

前回(6/23)のhogeフォルダがデスクトップに残っているかもしれないのでご注意ください。

農学生命情報科学 特論I 第3回

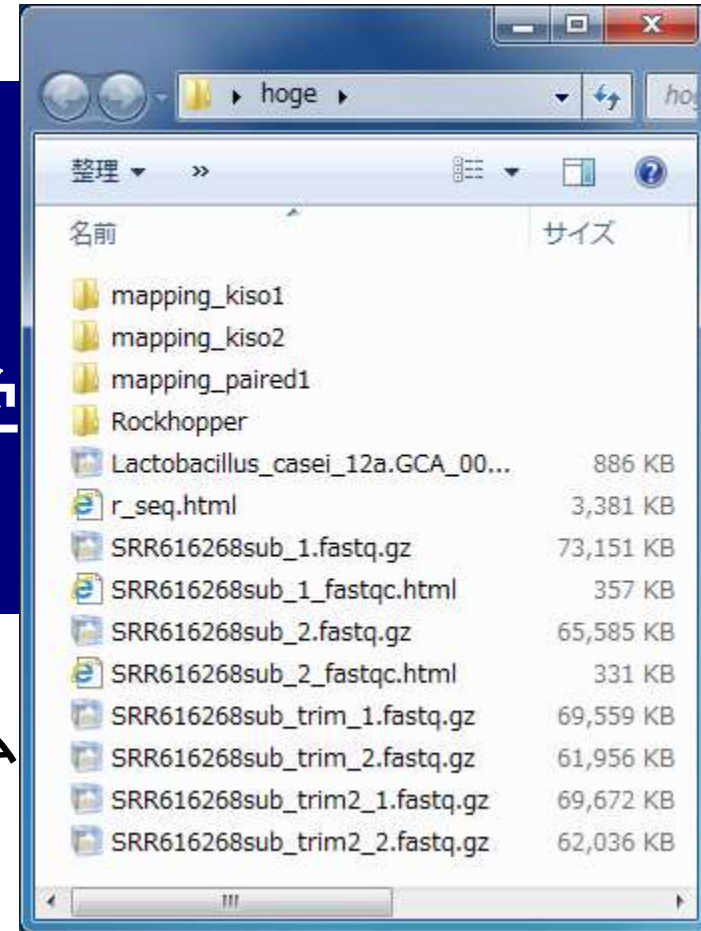
大学院農学生命科学研究科

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

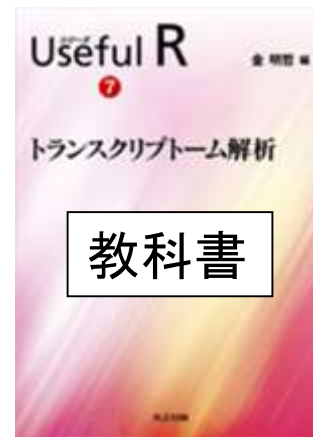
<http://www.iu.a.u-tokyo.ac.jp/~kadota/>



講義予定

NGSの普及により、以前は主にゲノム解析系で必要とされていた配列解析のためのスキルがトランスクリプトーム解析においても要求される時代になっています。本科目では、様々な局面で応用可能な配列解析系のスキルアップを目指し、RNA-Seqに基づくトランスクリプトーム解析を題材とした講義を行います。

- 第1回(2015年6月16日)
 - データベース、データ取得、ファイル形式、Quality Control
 - 教科書の1.3節周辺
- 第2回(2015年6月23日)
 - Quality Control、k-mer解析、トリミング(アダプター配列除去)
- 第3回(2015年6月30日)
 - フィルタリング、アセンブル、マッピング、カウント情報取得
 - 教科書の2.3節周辺
- 第4回(2015年7月7日)
 - クラスタリング、実験デザイン、分布(モデル)、発現変動解析
 - 教科書の3.3節周辺



Contents

■ 先週の課題やエラーについて

- 課題1の解説
- アダプター配列除去(QuasR)実行時のエラーは門田のミス
- Paired-endデータの取り扱い

■ フィルタリング(filtering)

- 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
- 発展形:リード数が異なる場合(ShortReadパッケージ)

■ アセンブル(assembly)

- ゲノム用とトランスクリプトーム用
- Rockhopper2(バクテリアのトランスクリプトーム用)

■ マッピング(mapping)

- 基礎、オプション、仮想データ説明
- マッピング本番、出力ファイル形式、オプションと結果の関係
- カウント情報取得

■ 実データ解析

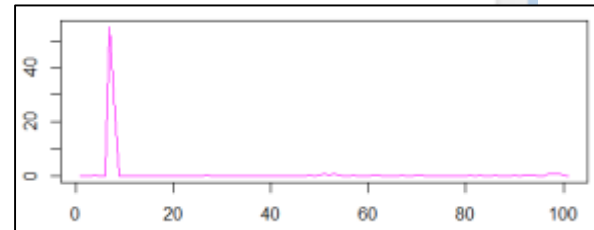
- QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

k-mer解析

「上級1」のコードをテンプレートにして、Position 7で特異的に出現しているAAGAGCAをRで確認してみよう。Head関数出力結果は最初の6個(つまりposition 1-6)までしか表示させていないので、ここではn=8としてposition 1-8まで表示させるようにしている。たしかにposition 7で期待値(Exp = 90.09)の55.15倍の出現回数になっていることが分かる。

```
#以下はおまけ(forループを用いて美しく...上級1)
param_len_ngs <- 107
param_obj <- "CGGGCCT"
Obs <- NULL
hoge <- param_len_ngs - nchar(param_obj) + 1
for(i in 1:hoge){
  Obs <- c(Obs, table(subseq(fasta,
}
head(Obs)
mean(Obs)
Exp <- mean(Obs)
head(Obs/Exp)
plot(Obs/Exp, type="l", col="red")
```

```
R Console
> param_len_ngs <- 107
> param_obj <- "AAGAGCA"
> Obs <- NULL
> hoge <- param_len_ngs - nchar(param_obj) + 1
> for(i in 1:hoge){
+   Obs <- c(Obs, table(subseq(fasta, start=i, width=
+ })
> head(Obs, n=8)
AAGAGCA AAGAGCA AAGAGCA AAGAGCA AAGAGCA AAGAGCA
      3      6      11      27      11      21
AAGAGCA AAGAGCA
      4968      2096
> mean(Obs)
[1] 90.08911
> Exp <- mean(Obs)
> head(Obs/Exp, n=8)
      AAGAGCA      AAGAGCA      AAGAGCA      AAGAGCA
0.03330036  0.06660073  0.12210133  0.29970326
      AAGAGCA      AAGAGCA      AAGAGCA      AAGAGCA
0.12210133  0.23310254 55.14540059 23.26585339
> plot(Obs/Exp, type="l", col="magenta")
> |
```

