #転写開始点上流の塩基配列$
> param_downstream <- 20 #転写開始点下流の塩基配列$
>
> #前処理(指定したパッケージ中のオブジェクト名をgenomeおよびtxdbに$
> library(param_bsgenome, character.only=T) #指定したパッケージの読$
要求されたパッケージ BSgenome をロード中です
要求されたパッケージ rtracklayer をロード中です
Error in loadNamespace(j <- i[[1L]], c(lib.loc, .libPaths()), versi$
  'XML' という名前のパッケージはありません ②
エラー: パッケージ 'rtracklayer' をロードできませんでした
> tmp <- ls(paste("package", param_bsgenome, sep=":")) #指定したパッ$
以下にエラー as.environment(pos) :
  検索リストに "package:BSgenome.Celegans.UCSC.ce6" という項目はあ$
> genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクト$
以下にエラー eval(parse(text = tmp)) :
  引数 'expr' の評価中にエラーが起きました (関数 'eval' に対するメ$
>
> library(param_txdb, character.only=T) #指定したパッケージの読み$
要求されたパッケージ GenomicFeatures をロード中です
Error in loadNamespace(j <- i[[1L]], c(lib.loc, .libPaths()), versi$
  'RCurl' という名前のパッケージはありません ③
エラー: パッケージ 'GenomicFeatures' をロードできませんでした

Jul 30 2015
```



# XMLのインストール

①XMLパッケージのインストール。XML自体はCRANから提供されているが、CRANパッケージ用のinstall.packages関数ではなく、Bioconductor用のbiocLite関数でもCRANパッケージのXMLをインストール可能であるというTipsを兼ねて紹介。

R Console

```
> source("http://bioconductor.org/biocLite.R")#
Bioconductor version 3.0 (BiocInstaller 1.16.5), ?biocLite for help
A new version of Bioconductor is available after installing the
most recent version of R; see http://bioconductor.org/install
> biocLite("XML", suppressUpdates=TRUE)
BioC_mirror: http://bioconductor.org
Using Bioconductor version 3.0 (BiocInstaller 1.16.5), R version
3.1.3.
Installing package(s) 'XML'
URL 'http://cran.fhcrc.org/bin/windows/contrib/3.1/XML_3.98-1.3.z$
Content type 'application/zip' length 4293309 bytes (4.1 MB)
開かれた URL
downloaded 4.1 MB

package 'XML' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:\Users\kadota\AppData\Local\Temp\RtmpER54eM\downloaded_p$
> |
```

①

# RCurlのインストール

①RCurlパッケージのインストール。  
RCurlもCRANから提供されているが、  
biocLite関数でインストール可能。

```
R Console
> source("http://bioconductor.org/biocLite.R") #おまじない
Bioconductor version 3.0 (BiocInstaller 1.16.5), ?biocLite for help
A new version of Bioconductor is available after installing the
most recent version of R; see http://bioconductor.org/install
> biocLite("XML", suppressUpdates=TRUE)
BioC_mirror: http://bioconductor.org
Using Bioconductor version 3.0 (BiocInstaller 1.16.5), R version
3.1.3.
Installing package(s) 'XML'
URL 'http://cran.fhcrc.org/bin/windows/contrib/3.1/XML_3.98-1.3.z$
Content type 'application/zip' length 4293309 bytes (4.1 MB)
開かれた URL
downloaded 4.1 MB

package 'XML' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:\Users\kadota\AppData\Local\Temp\RtmpER54eM\downloaded_p$
> biocLite("RCurl", suppressUpdates=TRUE)
```



# RCurlのインストール

①RCurlパッケージのインストール。  
RCurlもCRANから提供されているが、  
biocLite関数でインストール可能。

```
R Console
package 'XML' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
  C:\Users\kadota\AppData\Local\Temp\RtmpER54eM\downloaded_p$
> biocLite("RCurl", suppressUpdates=TRUE)
BioC_mirror: http://bioconductor.org
Using Bioconductor version 3.0 (BiocInstaller 1.16.5), R version
  3.1.3.
Installing package(s) 'RCurl'
  URL 'http://cran.fhcrc.org/bin/windows/contrib/3.1/RCurl_1.95-4.7$
Content type 'application/zip' length 2858866 bytes (2.7 MB)
開かれた URL
downloaded 2.7 MB

package 'RCurl' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
  C:\Users\kadota\AppData\Local\Temp\RtmpER54eM\downloaded_p$
> |
```



# R ver. 3.1.3で再挑戦

意気揚々と、再度例題3をコピー。しかし  
①今度はcheckmateというパッケージがないからダメだと言われていることがわかる。

```
R Console
> out_f <- "hoge3.fasta" #出力ファイル名を指定して$
> param_bsgenome <- "BSgenome.Celegans.UCSC.ce6" #パッケージ名を指$
> param_txdb <- "TxDb.Celegans.UCSC.ce6.ensGene" #パッケージ名を指$
> param_upstream <- 500 #転写開始点上流の塩基配列$
> param_downstream <- 20 #転写開始点下流の塩基配列$
>
> #前処理 (指定したパッケージ中のオブジェクト名をgenomeおよびtxdbに$
> library(param_bsgenome, character.only=T) #指定したパッケージの読$
要求されたパッケージ BSgenome をロード中です
要求されたパッケージ rtracklayer をロード中です
Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = $
 'checkmate' という名前のパッケージはありません ①
エラー: パッケージ 'rtracklayer' をロードできませんでした
> tmp <- ls(paste("package", param_bsgenome, sep=":")) #指定したパ$
以下にエラー as.environment(pos) :
  検索リストに "package:BSgenome.Celegans.UCSC.ce6" という項目は$
> genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクト$
以下にエラー eval(parse(text = tmp)) :
  引数 'expr' の評価中にエラーが起きました (関数 'eval' に対する$
>
> library(param_txdb, character.only=T) #指定したパッケージの読み$
要求されたパッケージ GenomicFeatures をロード中です
Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = $
```

# checkmateのインストール

①checkmateパッケージのインストール。  
checkmateもCRANから提供されているが、biocLite関数でインストール可能。

```
R Console
> fasta <- getSeq(genome, hoge) #指定した範囲の塩基配列情$
以下にエラー (function (classes, fdef, mtable) :
  unable to find an inherited method for function 'getSeq' for s$
> fasta #確認してるだけです
エラー: オブジェクト 'fasta' がありません
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fas$
以下にエラー is(x, "XStringSet") : オブジェクト 'fasta' がありま$
> biocLite("checkmate", suppressUpdates=TRUE) ①
BioC_mirror: http://bioconductor.org
Using Bioconductor version 3.0 (BiocInstaller 1.16.5), R version
  3.1.3.
Installing package(s) 'checkmate'
  URL 'http://cran.fhcrc.org/bin/windows/contrib/3.1/checkmate_1.6.$
Content type 'application/zip' length 393202 bytes (383 KB)
  開かれた URL
downloaded 383 KB

package 'checkmate' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
  C:\Users\kadota\AppData\Local\Temp\RtmpER54eM\downloaded_p$
> |
```

意気揚々と、再度例題3をコピー。  
無事パッケージのロードに成功!!

# R ver. 3.1.3で再々挑戦



```
R Console
> out_f <- "hoge3.fasta" #出力ファイル名を指定して$
> param_bsgenome <- "BSgenome.Celegans.UCSC.ce6" #パッケージ名を指$
> param_txdb <- "TxDb.Celegans.UCSC.ce6.ensGene" #パッケージ名を指$
> param_upstream <- 500 #転写開始点上流の塩基配列$
> param_downstream <- 20 #転写開始点下流の塩基配列$
>
```

```
> #前処理 (指定したパッケージ中のオブジェクト名を
> library(param_bsgenome, character.
要求されたパッケージ BSgenome をロード中です
要求されたパッケージ rtracklayer をロード中
> tmp <- ls(paste("package", par
> genome <- eval(parse(text=tmp)
>
> library(param_txdb, character.
要求されたパッケージ GenomicFeatures を
> tmp <- ls(paste("package", par
> txdb <- eval(parse(text=tmp))
>
> #本番
> gn <- sort(genes(txdb))
> hoge <- promoters(gn, upstream:
+ downstream=param_do
> fasta <- getSeq(genome, hoge)
> fasta
```

```
R Console
> fasta #確認してるだけ$
A DNASTringSet instance of length 27928
width seq names $
[1] 520 TTGCGCGTAAAATAT...TCTTGCCGCGTCAC Y74C9A.2
[2] 520 GCAGATAATTGAGGA...TCGAAACGATGGAA Y74C9A.1
[3] 520 AGAGGAATTCACCG...TGCGACGAGACAAA Y48G1C.12
[4] 520 TTACCTCCAGTGATT...AAATTCCCAGATCCC Y48G1C.4
[5] 520 CTTCAATCCCACACT...TGGAAGCAGCAGCA Y48G1C.2
... ..
[27924] 520 GGGTGTTACACTATG...TTAGTATAAAAAAAG MTCE.29
[27925] 520 ATATCTGCAGTATTT...TTAGTATAATTTCAG MTCE.30
[27926] 520 GTAGTATTTTTCCAA...ATAGTTTATATTAA MTCE.32
[27927] 520 TTTATATTATTACGG...TTACTTGTTTAAAT MTCE.33
[27928] 520 GATTAATATTTTTTG...TTAGTTTAAATAAAA MTCE.36
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", wid$
> |
```

# この程度は...

## バイオインフォマティクス人材育成カリキュラム(次世代シーケンサ) | NGSハンズオン講習会2015 NEW

NGSハンズオン講習会を2015年7月22日-8月6日の11日間で開催します。予定通り受講申込多数のため、予備日(8月26日、27日、28日)も実施することになっています。



はじめに(全)

- ・ 講習基本
- ・ 平成

### 2015年7月22日(2015.07.22):Bio-Linux 8とRのインストール状況確認(門田幸二、寺田 透)

書籍「[日本乳酸菌学会誌](#)」についてで示した通りのPC環境を構築しておきましょう。連載第1-3回、および第4回のウェブ資料W6-5までの予習は必須です。Rについても同様です。一週間程度はぎっちり時間をかけて予習しておきましょう。7/22は、以下に示すようなことができる(わかる)ようになっていることの確認を自分でしてもらう日です。門田自身全てを完全に把握しているわけではありませんし、ウェブ資料のページ数も膨大ですので、どこにどのようなことが書かれていたかの全体像の俯瞰やチェックリストという位置づけでもあります(頻りに更新しているのでときどきロードしましょう):

#### ・ Bio-Linux 8

#### ・ R

1. [インストール](#)についてをよく読み、ここに書いてある手順に従って2015年4月4日以降にインストールを行った
2. [インストール](#)についてで書いてある内容はBio-Linux8(ゲストOS)とは無関係であり、WindowsやMacintosh(つまりホストOS)上で行う作業である
3. ファイルの拡張子(.txtや.docxなど)はちゃんと表示されている
4. Rの起動と終了ができる。終了時に表示されるメッセージにうろたえない
5. R本体だけでなくRパッケージ群のインストールもちゃんと行った
6. Rパッケージ群のインストール確認も行い、エラーが出ないことを確認した
7. library関数を用いたRパッケージのロード中に、別のパッケージがないことに起因するエラーメッセージが出ることもあるが、必要なパッケージを個別にインストールするやり方を知っている
8. [基本的な利用法](#)をよく読み、予習を行った
9. 作業ディレクトリの変更ができる
10. 例題ファイルのダウンロード時に、拡張子が勝手に変わることがあるので注意する
11. 慣れないうちは、`getwd()`と`list.files()`を打ち込むことで、作業ディレクトリと入力ファイルの存在確認を行う
12. エラーに遭遇した際、「ありがちなミス1-4」に当てはまっていないかどうか自分で確認

R ver. 3.1.3での、ロードされているパッケージ群のバージョン情報に注目!

# sessionInfo (R 3.1.3)

```
R Console  
> sessionInfo()  
R version 3.1.3 (2015-03-09)  
Platform: x86_64-w64-mingw32/x64 (64-bit)  
Running under: Windows 7 x64 (build 7601) Service Pack 1  
  
locale:  
[1] LC_COLLATE=Japanese_Japan.932 LC_CTYPE=Japanese_Japan.932  
[3] LC_MONETARY=Japanese_Japan.932 LC_NUMERIC=C  
[5] LC_TIME=Japanese_Japan.932  
  
attached base packages:  
[1] stats4 parallel stats graphics grDevices utils datasets methods  
[9] base  
  
other attached packages:  
[1] TxDb.Celegans.UCSC.ce6.ensGene_3.0.0 GenomicFeatures_1.18.7  
[3] BSgenome.Celegans.UCSC.ce6_1.4.0 BSgenome_1.34.1  
[5] rtracklayer_1.26.3 BiocInstaller_1.16.5  
[7] Biostrings_2.34.1 XVector_0.6.0  
[9] AnnotationDbi_1.28.2 Biobase_2.26.0  
[11] GenomicRanges_1.18.4 GenomeInfoDb_1.2.4  
[13] IRanges_2.0.1 S4Vectors_0.4.0  
[15] BiocGenerics_0.12.1  
  
loaded via a namespace (and not attached):
```



# sessionInfo (R 3.1.3)

特に①～③のパッケージは、④2015年7月24日にインストールしたばかりであり、前述のR ver. 3.2.0実行結果よりも後にインストールしたものである。

other attached packages:

```
[1] TxDb.Celegans.UCSC.ce6.ensGene_3.0.0 GenomicFeatures_1.18.7
[3] BSgenome.Celegans.UCSC.ce6_1.4.0      BSgenome_1.34.1
[5] rtracklayer_1.26.3                    BiocInstaller_1.16.5
[7] Biostrings_2.34.1                     XVector_0.6.0
[9] AnnotationDbi_1.28.2                  Biobase_2.26.0
[11] GenomicRanges_1.18.4                  GenomeInfoDb_1.2.4
[13] IRanges_2.0.1                          S4Vectors_0.4.0
[15] BiocGenerics_0.12.1
```

loaded via a namespace (and not attached):

```
[1] base64enc_0.1-2          BatchJobs_1.6          BBmisc_1.9
[4] BiocParallel_1.0.3      biomaRt_2.22.0        bitops_1.0-6
[7] brew_1.0-6              checkmate_1.6.1       codetools_0.2-11
[10] crayon_1.3.0            DBI_0.3.1              digest_0.6.8
[13] fail_1.2                foreach_1.4.2          GenomicAlignments_1.2.2
[16] iterators_1.0.7         magrittr_1.5           memoise_0.2.1
[19] RCurl_1.95-4.7          Rsamtools_1.18.3      RSQLite_1.0.0
[22] sendmailR_1.2-1         stringi_0.4-1         stringr_1.0.0
[25] testthat_0.10.0         tools_3.1.3           XML_3.98-1.3
[28] zlibbioc_1.12.0
```

> date()

```
[1] "Fri Jul 24 21:21:59 2015"
```

> |

# 全体像

リポジトリ(パッケージ提供元)	パッケージ名	R本体	
		R ver. 3.1.3	R ver. 3.2.0
Bioconductor	<a href="#">TxDb.Celegans.UCSC.ce6.ensGene</a>	3.0.0	3.1.2
Bioconductor	<a href="#">BSgenome.Celegans.UCSC.ce6</a>	1.4.0	1.4.0
CRAN	XML	3.98-1.3	
CRAN	RCurl	1.95-4.7	1.95-4.6
CRAN	checkmate	1.6.1	-

# sessionInfo (R 3.2.0)

BSgenomeパッケージはver. 1.4.0と同じ。しかしTxDbパッケージは、ver. 3.1.2 (R ver. 3.2.0)とver. 3.0.0 (R ver. 3.1.3)で異なっていることがわかる。

```
R Console
> sessionInfo()
R version 3.2.0 (2015-04-16)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows 7 x64 (build 7601) Service Pack 1

locale:
[1] LC_COLLATE=Japanese_Japan.932  LC_CTYPE=Japanese_Japan.932
[3] LC_MONETARY=Japanese_Japan.932  LC_NUMERIC=C
[5] LC_TIME=Japanese_Japan.932

attached base packages:
[1] stats4      parallel    stats       graphics    grDevices   utils       datasets    methods
[9] base

other attached packages:
[1] TxDb.Celegans.UCSC.ce6.ensGene_3.1.2  BSgenome.Celegans.UCSC.ce6_1.4.0
[3] BSgenome.Hsapiens.UCSC.hg19_1.4.0    BSgenome_1.36.0
[5] rtracklayer_1.28.3                   Biostrings_2.36.1
[7] XVector_0.8.0                         RCurl_1.95-4.6
[9] bitops_1.0-6                          TxDb.Hsapiens.UCSC.hg19.knownGene_3.1.2
[11] GenomicFeatures_1.20.1                AnnotationDbi_1.30.1
[13] Biobase_2.28.0                        GenomicRanges_1.20.3
[15] GenomeInfoDb_1.4.0                   IRanges_2.2.1
[17] S4Vectors_0.6.0                       BiocGenerics_0.14.0

loaded via a namespace (and not attached):
```

# sessionInfo (R 3.2.0)

①XML、および②RCurlのバージョン。  
checkmateパッケージは、存在の有無は不明であるが、少なくともロードされておらず、R ver. 3.2.0では必要とされていないようだ。

R Console

other attached packages:

```
[1] TxDb.Celegans.UCSC.ce6.ensGene_3.1.2    BSgenome.Celegans.UCSC.ce6_1.4.0
[3] BSgenome.Hsapiens.UCSC.hg19_1.4.0      BSgenome_1.36.0
[5] rtracklayer_1.28.3                      Biostrings_2.36.1
[7] XVector_0.8.0                           RCurl_1.95-4.6
[9] bitops_1.0-6                             TxDb.Hsapiens.UCSC.hg19.knownGene_3.1.2
[11] GenomicFeatures_1.20.1                  AnnotationDbi_1.30.1
[13] Biobase_2.28.0                          GenomicRanges_1.20.3
[15] GenomeInfoDb_1.4.0                      IRanges_2.2.1
[17] S4Vectors_0.6.0                         BiocGenerics_0.14.0
```

loaded via a namespace (and not attached):

```
[1] zlibbioc_1.14.0          GenomicAlignments_1.4.1 BiocParallel_1.2.1
[4] tools_3.2.0              DBI_0.3.1                lambda.r_1.1.7
[7] futile.logger_1.4.1      futile.options_1.0.0     biomaRt_2.24.0
[10] RSQLite_1.0.0            Rsamtools_1.20.2         XML_3.98-1.1
```

> |

# 全体像を把握

①赤枠がごく最近(2015年7月24日)インストールしたCRAN提供パッケージたち。バージョン番号がR ver. 3.2.0のものに比べて大きいので、確かに最新版なのだろう。

リポジトリ(パッケージ提供元)	パッケージ名	R本体	
		R ver. 3.1.3	R ver. 3.2.0
Bioconductor	TxDb.Celegans.UCSC.ce6.ensGene	3.0.0	3.1.2
Bioconductor	BSgenome.Celegans.UCSC.ce6	1.4.0	1.4.0
CRAN	XML	3.98-1.3	3.98-1.1
CRAN	RCurl	1.95-4.7	1.95-4.6
CRAN	checkmate	1.6.1	-



# 全体像を把握

これらの結果から、2つの結論が導き出される。① CRAN提供パッケージは、R本体のバージョンに関係なく、常に最新版がインストールされる。② R本体のバージョンが上がると必要とされないパッケージが出ることもある(checkmate)

リポジトリ(パッケージ提供元)	パッケージ名	R本体	
		R ver. 3.1.3	R ver. 3.2.0
Bioconductor	TxDb.Celegans.UCSC.ce6.ensGene	3.0.0	3.1.2
Bioconductor	BSgenome.Celegans.UCSC.ce6	1.4.0	1.4.0
CRAN	XML	3.98-1.3	3.98-1.1
CRAN	RCurl	1.95-4.7	1.95-4.6
CRAN	checkmate	1.6.1	-



# R ver. 3.1.3で再挑戦

ここで、もう一度「R ver. 3.1.3で再挑戦」したときのエラーメッセージを読み返す。①  
黒枠は、BSgenome.Celegans.UCSC.ce6  
パッケージのロード部分。

R Console