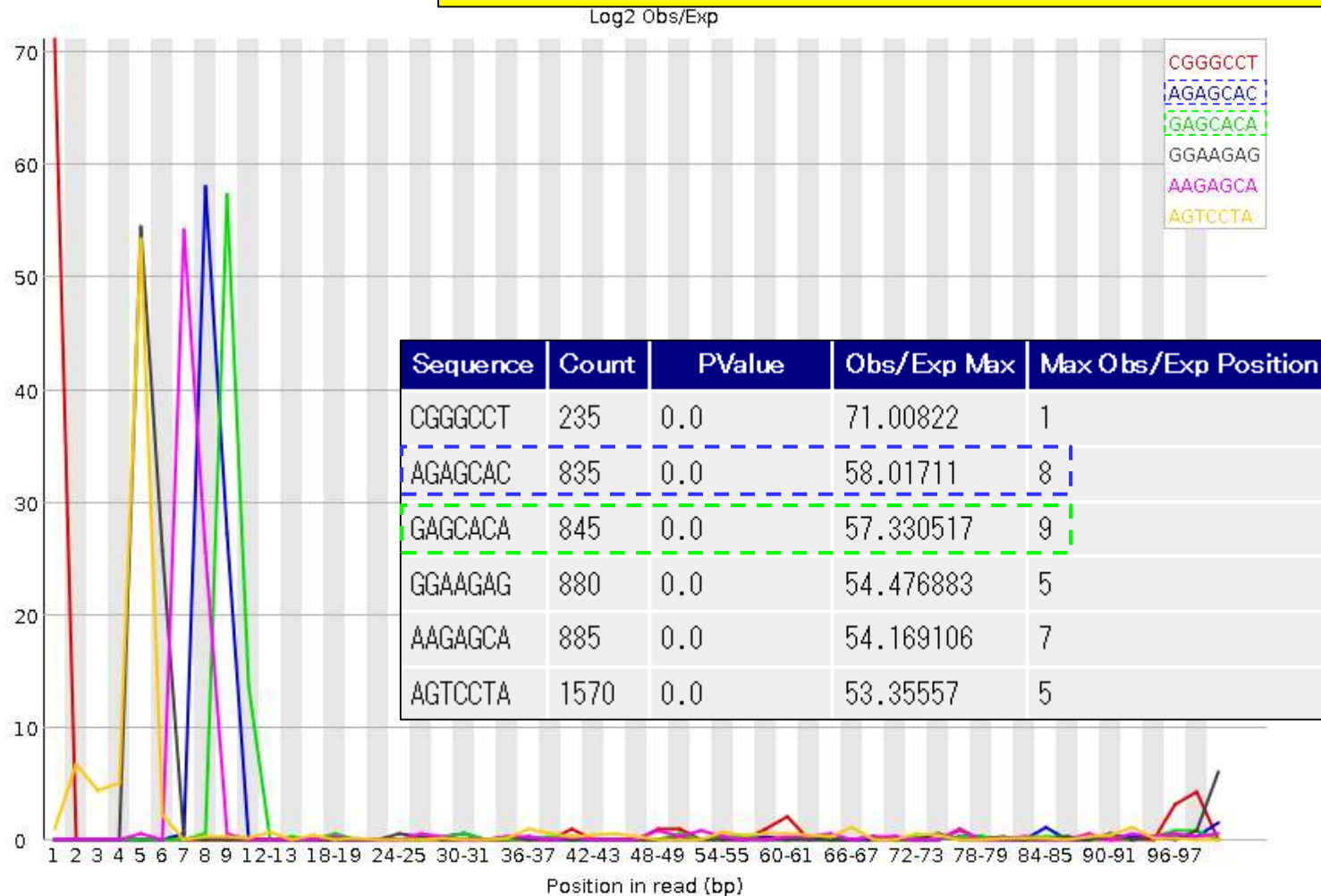


課題1

「上級1」のコードをテンプレートにして、① AGAGCAC、② GAGCACA、③長さを含めて任意、のk-mer解析を行い、得られた結果を示せ。

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**




課題1の解説

①AGAGCACのk-mer解析時にエラーが出たと思います。この原因は、AGAGCACが1つも存在しないポジションがあった場合に、そこがNAとなるからです。mean関数実行時に要素中にNAが1つでもあると計算できないという罠があったのです。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin
↓
#必要なパッケージをロード↓
library(Biostrings) #パッケージの読み込み↓
↓
#入力ファイルの読み込み↓
fasta <- readDNASTringSet(in_f, format="fastq")
fasta #確認
↓
#以下はおまけ(forループを用いて美しく...上)
param_len_ngs <- 107↓
param_obj <- "AGAGCAC"↓
Obs <- NULL↓
hoge <- param_len_ngs - nchar(param_obj)
for(i in 1:hoge){↓
  Obs <- c(Obs, table(subseq(fasta, start=i, end=i+param_len_ngs-1, step=1)))
}↓
head(Obs, n=10)↓
mean(Obs)↓
Exp <- mean(Obs)↓
head(Obs/Exp, n=10)↓
plot(Obs/Exp, type="l", col="blue")↓
```

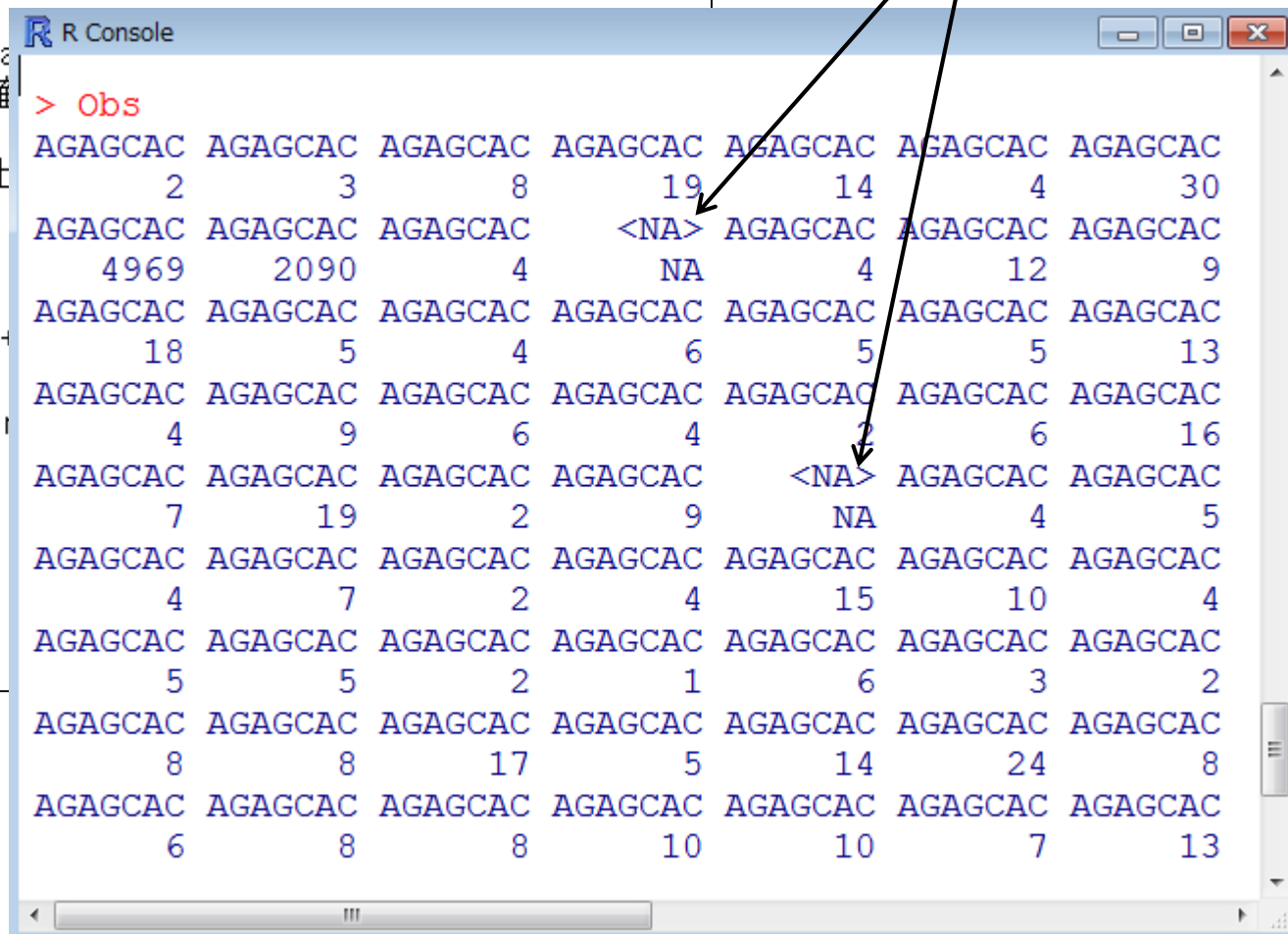
```
R Console
> head(Obs, n=10)
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC
      4969      2090      4
> mean(Obs)
[1] NA
> Exp <- mean(Obs)
> head(Obs/Exp, n=10)
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      NA      NA      NA      NA      NA      NA      NA
AGAGCAC AGAGCAC AGAGCAC
      NA      NA      NA
> plot(Obs/Exp, type="l", col="blue")
Error in plot.window(...) : 有限な 'ylim' の値が必要で$
追加情報: Warning messages:
1: In min(x) : min の引数に有限な値がありません: Inf $
2: In max(x) : max の引数に有限な値がありません: -Inf $
> |
```



課題1の解説

Obsと打ち込んで、何が起きているのかを悟る。画面に見えているだけでも、AGAGCACが存在しないpositionが少なくとも2ヶ所あることがわかる。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin_fに格納↓
↓
#必要なパッケージをロード↓
library(Biostrings) #パッケージの読み込み↓
↓
#入力ファイルの読み込み↓
fasta <- readDNASTringSet(in_f, format="fastq")
fasta #確認
↓
#以下はおまけ(forループを用いて美しく...上)
param_len_ngs <- 107↓
param_obj <- "AGAGCAC"↓
Obs <- NULL↓
hoge <- param_len_ngs - nchar(param_obj)
for(i in 1:hoge){↓
  Obs <- c(Obs, table(subseq(fasta, start=i, end=i+param_len_ngs-1, step=1)))
}↓
head(Obs, n=10)↓
mean(Obs)↓
Exp <- mean(Obs)↓
head(Obs/Exp, n=10)↓
plot(Obs/Exp, type="l", col="blue")↓
```



課題1の解説

最もシンプルな対処法は、mean関数中のna.rmオプションをTRUE (デフォルトはFALSE)にするやり方。NAの要素を読み飛ばすオプションなので、平均値の計算が若干不正確になるが、AGAGCACがposition 8で相対的に高頻度に出現しているということを調べるという点では全く問題ない。

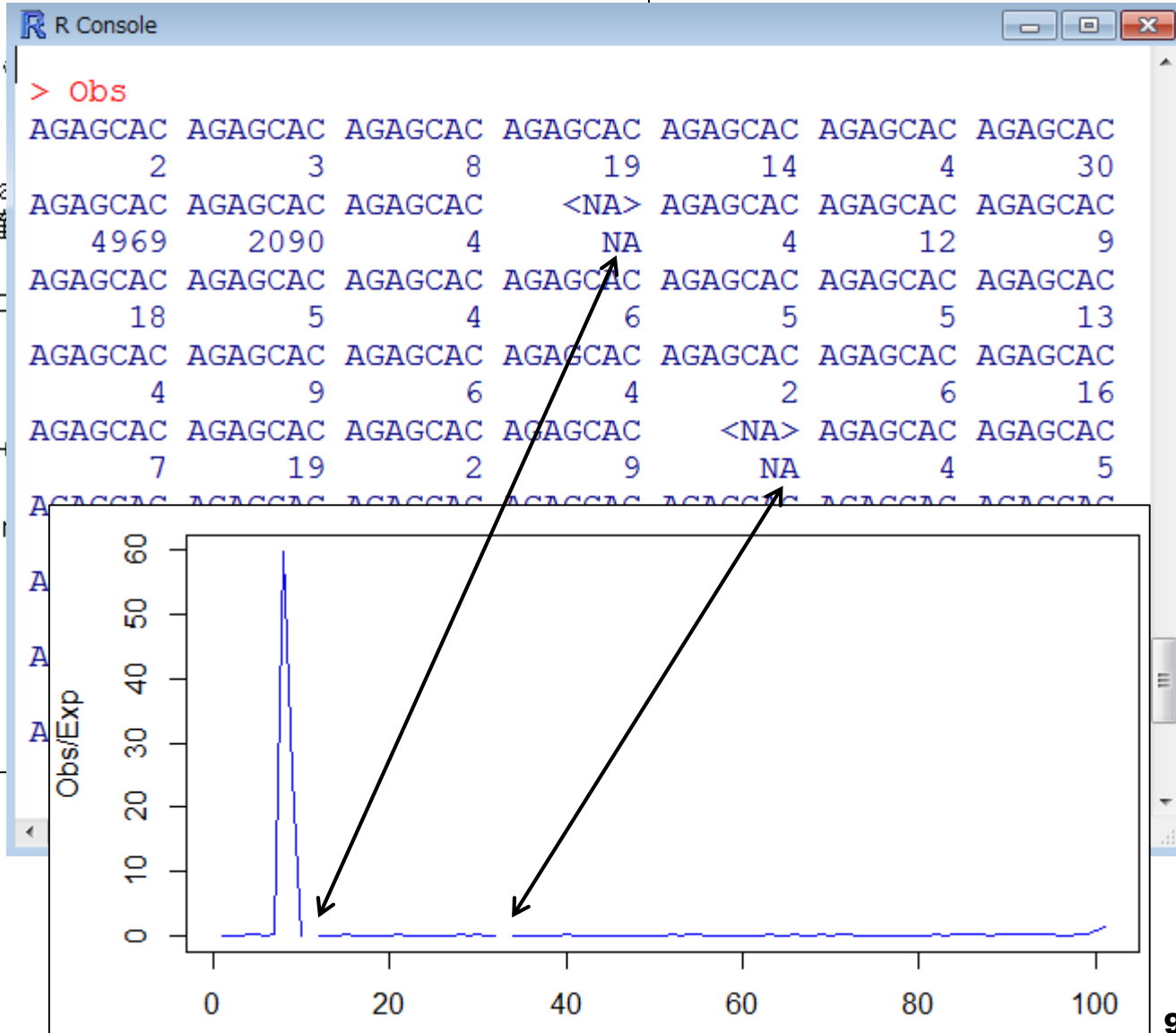
```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定して
↓
#必要なパッケージをロード↓
library(Biostrings) #パッケージの読み込み↓
↓
#入力ファイルの読み込み↓
fasta <- readDNASTringSet(in_f, format="fastq")
fasta #確認
↓
#以下はおまけ(forループを用いて美しく...上)
param_len_ngs <- 107↓
param_obj <- "AGAGCAC"↓
Obs <- NULL↓
hoge <- param_len_ngs - nchar(param_obj)
for(i in 1:hoge){↓
  Obs <- c(Obs, table(subseq(fasta, start=i, end=i+param_len_ngs, step=1)))
}↓
head(Obs, n=10)↓
mean(Obs, na.rm=TRUE)↓
Exp <- mean(Obs, na.rm=TRUE)↓
head(Obs/Exp, n=10)↓
plot(Obs/Exp, type="l", col="blue")↓
```

```
R Console
> head(Obs, n=10)
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC
      4969      2090      4
> mean(Obs, na.rm=TRUE)
[1] 83.08081
> Exp <- mean(Obs, na.rm=TRUE)
> head(Obs/Exp, n=10)
      AGAGCAC      AGAGCAC      AGAGCAC      AGAGCAC
0.02407295 0.03610942 0.09629179 0.22869301
      AGAGCAC      AGAGCAC      AGAGCAC      AGAGCAC
0.16851064 0.04814590 0.36109422 59.80924012
      AGAGCAC      AGAGCAC
25.15623100 0.04814590
> plot(Obs/Exp, type="l", col="blue")
> |
```


課題1の解説

平均値(この場合83.08081)さえ計算できれば、Obs/Expを計算可能。つまり、プロファイルを描画することができる。NAとなっている場所が2ヶ所あることがわかる。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin
↓
#必要なパッケージをロード↓
library(Biostrings) #パ
↓
#入力ファイルの読み込み↓
fasta <- readDNASTringSet(in_f, format="fa
fasta #確
↓
#以下はおまけ(forループを用いて美しく...上
param_len_ngs <- 107↓
param_obj <- "AGAGCAC"↓
Obs <- NULL↓
hoge <- param_len_ngs - nchar(param_obj) ↓
for(i in 1:hoge){↓
  Obs <- c(Obs, table(subseq(fasta, star
})↓
head(Obs, n=10)↓
mean(Obs)↓
Exp <- mean(Obs)↓
head(Obs/Exp, n=10)↓
plot(Obs/Exp, type="l", col="blue")↓
```



課題1の解説

もっとよりよい方法は、①NAとなった要素に0を入れること。is.na関数は、要素がNAのときにTRUE、そうでないときにFALSEを返す。

```
R Console
> mean(Obs)
[1] NA
> mean(Obs, na.rm=TRUE)
[1] 83.08081
> head(Obs, n=14)
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC <NA> AGAGCAC AGAGCAC AGAGCAC
      4969      2090      4      NA      4      12      9
> Obs[1:14]
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC <NA> AGAGCAC AGAGCAC AGAGCAC
      4969      2090      4      NA      4      12      9
> is.na(Obs)[1:14]
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      FALSE      FALSE      FALSE      FALSE      FALSE      FALSE      FALSE
AGAGCAC AGAGCAC AGAGCAC <NA> AGAGCAC AGAGCAC AGAGCAC
      FALSE      FALSE      FALSE      TRUE      FALSE      FALSE      FALSE
> |
```

課題1の解説

0代入前後で出現頻度の平均値が83.08081から81.43564に少し下がっているのがわかります。

```
R Console
> Obs[1:14]
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC <NA> AGAGCAC AGAGCAC AGAGCAC
      4969      2090      4      NA      4      12      9
> is.na(Obs)[1:14]
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      FALSE      FALSE      FALSE      FALSE      FALSE      FALSE      FALSE
AGAGCAC AGAGCAC AGAGCAC <NA> AGAGCAC AGAGCAC AGAGCAC
      FALSE      FALSE      FALSE      TRUE      FALSE      FALSE      FALSE
> Obs[is.na(Obs)] <- 0
> head(Obs, n=14)
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC <NA> AGAGCAC AGAGCAC AGAGCAC
      4969      2090      4      0      4      12      9
> mean(Obs)
[1] 81.43564
> mean(Obs, na.rm=TRUE)
[1] 81.43564
> |
```

課題1の解説

こんな感じになります。2015年6月23日の課題を該当部分のみ差し替えたいヒトは、回収箱に入れてかまいません。課題用紙は後ろにあります。先週提出分と突き合わせて評価しますので、考察部分など差し替える必要がない部分は空欄で構いません。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定して
↓
#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(in_f, format="fastq")
fasta #確認
↓
#以下はおまけ(forループを用いて美しく...上)
param_len_ngs <- 107↓
param_obj <- "AGAGCAC"↓
Obs <- NULL↓
hoge <- param_len_ngs - nchar(param_obj)
for(i in 1:hoge){↓
  Obs <- c(Obs, table(subseq(fasta, start=i, end=i+param_len_ngs-1, step=1)))
}↓
Obs[is.na(Obs)] <- 0 ①
head(Obs, n=10)↓
mean(Obs, na.rm=TRUE)↓
Exp <- mean(Obs, na.rm=TRUE)↓
head(Obs/Exp, n=10)↓
plot(Obs/Exp, type="l", col="blue")↓
```

```
R Console
> Obs[is.na(Obs)] <- 0
> head(Obs, n=10)
AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC AGAGCAC
      2      3      8      19      14      4      30
AGAGCAC AGAGCAC AGAGCAC
 4969    2090      4
> mean(Obs, na.rm=TRUE)
[1] 81.43564
> Exp <- mean(Obs, na.rm=TRUE)
> head(Obs/Exp, n=10)
      AGAGCAC      AGAGCAC      AGAGCAC      AGAGCAC
0.02455927 0.03683891 0.09823708 0.23331307
      AGAGCAC      AGAGCAC      AGAGCAC      AGAGCAC
0.17191489 0.04911854 0.36838906 61.01750760
      AGAGCAC      AGAGCAC
25.66443769 0.04911854
> plot(Obs/Exp, type="l", col="blue")
> |
```

バージョンアップ

ひっそりと修正しています。レポート中で問題点の指摘や解決法を示して頂いた方々に感謝m(_ _)m



- イントロ | NGS | 読み込み | FASTA形式 | [基本情報を取得](#) (last modified 2014/08/18)
- イントロ | NGS | 読み込み | FASTA形式 | [description行の記述を整形](#) (last modified 2014/04/05)
- イントロ | NGS | 読み込み | FASTQ形式 | **基礎** (last modified 2015/06/24) **1**
- イントロ | NGS | 読み込み | FASTQ形式 | **応用** (last modified 2015/06/18) **1**

イントロ | NGS | 読み込み | FASTQ形式 | 基礎 **NEW**

Sanger FASTQ形式ファイルを読み込むやり方を示します。「基礎」では、FASTQファイルの中身を全て読み込む手順を示します。入力・出力形式は、ともに非圧縮(fasta)・gzip圧縮(fasta.gz)ファイルが可能です。

- イントロ | ファイル **2**
- イントロ | 1. サンプルデータ
- イントロ | [SRR0374397](#) (et al., 2010)。
- イントロ | quality情報を
- イントロ | in_f <- "
- イントロ | #必要なパッケージをロ
- イントロ | library(Biostrings)
- イントロ | #入力ファイルの読み込
- イントロ | fasta <- readDNASTr
- イントロ | fasta
- イントロ | #以下はおまけ(部分配列
- イントロ | hoge <- subseq(fasta
- イントロ | hoge
- イントロ | head(table(hoge))
- イントロ | head(sort(table(hoge
- イントロ | table(hoge)["CGGGCC
- イントロ | hoge <- subseq(fasta
- イントロ | hoge
- イントロ | table(hoge)["CGGGCC

8. 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。NAを含む場合への各種対応策を2015年6月24日に追加しました(茂木朋貴氏、野間口達洋氏、他多くの受講生提供情報)。

```
in_f <- "SRR616268sub 1.fastq.gz"
#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 <- readDNASTring(fasta)

#以下はおまけ(部分配列)
hoge <- subseq(fasta, 1, 100000)
hoge
head(table(hoge))
head(sort(table(hoge)))
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, 1, 100000)
hoge
table(hoge)["CGGGCCT"]

#以下はおまけ(forループを用いて美しく...上級1)
param_len_ngs <- 107 #リード長を指定
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
hoge <- param_len_ngs - nchar(param_obj) + 1 #positionの右端の値を計算してhogeに格納
for(i in 1:hoge){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=nchar(param_obj)))[param_obj])
}
Obs[is.na(Obs)] <- 0 #NAの位置に0を代入
head(Obs) #最初の6個の要素を表示
mean(Obs, na.rm=TRUE) #平均値
Exp <- mean(Obs, na.rm=TRUE) #平均値をExpとして取り扱う
head(Obs/Exp)
plot(Obs/Exp, type="l", col="red")

#以下はおまけ(forループを用いて美しく...上級2)
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
hoge <- width(fasta)[1] - nchar(param_obj) + 1 #positionの右端の値を計算してhogeに格納
for(i in 1:hoge){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=nchar(param_obj)))[param_obj])
}
```

Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

アダプター配列除去

param_nrecのところをいくら変えてもエラーが出続ける問題に多くの受講生が遭遇しました。理由はほぼ間違いなく、単純に私がpreprocessReads関数中に明示的にparam_nrecを与えていなかったからです。結果として、それが適切に反映されずにデフォルトの1,000,000がいつまでも使われていたということでしょう。

前処理 | トリミング | アダプター配列除去(基礎) | QuasR

QuasRパッケージを用いたアダプター配列除去の基本形を示します。param_nrecオプションは、一度に処理するリード数を指定しているのですが、基本的に無視で構いません。デフォルトの1000000のときに、メモリ不足でフリーズしたの