```
> out_f <- "hoge3.fasta" #出力ファイル名を指定して$
> param_bsgenome <- "BSgenome.Celegans.UCSC.ce6" #パッケージ名を指$
> param_txdb <- "TxDb.Celegans.UCSC.ce6.ensGene" #パッケージ名を指$
> param_upstream <- 500 #転写開始点上流の塩基配列$
> param_downstream <- 20 #転写開始点下流の塩基配列$
>
> #前処理 (指定したパッケージ中のオブジェクト名をgenomeおよびtxdbに$
> library(param_bsgenome, character.only=T) #指定したパッケージの読$
要求されたパッケージ BSgenome をロード中です
要求されたパッケージ rtracklayer をロード中です
Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = $
  'checkmate' という名前のパッケージはありません
エラー: パッケージ 'rtracklayer' をロードできませんでした
> tmp <- ls(paste("package", param_bsgenome, sep=":")) #指定したパ$
以下にエラー as.environment(pos) :
  検索リストに "package:BSgenome.Celegans.UCSC.ce6" という項目は$
> genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクト$
以下にエラー eval(parse(text = tmp)) :
  引数 'expr' の評価中にエラーが起きました (関数 'eval' に対する$
>
> library(param_txdb, character.only=T) #指定したパッケージの読み$
要求されたパッケージ GenomicFeatures をロード中です
Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = $
```



①

# R ver. 3.1.3で再挑戦

文脈から、以下のように理解する。①  
BSgenome.Celegans.UCSC.ce6は内部的にrtracklayerを利用しているため、rtracklayerもロードしようとしている。しかし、②rtracklayerは内部的にcheckmateを利用している。このためcheckmateもロードしようとしたが、ckeckmateが未インストールのためこけた、というのが事の顛末。

```
R Console
> out_f <- "hoge3.fasta" #出力ファイル名を
> param_bsgenome <- "BSgenome.Celegans.UCSC.ce6" #パッケ
> param_txdb <- "TxDb.Celegans.UCSC.ce6.ensGene" #パッケ
> param_upstream <- 500 #転写開始点上流
> param_downstream <- 20 #転写開始点下流
>
> #前処理 (指定したパッケージ中のオブジェクト名をgenomeおよびtxdbに$
> library(param_bsgenome, character.only=T) #指定したパッケージの読$
要求されたパッケージ BSgenome をロード中です
要求されたパッケージ rtracklayer をロード中です
Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = $
'checkmate' という名前のパッケージはありません
エラー: パッケージ 'rtracklayer' をロードできませんでした
> tmp <- ls(paste("package", param_bsgenome, sep=":")) #指定したパ$
以下にエラー as.environment(pos) :
  検索リストに "package:BSgenome.Celegans.UCSC.ce6" という項目は$
> genome <- eval(parse(text=tmp)) #文字列tmpをRオブジェクト$
以下にエラー eval(parse(text = tmp)) :
  引数 'expr' の評価中にエラーが起きました (関数 'eval' に対する$
>
> library(param_txdb, character.only=T) #指定したパッケージの読み$
要求されたパッケージ GenomicFeatures をロード中です
Error in loadNamespace(i, c(lib.loc, .libPaths()), versionCheck = $
```

①

②



# 依存関係を含めると...

複雑。rtracklayer ver. 1.28.3ではcheckmateパッケージを内部的に使わなくなったのだろう、と理解する。たとえばBSgenomeとTxDbのバージョンがR本体の間で不変であったとしても、内部的に用いるパッケージのバージョンや、それに起因する実行結果の違いもありうる。

リポジトリ(パッケージ提供元)	パッケージ名	R本体	
		R ver. 3.1.3	R ver. 3.2.0
Bioconductor	TxDb.Celegans.UCSC.ce6.ensGene	3.0.0	3.1.2
CRAN	RCurl	1.95-4.7	1.95-4.6
Bioconductor	BSgenome.Celegans.UCSC.ce6	1.4.0	1.4.0
Bioconductor	rtracklayer	1.26.3	1.28.3
CRAN	XML	3.98-1.3	3.98-1.1
CRAN	checkmate	1.6.1	-



# 現実的な対策

私は少しの結果の違いは(突き詰めていっても不毛なので)気にしない。論文中に記載するのはR本体と主要なパッケージ(この場合BSgenomeとTxDb)のみが基本。

リポジトリ(パッケージ提供元)	パッケージ名	R本体	
		R ver. 3.1.3	R ver. 3.2.0
Bioconductor	TxDb.Celegans.UCSC.ce6.ensGene	3.0.0	3.1.2
CRAN	RCurl	1.95-4.7	1.95-4.6
Bioconductor	BSgenome.Celegans.UCSC.ce6	1.4.0	1.4.0
Bioconductor	rtracklayer	1.26.3	1.28.3
CRAN	XML	3.98-1.3	3.98-1.1
CRAN	checkmate	1.6.1	-

# R本体のバージョンは重要

## ■ 近年のリリース頻度

### □ R本体 (<http://www.r-project.org/>)

- 2015-06-18にver. 3.2.1をリリース
- 2015-04-16にver. 3.2.0をリリース
- 2015-03-09にver. 3.1.3をリリース
- 2014-10-31にver. 3.1.2をリリース
- ...
- 2012-03-30にver. 2.15.0をリリース
- ...

### □ Bioconductor (<http://bioconductor.org/>)は半年ごとにリリース

- 2015-04にver. 3.1をリリース (R ver. 3.2.1で動作確認)、提供パッケージ数: 1,024
- 2014-10にver. 3.0をリリース (R ver. 3.1.1で動作確認)、提供パッケージ数: 934
- 2014-04にver. 2.14をリリース (R ver. 3.1.0で動作確認)、提供パッケージ数: 824
- 2013-10にver. 2.13をリリース (R ver. 3.0で動作確認)、提供パッケージ数: 750
- 2013-04にver. 2.12をリリース (R ver. 3.0で動作確認)、提供パッケージ数: 672
- 2012-10にver. 2.11をリリース (R ver. 2.15.1で動作確認)、提供パッケージ数: 608
- 2012-04にver. 2.10をリリース (R ver. 2.15.0で動作確認)、提供パッケージ数: 553
- ...

R本体のバージョンがわかれば、Bioconductor提供パッケージのバージョンも概ね定まる。例えば、①R ver. 3.1.2では、どのタイミングで(例えば2015年8月に)パッケージのインストールを行おうが、Bioconductor ver. 3.0で提供されているパッケージのバージョンとなる。決して2015年4月に公開されたBioconductor ver. 3.1で提供されているパッケージがインストールされることはない。

①

# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# プロモーター配列取得

multi-FASTA形式のゲノム配列ファイルとGFF3形式などの一般的なアノテーションファイルがあれば、それを利用することもできる。②例題7の乳酸菌のプロモーター配列取得を示します。

- ・ [イントロ](#) | [一般](#) | [Tips](#) | [拡張子は同じで任意の文字を追加して保存](#) (last modified 2013/09/26)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [公共DBから](#) (last modified 2014/05/28)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [BSgenome](#) (last modified 2015/02/19) **NEW**
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [公共DBから](#) (last modified 2014/04/02)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [BSgenomeとTxDbから](#) (last modified 2015/02/20) **NEW**
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [GenomicFeatures\(Lawrence 2013\)](#) **①** (last modified 2015/02/20)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [トランスクリプトーム配列](#) | [公共DBから](#) (last modified 2014/04/02)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [トランスクリプトーム配列](#) | [biomaRt\(Drydenk 2009\)](#) (last modified 2015/02/20) **NEW**

## [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [GenomicFeatures\(Lawrence\\_2013\)](#) **NEW**

[GenomicFeatures](#)パッケージを主に用いてプロモーター配列(転写開始点近傍配列)を得るやり方を示します。ここでは、[イントロ | 一般 | 配列取得 | ゲノム配列\(BSgenomeから\)](#)で指定可能なゲノムと [イントロ | NGS | アノテーション情報取得 | TranscriptDb | GenomicFeatures\(Lawrence 2013\)](#)で作成可能なTranscriptDbオブジェクトを入力として、getPromoterSeqという関数を用いて、転写開始点から任意の[上流xxx塩基, 下流yyy塩基]分の塩基配列を取得して、FASTA形式ファイルで保存するやり方を示しています。multi-FASTA形式のゲノム配列ファイルとGFF3形式のアノテーションファイルのみからプロモーター配列を取得するやり方も示しています。

### 1. ヒト **②** 7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得する場合:

ヒトゲノム  
オブジェクトを

```
out_f
param_
param2
param3
param_
param_
#必要な
library
```

[Ensembl](#) (Flicek et al., Nucleic Acids Res., 2014)から提供されている [Lactobacillus casei 12A](#)の multi-FASTA形式ゲノム配列ファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome.Chromosome.fa](#))と GFF3形式のアノテーションファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.chromosome.Chromosome.gff3](#))を読み込むやり方です。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入力ファイル
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3"#入力ファイル名
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定

#必要なパッケージをロード
library(Rsamtools) #パッケージの読み込み
library(Biostrings) #パッケージの読み込み
library(GenomicFeatures) #パッケージの読み込み
```

2015年3月5-6日のHPCIハンズオン講習会時のスライド(2015.03.05)です。このときはR ver. 3.1.2でうまくいった。

# プロモーター配列取得

7. 乳酸菌ゲノム(*Lactobacillus casei* 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得する場合:

Ensembl (Flicek et al., *Nucleic Acids Res.*, 2014)から提供されている *Lactobacillus casei* 12Aの multi-FASTA形式ゲノム配列ファイル (*Lactobacillus casei* 12a.GCA\_000309565.2.25.chromosome.Chromosome.gff3)と GFF3形式のアノテーションファイル (*Lactobacillus casei* 12a.GCA\_000309565.2.25.dna.chromosome.Chromosome.fa)を読み込むやり方です。  
makeTranscriptDbFromGFF関数実行時に用いているuseGenesAsTranscripts=Tは、原核生物(prokaryotes)に代表される「転写物 = 遺伝子」の場合に指定します。デフォルトはFです。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入力ファイル名を指
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3"#入力ファイル名を指
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定
```

```
#必要なパッケージをロード
library(Rsamtools) #パッケージの
library(Biostrings) #パッケージの
library(GenomicFeatures) #パッケージの
```

```
#前処理(アノテーション情報を取得)
txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptD
format="gff3", useGenesAsTranscripts=T)#
txdb #確認してるだ
```

```
#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb)) #遺伝子の座標
hoge <- promoters(gn, upstream=param_upstream, #指定
downstream=param_downstream)#指定した範囲
hoge #確認してるだ
obj <- (ranges(hoge)@start >= 0) #条件を満たす
```

```
R Console
> #本番(配列取得)
> fasta <- getSeq(FaFile(in_f1), hoge) #指定した範囲の塩基配列情$
> names(fasta) <- names(hoge) #description情報を追加して$
> fasta #確認してるだけです
A DNAStringSet instance of length 2680
      width seq
[1] 520 CTTGAAAGCCTTGAAA...ATTTACGATCACCCG gene:LCA12A_0618
[2] 520 GAACGTGATAATGTGCG...AACAAATCGAAATCAC gene:LCA12A_0619
[3] 520 GTTTCACGATCAATCG...ACTGGATCACTTGGT gene:LCA12A_0620
[4] 520 CCACAGCATTGGCTGT...GGACAAGAAAGAAAC gene:LCA12A_0621
[5] 520 ACTGATCATTATGACC...TGATCGCCAAGAAAG gene:LCA12A_0622
...
[2676] 520 AAGTGACCGTTGCTTA...TGCATCTATCACTGC gene:LCA12A_0612
[2677] 520 ACACAGAAACTTTATG...AACCTTTCAAGGCAG gene:LCA12A_0613
[2678] 520 TGCGCAAGTCATATCG...AAATCGTCAAATCAA gene:LCA12A_0614
[2679] 520 CGAATTCCTGTGATT...TGCGCGTTTCGGCGG gene:LCA12A_0615
[2680] 520 GCGCCAATTCACGCTC...AACCAAACGGACTTT gene:LCA12A_0616
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width=50)#fast$
> |
```

# プロモーター配列取得

当時(2015年3月; R ver. 3.1.2)は、GFF3形式のアノテーションファイル読み込み時に①makeTranscriptDbFromGFF関数、および②useGenesAsTranscriptsオプション(デフォルトはF)を利用していた。

## 7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流20塩基]のプロモ

Ensembl (Flicek et al., Nucleic Acids Res., 2014)から提供されている [Lactobacillus](#)

([Lactobacillus casei 12a.GCA\\_000309565.2.25.chromosome.Chromosome.gff3](#))と

([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome.Chromosome.fa](#))を読み込むやり方です。

makeTranscriptDbFromGFF関数実行時に用いているuseGenesAsTranscripts=Tは、原核生物(prokaryotes)に代表される「転写物 = 遺伝子」の場合に指定します。デフォルトはFです。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入力ファイル名を指
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3"#入力ファイル名を指
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定
```

#必要なパッケージをロード

```
library(Rsamtools) #パッケージの読み込み
library(Biostrings) #パッケージの読み込み
library(GenomicFeatures) #パッケージの読み込み
```

#前処理(アノテーション情報を取得)

```
txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトを取得してtxdbに格納
                               format="gff3", useGenesAsTranscripts=T)#TranscriptDbオブジェクトを取得してtxdbに格納
txdb #確認してるだけです
```

#前処理(欲しい領域の座標情報取得)

```
gn <- sort(genes(txdb)) #遺伝子の座標情報を取得
hoge <- promoters(gn, upstream=param_upstream, #指定した範囲の座標情報を取得
                 downstream=param_downstream)#指定した範囲の座標情報を取得
hoge #確認してるだけです
obj <- (ranges(hoge)@start >= 0) #条件を満たすかどうかを判定した結果をobjに格納(座標が0よりも小さい)
```

# プロモーター配列取得

```

in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fasta" #入力ファイル
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3" #入力ファイル名
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定

#必要なパッケージをロード
library(Rsamtools) #パッケージの読み込み
library(Biostrings) #パッケージの読み込み
library(GenomicFeatures) #パッケージの読み込み

#前処理(アノテーション情報を取得)
txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトを取得してtxdbに格納
                               format="gff3", useGenesAsTranscripts=T) #TranscriptDbオブジェクトを取得してtxdbに格納
txdb #確認してるだけです

```

```

#前処理(欲しい領域の座標情報取
gn <- sort(genes(txdb))
hoge <- promoters(gn, upstr
                 downstream=param
hoge
obj <- (ranges(hoge))@start

```

```

R Console
> #前処理 (アノテーション情報を取得)
> txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトを$
+                               format="gff3", useGenesAsTranscripts=T) #TranscriptDbオブ$
Prepare the 'metadata' data frame ... metadata: OK
Warning messages:
1: 'makeTranscriptDbFromGFF' is deprecated.
   Use 'makeTxDbFromGFF' instead.
   See help("Deprecated")
2: 'useGenesAsTranscripts' is ignored and deprecated
3: In matchCircularity(seqlevels(gr), circ_seqs) :
   None of the strings in your circ_seqs argument match your seqnames.
> txdb #確認してるだけです
TxDb object:

```





# プロモーター配列取得

警告メッセージの中身。①

makeTranscriptDbFromGFF関数は削除予定だ。makeTxDbFromGFF関数を使用せよ。②useGenesAsTranscriptsオプションを与えても無視するし、このオプション自体も削除予定だよ。といった具合。R ver. 3.2.0では、今のところ③で示すようにうまく読み込めてはいる。

R Console

```
> #前処理 (アノテーション情報を取得)
> txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbFromGFF
+                               format="gff3", useGenesAsTranscripts=T) #TxDb
Prepare the 'metadata' data frame ... metadata: OK
Warning messages:
1: 'makeTranscriptDbFromGFF' is deprecated.
Use 'makeTxDbFromGFF' instead.
2: 'useGenesAsTranscripts' is ignored and deprecated
3: In matchCircularity(seqlevels(gr), circ_seqs) :
  None of the strings in your circ_seqs argument match your seqnames.
> txdb
```



#確認してるだけです

```
TxDb object:
# Db type: TxDb
# Supporting package: GenomicFeatures
# Data source: Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3
# Organism: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2799
# exon_nrow: 2800
# cds_nrow: 2681
# Db created by: GenomicFeatures package from Bioconductor
# Creation time: 2015-07-25 15:30:24 +0900 (Sat, 25 Jul 2015)
# GenomicFeatures version at creation time: 1.20.1
# RSQLite version at creation time: 1.0.0
# DBSCHEMAVERSION: 1.1
```



# プロモーター配列取得

2015年7月25日現在。③R ver. 3.1.2で正常動作したコードを(#のコメントアウトで)残しつつ、④R ver. 3.2.0で動作確認済みのコードがデフォルトとなるようにしている。

- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [ゲノム配列](#) | [BSgenome](#)(last modified 2015/04/22)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [公共DBから](#)(last modified 2014/04/02)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [BSgenomeとTxDbから](#)(last modified 2015/02/20)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [GenomicFeatures\(Lawrence 2013\)](#)(last modified 2015/05/09)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [トランスクリプトーム配列](#) | [公共DBから](#)(last modified 2015/05/09)
- ・ [イントロ](#) | [一般](#) | [配列取得](#) | [トランスクリプトーム配列](#) | [GenomicFeatures\(Lawrence 2013\)](#)(last modified 2015/05/09)



## イントロ | 一般 | 配列取得 | プロモーター配列 | GenomicFeatures(Lawrence\_2013) NEW

[GenomicFeatures](#)パッケージを主に用いてプロモーター配列(転写開始点近傍配列)を得るやり方を示します。ここでは、[イントロ | 一般 | 配列取得 | ゲノム配列\(BSgenomeから\)](#)で指定可能なゲノムと [イントロ | NGS | アノテーション情報取得 | TranscriptDb](#) [GenomicFeatures\(Lawrence 2013\)](#)で作成可能なTranscriptDbオブジェクトを入力して、転写開始点から任意の上流

xxx塩基, 下流yyy塩基]の配列を得るやり方を示します。ここでは、[イントロ | 一般 | 配列取得 | ゲノム配列\(BSgenomeから\)](#)で指定可能なゲノムと [イントロ | NGS | アノテーション情報取得 | TranscriptDb](#) [GenomicFeatures\(Lawrence 2013\)](#)で作成可能なTranscriptDbオブジェクトを入力して、転写開始点から任意の上流



### 7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得する場合:

[Ensembl](#) (Flicek et al., Nucleic Acids Res., 2014)から提供されている [Lactobacillus casei 12A](#)の multi-FASTA形式ゲノム配列ファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome.Chromosome.fa](#))と GFF3形式のアノテーションファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.chromosome.Chromosome.gff3](#))を読み込むやり方です。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入力ファイル
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3"#入力ファイル
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定
```

```
#必要なパッケージをロード
library(Rsamtools) #パッケージの読み込み
library(Biostrings) #パッケージの読み込み
library(GenomicFeatures) #パッケージの読み込み
```

```
#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトを取得してtxdbに格納
# format="gff3", useGenesAsTranscripts=T)#TranscriptDbオブジェクトを取得してtxdbに格納
txdb <- makeTxDbFromGFF(in_f2, format="gff3")#TranscriptDbオブジェクトを取得してtxdbに格納
txdb #確認してるだけです
```



1. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

2. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

3. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

4. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

5. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

6. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

7. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

8. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

9. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

10. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

11. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.

12. ヒト("BSgenomeHomoSapiens")ヒトゲノム ver.


# R ver. 3.2.0で実行

makeTranscriptDbFromGFF関数  
やuseGenesAsTranscriptsオプション  
を使っていないので、①それ系  
の警告は出ていないことがわかる。

```
> #前処理 (アノテーション情報を取得)
> #txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブ$
> #
> #           format="gff3", useGenesAsTranscripts=T) #Trans$
> txdb <- makeTxDbFromGFF(in_f2, format="gff3") #Transcript$
Prepare the 'metadata' data frame ... metadata: OK
```

Warning message:

```
In matchCircularity(seqlevels(c
  None of the strings in your c
```

```
> txdb
TxDb object: 
# Db type: TxDb
# Supporting package: GenomicFe
# Data source: Lactobacillus_ca
# Organism: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2799
# exon_nrow: 2800
```