4. gzip圧縮状態のFASTQ形式ファイル(SRR609266.fastq.gz)の場合:

small RNA-seqデータ(ファイルサイズは400Mb弱、11928428リード)です。このファイルに対するFastQC実行結果として「RNA PCR Primer, Index 1」(RPI1)が含まれているとレポートされた49 bpをアダプター配列として入力しています。ここでは、アダプター配列以外はデフォルトで実行しています。アダプター配列の位置は5'側(左側)ではなく3'側(右側)にあるという前提であり、右側のアダプター配列しかトリミングしないやり方です。それが、preprocessReads関数実行時にRpatternのみ記載している理由です。約5分。

```
in_f <- "SRR609266.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_adapter <- "TGGAATTCTCGGGTCCAAGGAAGTCCAGTCACATCACGATCTCGTATG" #アダプター配列
param_nrec <- 500000 #一度に処理するリード数を指定

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #アダプター配列除去を実行
                       outputFilename=out_f, #アダプター配列除去を実行
                       Rpattern=param_adapter) #アダプター配列除去を実行

res #確認してるだけです
file.size(in_f) #入力ファイルのサイズを表示
file.size(out_f) #出力ファイルのサイズを表示
```

1. gzip圧縮状態のFASTQ形式ファイル(SRR609266.fastq.gz)の場合:

small RNA-seqデータ(ファイルサイズは400Mb弱、11928428リード)です。このファイルに対するFastQC実行結果として「RNA PCR Primer, Index 1」(RPI1)が含まれているとレポートされた49 bpをアダプター配列として入力しています。ここでは、アダプター配列以外はデフォルトで実行しています。アダプター配列の位置は5'側(左側)ではなく3'側(右側)にあるという前提であり、右側のアダプター配列しかトリミングしないやり方です。それが、preprocessReads関数実行時にRpatternのみ記載している理由です。約5分。

```
in_f <- "SRR609266.fastq.gz"
out_f <- "hoge4.fastq.gz"
param_adapter <- "TGGAATTCTCGGGTCCAAGGAAGTCCAGTCACATCACGATCTCGTATG"
param_nrec <- 500000

#必要なパッケージをロード
library(QuasR)

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f,
                       outputFilename=out_f,
                       Rpattern=param_adapter)
```

休み時間にも動作確認してみてください。きつとうまくいくことでしょう。

アダプター配列除去

- 前処理 | フィルタリング | [GFF/GTF形式ファイル](#) (last modified 2013/10/10)
- 前処理 | フィルタリング | 組合せ | [ACGTのみ & 指定した長さの範囲の配列](#) (last modified 2014/01/01)
- 前処理 | トリミング | ポリア配列除去 | [ShortRead\(Morgan 2009\)](#) (last modified 2014/01/01)
- 前処理 | トリミング | アダプター配列除去(基礎) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/01/01) **1**
- 前処理 | トリミング | アダプター配列除去(基礎) | [girafe\(Toedling 2010\)](#) (last modified 2010/01/01)

前処理 | トリミング | アダプター配列除去(基礎) | QuasR(Gaidatzis 2015) NEW

QuasRパッケージを用いたアダプター配列除去の基本形を示します。param_nrecオプションは、一度に処理するリード数を指定しているのですが、基本的に無視で構いません。デフォルトの1000000のときに、メモリ不足でフリーズしたの

で、デフォルトの2000000に増やして実行したところ、メモリ不足でフリーズしたのと同じです。それでも「ファイル」-「ディレクトリ」

4. gzip圧縮FASTQ形式ファイル(SRR609266.fastq.gz)の場合:

small RNA-seqデータ(ファイルサイズは400Mb弱、11928428リード)です。このファイルに対するFastQC実行結果として「RNA PCR Primer, Index 1」(RPI1)が含まれているとレポートされた49 bpをアダプター配列として入力しています。ここでは、アダプター配列以外はデフォルトで実行しています。アダプター配列の位置は5'側(左側)ではなく3'側(右側)にあるという前提であり、右側のアダプター配列しかトリミングしないやり方です。それが、preprocessReads関数実行時にRpatternのみ記載している理由です。約5分。

1. gzip圧縮状態のFASTQ形式ファイル

small RNA-seqデータ中の記述からGSE41410(Zhu 2013)の7を参照すると、アダプター配列は"TGGAATTCTCGGGTGCCTCAAGGAACTCCAGTCACATCACGATCTCGTATG"以外はデフォルトであり、右側のアダプター配列しかトリミングしない理由です。約5分。

```

in_f <- "SRR609266.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 adapter <- "TGGAATTCTCGGGTGCCTCAAGGAACTCCAGTCACATCACGATCTCGTATG" #アダプター配列
param_nrec <- 500000 #一度に処理するリード数を指定

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                        outputFilename=out_f, #out_fを出力とする
                        Rpattern=param_adapter, #リードの右側(Right)のアダプターを除去
                        nrec=param_nrec) #一度に処理するリード数

res #確認してるだけです
file.size(in_f) #入力ファイルのサイズを表示
file.size(out_f) #出力ファイルのサイズを表示
    
```


Paired-endの取扱い

- 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 2015/06/22) 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 2015/06/23) 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 | 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 | アダプター配列除去(応用) | QuasR(Gaidatzis_2015) NEW

- 前処理 | トリミング | ...
- アセンブル | ...
- アセンブル | ...
- マッピング | ...
- マッピング | ...
- マッピング | ...
- マッピング | ...
- マッピング | ...
- マッピング | ...
- マッピング | ...
- マッピング | ...

QuasRパッケージをインストールしてR Guiが終了した後に、6月21日にparam_nrecを指定して実行してください。デフォルトは「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は「GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG」と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリミングしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。残念ながら、QuasR (ver. 1.8.2)では「Removing adapters from paired-end samples is not yet supported」と出ます。2015年6月23日に開発者に実装を心待ちにしているとメールをしておきました。

```

in_f <- "SRR6092
out_f <- "hoge1.
param_adapter <-

```

2. gzip圧縮FASTQ形式ファイル(SRR616268sub_1.fastq.gzとSRR616268sub_2.fastq.gz)の場合:
 乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。paired-endデータです。SRR616268sub_1.fastq.gzは、約75MB、全リード107 bpです。SRR616268sub_2.fastq.gzは、約67MB、全リード93 bpです。FastQC実行結果として「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は「GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG」と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリミングしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。残念ながら、QuasR (ver. 1.8.2)では「Removing adapters from paired-end samples is not yet supported」と出ます。2015年6月23日に開発者に実装を心待ちにしているとメールをしておきました。

```

in_f1 <- "SRR616268sub_1.fastq.gz"
in_f2 <- "SRR616268sub_2.fastq.gz"
out_f1 <- "hoge1_1.fastq.gz"
out_f2 <- "hoge1_2.fastq.gz"
param_adapter <- "GATCGGAAGAGCACACG
param_nrec <- 500000

#必要なパッケージをロード
library(QuasR)

#本番(前処理)
res <- preprocessReads(filename=in_
  filenameMate=in_f2,
  outputFilename=out_f1,
  outputFilenameMate=out_f2,
  Lpattern=param_adapter)
res

```

```

R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> list.files(pattern="sub_..fastq")
[1] "SRR616268sub_1.fastq.gz"
[2] "SRR616268sub_1_fastqc.html"
[3] "SRR616268sub_2.fastq.gz"
[4] "SRR616268sub_2_fastqc.html"
> |

```

#前処理を実行
 #前処理を実行
 #確認してるだけです

Paired-endの取扱い

QuasR (ver. 1.8.2)では、まだpaired-endデータのアダプター配列除去には対応できていないようだといいましたが、その後すぐに開発者が回避策(workaround)を教えてくださいました。

2. gzip圧縮FASTQ形式ファイル(SRR616268sub_1.fastq.gzとSRR616268sub_2.fastq.gz)の場合

乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。paired-endデータです。SRR616268sub_1.fastq.gzは、約75MB、全リード107 bpです。SRR616268sub_2.fastq.gzは、約67MB、全リード93 bpです。

TruSeq Adapter, Index 3が含まれているとレポートされました。これでググると塩基配列情報は

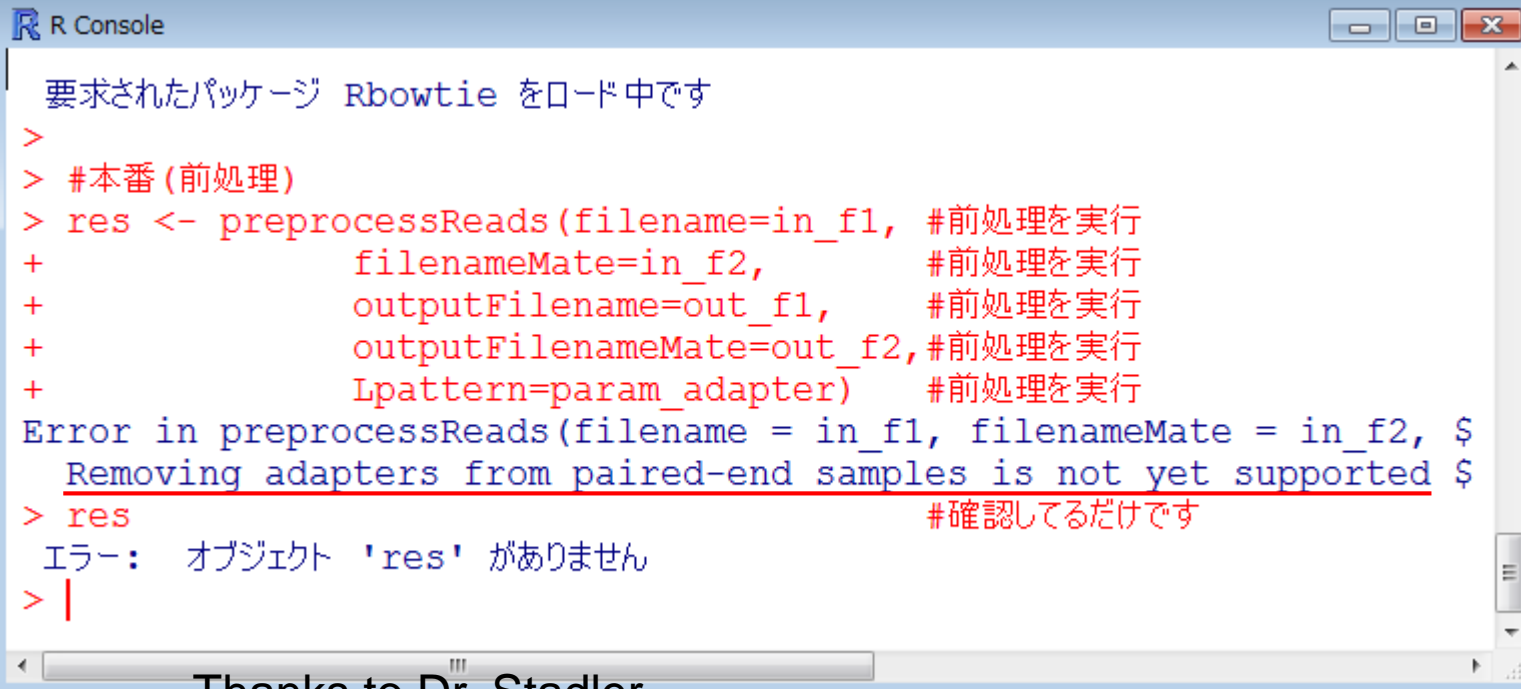
"GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG"と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリミングしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。残念ながら、QuasR (ver. 1.8.2)では「Removing adapters from paired-end samples is not yet supported」と出ます。2015年6月23日に開発者に実装を心待ちにしているとメールをしておきました。

```
in_f1 <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin_f1に格納
in_f2 <- "SRR616268sub_2.fastq.gz" #入力ファイル名を指定してin_f2に格納
out_f1 <- "hoge1_1.fastq.gz" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge1_2.fastq.gz" #出力ファイル名を指定してout_f2に格納
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG"
param_nrec <- 500000
```

```
#必要なパッケージをロード
library(QuasR)
```

```
#本番(前処理)
```

```
res <- preprocessReads(filename=in_f1, filenameMate=in_f2,
                       outputFilename=out_f1, outputFilenameMate=out_f2,
                       outputFilenameLpattern=out_f1, outputFilenameMateLpattern=out_f2,
                       Lpattern=param_adapter)
res
```



Thanks to Dr. Stadler

回避策

3. 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100万リード分です。paired-endデータのforward側、約75MB、全リード107 bpです。FastQC実行結果(SRR616268sub_1_fastqc.html)として「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は
 "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGCTTCTGCTTG"と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリムしないやり方です。この例題では、アダプター配列除去のみしか行わず出力ファイル中のリード数が入力ファイルと同じになるようにしています。これを確実に実現するために、許容するNの数としてトテツモナイ大きさの1000をparam_nBasesで与え、アダプター配列除去後の許容する最低配列長として全てが残る0をparam_minLengthで与えています。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。結局、出力ファイルでもう一度FastQCを実行すると、「TruSeq Adapter, Index 3」がOverrepresented sequencesの項目から消えていることまでは確認しました。しかし、まだ「TruSeq Adapter, Index 2」が残っています。ファイルサイズが、74,906,576 bytesから71,542,907 bytesに若干減っていることがわかります。

基本戦略は「single-endとして別々にアダプター配列除去を行うが、出力ファイルのリード数が変わらないようにする」です。乳酸菌paired-end RNA-seqデータ(SRR626268)のアダプター配列除去を行う一連の手順の基本形は、例題3-5です。例題3では、Forward側のリードファイルを入力として、FastQCでレポートされた「TruSeq Adapter, Index 3」の除去を行っています。約1分

```

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 <- "hoge3.fastq.gz" #出力ファイル名を指定
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGC" #アダプター配列
param_nBases <- 1000 #許容するNの数を指定
param_minLength <- 0 #アダプター配列除去後
param_nrec <- 500000 #一度に処理するリード数

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                       outputFilename=out_f, #out_fを出力とする
                       Lpattern=param_adapter, #リードの左側(Left)
                       nBases=param_nBases, #許容するNの数
                       minLength=param_minLength, #アダプター配列除去後
                       nrec=param_nrec) #一度に処理するリード数

res #確認してるだけです
file.size(in_f) #入力ファイルのサイズ
file.size(out_f) #出力ファイルのサイズ
    
```

```

R Console
> res # $
totalSequences SRR616268sub_1.fastq.gz 1000000
matchTo5pAdapter 456561
matchTo3pAdapter 0
tooShort 0
tooManyN 0
lowComplexity 0
totalPassed 1000000
> file.size(in_f) # $
[1] 74906576
> file.size(out_f) # $
[1] 71542907
> |
    
```

回避策

例題4では、例題3の出力ファイルを入力として、もう1つレポートされていた「TruSeq Adapter, Index 2」の除去を行っています。FastQC実行環境にない場合でも、用いた実験プロトコルが分かれば、候補となるadapterやprimer配列の一部で文字列検索すればある程度わかります。もちろん、この場合はIndex配列部分を含めたりしたほうが良いとは思いますが。約1分

4. gzip圧縮FASTQ形式ファイル(hoge3.fastq.gz)の場合:

3.で得られたhoge3.fastq.gzファイル中には、まだ「TruSeq Adapter, Index 2」が残っています。「TruSeq Adapter, Index 2」の塩基配列情報は "GATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTCGTATGCCGTCTTCTGCTTG"と書いてあったので、これを入力として与えます。約1分。ファイルサイズが、71,542,907 bytesから71,343,605 bytesに若干減っていることがわかります。この出力ファイル(hoge4.fastq.gz)と同じものが [SRR616268sub_trim2_1.fastq.gz](#) (1,000,000リード、71,343,605 bytes)です。FastQC実行結果は [SRR616268sub_trim2_1_fastqc.html](#)です。

```

in_f <- "hoge3.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_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTCGTATGCCGTCTTCTGCTTG" #アダ
param_nBases <- 1000             #許容するNの数を指定
param_minLength <- 0             #アダプター配列除去後の許容する最低配列長を指定
param_nrec <- 500000             #一度に処理するリード数を指定

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                        outputFilename=out_f, #out_fを出力とする
                        lpattern=param_adapter, #リードの左側(Left)
                        nBases=param_nBases, #許容するNの数
                        minLength=param_minLength, #アダプター配列除去後
                        nrec=param_nrec) #一度に処理するリード数
res #確認してるだけです
file.size(in_f) #入力ファイルのサイズ
file.size(out_f) #出力ファイルのサイズ

```

```

R Console
> res
totalSequences      hoge3.fastq.gz      #5
                    1000000
matchTo5pAdapter    170381
matchTo3pAdapter    0
tooShort             0
tooManyN             0
lowComplexity        0
totalPassed         1000000
> file.size(in_f)   #5
[1] 71542907
> file.size(out_f)  #5
[1] 71343605
> |

```

回避策

5. gzip圧縮FASTQ形式ファイル(SRR616268sub_2.fastq.gz)の場合:

乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。paired-endデータのreverse側、約67MB、93 bpです。FastQC実行結果(SRR616268sub_2_fastqc.html)として「Illumina Single End PCR Primer 1」があるとレポートされました。これでググると塩基配列情報は "AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT" と書いて、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプターしかトリミングしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。ファイルサイズが、67,158,462 bytesから63,524,791 bytesに若干減っていることがわかります。この出力(hoge5.fastq.gz)と同じものが SRR616268sub_trim2_2.fastq.gz(1,000,000リード、63,524,791 bytes)です。行結果は SRR616268sub_trim2_2_fastqc.html です。

例題5では、Reverse側のリードファイルを入力として、FastQCでレポートされた「Illumina Single End PCR Primer 1」の除去を行っています。Primer配列の場合は、逆相補鎖(reverse complement)を作成してから与えないとうまく除去できなかったことから、このようにしています。FastQCを実行できる環境にあれば、「overrepresented sequences」項目のpossible sourceのところNo Hitになったかどうかでうまくいったかどうかの判定ができます。約1分

```

in_f <- "SRR616268sub_2.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge5.fastq.gz" #出力ファイル名を指定してout_fに格納
param_adapter <- "AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT" #
param_nBases <- 1000 #許容するNの数を指定
param_minLength <- 0 #アダプター配列除去後
param_nrec <- 500000 #一度に処理するリード数

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

#前処理(reverse complementの作成)
hoge <- reverseComplement(DNAString(param_adapter)) #逆相補

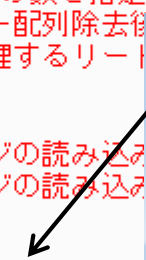
#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                        outputFilename=out_f, #out_fを出力とする
                        Lpattern=as.character(hoge), #リードの左側(Left)
                        nBases=param_nBases, #許容するNの数
                        minLength=param_minLength, #アダプター配列除去後
                        nrec=param_nrec) #一度に処理するリード数
res #確認してるだけです

```

```

R Console
> res # $
SRR616268sub_2.fastq.gz
totalSequences 1000000
matchTo5pAdapter 273793
matchTo3pAdapter 0
tooShort 0
tooManyN 0
lowComplexity 0
totalPassed 1000000
> file.size(in_f) # $
[1] 67158462
> file.size(out_f) # $
[1] 63524791
> |

```



Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

フィルタリング: 基本形

リード数の揃ったpaired-endデータを入力として、許容するNの数と最短配列長でフィルタリングするやり方を示します。

- 前処理 | フィルタリング | [Illuminaの pass filtering](#) (last modified 2013/06/19)
- 前処理 | フィルタリング | [GFF/GTF形式ファイル](#) (last modified 2013/10/10)
- 前処理 | フィルタリング | 組合せ | [ACGTのみ & 指定した長さの範囲の配列](#) (last modified 2014/06/11)
- 前処理 | フィルタリング | paired-end | 配列長とN数 | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/26) **NEW**
- 前処理 | フィルタリング | paired-end | 共通リード抽出 | [ShortRead\(Morgan 2007\)](#) (last modified 2015/06/26)

前処理 | フィルタリング | paired-end | 配列長とN数 | QuasR(Gaidatzis_2015)

NEW

- アセンブル | [について](#) (last modified 2015/06/26)
- アセンブル | [ゲノム用](#) (last modified 2015/06/26)
- アセンブル | [トランスクリプト](#) (last modified 2015/06/26)
- マッピング | [について](#) (last modified 2015/06/26)
- マッピング | [basic aligner](#) (last modified 2015/06/26)
- マッピング | [splice-aware](#) (last modified 2015/06/26)
- マッピング | [Bisulfite seq](#) (last modified 2015/06/26)
- マッピング | [\(ESTレベル\)](#) (last modified 2015/06/26)
- マッピング | [基礎](#) (last modified 2015/06/26)
- マッピング | [single-end](#) (last modified 2015/06/26)
- マッピング | [single-end](#) (last modified 2015/06/26)
- マッピング | [single-end](#) (last modified 2015/06/26)
- マッピング | [single-end](#) (last modified 2015/06/26)
- マッピング | [マップ後](#) (last modified 2015/06/26)