```
R Console
> fasta #確認してるだけで$
A DNAStringSet instance of length 2726
      width seq
[1] 520 CTTGAAAGCCTTGA...TTTACGATCACCCG LCA12A_0618
[2] 520 GAACTGATAATGTG...ACAATCGAAATCAC LCA12A_0619
[3] 520 GTTTCACGATCAAT...CTGGATCACTTGGT LCA12A_0620
[4] 520 CCACAGCATTGGCT...GACAAGAAAGAAAC LCA12A_0621
[5] 520 ACTGATCATTATGA...GATCGCCAAGAAAG LCA12A_0622
...
[2722] 520 AAGTGACCGTTGCT...GCATCTATCACTGC LCA12A_0612
[2723] 520 ACACAGAAACTTTA...ACCTTTCAAGGCAG LCA12A_0613
[2724] 520 TGCGCAAGTCATAT...AATCGTCAAATCAA LCA12A_0614
[2725] 520 CGAATTCCTGTGA...GCGCGTTTCGGCGG LCA12A_0615
[2726] 520 GCGCCAATTCACGC...ACCAAACGGACTTT LCA12A_0616
>
> #ファイルに保存
> writeXStringSet(fasta, file=out_f, format="fasta", width$
> |
```

# R ver. 3.1.2で実行

R ver. 3.1.2(正確にはGenomicFeatures ver. 1.18.7)当時は、まだmakeTxDbFromGFF関数が実装されていないので、①makeTxDbFromGFF関数を見つけることができないのは当たり前。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chr1.gff3"
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.gff3"
out_f <- "hoge7.fasta"
param_upstream <- 500
param_downstream <- 20
```

#出力ファイル名を指定してout\_fに格納  
#転写開始点上流の塩基配列数を指定  
#転写開始点下流の塩基配列数を指定

#必要なパッケージをロード

```
library(Rsamtools)
library(Biostrings)
library(GenomicFeatures)
```

#前処理(アノテーション情報を取得)

```
#txdb <- makeTranscriptDbFromGFF(in_f1,
#                               format="gff3", useGenesAsTranscripts=T)
txdb <- makeTxDbFromGFF(in_f2, format="gff3", useGenesAsTranscripts=T)
```

#前処理(欲しい領域の座標情報取得)

```
gn <- sort(genes(txdb))
hoge <- promoters(gn, upstream=param_upstream,
                 downstream=param_downstream)
hoge
```

```
R Console
The following object is masked from 'package:GenomeInfoDb':
  species

警告メッセージ:
1: パッケージ 'GenomicFeatures' はバージョン 3.1.3 の R の下で$
2: パッケージ 'AnnotationDbi' はバージョン 3.1.3 の R の下で造$
>
> #前処理 (アノテーション情報を取得)
> #txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクト$
> #                               format="gff3", useGenesAsTranscripts=T) #Transcript
> txdb <- makeTxDbFromGFF(in_f2, format="gff3") #TranscriptDbオブ
① エラー: 関数 "makeTxDbFromGFF" を見つけることができませんでした
> txdb
エラー: オブジェクト 'txdb' がありません
> package.version("GenomicFeatures")
[1] "1.18.7"
> packageVersion("GenomicFeatures")
[1] '1.18.7'
> |
```

# R ver. 3.1.3で実行

R ver. 3.1.3(正確にはGenomicFeatures ver. 1.18.7)当時は、まだmakeTxDbFromGFF関数が実装されていないので、①makeTxDbFromGFF関数を見つけることができないのは当たり前。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chr1.chromosomes.gff3"
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosomes.gff3"
out_f <- "hoge7.fasta"
param_upstream <- 500
param_downstream <- 20
```

#出力ファイル名を指定してout\_fに格納  
#転写開始点上流の塩基配列数を指定  
#転写開始点下流の塩基配列数を指定

#必要なパッケージをロード

```
library(Rsamtools)
library(Biostrings)
library(GenomicFeatures)
```

#前処理(アノテーション情報を取得)

```
#txdb <- makeTranscriptDbFromGFF(in_f1,
#                               format="gff3", useGenesAsTranscripts=T)
txdb <- makeTxDbFromGFF(in_f2, format="gff3", useGenesAsTranscripts=T)
txdb
```

#前処理(欲しい領域の座標情報取得)

```
gn <- sort(genes(txdb))
hoge <- promoters(gn, upstream=param_upstream, downstream=param_downstream)
hoge
```

```
R Console

次のパッケージを付け加えます: 'AnnotationDbi'

The following object is masked from 'package:GenomeInfoDb':
  species

>
> #前処理(アノテーション情報を取得)
> #txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクト$
> #                               format="gff3", useGenesAsTranscripts=T) #TranscriptDbオブジェクト
> txdb <- makeTxDbFromGFF(in_f2, format="gff3") #TranscriptDbオブジェクト
① エラー: 関数 "makeTxDbFromGFF" を見つけることができませんでした
> txdb
エラー: オブジェクト 'txdb' がありません
> package.version("GenomicFeatures")
[1] "1.18.7"
> packageVersion("GenomicFeatures")
[1] '1.18.7'
> |
```

# 定期的にバージョンアップ

バグの修正や新たな機能がどんどん追加されている。最新版の利用をお勧め。推奨の毎年5月と11月ごろにバージョンアップをちゃんとやったヒトは、R ver. 3.2.0以上なはず。

## ■ 近年のリリース頻度

### □ R本体 (<http://www.r-project.org/>)

- 2015-06-18にver. 3.2.1をリリース
- 2015-04-16にver. 3.2.0をリリース
- 2015-03-09にver. 3.1.3をリリース
- 2014-10-31にver. 3.1.2をリリース
- ...
- 2012-03-30にver. 2.15.0をリリース
- ...

### □ Bioconductor (<http://bioconductor.org/>)は半年ごとにリリース

- 2015-04にver. 3.1をリリース (R ver. 3.2.1で動作確認)、提供パッケージ数: 1,024
- 2014-10にver. 3.0をリリース (R ver. 3.1.1で動作確認)、提供パッケージ数: 934
- 2014-04にver. 2.14をリリース (R ver. 3.1.0で動作確認)、提供パッケージ数: 824
- 2013-10にver. 2.13をリリース (R ver. 3.0で動作確認)、提供パッケージ数: 750
- 2013-04にver. 2.12をリリース (R ver. 3.0で動作確認)、提供パッケージ数: 672
- 2012-10にver. 2.11をリリース (R ver. 2.15.1で動作確認)、提供パッケージ数: 608
- 2012-04にver. 2.10をリリース (R ver. 2.15.0で動作確認)、提供パッケージ数: 553
- ...

# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# データの型

簡単にいうと、「型 = 見た目」です。例えば、①  
 makeTxDbFromGFF関数は、(GFF3形式の)アノテーション  
 ファイルを読み込んだものを②TxDbという形式のオブジェク  
 トとして格納しています。重要なのは、細かい中身の理解で  
 はなく、ちゃんと読み込めているっぽいことを確認すること

7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流200塩基]の配列

Ensembl (Flicek et al., Nucleic Acids Res., 2014)から提供されている  
 配列ファイル(Lactobacillus casei 12a.GCA\_000309565.2.25.dna.gz)とアノ  
 テーションファイル(Lactobacillus casei 12a.GCA\_000309565.2.25.gff3)

```

in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.gz"
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.gff3"
out_f <- "hoge7.fasta" #出力ファイル
param_upstream <- 500 #転写開始位置の上流塩基数
param_downstream <- 200 #転写開始位置の下流塩基数

#必要なパッケージをロード
library(Rsamtools) #パッケージ
library(Biostrings) #パッケージ
library(GenomicFeatures) #パッケージ

#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクト
#                               format="gff3", useGenesAsTranscripts=TRUE)
txdb <- makeTxDbFromGFF(in_f2, format="gff3") #確認して

#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb)) #遺伝子座標
hoge <- promoters(gn, upstream=param_upstream, downstream=param_downstream) #指定した領域の座標
hoge
    
```

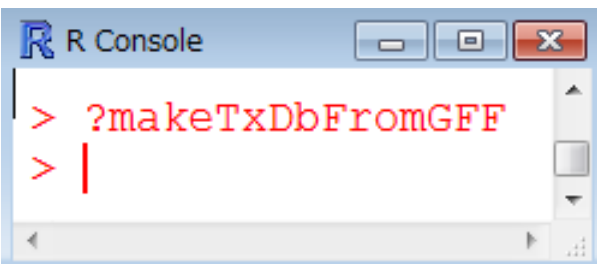
```

R Console
> txdb <- makeTxDbFromGFF(in_f2, format="gff3") #TranscriptDbオブジェクトを作成
Prepare the 'metadata' data frame ... metadata: OK
Warning message:
In matchCircularity(seqlevels(gr), circ_seqs) :
  None of the strings in your circ_seqs argument matched any of the
  seqlevels
> txdb #確認してる$
TxDb object:
# Db type: TxDb
# Supporting package: GenomicFeatures
# Data source: Lactobacillus_casei_12a.GCA_000309565.2.25.gff3
# Organism: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2799
# exon_nrow: 2800
# cds_nrow: 2681
# Db created by: GenomicFeatures package from Bioconductor
# Creation time: 2015-07-25 17:25:09 +0900 (Sat, 25 Jul 2015)
# GenomicFeatures version at creation time: 1.20.1
# RSQLite version at creation time: 1.0.0
# DBSCHEMAVERSION: 1.1
>
    
```





# データの型



```
R Console  
> ?makeTxDbFromGFF  
> |
```

`makeTxDbFromGFF {GenomicFeatures}` R Documentation

Make a TxDb object from annotations available as a GFF3 or GTF file

## Description

The `makeTxDbFromGFF` function allows the user to make a [TxDb](#) object from transcript annotations available as a GFF3 or GTF file.

## Usage

```
makeTxDbFromGFF(file,  
                 format=c("gff3", "gtf"),  
                 exonRankAttributeName=NA,  
                 gffGeneIdAttributeName=NA,  
                 chrominfo=NA,  
                 dataSource=NA,  
                 organism=NA,  
                 circ_seqs=DEFAULT_CIRC_SEQS,  
                 mirBaseBuild=NA,  
                 useGenesAsTranscripts=FALSE,  
                 gffTxName="mRNA",  
                 species=NA)
```

# データの型

①これはGenomicFeaturesパッケージが提供している関数。②GFF3形式以外にGTF形式も読み込めるようだ。③GTF形式ファイルの場合はformat="gtf"と書けばよさそう

```
R Console  
> ?makeTxDbFromGFF  
> |
```

```
makeTxDbFromGFF {GenomicFeatures} R Documentation
```

Make a TxDb object from annotations available as a GFF3 or GTF file

## Description

The `makeTxDbFromGFF` function allows the user to make a [TxDb](#) object from transcript annotations available as a GFF3 or GTF file.

## Usage

```
makeTxDbFromGFF(file,  
                 format=c("gff3", "gtf"),  
                 exonRankAttributeName=NA,  
                 gffGeneIdAttributeName=NA,  
                 chrominfo=NA,  
                 dataSource=NA,  
                 organism=NA,  
                 circ_seqs=DEFAULT_CIRC_SEQS,  
                 miRBaseBuild=NA,  
                 useGenesAsTranscripts=FALSE,  
                 gffTxName="mRNA",  
                 species=NA)
```

# データの型

`makeTxDbFromGFF` is a convenience function that reads a GFF file to the [makeTxDbFromGRanges](#) function.

## Value

A [TxDb](#) object.



## Author (s)

M. Carlson and H. Pages

## See Also

- [makeTxDbFromGRanges](#), which `makeTxDbFromGFF` is based on, for making a [TxDb](#) object from a [GRanges](#) object.
- The [import](#) function in the `rtracklayer` package (also used by `makeTxDbFromGFF` internally).
- [makeTxDbFromUCSC](#) and [makeTxDbFromBiomart](#) for convenient ways to make a [TxDb](#) object from UCSC or BioMart online resources.



htmlマニュアルの下部に移動。①`makeTxDbFromGFF`関数実行結果として得られるものはTxDbという形式のオブジェクト。②もちろんR上でアノテーション情報を格納する形式はTxDbだけではなく、GRangesという形式も存在する。GRanges → TxDbへの型変換は`makeTxDbFromGRanges`関数で可能であることがわかる。アノテーション情報源は、③UCSCや④BioMartが有名だが、そこからうまく情報を抽出してTxDbオブジェクトを得るための関数も存在することが分かる。こんな感じで関数の知識や利用の幅を広げていく。

# library(help=パッケージ名)

①指定したパッケージ中で利用可能な関数を概観できます。②下のほうに移動すると、③など(Rで)塩基配列解析で紹介していない様々な関数があることがわかります。

R Console

```
> library(help=GenomicFeatures)
> |
```

パッケージ 'GenomicFeatures' のドキュメント

パッケージ 'GenomicFeatures' の情報

記述:

Package:	DEFAULT_CIRC_SEQS	character vector: strings that are usually circular chromosomes
Title:	FeatureDb-class	FeatureDb objects
	TxDb-class	TxDb objects
	asBED, TxDb-method	Coerce to file format structures
Version:	extractTranscriptSeqs	Extract transcript sequences from chromosomes
Author:	extractUpstreamSeqs	Extract sequences upstream of a set of genes or transcripts
	features	Extract simple features from a FeatureDb object
	getPromoterSeq	Get gene promoter sequences
	id2name	Map internal ids to external names for a given feature type
	makeFeatureDbFromUCSC	Making a FeatureDb object from annotations available at the UCSC Genome Browser
	makeTxDb	Making a TxDb object from user supplied annotations

# パッケージ概観


②Bioconductor (ver. 3.1)上での GenomicFeaturesのトップページ。③ダウンロード数はトップ5%であることがわかる

- インポート | 一般 | 配列取得 | ゲノム配列 | [BSgenome](#)(last modified 2015/04/22)
- インポート | 一般 | 配列取得 | プロモーター配列 | [公共DBから](#)(last modified 2014/04/02)
- インポート | 一般 | 配列取得 | プロモーター配列 | [BSgenomeとTxDbから](#)(last modified 2015/02/20)
- インポート | 一般 | 配列取得 | プロモーター配列 | [GenomicFeatures\(Lawrence 2013\)](#) **①**
- インポート | 一般 | 配列取得 | トランスクリプトーム配列 | [公共DBから](#)(last modified 2015/05/09)
- インポート | 一般 | 配列取得 | トランスクリプトーム配列 | [GenomicFeatures\(Lawrence 2013\)](#)(last modified 2015/05/09)

**イントロ | 一般 | 配列取得 | プロモーター配列 | GenomicFeatures(Lawrence\_2013) NEW**

GenomicFeaturesパッケージを主に用いてプロモーター配列(転写開始点近傍配列)を得るやり方を示します。  
 ここでは、[イントロ | 一般 | 配列取得 | ゲノム配列\(BSgenomeから\)](#)で指定可能なゲノムと [イントロ | NGS | アノテーション情報取得 | TranscriptionStartSites](#)を組み合わせ、プロモーター配列を得ることが出来ます。

②  
 xxx塩基、下は、どの生状の指定が Chr4, Chr5, genome@se形式のアノ  
 1. ヒト("BSgenome.Hsapiens")  
 ヒトゲノム



Search:   
[Home](#)   [Install](#)   [Help](#)   [Developers](#)   [About](#)

Home » [Bioconductor 3.1](#) » [Software Packages](#) » GenomicFeatures

## GenomicFeatures

platforms **all**

downloads **top 5%**

posts **32 / 2 / 3 / 8**

in Bioc **5.5 years**

build **ok**

commits **15.17**

test coverage **70%**

### Tools for making and manipulating transcript centric annotations

Bioconductor version: Release (3.1)

A set of tools and methods for making and manipulating transcript centric annotations. With these tools

#### Documentation »

*Bioconductor*

- Package [vignettes](#) and manuals.
- [Workflows](#) for learning and use.
- [Course and conference](#) material.
- [Videos](#).
- Community [resources](#) and [tutorials](#).

R / [CRAN](#) packages and [documentation](#)

# パッケージ概観

少し下のほうに移動。①原著論文はきっちり引用。②パッケージのインストール法。③パッケージのドキュメント(マニュアルではない)をローカル環境で開く方法。④具体的なやり方。④の中身は⑤と同じだが、ネットワークかローカルかの違いがある。⑥利用可能な関数マニュアルはここでも見られる。

Author: M. Carlson, H. Pages, P. Aboyoun, S. Falcon, M.

Maintainer: Bioconductor Package Maintainer <maintainer@bioconductor.org>

Citation (from within R, enter `citation("GenomicFeatures")`):

Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, Morgan M and Carey V (2013). "Software for Computing and Annotating Genomic Ranges." *PLoS Computational Biology*, 9. <http://dx.doi.org/10.1371/journal.pcbi.1003118>, <http://www.ploscompbiol.org/article/info%3Adoi%2F10.1371%2Fjournal.pcbi.1003118>



## Installation

To install this package, start R and enter:

```
source("http://bioconductor.org/biocLite.R")
biocLite("GenomicFeatures")
```



## Documentation

To view documentation for the version of this package installed in your system, start R and enter:

```
browseVignettes("GenomicFeatures")
```



```
R Console
> browseVignettes("GenomicFeatures")
starting httpd help server ... done
> |
```



- [PDF](#) [R Script](#) Making and Utilizing TxDb Objects
- [PDF](#) Reference Manual
- [Text](#) NEWS



# 目的をおさらい

乳酸菌 *L. casei* 12A株のゲノム配列とアノテーションファイルを読み込んで、全遺伝子の上流500塩基、および下流20塩基分の塩基配列を抽出してプロモーター解析(するための入力データを取得)したい! 黒枠までで行ったことは、アノテーションファイルを読み込んでtxdbオブジェクトに格納したところまで。

## 7. 乳酸菌ゲノム(*Lactobacillus casei* 12A)の[上流500塩基, 下流20塩基]のプロ

[Ensembl](#) (Flicek et al., *Nucleic Acids Res.*, 2014)から提供されている *Lactobacillus casei* 12a.GCA\_000309565.2.25.dna.chromosome配列ファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome](#))とアノテーションファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.gff3](#))

```

in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome"
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.gff3"
out_f <- "hoge7.fasta" #出力ファイル
param_upstream <- 500 #転写開始位置の上流塩基数
param_downstream <- 20 #転写開始位置の下流塩基数

#必要なパッケージをロード
library(Rsamtools) #パッケージ
library(Biostrings) #パッケージ
library(GenomicFeatures) #パッケージ

#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクト
#                               format="gff3", useGenesAsTranscripts=TRUE)
txdb <- makeTxDbFromGFF(in_f2, format="gff3") #確認して

#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb)) #遺伝子の座標情報
hoge <- promoters(gn, upstream=param_upstream, downstream=param_downstream) #指定した領域のプロモーター情報
hoge
    
```

```

> txdb <- makeTxDbFromGFF(in_f2, format="gff3") #TranscriptDbオブジェクトを作成
Prepare the 'metadata' data frame ... metadata: OK
Warning message:
In matchCircularity(seqlevels(gr), circ_seqs) :
  None of the strings in your circ_seqs argument matched
> txdb #確認してる$
TxDb object:
# Db type: TxDb
# Supporting package: GenomicFeatures
# Data source: Lactobacillus_casei_12a.GCA_000309565.$
# Organism: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2799
# exon_nrow: 2800
# cds_nrow: 2681
# Db created by: GenomicFeatures package from Bioconductor
# Creation time: 2015-07-25 17:25:09 +0900 (Sat, 25 Jul 2015)
# GenomicFeatures version at creation time: 1.20.1
# RSQLite version at creation time: 1.0.0
# DBSCHEMAVERSION: 1.1
> |
    
```

# 残りのコードを概説

```
#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトとして作成
# format="gff3", useGenesAsTranscripts=T) #TranscriptDbオブジェクトとして作成
txdb <- makeTxDbFromGFF(in_f2, format="gff3") #TranscriptDbオブジェクトとして作成
txdb #確認してるだけです
```



```
#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb)) #遺伝子の座標情報を取得
hoge <- promoters(gn, upstream=param_upstream, #指定した範囲の座標情報を取得
downstream=param_downstream) #指定した範囲の座標情報を取得
hoge #確認してるだけです
obj <- (ranges(hoge)@start >= 0) #条件を満たすかどうかを判定した結果をobjに格納(座標が0よりも小さい)
hoge <- hoge[obj] #objがTRUEとなる要素のみ抽出した結果をhogeに格納
hoge #確認してるだけです
```



```
#本番(配列取得)
fasta <- getSeq(FaFile(in_f1), hoge) #指定した範囲の塩基配列情報を取得
names(fasta) <- names(hoge) #description情報を追加している
fasta #確認してるだけです
```



```
#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50) #fastaの中身を指定したファイル名で保存
```

①TxDbオブジェクト(オブジェクト名txdb)を出発点として、欲しい領域(上流500塩基、下流20塩基)情報を取得。②ゲノム配列情報FaFile(in\_f1)と欲しい領域情報hogeを入力として塩基配列を取得。③FASTA形式ファイルで保存。このあたりは、H26年度受講生追跡調査の「配列情報を使った次の解析へのinput fileの加工が難関だと感じている」に対する回答。特にプロモーター解析などは「…で、その入力ファイルをどうやって作るの?」かが実務担当者の直面する課題。



# データの型

①黒枠部分は感覚的にpromotersという関数を用いて転写開始点上流と下流の範囲を取得した結果をhogeに格納してるんだらうな、くらいは分かりますよね。そしてなぜ②のところではtxdbを使っていないのか?という素朴な疑問があるでしょう。理由はpromoters関数が入力としてTxDbオブジェクトを想定していないからです(おそらく2014年10月リリースのBioconductor ver. 3.0までは正解)。2015年4月リリースのBioconductor ver. 3.1では、GenomicFeaturesパッケージから提供されているpromoters関数はTxDbに対応済み。

```
#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f1, format="gff3", useGenesAsFeatures=TRUE)
#txdb <- makeTxDbFromGFF(in_f2, format="gff3")
txdb
```

```
#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb))
```

```
hoge <- promoters(gn, upstream=param_upstream, downstream=param_downstream)
```

```
hoge
obj <- (ranges(hoge)@start >= 0)
hoge <- hoge[obj]
hoge
```

```
#本番(配列取得)
fasta <- getSeq(FaFile(in_f1), hoge)
names(fasta) <- names(hoge)
fasta
```

```
#ファイルに保存
writeXStringSet(fasta, file=out_f, format="fasta", width=50)
```

##指定した範囲の座標情報を取得

#指定した範囲の座標情報を取得

#確認してるだけです

#条件を満たすかどうかを判定した結果をobjに格納(座標が0よりも)

#objがTRUEとなる要素のみ抽出した結果をhogeに格納

#確認してるだけです

#指定した範囲の塩基配列情報を取得

#description情報を追加している

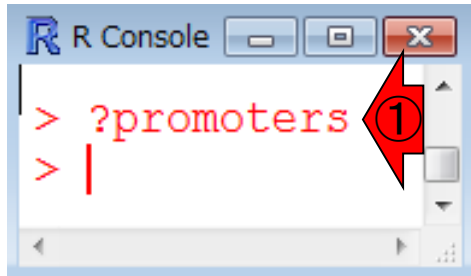
#確認してるだけです

#fastaの中身を指定したファイル名で保存



# データの型

①「?promoters」。門田のR ver. 3.2.0環境では、promotersという同じ名前の関数が3つのパッケージ(IRanges, GenomicFeatures, and GenomicRanges)から提供されていることがわかる。② GenomicRangesパッケージから提供されているpromoters関数は入力としてTxDbオブジェクトを受け付けないことがわかる。しかし、③GenomicFeaturesパッケージから提供されているpromoters関数は、TxDbを入力としていることが2015年7月26日に判明(爆)。



```
R Console  
> ?promoters  
> |
```

Help on topic 'promoters' was found in the following packages:

[Intra range transformations of a Ranges, Views, RangesList, or MaskCollection object](#)

(in package [IRanges](#) in library C:/Users/kadota/Documents/R/win-library/3.2)

[Extract genomic features from an object](#)

(in package [GenomicFeatures](#) in library C:/Users/kadota/Documents/R/win-library/3.2)

[Intra range transformations of a GRanges or GRangesList object](#)

(in package [GenomicRanges](#) in library C:/Users/kadota/Documents/R/win-library/3.2)

# データの型

③ GenomicFeatures パッケージから提供されている promoters 関数は、④ TxDb を入力としていることが 2015 年 7 月 26 日に判明(爆)。

Help on topic 'promoters' was found in the following packages:

[Intra range transformations of a Ranges, Views, RangesList, or MaskCollection object](#)

(in package [IRanges](#) in library C:/Users/kadota/Documents/R/win-library/3.2)

[Extract genomic features from an object](#)

(in package [GenomicFeatures](#) in library C:/Users/kadota/Documents/R/win-library/3.2)

[Intra range transformations of a GRanges or GRangesList object](#)

(in package [Genom](#)

transcripts {GenomicFeatures}

R Documentation

## Extract genomic features from an object

### Description

Generic functions to extract genomic features from an object. This page documents the methods for [TxDb](#) objects only.

### Usage

```
transcripts(x, ...)
## S4 method for signature 'TxDb'
transcripts(x, vals=NULL, columns=c("tx_id", "tx_name"))

exons(x, ...)
## S4 method for signature 'TxDb'
exons(x, vals=NULL, columns="exon_id")

cds(x, ...)
## S4 method for signature 'TxDb'
cds(x, vals=NULL, columns="cds_id")

genes(x, ...)
## S4 method for signature 'TxDb'
genes(x, vals=NULL, columns="gene_id", single.strand.genes.only=TRUE)

## S4 method for signature 'TxDb'
promoters(x, upstream=2000, downstream=200, ...)
```



# プロモーター配列取得

②例題8に、③TxDbオブジェクトを入力としてpromoters関数を実行するコードを追加。1つ上の行は過去の遺物(legacy)として残しているだけ。(Rで)塩基配列解析自体がデカすぎて、もはや全体の動作確認もままならないし、修正・変更もしきれていないのが現状

- インタロ | 一般 | 配列取得 | プロモーター配列 | [BSgenomeとTxDbから](#)(last modified 2015...
- インタロ | 一般 | 配列取得 | プロモーター配列 | [GenomicFeatures\(Lawrence 2013\)](#)
- インタロ | 一般 | 配列取得 | トランスクリプトーム配列 | [公共DBから](#)(last modified 2015-05...
- インタロ | 一般 | 配列取得 | ...
- インタロ | 一般 | 配列取得 | ...
- インタロ | 一般 | 配列取得 | ...

**①** [イントロ](#) | [一般](#) | [配列取得](#) | [プロモーター配列](#) | [Genomi](#)

[GenomicFeatures](#) パッケージを主に用いてプロモーター配列(転写開始点近傍配列)を得るやり方を示します。

**②** [GenomicFeatures](#) | [BSgenome](#) | [TranscriptDb](#) | [TxDb](#) | [GenomicFeatures](#)

ここでは、[TxDb](#) から、xxx塩基、下流xxx塩基、上流xxx塩基の範囲で、どの生物種にも適用可能な形式のアンノテーションを取得します。

**③** [1. ヒト\("BSgenome.Hsapiens.UCSC.hg19"\)](#)

ヒトゲノム ver. 19

## **8. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流100塩基、下流0塩基]のプロモーター配列を取得する場合:**

7と基本的に同じですが、少なくとも2015年4月リリースのBioconductor ver. 3.1以降で、[GenomicFeatures](#) パッケージが提供するpromoters関数がTxDbオブジェクトをそのまま読み込めることが判明したので、それを反映させています。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa" #入力ファ
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3" #入力ファ
out_f <- "hoge8.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 100 #転写開始点上流の塩基配列数を指定
param_downstream <- 0 #転写開始点下流の塩基配列数を指定
```

```
#必要なパッケージをロード
library(Rsamtools) #パッケージの読み込み
library(Biostrings) #パッケージの読み込み
library(GenomicFeatures) #パッケージの読み込み
```

```
#前処理(アンノテーション情報を取得)
txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトを取得してtxdbに格納
# format="gff3", useGenesAsTranscripts=T) #TranscriptDbオブジェクトを取得してtxdbに
txdb <- makeTxDbFromGFF(in_f2, format="gff3") #TranscriptDbオブジェクトを取得してtxdbに格納
txdb #確認してるだけです
```

```
#前処理(欲しい領域の座標情報取得)
#gn <- sort(genes(txdb)) #遺伝子の座標情報を取得
hoge <- promoters(txdb, upstream=param_upstream, #指定した範囲の座標情報を取得
downstream=param_downstream) #指定した範囲の座標情報を取得
hoge #確認してるだけです
```

# 遺言1

- R本体とパッケージのバージョンアップを毎年5月と11月に行うべし
  - 依存関係など様々な問題はあるものの、落ち着いて必要なパッケージを個別にインストールしていけば大丈夫
  - 第一義的に重要なのは、Linuxでプログラムのインストールができることと同様、library(XXX)の部分でこけないようにすること。

- [はじめに](#) (last modified 2015/03/31)
- [参考資料\(講義、講習会、本など\)](#) (last modified 2015/07/07) **NEW**
- [過去のお知らせ](#) (last modified 2015/07/08) **NEW**
- [インストール | について](#) (last modified 2015/04/04)
- インストール | R本体 | 最新版 | [Win用](#) (last modified 2015/03/22) 推奨
- インストール | R本体 | 最新版 | [Mac用](#) (last modified 2015/04/22) 推奨
- インストール | R本体 | 過去版 | [Win用](#) (last modified 2015/03/22)
- インストール | R本体 | 過去版 | [Mac用](#) (last modified 2015/03/22)
- インストール | Rパッケージ | [ほぼ全て\(20GB以上?!\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [必要最小限プラスアルファ\(数GB?!\)](#) ① modified 2015/07/24) 推奨 **NEW**
- インストール | Rパッケージ | [必要最小限プラスアルファ\(アグリバイター 居室のみ\)](#) (last modified 2015/06/16)
- インストール | Rパッケージ | [必要最小限\(数GB?!\)](#) (last modified 2015/05/25)
- インストール | Rパッケージ | [個別](#) (last modified 2015/06/10)

②

①

# 遺言2

## ■ R最新版でうまくいったなら、過去版のエラー原因究明は時間の無駄!

- ① R ver. 3.1.2のときはmakeTranscriptDbFromGFF関数でうまくいくはずだが、エラーが出る現象に遭遇した(2015.07.07のアグリバイオ大学院講義)。しかし、② R ver. 3.2.1でmakeTxDbFromGFF関数で正常動作することがわかり、その時点で思考停止

### 7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得する場合:

Ensembl (Flicek et al., Nucleic Acids Res., 2014)から提供されている [Lactobacillus casei 12A](#)の multi-FASTA形式ゲノム配列ファイル(Lactobacillus casei 12a.GCA\_000309565.2.25.dna.chromosome.Chromosome.fa)と GFF3形式のアノテーションファイル(Lactobacillus casei 12a.GCA\_000309565.2.25.chromosome.Chromosome.gff3)を読み込むやり方です。

```

in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入力ファ
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3"#入力ファイ
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定

#必要なパッケージをロード
library(Rsamtools) #パッケージの読み込み
library(Biostrings) #パッケージの読み込み
library(GenomicFeatures) #パッケージの読み込み

#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDbオブジェクトを取得してtxdbに格納
# format="gff3", useGenesAsTranscripts=T)#TranscriptDbオブジェクトを取得してtxdbに格納
txdb <- makeTxDbFromGFF(in_f2, format="gff3")#TranscriptDbオブジェクトを取得してtxdbに格納
txdb #確認してるだけです

#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb)) #遺伝子の座標情報を取得
hoge <- promoters(gn, upstream=param_upstream, #指定した範囲の座標情報を取得
downstream=param_downstream)#指定した範囲の座標情報を取得
hoge #確認してるだけです
    
```



# 遺言3

- 得られた結果を無条件に受け入れるな！ 様々な角度で検証せよ。
  - 実はプログラムのバグで、GFF3形式ファイルの読み込みに失敗していただけ、みたいなオチもよくあるだろう。フリーソフトの利用は自己責任である。

## 7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得する場合:

Ensembl (Flicek et al., Nucleic Acids Res., 2014)から提供されている [Lactobacillus casei 12A](#)の multi-FASTA形式ゲノム配列ファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome.Chromosome.f](#))と GFF3形式のアノテーションファイル([Lactobacillus casei 12a.GCA\\_000309565](#))

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.f"
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.gff3"
out_f <- "hoge7.fasta"
param_upstream <- 500
param_downstream <- 20

#必要なパッケージをロード
library(Rsamtools)
library(Biostrings)
library(GenomicFeatures)

#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #TranscriptDb object
#                               format="gff3", useGenesAsTranscripts=TRUE)
txdb <- makeTxDbFromGFF(in_f2, format="gff3")
txdb

#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb))
hoge <- promoters(gn, upstream=param_upstream, downstream=param_downstream)
hoge
```

```
> txdb
TxDb object:
# Db type: TxDb
# Supporting package: GenomicFeatures
# Data source: Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.f
# Organism: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2799
# exon_nrow: 2800
# cds_nrow: 2681
# Db created by: GenomicFeatures package from Bioconductor
# Creation time: 2015-07-26 11:26:55 +0900 (Sun, 26 Jul 2015)
# GenomicFeatures version at creation time: 1.20.1
# RSQLite version at creation time: 1.0.0
# DBSCHEMAVERSION: 1.1
> |
```

# Rを用いた検証例

txdbオブジェクトの元データは、①のGFF3ファイル。②乳酸菌L. casei 12Aの遺伝子数(=転写物数)は、2,799個となっている。これが正しいかを検証する。

## 7. 乳酸菌ゲノム(Lactobacillus casei 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得

Ensembl (Flicek et al., Nucleic Acids Res., 2014)から提供されている [Lactobacillus casei 12A](#)の multi-FASTA形式ゲノム配列ファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome.Chromosome.fa](#))と GFF3形式のアノテーションファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.chromosome.Chromosome.gff3](#))を読み込むやり方です。

```
in_f1 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa" #入力ファ
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3" #GFF3ファイル
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
param_upstream <- 500 #転写開始点上流の塩基配列数を指定
param_downstream <- 20 #転写開始点下流の塩基配列数を指定
```

```
#必要なパッケージをロード
library(Rsamtools) #バ
library(Biostrings) #バ
library(GenomicFeatures) #バ

#前処理(アノテーション情報を取得)
#txdb <- makeTranscriptDbFromGFF(in_f2, #Tr
# format="gff3", useGenesAsTransc
txdb <- makeTxDbFromGFF(in_f2, format="gff3") #確
txdb

#前処理(欲しい領域の座標情報取得)
gn <- sort(genes(txdb)) #追
hoge <- promoters(gn, upstream=param_upstr #指
downstream=param_downstream) #確
hoge
```

```
R Console
> txdb
TxDb object:
# Db type: TxDb
# Supporting package: GenomicFeatures
# Data source: Lactobacillus_casei_12a.GCA_000309565.2.25
# Organism: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2799
# exon_nrow: 2800
# cds_nrow: 2681
# Db created by: GenomicFeatures package from Bioconductor
# Creation time: 2015-07-26 11:26:55 +0900 (Sun, 26 Jul 2015)
# GenomicFeatures version at creation time: 1.20.1
# RSQLite version at creation time: 1.0.0
# DBSCHEMAVERSION: 1.1
> |
```



# Rを用いた検証例

(9列目の) "ID=gene" という文字列を含む行数を調べてみるというのでは、という思想のもので検証する。

## 7. 乳酸菌ゲノム(*Lactobacillus casei* 12A)の[上流500塩基, 下流20塩基]のプロモーター配列を取得する場合:

Ensembl (Flicek et al., *Nucleic Acids Res.*, 2014)から提供されている *Lactobacillus casei* 12Aの multi-FASTA形式ゲノム配列ファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.dna.chromosome.Chromosome.fa](#))と GFF3形式のアノテーションファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.chromosome.Chromosome.gff3](#))を読み込むやり方です。

```
in_f1 <- "Lactobacillus casei 12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入力ファ
in_f2 <- "Lactobacillus casei 12a.GCA_000309565.2.25.chromosome.Chromosome.gff3"#入力ファイ
out_f <- "hoge7.fasta" #出力ファイル名を指定してout_fに格納
```

```
paran ##gff-version 3
paran ##sequence-region Chromosome 1 2907892
```

#必要	Chromosome	ena	gene	1	1350	.	+	.	ID=gene:LCA12A_0617;assembly_name=GCA_000309565.2;assembly_r
libra	Chromosome	ena	gene	1523	2662	.	+	.	ID=gene:LCA12A_0618;assembly_name=GCA_000309565.2;assembly_r
libra	Chromosome	ena	gene	3240	3452	.	+	.	ID=gene:LCA12A_0619;assembly_name=GCA_000309565.2;assembly_r
libra	Chromosome	ena	gene	3449	4564	.	+	.	ID=gene:LCA12A_0620;assembly_name=GCA_000309565.2;assembly_r
#前処	Chromosome	ena	gene	4817	6778	.	+	.	ID=gene:LCA12A_0621;assembly_name=GCA_000309565.2;assembly_r
#txdb	Chromosome	ena	gene	6840	9461	.	+	.	ID=gene:LCA12A_0622;assembly_name=GCA_000309565.2;assembly_r
#	Chromosome	ena	gene	9566	10270	.	-	.	ID=gene:LCA12A_0623;assembly_name=GCA_000309565.2;assembly_r
txdb	Chromosome	ensembl	CDS	1	1350	.	+	0	ID=CDS:EKP96483;Parent=transcript:EKP96483;assembly_r
txdb	Chromosome	ensembl	exon	1	1350	.	+	.	Name=EKP96483-1;Parent=transcript:EKP96483;assembly_r
#前処	Chromosome	ensembl	CDS	1523	2662	.	+	0	ID=CDS:EKP96484;Parent=transcript:EKP96484;assembly_r
gn <	Chromosome	ensembl	exon	1523	2662	.	+	.	Name=EKP96484-1;Parent=transcript:EKP96484;assembly_r
hoge	Chromosome	ensembl	CDS	3240	3452	.	+	0	ID=CDS:EKP96485;Parent=transcript:EKP96485;assembly_r
hoge	Chromosome	ensembl	exon	3240	3452	.	+	.	Name=EKP96485-1;Parent=transcript:EKP96485;assembly_r
<	Chromosome	ensembl	CDS	3449	4564	.	+	0	ID=CDS:EKP96486;Parent=transcript:EKP96486;assembly_r
	Chromosome	ensembl	exon	3449	4564	.	+	.	Name=EKP96486-1;Parent=transcript:EKP96486;assembly_r
	Chromosome	ensembl	CDS	4817	6778	.	+	0	ID=CDS:EKP96487;Parent=transcript:EKP96487;assembly_r
	Chromosome	ensembl	exon	4817	6778	.	+	.	Name=EKP96487-1;Parent=transcript:EKP96487;assembly_r
	Chromosome	ensembl	CDS	6840	9461	.	+	0	ID=CDS:EKP96488;Parent=transcript:EKP96488;assembly_r



# Rを用いた検証例

txdbオブジェクトの元データは、①のGFF3ファイル。乳酸菌L. casei 12Aの遺伝子数(=転写物数)は、2,799個となっている。Rで③"ID=gene"を含む行数も④2,799個だったので、大丈夫だろうと判断する。

- 書籍 | 日本乳酸菌学会誌 | [第4回クオリティコントロールとプログラムのインストール](#)
- イントロ | 一般 | [ランダムに行を抽出](#) (last modified 2014/07/17)
- イントロ | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2014/07/17)
- イントロ | 一般 | [任意のキーワードを含む行を抽出\(基礎\)](#) ① (last modified 2014/04/11)
- イントロ | 一般 | [ランダムな塩基配列を生成](#) (last modified 2014/06/16)
- イントロ | 一般 | [任意の長さの可能な全ての塩基配列を作成](#) (last modified 2015/02/02)

② 14. GFF3形式のタブ区切りテキストファイル([Lactobacillus casei 12a.GCA\\_000309565.2.25.chromosome.Chromosome.gff3](#))に対して、"ID=gene"という文字列が含まれる行全体を出力したい場合:

[乳酸菌ゲノム\(Lactobacillus casei 12A\)](#)のアノテーションファイルです。4をベースに作成。