QuasRパッケージを用いたフィルタリング例を示します。「ファイル」→「ディレクトリ」の選択

1. gzip圧縮FASTQ形式ファイル

前処理 | トリミング | アダプター配列除去 | paired-endデータです。SRR1616268sub_trim2_1.fastq.gzとSRR1616268sub_trim2_2.fastq.gzのtooManyNが1,551リードあります。尚、998,431 + 107 + 1,551 (union)は(1,000,000 - 998,431)

```
in_f1 <- "SRR1616268sub_trim2_1.fastq.gz"
in_f2 <- "SRR1616268sub_trim2_2.fastq.gz"
out_f1 <- "hoge1_1.fastq.gz"
out_f2 <- "hoge1_2.fastq.gz"
param_nrec <- 500000
```

```
#必要なパッケージをロード
library(QuasR)
```

#本番(前処理)

```
res <- preprocessReads(filename=in_f1,
```

2. gzip圧縮FASTQ形式ファイル(SRR616268sub_trim2_1.fastq.gzとSRR616268sub_trim2_2.fastq.gz)の場合:

許容するN数(param_nBases)と最短配列長(param_minLength)を明示的に与えるやり方です。tooShortが119、tooManyNが1,762リードあったことがわかります。結果として、998,208リードからなるファイルが出力されています。例題6に比べて若干生き残るリード数が減るのは、param_nBases(デフォルトは2)と param_minLength(デフォルトは14)での条件がデフォルトよりも若干厳しめ(1と18)だからです。

```
in_f1 <- "SRR616268sub_trim2_1.fastq.gz" #入力ファイル名を指定してin_f1に格納
in_f2 <- "SRR616268sub_trim2_2.fastq.gz" #入力ファイル名を指定してin_f2に格納
out_f1 <- "hoge2_1.fastq.gz" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge2_2.fastq.gz" #出力ファイル名を指定してout_f2に格納
param_nBases <- 1 #許容するNの数を指定
param_minLength <- 18 #アダプター配列除去後の許容する最短配列長を指定
param_nrec <- 500000 #一度に処理するリード数を指定
```

#必要なパッケージをロード

```
library(QuasR) #パッケージの読み込み
```

#本番(前処理)

```
res <- preprocessReads(filename=in_f1, #前処理を実行
                       filenameMate=in_f2, #前処理を実行
                       outputFilename=out_f1, #前処理を実行
                       outputFilenameMate=out_f2, #前処理を実行
                       nBases=param_nBases, #許容するNの数
                       minLength=param_minLength, #アダプター配列除去後の許容する最短配列長
                       nrec=param_nrec) #一度に処理するリード数
#確認してるだけです
```

res

フィルタリング: 基本形

この例題で用いているファイルは hoge4.fastq.gz と hoge5.fastq.gz でも構いません。中身は同じです。

2. gzip圧縮FASTQ形式ファイル(SRR616268sub_trim2_1.fastq.gzとSRR616268sub_trim2_2.fastq.gz)の場合:

許容するN数(param_nBases)と最短配列長(param_minLength)を明示的に与えるやり方です。tooShortが119, tooManyNが1,762リードあったことがわかります。結果として、998,208リードからなるファイルが出力されています。例題6に比べて若干生き残るリード数が減るのは、param_nBases(デフォルトは2)と param_minLength(デフォルトは14)での条件がデフォルトよりも若干厳しめ(1と18)だからです。

```
in_f1 <- "SRR616268sub_trim2_1.fastq.gz" #入力ファイル名を指定してin_f1に格納
in_f2 <- "SRR616268sub_trim2_2.fastq.gz" #入力ファイル名を指定してin_f2に格納
out_f1 <- "hoge2_1.fastq.gz" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge2_2.fastq.gz" #出力ファイル名を指定してout_f2に格納
param_nBases <- 1 #許容するNの数を指定
param_minLength <- 18 #アダプター配列除去後の許容する最短配列長を指定
param_nrec <- 500000 #一度に処理するリード数を指定
```

#必要なパッケージをロード

```
library(QuasR) #パッケージの読み込み
```

#本番(前処理)

```
res <- preprocessReads(filename=in_f1, #前処理を実行
                       filenameMate=in_f2, #前処理を実行
                       outputFilename=out_f1, #前処理を実行
                       outputFilenameMate=out_f2, #前処理を実行
                       nBases=param_nBases, #許容するNの数
                       minLength=param_minLength, #アダプター配列除去後
                       nrec=param_nrec) #一度に処理するリード
res #確認してるだけです
```

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> file.info(list.files(pattern="fastq.gz"))
              size
hoge3.fastq.gz 71542907
hoge4.fastq.gz 71343605
hoge5.fastq.gz 63524791
SRR609266.fastq.gz 383840833
SRR616268sub_1.fastq.gz 74906576
SRR616268sub_2.fastq.gz 67158462
SRR616268sub_trim_1.fastq.gz 71227695
SRR616268sub_trim_2.fastq.gz 63442904
SRR616268sub_trim2_1.fastq.gz 71343605
SRR616268sub_trim2_2.fastq.gz 63524791
```


フィルタリング：基本形

リード長が18塩基未満を捨て、1塩基だけNを許容するという閾値でフィルタリングをした結果、998,208リードが残ったということ。約2分

2. gzip圧縮FASTQ形式ファイル(SRR616268sub_trim2_1.fastq.gzとSRR616268sub_trim2_2.fastq.gz)の

許容するN数(param_nBases)と最短配列長(param_minLength)を明示的に与えるやり方です。tooShortが119、tooManyNが1,762リードあったことがわかります。結果として、998,208リードからなるファイルが出力されています。例題6に比べて若干生き残るリード数が減るのは、param_nBases(デフォルトは2)とparam_minLength(デフォルトは14)での条件がデフォルトよりも若干厳しめ(1と18)だからです。

```
in_f1 <- "SRR616268sub_trim2_1.fastq.gz"#入力ファイル名を指定してin_f1に格納
in_f2 <- "SRR616268sub_trim2_2.fastq.gz"#入力ファイル名を指定してin_f2に格納
out_f1 <- "hoge2_1.fastq.gz" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge2_2.fastq.gz" #出力ファイル名を指定してout_f2に格納
param_nBases <- 1 #許容するNの数を指定
param_minLength <- 18 #アダプター配列除去後の許容する最短配列長を指定
param_nrec <- 500000 #一度に処理するリード数を指定
```

#必要なパッケージ

```
library(QuasR)
```

#本番(前処理)

```
res <- preprocess
```

```
fi
ou
ou
nB
mi
nr
```

```
res
```

```
R Console
+ minLength=param_minLength, #アダプター配列除去後の許容する最短$
+ nrec=param_nrec) #一度に処理するリード数
filtering .../SRR616268sub_trim2_1.fastq.gz and
SRR616268sub_trim2_2.fastq.gz
> res #確認してるだけです
SRR616268sub_trim2_1.fastq.gz:SRR616268sub_trim2_2.fastq.gz
totalSequences 1000000
matchTo5pAdapter NA
matchTo3pAdapter NA
tooShort 119
tooManyN 1762
lowComplexity 0
totalPassed 998208
> |
```

フィルタリング：発展形

- 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 2015/06/24) 推奨 **①**
- 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 2013/06/15)
- 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](#) **②** (last modified 2014/08/27)
- 前処理 | フィルタリング | [ACGTのみからなる配列を抽出](#) (last modified 2014/08/04)
- 前処理 | フィルタリング | [ACGT以外のcharacter "-"をNに変換](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [ACGT以外の文字数が閾値以下の配列を抽出](#) (last modified 2013/09/27)
- 前処理 | フィルタリング | [重複のない配列セットを作成](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [指定した長さ以上の配列を抽出](#) **③** (last modified 2014/02/07)
- 前処理 | フィルタリング | [任意のリード\(サブセット\)を抽出](#) (last modified 2014/08/21)
- 前処理 | フィルタリング | [指定した長さの範囲の配列を抽出](#) (last modified 2015/02/26)
- 前処理 | フィルタリング | [任意のIDを含む配列を抽出](#) (last modified 2013/06/18)
- 前処理 | フィルタリング | [Illuminaのpass filtering](#) (last modified 2013/06/19)
- 前処理 | フィルタリング | [GFF/GTF形式ファイル](#) (last modified 2013/10/10)
- 前処理 | フィルタリング | 組合せ | [ACGTのみ & 指定した長さの範囲の配列](#) (last modified 2014/06/11)
- 前処理 | フィルタリング | paired-end | 配列長とN数 | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/26) 1
- 前処理 | フィルタリング | paired-end | 共通リード抽出 | [ShortRead\(Morgan 2009\)](#) **④** (last modified 2015/06/2)

基本戦略は、paired-endをsingle-endとして独立に取り扱い、最後に共通リードを抽出。例えば①～③あたりの任意のフィルタリングをsingle-endとして独立に実行。最後に④でリード数の異なるpaired-endのファイルを入力として、共通リードを出力。ここでは、①の例題6-8を実行して得られたリード数の異なるpaired-endファイルを入力として、④を実行するやり方を示します。個人見解としては、④の入力ファイルの段階ではFASTQではなくFASTA形式でいいと思います。一連のフィルタリングでそここのクオリティが保証されているので、わざわざクオリティ情報を保持しておく必要はないだろうという思想

フィルタリング: 発展形

Small RNA-seqデータの場合はリードの右側にあるアダプター配列除去でしたが、乳酸菌RNA-seqデータの場合は左側にアダプターがついているのでRpatternではなくLpatternになっている点に注意。

- 前処理 | トリミング | ボリア配列除去 | 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 2015/06/24)
- 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(基礎) | QuasR(Gaidatzis_2015)

NEW

QuasRパッケージ

②

5. gzip圧縮FASTQ形式ファイル(SRR616268sub_1.fastq.gz)の場合:

一度に処理するときに、メモリ不足でフ長くなりますが、エラーさらに小さくして対処し村浩正 氏提供情報)。「ファイル」-「ディレクトリ」を選択し、「gzip圧縮FASTQ形式」を選択し、small RNA-seqデータ(Genomics, 2013)中のFASTQ or SRA | SRAから、おそらくアダプター配列を削除できます。ここでは(左側)ではなく3'側(右側)にアダプター配列があることがありますが、preprocessR

乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。paired-endデータのforward側、約75MB(74,906,576 bytes)、全リード107 bpです。FastQC実行結果(SRR616268sub_1_fastqc.html)として「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGGATCTCGTATGCCGTCTTCTGCTTG"と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリミングしないやり方です。それが、preprocessR関数実行時にLpatternのみ記載している理由です。約2分。出力ファイル(998,664リード、71,426,499 bytes)でもう一度FastQCを実行すると、「TruSeq Adapter, Index 3」がOverrepresented sequencesの項目から消えていることまでは確認しました。しかし、まだ「TruSeq Adapter, Index 2」が残っています。