```

in_f <- "Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3" #入力ファイル名を
out_f <- "hoge14.txt" #出力ファイル名を指定してout_fに格納
param <- "ID=gene" #検索したい文字列を指定

#入力ファイルの読み込み
data <- readLines(in_f)
length(data)

#本番(リストファイル中の要素一つ一つに対して)
hoge <- sapply(param, grep, x=data)
hoge <- unique(hoge)
out <- data[hoge]
length(out)

#ファイルに保存
writeLines(out, out_f)

```

```

R Console
# > data <- readLines(in_f) #in_fで指$
# > length(data) #オブジェ$
[1] 11081
# > #本番(リストファイル中の要素一つ一つに対して、要$
# > hoge <- sapply(param, grep, x=data) #リストフ$
# > hoge <- unique(hoge) #得られるh$
# > out <- data[hoge] #hogeで指$
# > length(out) #オブジェ$
[1] 2799
# > #ファイルに保存
# > writeLines(out, out_f) #outの中身$
# > |

```

# Linuxを用いた検証例

Linuxをある程度使えるヒトは、普通grepコマンドを使います。①のmvコマンドは、ファイル名が長いので、hoge.gff3にrenameしているだけです。②の-cは、③で指定した文字列と一致した行数を表示させるオプション。

```
iu@bielinux[~/Desktop/mac_share]
iu@bielinux[mac_share] pwd
/home/iu/Desktop/mac_share
iu@bielinux[mac_share] ls
Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3
iu@bielinux[mac_share] mv Lactobacillus_casei_12a.GCA_000309565.2.25.ch
romosome.Chromosome.gff3 hoge.gff3
iu@bielinux[mac_share] ls
hoge.gff3
iu@bielinux[mac_share] grep -c ID=gene hoge.gff3
2799
iu@bielinux[mac_share] █
```

[ 3:08午後 ]  
[ 3:08午後 ]  
[ 3:08午後 ]  
[ 3:08午後 ]



# Linuxを用いた検証例

①-vは、一致しないものを表示するオプション。  
-cvと併用することで、一致しない行数を表示。  
②Linuxだと他の気になるキーワードも簡単に調査可能。  
③^をつけることで「Chromosomeを含む行数」ではなく「行頭にChromosomeを含む行数」をカウントすることもできる。

```
iu@bielinux[~/Desktop/mac_share]
iu@bielinux[mac_share] pwd
/home/iu/Desktop/mac_share
iu@bielinux[mac_share] ls
Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3
iu@bielinux[mac_share] mv Lactobacillus_casei_12a.GCA_000309565.2.25.chromosome.Chromosome.gff3 hoge.gff3
iu@bielinux[mac_share] ls
hoge.gff3
iu@bielinux[mac_share] grep -c ID=gene hoge.gff3
2799
① iu@bielinux[mac_share] grep -cv ID=gene hoge.gff3
8282
iu@bielinux[mac_share] grep -c gene hoge.gff3
5598
iu@bielinux[mac_share] grep -c CDS hoge.gff3
2681
iu@bielinux[mac_share] grep -c exon hoge.gff3
2807
③ iu@bielinux[mac_share] grep -c ^Chromosome hoge.gff3
11079
iu@bielinux[mac_share] █
```



# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# FASTQファイル解析

①ファイル形式の変換(FASTQ → FASTA)や、FastQCのいくつかの項目と同じ解析をRでも実行可能である: ②Overrepresented sequences、③Per base sequence content (リードのpositionごとの出現確率)、④Kmer Content。













- [イントロ | NGS | 読み込み | FASTQ形式 | 基礎](#) **④** (last modified 2015/06/24)
- [イントロ | NGS | 読み込み | FASTQ形式 | 応用](#) (last modified 2015/06/18)
- [イントロ | NGS | 読み込み | FASTQ形式 | description行の記述を整形](#) (last modified 2015/06/18)
- [イントロ | NGS | 読み込み | Illuminaの \\* seq.txt](#) (last modified 2013/06/13)
- [イントロ | NGS | 読み込み | Illuminaの \\* qseq.txt](#) (last modified 2013/06/17)
- [イントロ | ファイル形式の変換 | について](#) (last modified 2014/06/09)
- [イントロ | ファイル形式の変換 | BAM → BED](#) (last modified 2014/06/21)
- [イントロ | ファイル形式の変換 | FASTQ → FASTA](#) **①** (last modified 2015/06/15)
- [イントロ | ファイル形式の変換 | Genbank → FASTA](#) (last modified 2014/03/10)
- [イントロ | ファイル形式の変換 | qseq → FASTA](#) (last modified 2013/06/17)
- [イントロ | ファイル形式の変換 | qseq → Illumina FASTQ](#) (last modified 2013/06/17)
- [イントロ | ファイル形式の変換 | qseq → Sanger FASTQ](#) (last modified 2013/08/19)
- [前処理 | クオリティコントロール | について](#) (last modified 2015/06/25)
- [前処理 | クオリティチェック | QuasR\(Gaidatzis 2015\)](#) **③** (last modified 2015/06/15)
- [前処理 | クオリティチェック | qrc](#) (last modified 2014/06/17)
- [前処理 | クオリティチェック | PHREDスコアに変換](#) (last modified 2013/06/18)
- [前処理 | クオリティチェック | 配列長分布を調べる](#) (last modified 2015/06/22)
- [前処理 | クオリティチェック | Overrepresented sequences | ShortRead\(Morgan 2009\)](#) **②** (last modified 2015/07/21)
- [前処理 | トリミング | ポリA配列除去 | ShortRead\(Morgan 2009\)](#) (last modified 2014/06/11)
- [前処理 | トリミング | アダプター配列除去\(基礎\) | QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/26) 推奨
- [前処理 | トリミング | アダプター配列除去\(基礎\) | girafe\(Toedling 2010\)](#) (last modified 2014/06/11)
- [前処理 | トリミング | アダプター配列除去\(基礎\) | ShortRead\(Morgan 2009\)](#) (last modified 2014/06/21)
- [前処理 | トリミング | アダプター配列除去\(応用\) | QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/26) 推奨
- [前処理 | トリミング | アダプター配列除去\(応用\) | ShortRead\(Morgan 2009\)](#) (last modified 2014/06/18)
- [前処理 | トリミング | 指定した末端塩基数だけ除去](#) (last modified 2015/06/29) **NEW**

# FastQC実行結果

①頻出する配列をリストアップ。②トップは「CCCCGGTATA…」という50塩基の配列で14,383回出現。Percentageは1.4383%。全部で100万リードなので妥当。オリジナル107bpのうち最初の50bpで解析している。復習

## FastQC Report

### Summary

-  [Basic Statistics](#)
-  [Per base sequence quality](#)
-  [Per tile sequence quality](#)
-  [Per sequence quality scores](#)
-  [Per base sequence content](#)
-  [Per sequence GC content](#)
-  [Per base N content](#)
-  [Sequence Length Distribution](#)
-  [Sequence Duplication Levels](#)
-  [Overrepresented sequences](#)
-  [Adapter Content](#)
-  [Kmer Content](#)

### Overrepresented sequences

Sequence	Count	Percentage	Possible Source
CCCCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCA	14383	1.4383	No Hit
GGCCTATTCAGTGCAGTGCACCTTGGCGTACGACCCCTTCTCCGAAGT	11044	1.1044	No Hit
GTGCTTTTCACCTTCCCTCACGGTACTGGTTCCTACTATCGGTCCTAGGG	8892	0.8892000000000001	No Hit
CCCAGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCAC	8474	0.8474	No Hit
GTCAGTGGAGTATTTAGCCTGGGAGATGGTCTCCCGGATCCGACG	8189	0.8189	No Hit
GCCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCAGATCGACCTTC	8132	0.8132	No Hit
GTCCAGTCTACAACCCGAGAAGCAAGCTTCTCGTTTGGGCTCTTCCC	6663	0.6663	No Hit
GTGGTTTGGGTACGGGTAGTTTATTTCTCACTAGAAGCTTTTCTTGGC	6411	0.6411	No Hit
GGTCACTTGGTTTCGGGTCTACATCTGCTTACTCATTGCCCCTGTTTCTCAGA	5502	0.5502	No Hit
GCCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCAGATCAACCTTC	4845	0.48450000000000004	No Hit
CCCTCCATCGCTTAAACAAAATAAACTAGTGCAGGAATCTCAACCTGCTT	4395	0.43949999999999995	No Hit
CCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCACT	4385	0.4385	No Hit
CCCAGTCTGCCCGCCAGCTATGTATTCAGTACAAGCAATACACTG	4366	0.4366	No Hit
CCACAGTTTCGGTATTATGCTTAGCCCCGGTATATTTTCGGCGCAGTGCC	4314	0.4314	No Hit
CTGGGCTGTTCCCTTTTCGACAATGGACCTATCGCTCACTGTCTGACTC	4113	0.41130000000000005	No Hit
CCGCCGACTCAGGATCCTGGACGGAGGGTTCGACGTTTCGTTACAGGG	4081	0.4081	No Hit
CCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCAGATCGACCTTCA	3846	0.3846	No Hit
GTAGGTCAGTGGTTTCGGGTCTACATCTGCTTACTCATTGCCCCTGTTT	3823	0.3823	No Hit

subseqとtable関数の併用で同じ結果が得られる例として提示。これもただの復習。

# ShortRead実行結果

- 前処理 | クオリティチェック | [配列長分布を調べる](#) (last modified 2015/06/22)
- 前処理 | クオリティチェック | Overrepresented sequences | [ShortRead\(Morgan 2009\)](#)
- 前処理 | トリミング | ポリア配列除去 | [ShortRead\(Morgan 2009\)](#) (last modified 2014/07/11)
- 前処理 | トリミング | マダボカー配列除去(其礎) | [OverR\(Gaidatzis 2015\)](#) (last modified 2015/06/22)
- 前処理 | **前処理 | クオリティチェック | Overrepresented sequences | ShortRead**
- 前処理 | [ShortRead](#) パッケージを用いてリードの種類ごとの出現回数を得るやり方を示します。FastQCの Over
- 前処理 | 一般的な
- 前処理 | 「ファイル



## 前処理 | クオリティチェック | Overrepresented sequences | ShortRead

ShortRead パッケージを用いてリードの種類ごとの出現回数を得るやり方を示します。FastQCの Over

### 4. gzip圧縮FASTQ形式ファイル(SRR616268sub\_1.fastq.gz)の場合:

乳酸菌RNA-seqデータSRR616268の最初の100万リード分(約73MB)です。長さは全て107 bpです。3と基本的に同じですが、FastQCのデフォルトと同じく、最初の50bp分のみで解析するやり方です。「イントロ」一般「指定した範囲の配列を取得」のテクニ

```
1. gzip圧縮FASTQ形式ファイル
small RNA
in_f <- "SRR616268sub_1.fastq.gz"
out_f <- "hoge4.txt"
#必要なパッケージをロード
library(ShortRead)
#入力ファイルの読み込み
fastq <- readFastq(in_f)
sread(fastq)
```

```
in_f <- "SRR616268sub_1.fastq.gz"
out_f <- "hoge4.txt"
param <- c(1, 50)
#必要なパッケージをロード
library(ShortRead)
#入力ファイルの読み込み
fastq <- readFastq(in_f)
fasta <- sread(fastq)
#前処理(指定した範囲の配列を取得)
fasta <- subseq(fasta, param)
```

```
R Console
> head(out) #確認してるだけ$
fasta
CCCCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCA
14383
GGCCTATTCAGTGCAGGCTGACCTTGCAGGTCAGCACCCCTTCTTCCGAAGT
11044
GTGCTTTTACCTTTCCCTCACGGTACTGGTTCAGTACTATCGGTCAGTGGG
8892
CCCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCAC
8474
GTCAGTGGGAGTATTTAGCCTTGGGAGATGGTCCTCCCGGATTCCGACG
8189
GCCGGCATTCTCACTTCTAAGCGCTCCAGCCGTCCTCACGATCGACCTTC
8132
>
> #ファイルに保存
> tmp <- cbind(names(out), out) #保存したい情報$
> write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F)
```



# FastQC実行結果

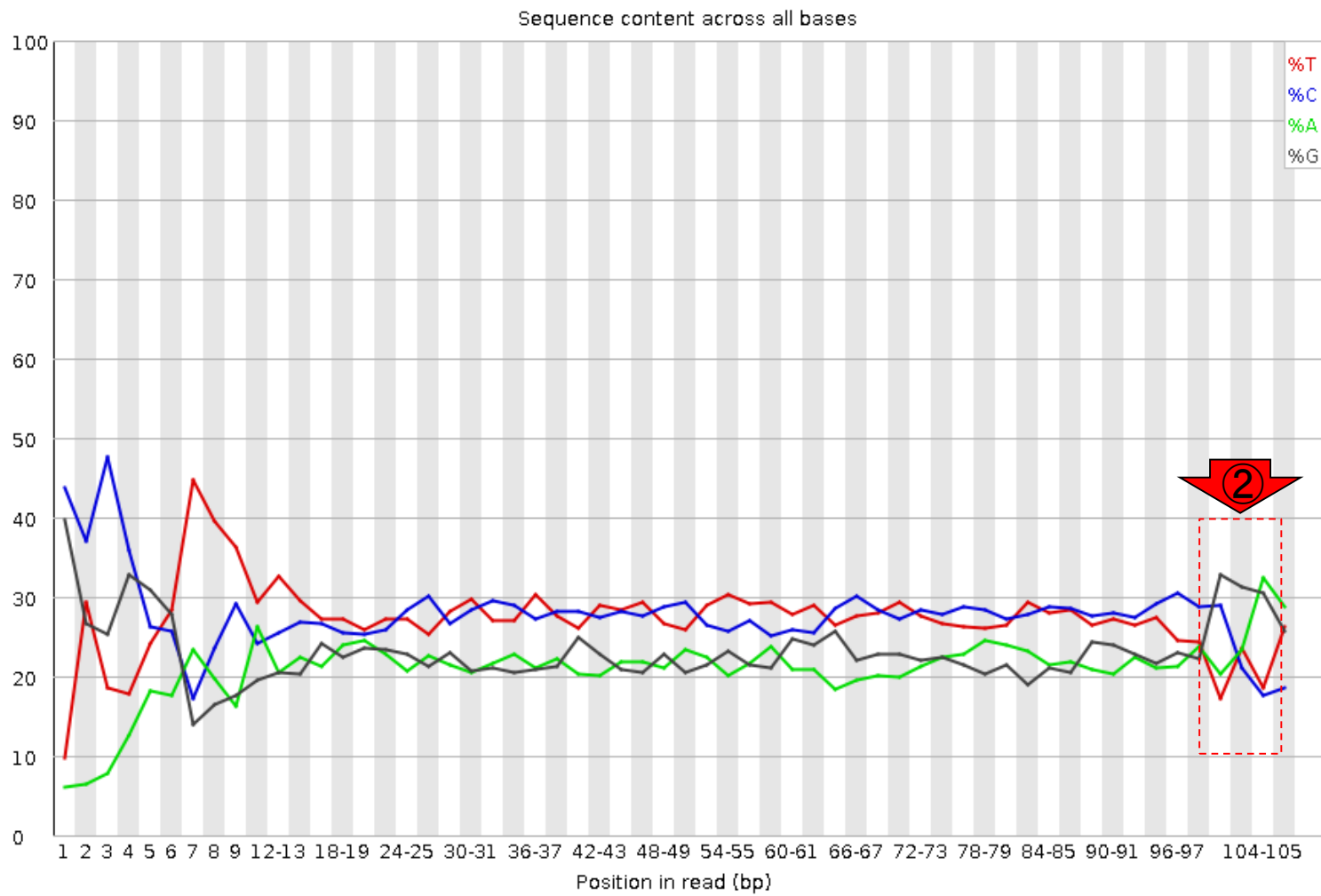
①ポジションごとの塩基の出現確率。FastQC (ver. 0.11.3)のデフォルトオプションでの実行結果は、最初の10塩基までは塩基ごとに、それ以降は適当に数塩基分を平均化した出現確率を表示(していることにQuasR実行結果と比較することで後に気づいた)。②赤枠部分に注目!

## FastQC Report

### Summary

- ✓ [Basic Statistics](#)
- ✓ [Per base sequence quality](#)
- ✓ [Per tile sequence quality](#)
- ✓ [Per sequence quality scores](#)
- ✗ [Per base sequence content](#) ①
- ! [Per sequence GC content](#)
- ✓ [Per base N content](#)
- ✓ [Sequence Length Distribution](#)
- ✗ [Sequence Duplication Levels](#)
- ✗ [Overrepresented sequences](#)
- ✓ [Adapter Content](#)
- ✗ [Kmer Content](#)

### ✗ Per base sequence content



# QuasR実行結果

①QuasRというパッケージは、Quality Control部分だけで比較すると、FastQCの簡易版のような位置づけ。②例題4実行結果のhoge4.pdfの一部が右下プロットで、FastQCのPer base sequence contentと同じもの。③の赤枠部分に注目!

- インポート | ファイル形式の変換 | [qseq --> Sanger FASTQ](#) (last modified 2013/08/19)
- 前処理 | [クオリティコントロール](#) | [について](#) (last modified 2015/06/25)
- 前処理 | クオリティチェック | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/15)
- 前処理 | クオリティチェック | [qrc](#) (last modified 2014/06/17)
- 前処理 | クオリティチェック | [PHREDスコアに変換](#) (last modified 2013/06/18)

## 前処理 | クオリティチェック | QuasR(Gaidatzis\_2015)

QuasRパッケージを用いてQCレポートファイルを出力するやり方を示します。FastQCのR版のようなqrcよりも相当ストイックな出力結果です。

### 4. FASTQ形式ファイル(SRR616268sub\_1.fastq.gz)の場合:

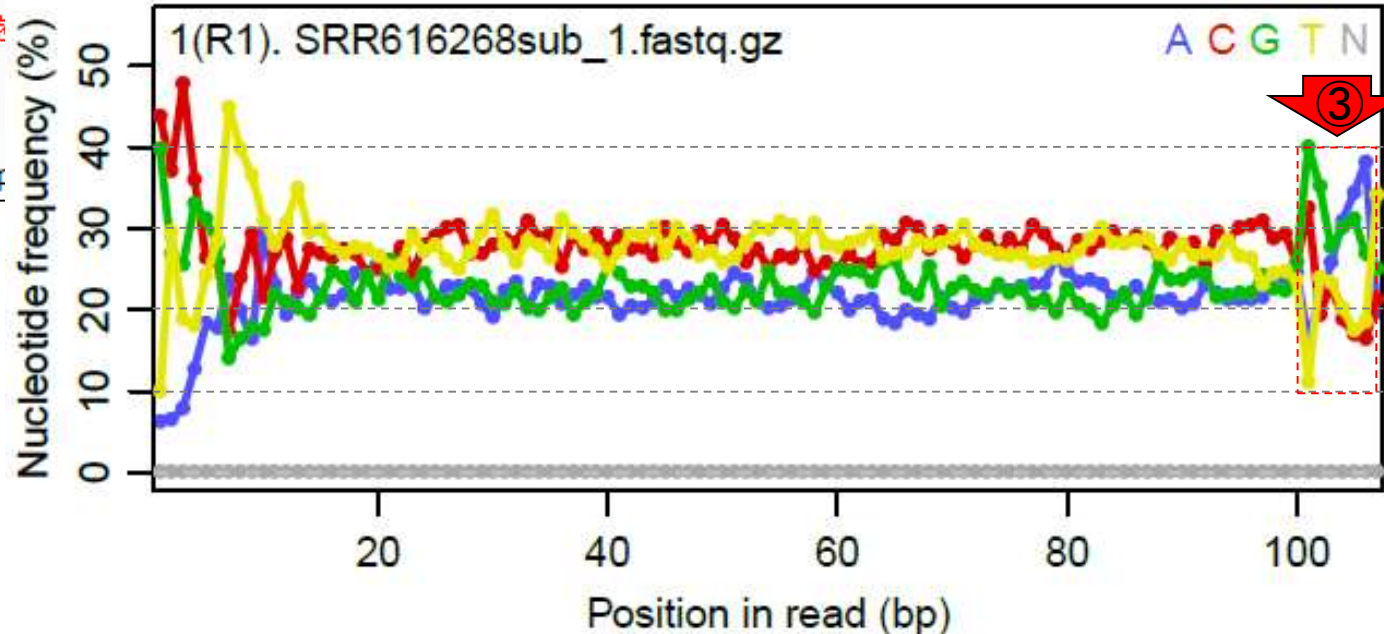
乳酸菌RNA-seqデータSRR616268の最初の100万リード分(約73MB)です。

- 1. サンプル名 (SRR037439 al., 2010)。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge4.pdf" #出力ファイル名を指定してout_fに格納

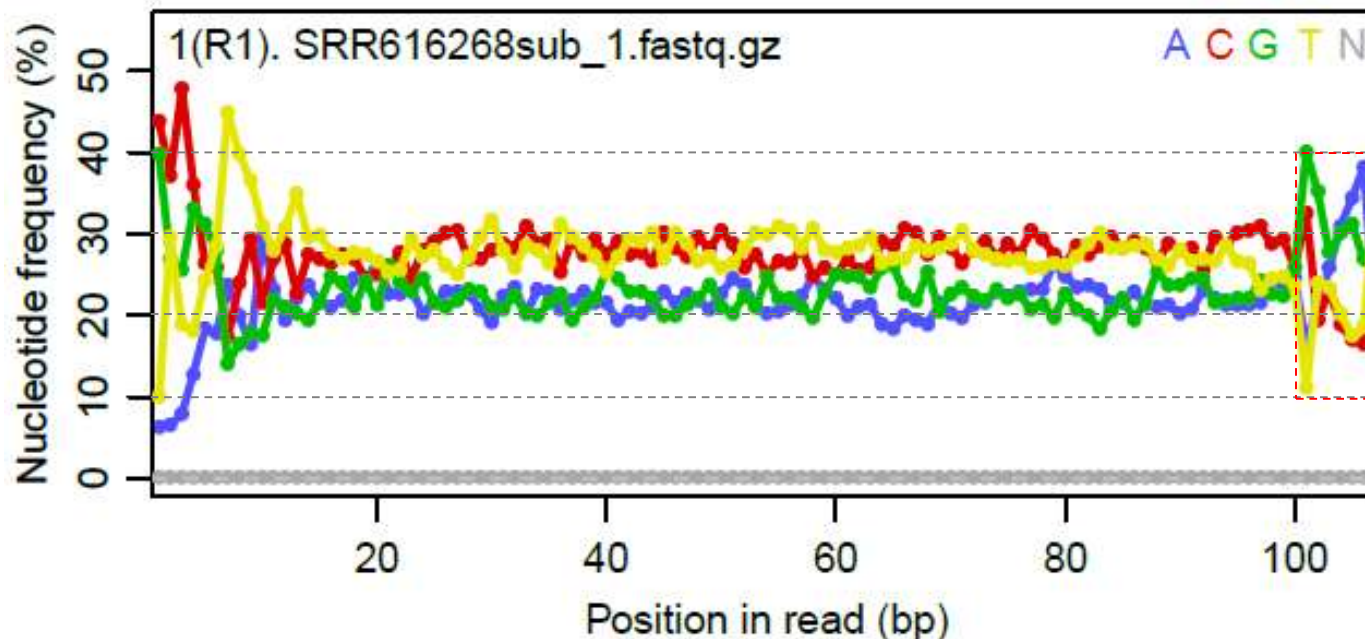
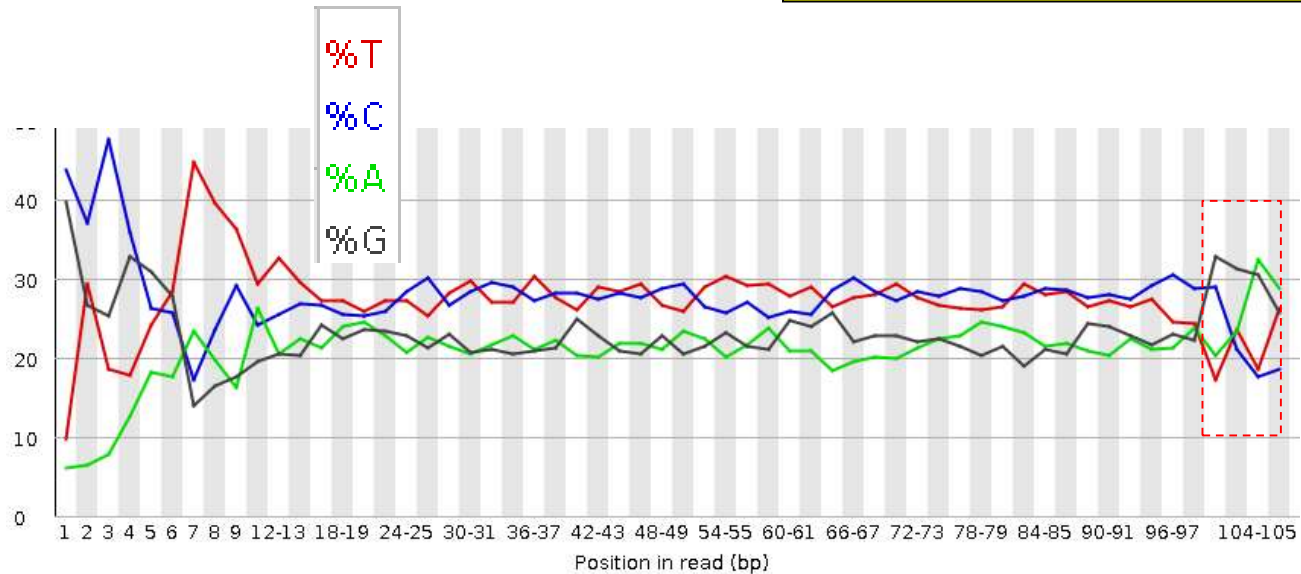
#必要なパッケージを
library(QuasR)

#本番
qQCReport(in_f, out_f)
```



# FastQC vs. QuasR

①FastQCのほうの横軸の数値をよく眺めると、こちらは2塩基分ずつ平均化したものをプロットしているのではないかとと思うに至る



# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# FastQCマニュアル

①fastqc2コマンドのマニュアルを表示。予めマニュアル全体を眺めて、`--nogroup`オプションを用いれば平均化(grouping)せず塩基ごとに表示できることが分かったうえで、②「`grep -A 5 nogroup`」とパイプでつなげて必要最小限の情報を表示

```
iu@bielinux[~]
iu@bielinux[iu] fastqc2 -v                               [10:46午後]
FastQC v0.11.3
iu@bielinux[iu] fastqc2 -h | grep -A 5 nogroup           [10:46午後]
  --nogroup      Disable grouping of bases for reads >50bp. All reports will
                  show data for every base in the read. WARNING: Using this
                  option will cause fastqc to crash and burn if you use it on
                  really long reads, and your plots may end up a ridiculous size.
                  You have been warned!
iu@bielinux[iu] █                                       [10:46午後]
```

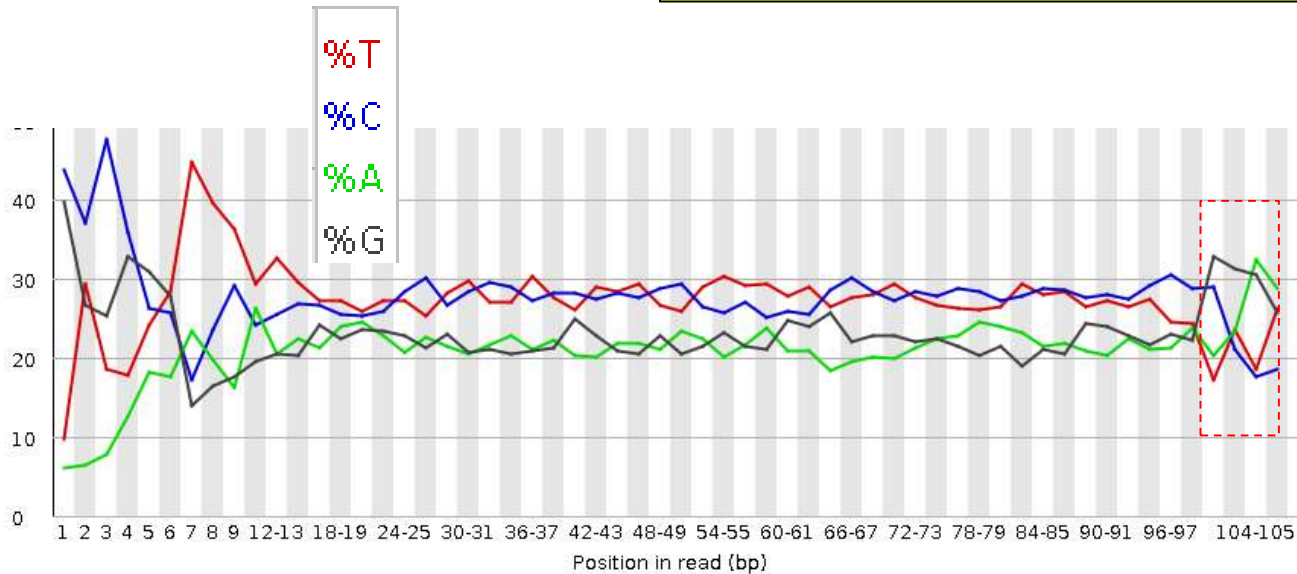
①

②



# FastQC

①このFastQCレポートhtmlは、②FastQCをデフォルトオプションで実行した結果。③得られたhtmlファイル名をresult\_without\_nogroup.htmlに変更。



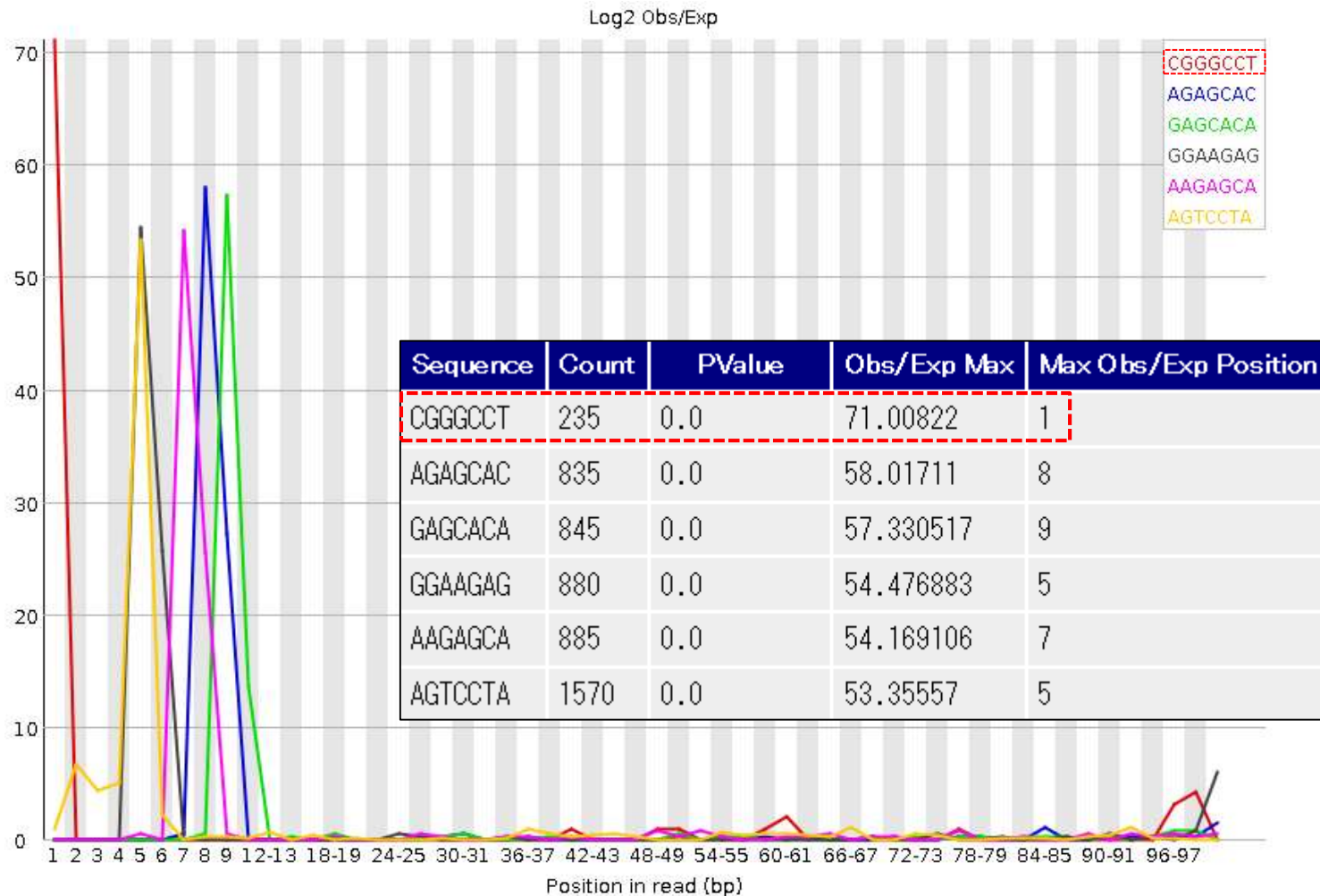
```

iu@bielinux[~/Desktop/mac_share]
iu@bielinux[mac_share] pwd [11:06午後]
/home/iu/Desktop/mac_share
iu@bielinux[mac_share] ls -la SRR616268sub_1.fastq.gz [11:06午後]
-rwxrwxrwx 1 root root 74906576 7月 13 17:06 SRR616268sub_1.fastq.gz
iu@bielinux[mac_share] fastqc2 -q SRR616268sub_1.fastq.gz [11:06午後]
iu@bielinux[mac_share] ls -la *.html [11:06午後]
-rwxrwxrwx 1 root root 365019 7月 26 23:06 SRR616268sub_1_fastqc.html
iu@bielinux[mac_share] mv SRR616268sub_1_fastqc.html result_without_nogroup.html
iu@bielinux[mac_share] date [11:07午後]
2015年 7月 26日 日曜日 23:07:50 JST
iu@bielinux[mac_share] [11:07午後]
  
```

# FastQC

## Summary

- ✔ [Basic Statistics](#)
- ✔ [Per base sequence quality](#)
- ✔ [Per tile sequence quality](#)
- ✔ [Per sequence quality scores](#)
- ✘ [Per base sequence content](#)
- ! [Per sequence GC content](#)
- ✔ [Per base N content](#)
- ✔ [Sequence Length Distribution](#)
- ✘ [Sequence Duplication Levels](#)
- ✘ [Overrepresented sequences](#)
- ✔ [Adapter Content](#)
- ✘ [Kmer Content](#) ①



# FastQC

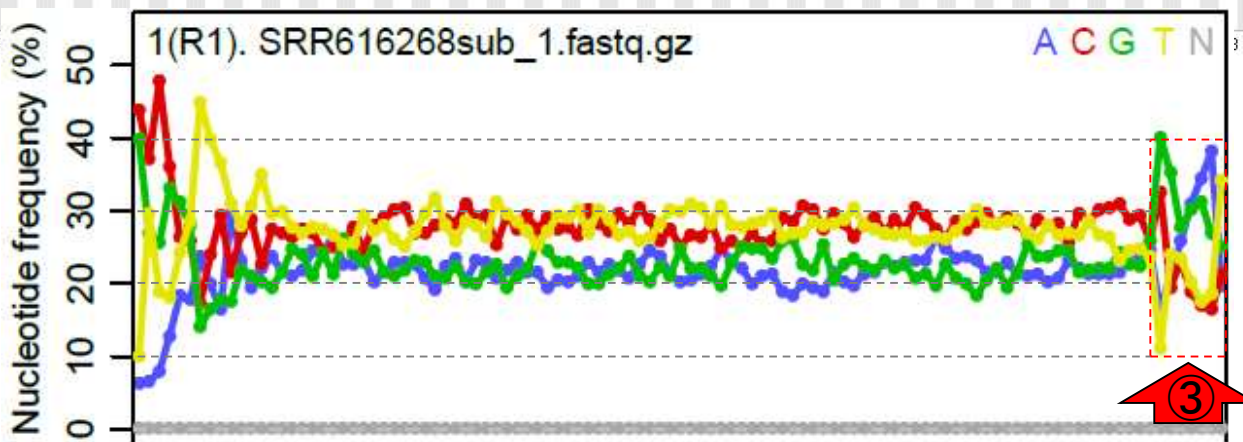
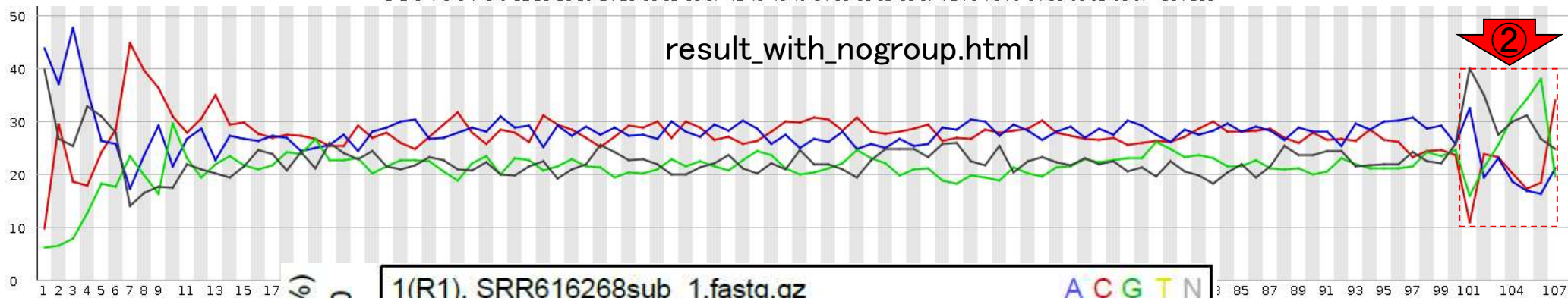
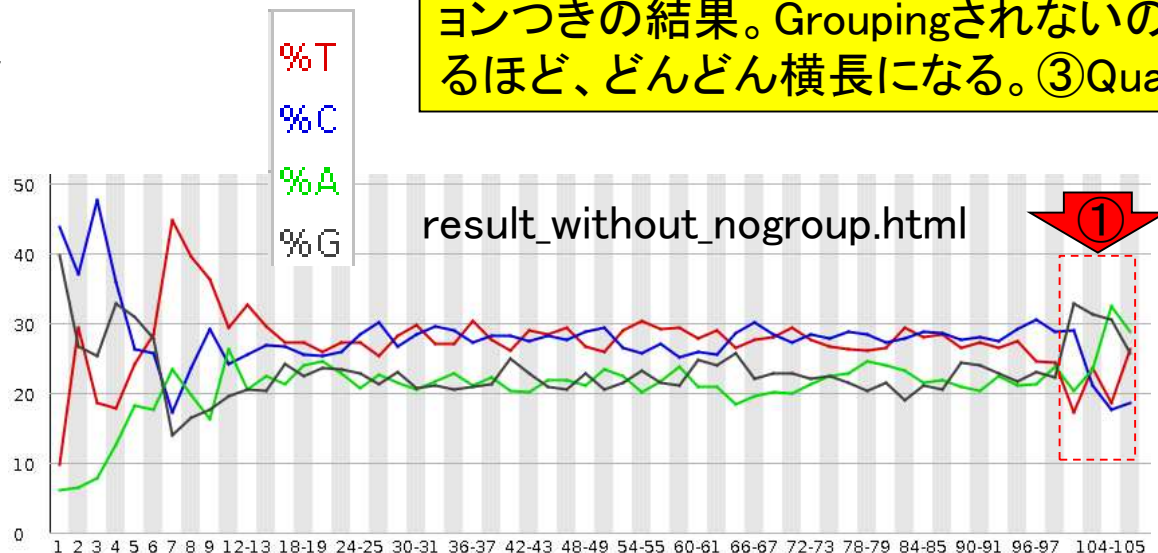
FastQC (ver. 0.11.3)を①デフォルトで実行、②-  
-nogroupオプションつきで実行。③得られたhtml  
ファイル名をresult\_with\_nogroup.htmlに変更。

```
iu@bielinux[~/Desktop/mac_share]
iu@bielinux[mac_share] pwd [11:06午後]
/home/iu/Desktop/mac_share
iu@bielinux[mac_share] ls -la SRR616268sub_1.fastq.gz [11:06午後]
-rwxrwxrwx 1 root root 74906576 7月 13 17:06 SRR616268sub_1.fastq.gz
① iu@bielinux[mac_share] fastqc2 -q SRR616268sub_1.fastq.gz [11:06午後]
iu@bielinux[mac_share] ls -la *.html [11:06午後]
-rwxrwxrwx 1 root root 365019 7月 26 23:06 SRR616268sub_1_fastqc.html
iu@bielinux[mac_share] mv SRR616268sub_1_fastqc.html result_without_nogroup.html
iu@bielinux[mac_share] date [11:07午後]
2015年 7月 26日 日曜日 23:07:50 JST
② iu@bielinux[mac_share] fastqc2 -q --nogroup SRR616268sub_1.fastq.gz [11:07午後]
iu@bielinux[mac_share] mv SRR616268sub_1_fastqc.html result_with_nogroup.html ③
iu@bielinux[mac_share] ls -la result_with* [10:43午前]
-rwxrwxrwx 1 root root 413246 7月 27 10:42 result_with_nogroup.html
-rwxrwxrwx 1 root root 365019 7月 26 23:06 result_without_nogroup.html
iu@bielinux[mac_share] date [10:43午前]
2015年 7月 27日 月曜日 10:43:19 JST
iu@bielinux[mac_share] fastqc2 -v [10:43午前]
FastQC v0.11.3
iu@bielinux[mac_share] █ [10:43午前]
```