```

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 <- "hoge5.fastq.gz" #出力ファイル名を指定してout_fに格納
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGGATCTCGTATGCCGTCTTCTGCTTG" #アダプター配列
param_nrec <- 500000 #一度に処理するリード数を指定

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

#本番(アダプター配列除去)
res <- preprocessR(filename=in_f, #in_fを入力とする
                    outputfilename=out_f, #out_fを出力とする
                    Lpattern=param_adapter, #リードの左側(Left)のアダプターを除去
                    nrec=param_nrec) #一度に処理するリード数

res #確認してるだけです
file.size(in_f) #入力ファイルのサイズを表示
file.size(out_f) #出力ファイルのサイズを表示
    
```

出力予定ファイルが存在する状態で実行するとエラーが出るので注意。

フィルタリング: 発展形

5. 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100万リード分です。paired-endデータのforward側、約75MB(74,906,576 bytes)、全リード107 bpです。FastQC実行結果([SRR616268sub_1_fastqc.html](#))として「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は

"GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG"と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリムしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。約2分。出力ファイル(998,664リード、71,426,499 bytes)でもう一度FastQCを実行すると、「TruSeq Adapter, Index 3」がOverrepresented sequencesの項目から消えていることまでは確認しました。しかし、まだ「TruSeq Adapter, Index 2」が残っています。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定してin_f
out_f <- "hoge5.fastq.gz" #出力ファイル名を指定してout_f
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGT"
param_nrec <- 500000 #一度に処理するリード数を指定

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                        outputFilename=out_f, #out_fを出力とする
                        Lpattern=param_adapter, #リードの左側(Left)のアダプ
                        nrec=param_nrec) #一度に処理するリード数
#確認してるだけです
res
file.size(in_f) #入力ファイルのサイズを表示
file.size(out_f) #出力ファイルのサイズを表示
```

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> list.files(pattern="fastq.gz")
[1] "hoge2_1.fastq.gz"
[2] "hoge2_2.fastq.gz"
[3] "hoge3.fastq.gz"
[4] "hoge4.fastq.gz"
[5] "hoge5.fastq.gz"
[6] "SRR609266.fastq.gz"
[7] "SRR616268sub_1.fastq.gz"
[8] "SRR616268sub_2.fastq.gz"
[9] "SRR616268sub_trim_1.fastq.gz"
[10] "SRR616268sub_trim_2.fastq.gz"
[11] "SRR616268sub_trim2_1.fastq.gz"
[12] "SRR616268sub_trim2_2.fastq.gz"
> |
```

フィルタリング: 発展形

①こんな感じになります。②ここで見えている出力ファイルのサイズは、既に存在していたhoge5.fastq.gzのファイルサイズであり、ここでのコピー実行結果で得られたものではない点に注意。

5. 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100万リード分です。paired-endデータのforward側、約75 bytes、全リード107 bpです。FastQC実行結果(SRR616268sub_1_fastqc.html)として「TruSeq Adapter

れているとレポートされました。これでググると塩基配列情報は "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGCTTCTCTGCTTG"と書いて

あったので、これを入力として与えます。ター配列しかトリムしないやり方です。2分。出力ファイル(998,664リード、71,42 Overrepresented sequencesの項目から消残っています。

```
in_f <- "SRR616268sub_1.fastq.gz"
out_f <- "hoge5.fastq.gz"
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGCTTCTCTGCTTG"
param_nrec <- 500000
```

#必要なパッケージをロード
library(QuasR)

#本番(アダプター配列除去)

```
res <- preprocessReads(filename=in_f,
                        outputFilename=out_f,
                        Lpattern=param_adapter,
                        nrec=param_nrec)
```

```
res
file.size(in_f)
file.size(out_f)
```

```
R Console
> #本番 (アダプター配列除去)
> res <- preprocessReads(filename=in_f, #in_fを入力とする
+                         outputFilename=out_f, #out_fを出力とする
+                         Lpattern=param_adapter, #リードの左側(Left)のアダ$
+                         nrec=param_nrec) #一度に処理するリード数
Error in preprocessReads(filename = in_f, outputFilename = out_f, $
existing output files: hoge5.fastq.gz
> res #確認してるだけです
SRR616268sub_trim2_1.fastq.gz:SRR616268sub_trim2_1
totalSequences $
matchTo5pAdapter $
matchTo3pAdapter $
tooShort $
tooManyN $
lowComplexity $
totalPassed $
> file.size(in_f) #入力ファイルのサイズを表$
[1] 74906576
> file.size(out_f) #出力ファイルのサイズを表$
[1] 63524791
> |
```



フィルタリング: 発展形

5. 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100万リード分です。paired-endデータのforward側、約75MB(74,906,576 bytes)、全リード107 bpです。FastQC実行結果([SRR616268sub_1_fastqc.html](#))として「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は

"GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG"と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリムしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。約2分。出力ファイル(998,664リード、71,426,499 bytes)でもう一度FastQCを実行すると、「TruSeq Adapter, Index 3」がOverrepresented sequencesの項目から消えていることまでは確認しました。しかし、まだ「TruSeq Adapter, Index 2」が残っています。

```
in_f <- "SRR616268sub_1.fastq.gz" #入力ファイル名を指定して
out_f <- "hoge5.fastq.gz" #出力ファイル名を指定して
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTC"
param_nrec <- 500000 #一度に処理するリード数を

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                       outputFilename=out_f, #out_fを出力とする
                       Lpattern=param_adapter, #リードの左側(Left)のアダプター配列を指定
                       nrec=param_nrec) #一度に処理するリード数

res #確認してるだけです
file.size(in_f) #入力ファイルのサイズを表
file.size(out_f) #出力ファイルのサイズを表
```

```
R Console
+ Lpattern=param_adapter, # $
+ nrec=param_nrec) # $
filtering SRR616268sub_1.fastq.gz
> res # $
SRR616268sub_1.fastq.gz
totalSequences 1000000
matchTo5pAdapter 456561
matchTo3pAdapter 0
tooShort 89
tooManyN 1247
lowComplexity 0
totalPassed 998664
> file.size(in_f) # $
[1] 74906576
> file.size(out_f) # $
[1] 71426499
> |
```

フィルタリング：発展形

7. gzip圧縮FASTQ形式ファイル(hoge5.fastq.gz)の場合:

5.で得られたhoge5.fastq.gzファイル中には、まだ「TruSeq Adapter, Index 2」が残っています。「TruSeq Adapter, Index 2」の塩基配列情報は
 "GATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTCGTATGCCGTCTTCTGCTTG"と書いてあったので、これを入力として与えます。約2分。この出力ファイル(hoge7.fastq.gz)と同じものが
[SRR616268sub_trim_1.fastq.gz](#)(998,658リード、71,227,695 bytes)です。[FastQC](#)実行結果は
[SRR616268sub_trim_1_fastqc.html](#)です。

```
in_f <- "hoge5.fastq.gz"           #入力ファイル名を指定してin_fに格納
out_f <- "hoge7.fastq.gz"         #出力ファイル名を指定してout_fに格納
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTC"
param_nrec <- 500000              #一度に処理するリード数を指定

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

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                        outputFilename=out_f, #out_fを出力とする
                        Lpattern=param_adapter, #リードの左側(Left)のアダプター配列を指定
                        nrec=param_nrec) #一度に処理するリード数を指定

res                                #確認してるだけです
file.size(in_f)                   #入力ファイルのサイズを表す
file.size(out_f)                  #出力ファイルのサイズを表す
```

```
R Console
+      Lpattern=param_adapter,    # $
+      nrec=param_nrec)          # $
  filtering hoge5.fastq.gz
> res                             # $
                                     hoge5.fastq.gz
totalSequences                       998664
matchTo5pAdapter                     170144
matchTo3pAdapter                      0
tooShort                              6
tooManyN                              0
lowComplexity                         0
totalPassed                           998658
> file.size(in_f)                  # $
[1] 71426499
> file.size(out_f)                 # $
[1] 71227695
> |
```

例題8のpaired-end reverse側も実行。出力ファイルのサイズは63,442,904 bytes、999,136リード

フィルタリング：発展形

8. gzip圧縮FASTQ形式ファイル(SRR616268sub_2.fastq.gz)の場合:

乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。paired-endデータのreverse側、約67MB、全リード93 bpです。FastQC実行結果([SRR616268sub_2_fastqc.html](#))として「Illumina Single End PCR Primer 1」が含まれているとレポートされました。これをググると塩基配列情報は "AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT" と書いてあったので、これを入力として与えます。アダプター配列の位置は5'側(左側)にあるという前提であり、左側のアダプター配列しかトリムしないやり方です。それが、preprocessReads関数実行時にLpatternのみ記載している理由です。約2分。この出力ファイル([hoge8.fastq.gz](#))と同じものが [SRR616268sub_trim_2.fastq.gz](#) (999,136リード、63,442,904 bytes)です。FastQC実行結果は [SRR616268sub_trim_2_fastqc.html](#) です。

```
in_f <- "SRR616268sub_2.fastq.gz" #入力ファイル名を指定して
out_f <- "hoge8.fastq.gz" #出力ファイル名を指定して
param_adapter <- "AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCT" #アダプター配列
param_nrec <- 500000 #一度に処理するリード数を指定する

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

#前処理(reverse complementの作成)
hoge <- reverseComplement(DNAString(param_adapter)) #逆相補鎖(nrec)を作成