# FastQC

①FastQCデフォルトの結果。②FastQC --nogroupオプション付きの結果。Groupingされないで、配列長が長くなるほど、どんどん横長になる。③QuasRデフォルトの結果



# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

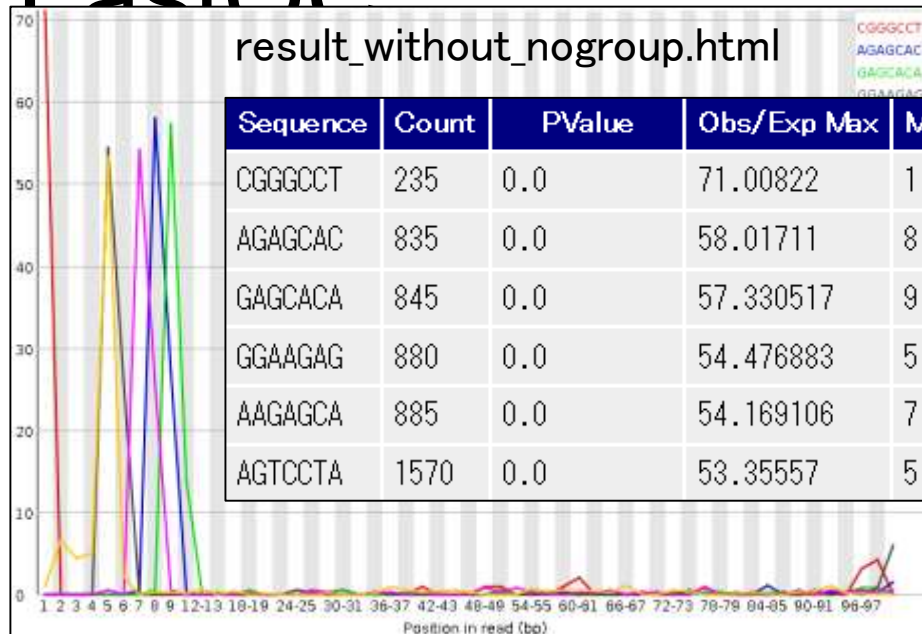
- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

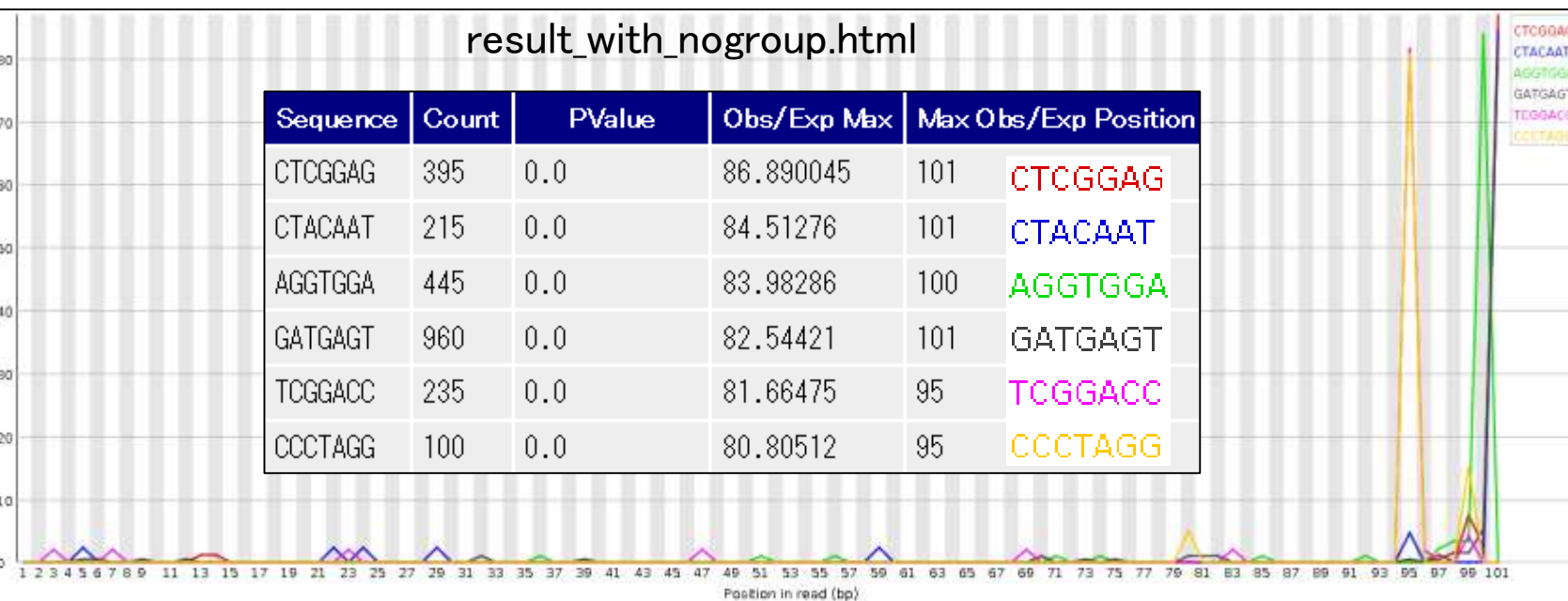
# FastQC

Kmer Contentの項目を比較。違いは--nogroupの有無のみ…何か変。結論としては、--nogroupをつけないときの結果は(門田の感覚では)バグ

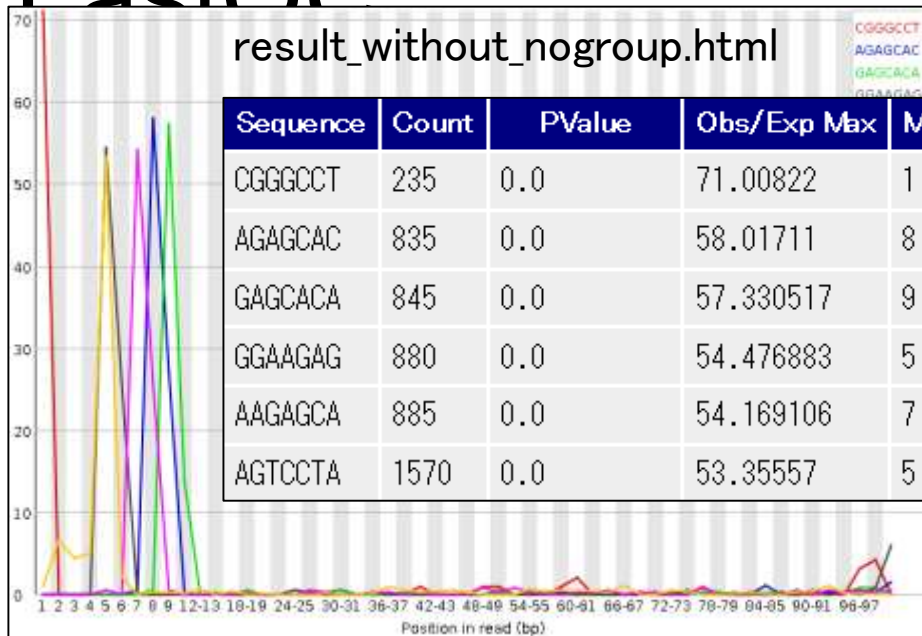


result\_with\_nogroup.html

Sequence	Count	PValue	Obs/Exp Max	Max Obs/Exp Position
CTCGGAG	395	0.0	86.890045	101
CTACAAT	215	0.0	84.51276	101
AGGTGGA	445	0.0	83.98286	100
GATGAGT	960	0.0	82.54421	101
TCGGACC	235	0.0	81.66475	95
CCCTAGG	100	0.0	80.80512	95



# FastQC



開発者側の論理としては、デフォルト(--nogroupオプションをつけない)だとKmer Content解析結果も塩基をgroupingする、という思想なのかもしれない。ただしこれは実質的に5'側が上位に来る傾向となり、3'側に除去すべきアダプターやプライマー配列が存在していてもそれに気づけない(実際門田がこれにハマった)。

```
iu@bielinux[~]
iu@bielinux[iu] fastqc2 -v
FastQC v0.11.3
iu@bielinux[iu] fastqc2 -h | grep -A 5 nogroup
--nogroup      Disable grouping of bases for reads >50bp. All reports will
                show data for every base in the read. WARNING: Using this
                option will cause fastqc to crash and burn if you use it on
                really long reads, and your plots may end up a ridiculous size.
                You have been warned!
iu@bielinux[iu]
```

# 様々な状況証拠を収集

## ■ リファレンス配列がある場合 (マッピング)

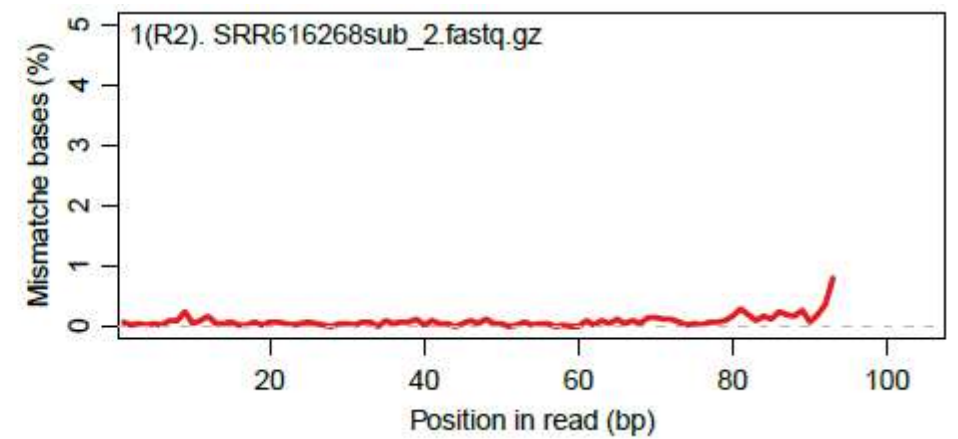
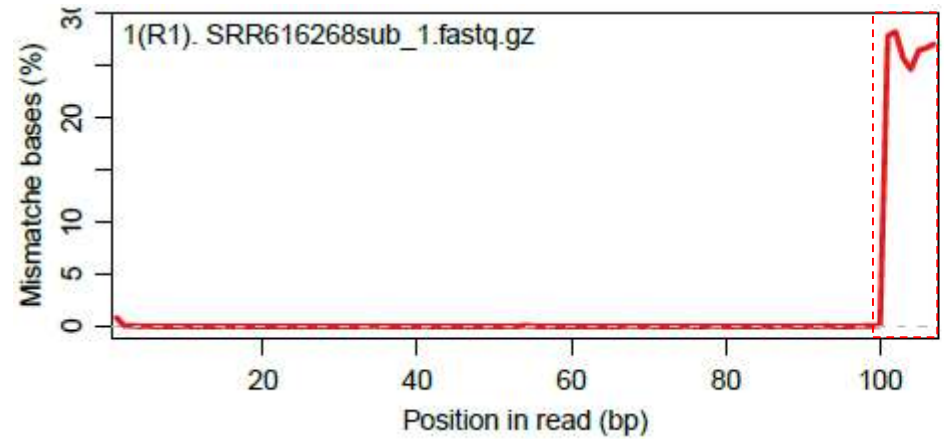
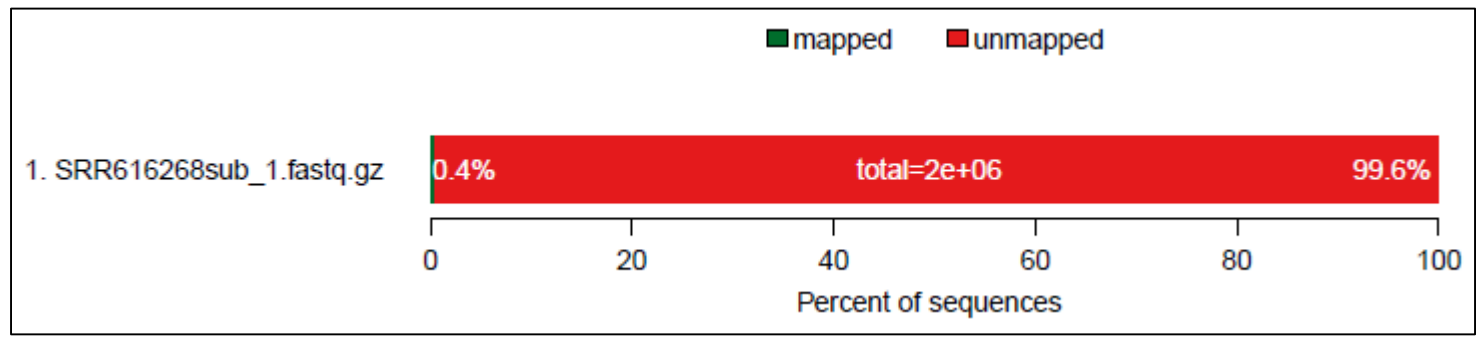
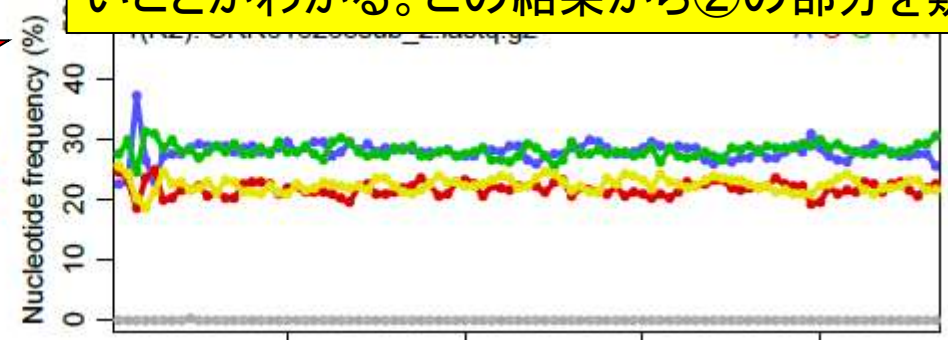
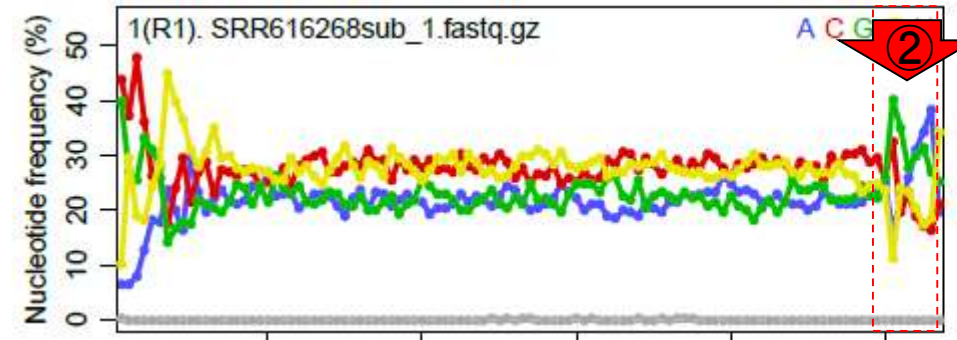
1. 診断: (比較的多めのミスマッチ数で) マッピングを行い、マップされたものの中でミスマッチがあったポジションの分布を眺める。もし3'側にミスマッチの偏りがあれば、それはプライマー/アダプター配列由来塩基という判断を下せる。
2. 対策: 該当部分の塩基を一定数トリム
3. 検証: マップ率が向上するかどうかを調べる。

## ■ リファレンス配列がない場合 (アSEMBル)

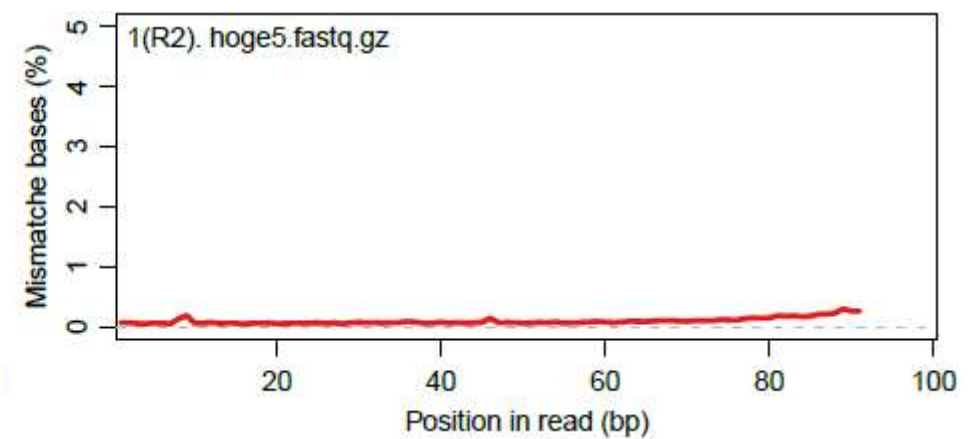
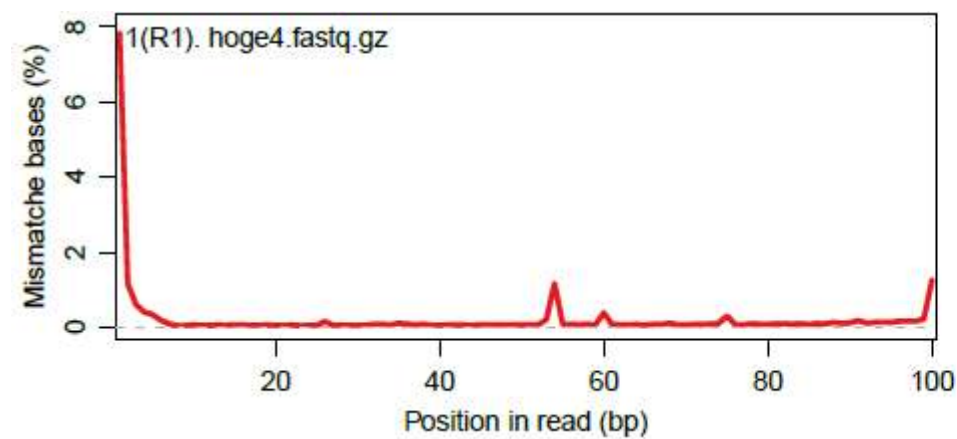
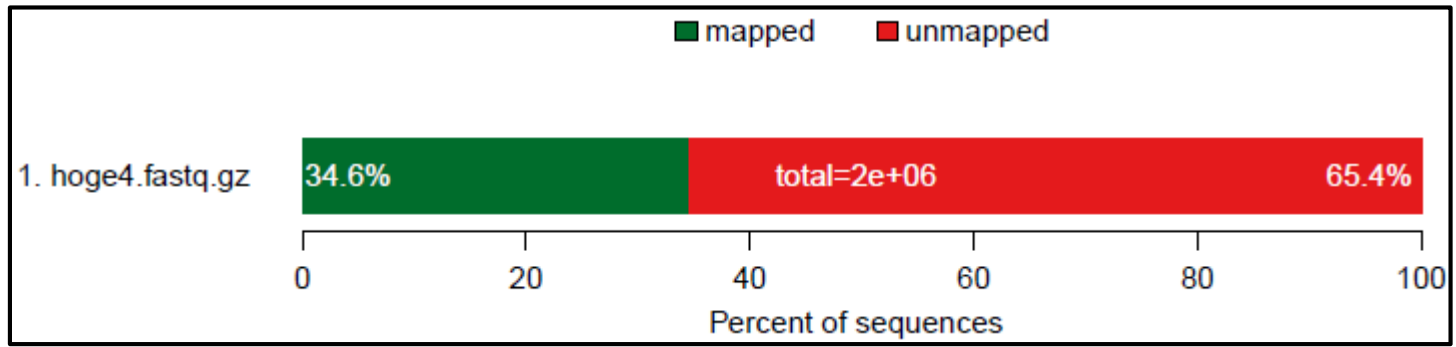
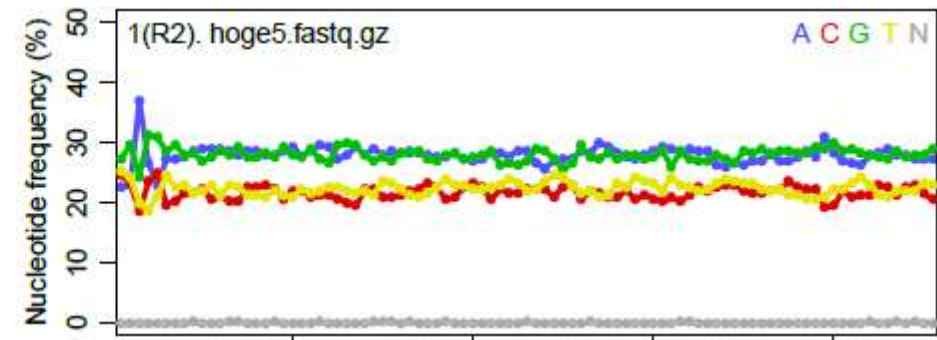
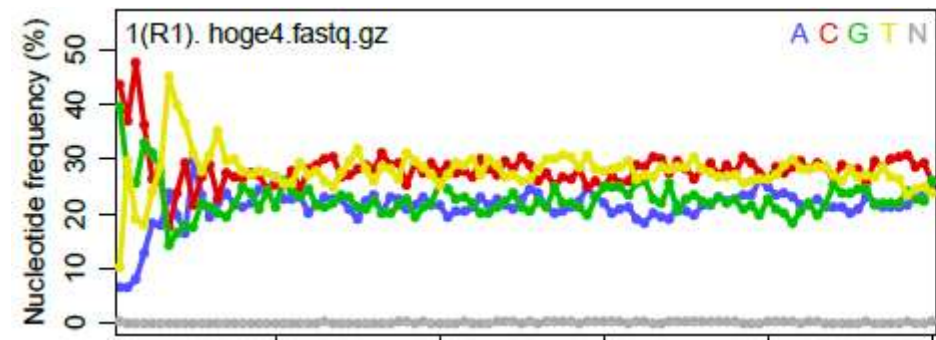
1. 診断: paired-endデータの場合は、single-endデータとしてde novo assembleするとある程度アSEMBルされる(一定数のコンティグが得られる)一方、paired-endデータで実行すると途端にアSEMBルされなくなる。
2. 対策: 1塩基ずつトリムしてはアSEMBル、みたいなことを繰り返して、一定範囲のデータを得る。
3. 検証: アSEMBルが劇的に改善されるところが出てくる(はず)。

# QuasR(トリム前)

QuasRはマッピングもできる。乳酸菌ゲノム配列にマッピングした結果のPDFファイルの一部を表示。  
①ほとんどマップされていない。100万リード中、わずか0.4% (約4,000リード)しかマップされていないことがわかる。この結果から②の部分に疑う。



# QuasR(トリム後)



# 末端塩基のトリム

- 前処理 | トリミング | アダプター配列除去(応用) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/29) NEW
- 前処理 | トリミング | アダプター配列除去(応用) | [ShortRead\(Morgan 2009\)](#) (last modified 2015/06/29) NEW
- 前処理 | トリミング | 指定した末端塩基数だけ除去 **①** (last modified 2015/06/29) NEW
- 前処理 | フィルタリング | PHREDスコアが低い塩基をNに置換 (last modified 2014/03/03)
- 前処理 | フィルタリング | PHREDスコアが低い配列(リード)を除去 (last modified 2014/08)

## 前処理 | トリミング | 指定した末端塩基数だけ除去 NEW

- 前処理 | フォーマット | 3'末端を指定塩基数分だけトリムするやり方を示します。
- 前処理 | フォーマット | 「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

### 1. multi-FASTQ形式ファイルの場合: **②** 4. FASTQ形式ファイル(SRR616268sub\_1.fastq.gz)の場合:

イントロ

```
in_f <- "hoge1.fastq.gz"
out_f <- "hoge2.fastq.gz"
param_trim <- 7

#必要なパッケージをロード
library(Biostrings)
library(ShortRead)

#入力ファイルの読み込み
fasta <- readFastq(in_f)

#本番
hoge <- width(sread(fasta)) - param_trim
```

乳酸菌RNA-seqデータSRR616268の最初の100万リード分(約73MB; 74,906,576 bytes)です。paired-endのforward側です。長さは全て107 bpです。3'側の7 bp分をトリムするので、出力ファイル中の長さは100 bpになります。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge4.fastq.gz" #出力ファイル名を指定してout_fに格納
param_trim <- 7 #3'末端のトリムしたい塩基数を指定

#必要なパッケージをロード
library(Biostrings) #パッケージの読み込み
library(ShortRead) #パッケージの読み込み

#入力ファイルの読み込み
fastq <- readFastq(in_f) #in_fで指定したファイルの読み込み
sread(fastq) #配列情報を表示

#本番
hoge <- width(sread(fastq)) - param_trim #トリム後のend位置情報を取得
hoge[hoge < 1] <- 1 #トリム後のend位置が1未満の場合には1にしている
hoge1 <- DNASTringSet(sread(fastq), start=1, end=hoge) #sread(fastq)から指定した範囲を抽出
hoge2 <- BStringSet(quality(quality(fastq)), start=1, end=hoge) #quality(quality(fastq))から指定した範囲を抽出
fastq <- ShortReadQ(hoge1, hoge2, id(fastq)) #ShortReadQというクラスオブジェクトを作成して
```



①トリム前のリード情報。②7塩基をトリムするので、③赤枠部分をトリムした結果が返される(はず)。

# 末端塩基のトリム

## 4. FASTQ形式ファイル(SRR616268sub\_1.fastq.gz)の場合:

乳酸菌RNA-seqデータSRR616268の最初の100万リード分(約73MB; 74,906,576 bytes)です。paired-endのforward側です。長さは全て107 bpです。3'側の7 bp分をトリムするので、出力ファイル中の長さは100 bpになります。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge4.fastq.gz" #出力ファイル名を指定してout_fに格納
param_trim <- 7 #3'末端のトリムしたい塩基数を指定
```

```
#必要なパッケージをロード
library(Biostrings)
library(ShortRead)
```

```
#入力ファイルの読み込み
fastq <- readFastq(in_f)
sread(fastq)
```

```
#本番
hoge <- width(sread(fastq))
hoge[hoge < 1] <- 1
hoge1 <- DNASTringSet(sread(fastq)[hoge])
hoge2 <- BStringSet(quality(sread(fastq)[hoge]))
fastq <- ShortRead0(hoge1, hoge2)
```

```
R Console
> fastq <- readFastq(in_f)
> sread(fastq)
#in_fで指定したファイルの配列情報を表示
A DNASTringSet instance of length 100000
width seq
[1] 107 AGCCCGACTTTCGTC CCTGCTC...TTCTGAGGGAACCTTCTCTAAC
[2] 107 GATCTGGGCTGTTCCCCTTTCG...CTGAATTCAGTAACCTCCGAAA
[3] 107 CCGGTATATTTTCGGCGCAGTG...GTTGTCTGTGCAACGATCTTAC
[4] 107 CTTGATACCGCCCAAGAACTT...GCCATCTGTGGCCATAAAGACC
[5] 107 CCCCGGTATATTTTCGGCGCAG...TAGTTGTCTGTGCAAGGTCCAT
...
[999996] 107 CCCCGGTATATTTTCGGCGCAG...TAGTTGTCTGTGCAAGGTCCAT
[999997] 107 TTCGGGTCTACATCTGCTTACT...AACGTAACCTCGCCGGGTGAAAT
[999998] 107 CGTCCATCCCGGTCCTCTCGTA...CGTTCTGAACCCAGCCAAGAGT
[999999] 107 CTAGGGAGTATTTAGCCTTGGG...AGGGTTCGACGTTTCCTGGACA
[1000000] 107 GCCTTGTCAATCAAGGTGAGCA...TCTTTTCACCTTGACAGCAGCT
> |
```

# 末端塩基のトリム

①トリム後のリード情報。赤枠部分の塩基がなくなっているのが分かる。また、②配列長も100 bpになっていることがわかる。

```
#入力ファイルの読み込み
fastq <- readFastq(in_f)
sread(fastq)

#本番
hoge <- width(sread(fastq)) - param_trim #トリム後のend位置情報を取得
hoge[hoge < 1] <- 1 #トリム後のend位置が1未満の場合には1にしている
hoge1 <- DNAStringSet(sread(fastq), start=1, end=hoge) #sread(fastq)から指定した範囲を抽出
hoge2 <- BStringSet(quality(quality(fastq)), start=1, end=hoge) #quality(quality(fastq))から指定した範囲を抽出
fastq <- ShortReadQ(hoge1, hoge2, id(fastq)) #ShortReadQというクラスオブジェクトを作成し、hoge1, hoge2, id(fastq)を渡す
sread(fastq)
quality(fastq)

#ファイルに保存
writeFastq(fastq, out_f, compress=FALSE)
```

```
R Console
> sread(fastq) #配列情報を表示
A DNAStringSet instance of length 1000000
      width seq
[1] 100 AGCCCGACTTTCGTCCCTGCTC...CCAACCATTCTGAGGGAACCTT
[2] 100 GATCTGGGCTGTTCCCCTTTCG...AGTTTATCTGAATTCAGTAACC
[3] 100 CCGGTATATTTTCGGCGCAGTG...CATCCTAGTTGTCTGTGCAACG
[4] 100 CTTCGATACCGCCAAGAACTT...CAAAGGTGCCATCTGTGGCCAT
[5] 100 CCCCGGTATATTTTCGGCGCAG...AACATCCTAGTTGTCTGTGCAA
...
[999996] 100 CCCCGGTATATTTTCGGCGCAG...AACATCCTAGTTGTCTGTGCAA
[999997] 100 TTCGGGTCTACATCTGCTTACT...GCAAGCAAACGTAACCTCGCCGG
[999998] 100 CGTCCATCCCGGTCCTCTCGTA...CTCACGACGTTCTGAACCCAGC
[999999] 100 CTAGGGAGTATTTAGCCTTGGG...TGGACGGAGGGTTCGACGTTTC
[1000000] 100 GCCTTGTCAATCAAGGTGAGCA...GCATCCATCTTTTCACCTTGAC
> |
```



# Tips: 例外処理

実用上は、①このあたりの例外処理が重要。このプログラムはリードごとにトリム後の配列長情報をhogeに予め格納している。そしてトリム後の配列長がもし1 bp未満になったら、1 bpは残すような処理をしている。このあたりの想定外への対処が、いわゆる「bug fix」。

```
#入力ファイルの読み込み
fastq <- readFastq(in_f) #in_fで指定したファイル名
sread(fastq) #配列情報を表示

#本番
hoge <- width(sread(fastq)) - param_trim #トリム後のend位置情報を取得
hoge[hoge < 1] <- 1 #トリム後のend位置が1未満の場合には1にしている
hoge1 <- DNAStringSet(sread(fastq), start=1, end=hoge) #sread(fastq)から指定した範囲を抽出
hoge2 <- BStringSet(quality(quality(fastq)), start=1, end=hoge) #quality(quality(fastq))から指定した範囲を抽出
fastq <- ShortReadQ(hoge1, hoge2, id(fastq)) #ShortReadQというクラスオブジェクトを作成
sread(fastq) #配列情報を表示
quality(fastq) #quality情報を表示

#ファイルに保存
writeFastq(fastq, out_f, compress=T) #fastqの中身を指定したファイル名で保存
```

```
R Console
> length(hoge)
[1] 1000000
> hoge[1:10]
[1] 100 100 100 100 100 100 100 100 100 100
> |
```

# Contents

## ■ パッケージ

- CRANとBioconductor
- 推奨パッケージインストール手順のおさらい
- Bioconductor概観 → ゲノム配列パッケージ(BSgenome)
- ヒトゲノム情報パッケージの解析

## ■ プロモーター配列取得(sessionInfo、バージョンの違い)

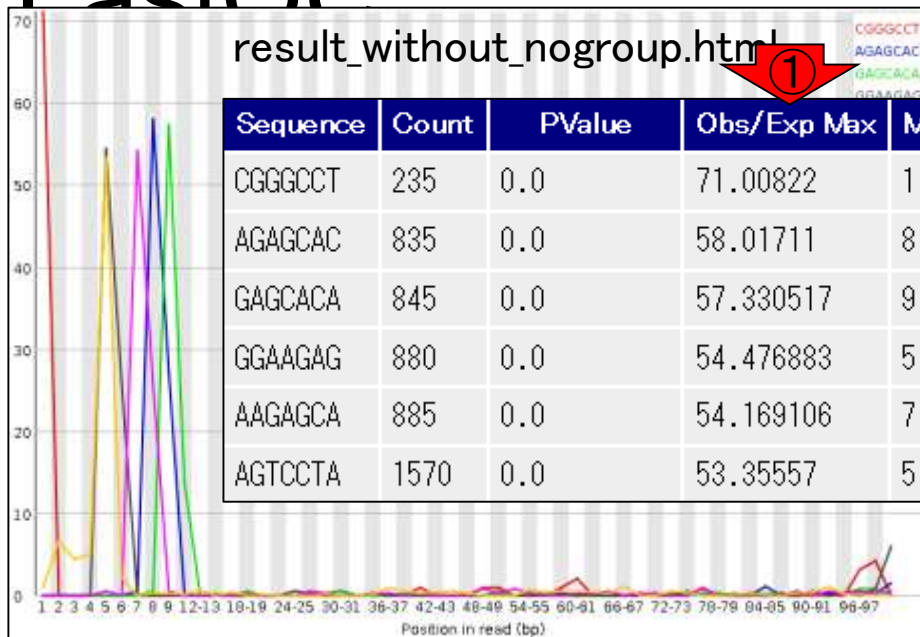
- Rパッケージ(ゲノムとアノテーションパッケージの併用)
- FASTA形式ファイルとGFF3形式ファイル
- データの型

## ■ FASTQファイルの各種解析(LinuxとRを相補的に活用)

- Linux (FastQC)とR (ShortRead)で同じ: Overrepresented sequences項目
- Linux (FastQC)とR (QuasR)で異なる見栄え: Per base sequence content項目。  
FastQCに--nogroupというオプションがあることを知る。
- FastQCのオプション(デフォルトと--nogroupあり)の違いによるKmer Content項目の結果の違い → Rで検証

# FastQC

Kmer Contentの項目で、①Obs/Exp Max  
やプロファイルの図を(若干数値が代わり  
ますが)Rで再現することができます。



Sequence	Count	PValue	Obs/Exp Max	Max Obs/Exp Position
CGGGCCT	235	0.0	71.00822	1
AGAGCAC	835	0.0	58.01711	8
GAGCACA	845	0.0	57.330517	9
GGAAGAG	880	0.0	54.476883	5
AAGAGCA	885	0.0	54.169106	7
AGTCCTA	1570	0.0	53.35557	5

result\_with\_nogroup.html

Sequence	Count	PValue	Obs/Exp Max	Max Obs/Exp Position
CTCGGAG	395	0.0	86.890045	101
CTACAAT	215	0.0	84.51276	101
AGGTGGA	445	0.0	83.98286	100
GATGAGT	960	0.0	82.54421	101
TCGGACC	235	0.0	81.66475	95
CCCTAGG	100	0.0	80.80512	95

# 参考資料

アグリバイオ大学院講義の枠組みで、FastQC、マッピング、アSEMBル、Kmer ContentのRでの検証などについて講義しています。2015.07.07, 2015.06.30, 2015.06.23などでキーワード検索。

## (Rで)塩基配列解析

～NGS, RNA-seq, ゲノム, トランスクリプトーム, 正規化, 発現変動, 統計, (last modified 2015/07/26 since 2011)

- [はじめに](#) (last modified 2015/03/31)
- [参考資料\(講義, 講習会, 本など\)](#) (last modified 2015/07/07) **NEW**
- [過去のお知らせ](#) (last modified 2015/07/08) **NEW**

### What's new

- このウェブサイトでフリー的な利用とめた

## 参考資料(講義, 講習会, 本など) **NEW**

基本的に私個人の個人ページに記載してあるものです。かなり古い講演資料などの情報をもとに勉強されている方

### 講習会, 講義, 講演資料

### 書籍, 学会誌

- 孫建強, [らNGSで](#)  
内容: 日  
違い」と

- 門田幸二「[講義資料](#)」, [アグリバイオインフォマティクス教育研究プログラム](#)の大学院講義科目: [農学生命情報科学特論](#)第4回, 東京大学(東京), 2015.07.07  
内容: 教科書の3.3節周辺。FastQC (ver. 0.11.3)の--nogroupオプションの有無とKmer Contentの項目の挙動の違い。様々な角度で動作確認および検証することの重要性。adapters/primers除去後の乳酸菌paired-end RNA-seqデータのマッピング。アンテーション情報がないときとアンテーションファイル(GFF3)を利用したときのカウント情報取得実例。RPKMの基本的な考え方。配列長とカウント数の関係。原著論文(Blekhman et al., 2010)の公共カウントデータを利用した各種解析。エクセルファイルでRで読み込めること、サブセット抽出や整形のテクニック。サンプル間クラスタリング結果での実験デザインの説明や発現変動解析結果の予想。TCCパッケージを用いた発現変動解析と結果の解釈。同一群同士と異なる群間の発現変動解析結果を示して、分布やモデルの意味を理解。倍率変化の危険性やFDRの結果との比較を通じて、なぜ先人が倍率変化のみで信頼性の高い結果を得てきたのかなど。2コマ(2×90 min)分。
- 門田幸二「[講義資料](#) (2015.07.01版)」, [アグリバイオインフォマティクス教育研究プログラム](#)の大学院講義科目: [農学生命情報科学特論](#)第3回, 東京大学(東京), 2015.06.30  
内容: 教科書の2.3節周辺。NAの取り扱い。paired-endの取り扱いとして、アダプタープライマー(adapters/primers)配列の除去や、Nの数およびリード長でのフィルタリングの基礎(QuasR)と応用(ShortRead)。アSEMBル(ゲノム用とトランスクリプトーム用)。Rockhopperを用いた乳酸菌RNA-seqデータのアSEMBル(single-end (SE)とpaired-end(PE)両方)。マッピングの基礎。仮想データ作成とQuasRを用いたマッピング。出力ファイル形式(SAM/BAMとBED)。bowtieを用いたマッピング時のオプション変更と結果の関係。カウント情報取得。マップされた領域情報取得などを単一・複数リードファイルでデモ。QuasRを用いた乳酸菌RNA-seq paired-endデータのマッピング。QCLレポートの概観とMismatch basesのポジション分布を眺めて、アSEMBルとマッピングの精度向上へのヒントを得る。2コマ(2×90 min)分。
- 門田幸二「[講義資料](#)」, [アグリバイオインフォマティクス教育研究プログラム](#)の大学院講義科目: [農学生命情報科学特論](#)第2回, 東京大学(東京), 2015.06.23