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #in_fを入力とする
                        outputFilename=out_f, #out_fを出力とする
                        Lpattern=as.character(hoge), #リードの左側(Left)の配列を指定する
                        nrec=param_nrec) #一度に処理するリード数

res #確認してるだけです
file.size(in_f) #入力ファイルのサイズを表す
file.size(out_f) #出力ファイルのサイズを表す
```

```
R Console
+ Lpattern=as.character(hoge) $
+ nrec=param_nrec) # $
filtering SRR616268sub_2.fastq.gz # $
> res # $
totalSequences 1000000
matchTo5pAdapter 273793
matchTo3pAdapter 0
tooShort 17
tooManyN 847
lowComplexity 0
totalPassed 999136
> file.size(in_f) # $
[1] 67158462
> file.size(out_f) # $
[1] 63442904
> |
```


乳酸菌RNA-seqデータ

①リード数の異なるpaired-endのファイルを入力として、共通リードを出力するやり方を示します。

■ 100万リードのオリジナル?!ファイル

- SRR616268sub_1.fastq.gz: 107 bp, 74,906,576 bytes
- SRR616268sub_2.fastq.gz: 93 bp, 67,158,462 bytes



■ 100万リードのアダプター除去後のファイル

- SRR616268sub_trim2_1.fastq.gz(hoge4.fastq.gz): 71,343,605 bytes
- SRR616268sub_trim2_2.fastq.gz : 63,524,791 bytes
- hoge2_1.fastq.gz: 998,208リード、71,194,586 bytes
- hoge2_2.fastq.gz: 998,208リード、63,375,473 bytes

■ アダプター除去およびフィルタリング後のファイル

- SRR616268sub_trim_1.fastq.gz(hoge7.fastq.gz): 998,658リード、71,227,695 bytes
- SRR616268sub_trim_2.fastq.gz(hoge8.fastq.gz): 999,136リード、63,442,904 bytes



file.info関数でもファイル名、サイズなど様々な情報をみることができます。

Tips : file.info関数

100万リードのオリジナル?!ファイル

- SRR616268sub_1.fastq.gz : 107 bp, 74,906,576 bytes
- SRR616268sub_2.fastq.gz : 93 bp, 67,158,462 bytes



100万リードのアダプター除去後のファイル

- SRR616268sub_trim2_1.fastq.gz
- SRR616268sub_trim2_2.fastq.gz
- hoge2_1.fastq.gz : 998,208リード
- hoge2_2.fastq.gz : 998,208リード

アダプター除去およびフィルタ

- SRR616268sub_trim_1.fastq.gz
- SRR616268sub_trim_2.fastq.gz

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> file.info(list.files(pattern="fastq.gz"))

      size isdir mode
hoge2_1.fastq.gz 71194586 FALSE 666
hoge2_2.fastq.gz 63375473 FALSE 666
hoge3.fastq.gz   71542907 FALSE 666
hoge4.fastq.gz   71343605 FALSE 666
hoge5.fastq.gz   71426499 FALSE 666
hoge7.fastq.gz   71227695 FALSE 666
hoge8.fastq.gz   63442904 FALSE 666
SRR609266.fastq.gz 383840833 FALSE 666
SRR616268sub_1.fastq.gz 74906576 FALSE 666
SRR616268sub_2.fastq.gz 67158462 FALSE 666
SRR616268sub_trim_1.fastq.gz 71227695 FALSE 666
SRR616268sub_trim_2.fastq.gz 63442904 FALSE 666
SRR616268sub_trim2_1.fastq.gz 71343605 FALSE 666
SRR616268sub_trim2_2.fastq.gz 63524791 FALSE 666
```

フィルタリング: 発展形

入力はリード数の異なるpaired-endファイル。出力はリード数が同じpaired-endファイルです。description行の記述形式次第ではエラーが出ると思いますが、適宜変更してください。



- 前処理 | フィルタリング | [GFF/GTF形式ファイル](#) (last modified 2013/10/10)
- 前処理 | フィルタリング | 組合せ | [ACGTのみ & 指定した長さの範囲の配列](#) (last modified 2014/06/20)
- 前処理 | フィルタリング | paired-end | 配列長とN数 | [QuasR\(Gaidatzis_2015\)](#) (last modified 2014/06/20)
- 前処理 | フィルタリング | paired-end | 共通リード抽出 | [ShortRead\(Morgan_2009\)](#) (last modified 2014/06/20)
- アセンブル | [について](#) (last modified 2014/06/20)
- アセンブル | [ゲノム用](#) (last modified 2015/05/15)
- アセンブル | [トランスクリプトーム\(転写\)](#) (last modified 2015/05/15)
- マッピング | [について](#) (last modified 2015/05/15)
- マッピング | [basic aligner](#) (last modified 2015/05/15)
- マッピング | [splice-aware aligner](#) (last modified 2015/05/15)
- マッピング | [Bisulfite sequencing用](#) (last modified 2015/05/15)
- マッピング | [\(ESTレベルの長さのcc\)](#) (last modified 2015/05/15)
- マッピング | [基礎](#) (last modified 2015/05/15)
- マッピング | single-end | [ゲノム | basic](#) (last modified 2015/05/15)
- マッピング | single-end | [ゲノム | basic](#) (last modified 2015/05/15)
- マッピング | single-end | [ゲノム | splice-aware](#) (last modified 2015/05/15)
- マップ後 | [について](#) (last modified 2015/05/15)
- マップ後 | [出力ファイル形式について](#) (last modified 2015/05/15)
- マップ後 | [出力ファイルの読み込み](#) (last modified 2015/05/15)

前処理 | フィルタリング | paired-end | 共通リード抽出 | ShortRead(Morgan_2009)

NEW

ShortReadパッケージを用いて、ファイル内での配列長とリード数が異なるpaired-endデータの2つのファイルを入力として、両方に共通して存在するリードのみ抽出するやり方を示します。QuasRは入力ファイルのリード数が揃ってないと受けつけないので、ShortReadパッケージで別々に読み込む戦略を採用しています。

1. gzip圧縮FASTQ形式ファイル(SRR616268sub_trim_1.fastq.gzとSRR616268sub_trim_2.fastq.gz)の場合:

前処理 | [トリミング](#) | [アダプター配列除去\(基礎\)](#) | [QuasR\(Gaidatzis_2015\)](#) の例題7と8の実行結果ファイルで、100万リード弱のpaired-endデータです。SRR616268sub_trim_1.fastq.gzは、998,658リード、71,227,695 bytes (約71MB)です。SRR616268sub_trim_2.fastq.gzは、999,136リード、63,442,904 bytes (約63MB)です。前処理のところで長ったらしい作業をしているのは、FASTQファイル中のdescription行でリードのIDの対応付けがそのままではできないので、対応付け可能な部分文字列を抽出しているためです。イントロ | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [description行の記述を整形](#)のコードをテンプレートにしています。共通リードは998,428個です。リード数がほとんど減っていないにもかかわらず、ファイルサイズが(71,227,695 bytes -> 63,446,240 bytes; 63,442,904 bytes -> 56,019,410 bytes)と大幅に減っているのは、description部分の記述内容が減っていることに起因します。

```
in_f1 <- "SRR616268sub_trim_1.fastq.gz" #入力ファイル名を指定してin_f1に格納
in_f2 <- "SRR616268sub_trim_2.fastq.gz" #入力ファイル名を指定してin_f2に格納
out_f1 <- "hoge1_1.fastq.gz" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge1_2.fastq.gz" #出力ファイル名を指定してout_f2に格納

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

#入力ファイルの読み込み
fastq1 <- readFastq(in_f1) #in_f1で指定したファイルの読み込み
fastq2 <- readFastq(in_f2) #in_f2で指定したファイルの読み込み
head(id(fastq1), n=4) #description部分を表示
head(id(fastq2), n=4) #description部分を表示
```

フィルタリング：発展形

最初の4リード分のdescription部分を表示。ここで見えているもので判断すると、SRR616268.7は2つの入力ファイル中に存在する。そして、少なくともforward側のファイル(*_1.fastq.gz)中の3リード(SRR616268.20-22)は、reverse側のファイル(*_2.fastq.gz)中には存在しないことがわかる。

1. gzip圧縮FASTQ形式ファイル(SRR616268sub_trim_1.fastq.gzとSRR616268sub_trim_2.fastq.gz)

前処理 | トリミング | アダプター配列除去(基礎) | QuasR(Gaidatzis 2015) の例題7と8の実行100万リード弱のpaired-endデータです。SRR616268sub_trim_1.fastq.gzは、998,658リード、(約71MB)です。SRR616268sub_trim_2.fastq.gzは、999,136リード、63,442,904 bytes(約63MB)のところで長ったらしい作業をしているのは、FASTQファイル中のdescription行でリードのIDのままではできないので、対応付け可能な部分文字列を抽出しているためです。イントロFASTQ形式 | description行の記述を整形のコードをテンプレートにしています。共通リードはリード数がほとんど減っていないにもかかわらず、ファイルサイズが(71,227,695 bytes -> 63,446,240 bytes; 63,442,904 bytes -> 56,019,410 bytes)と大幅に減っているのは、description行の抽出に起因します。

```
in_f1 <- "SRR616268sub_trim_1.fastq.gz" #入力ファイル
in_f2 <- "SRR616268sub_trim_2.fastq.gz" #入力ファイル
out_f1 <- "hoge1_1.fastq.gz" #出力ファイル
out_f2 <- "hoge1_2.fastq.gz" #出力ファイル

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

#入力ファイルの読み込み
fastq1 <- readFastq(in_f1) #in_f1で指$
fastq2 <- readFastq(in_f2) #in_f2で指$
head(id(fastq1), n=4) #descripti$
head(id(fastq2), n=4) #descripti$

#前処理(FASTQのdescription行を整形; in_f1)
fastq <- fastq1 #fastqとし
hoge <- strsplit(as.character(id(fastq)), " ", f)
description <- BStringSet(sapply(hoge, "[[", 1)) #
fastq <- ShortReadQ(sread(fastq), quality(fastq))
head(id(fastq), n=4) #descripti$
fastq1 <- fastq #元に戻す
```

```
R Console
> #入力ファイルの読み込み
> fastq1 <- readFastq(in_f1) #in_f1で指$
> fastq2 <- readFastq(in_f2) #in_f2で指$
> head(id(fastq1), n=4) #descripti$
A BStringSet instance of length 4
width seq
[1] 44 SRR616268.7 2291:6:...412:2249 length=107
[2] 45 SRR616268.20 2291:6...720:2221 length=107
[3] 45 SRR616268.21 2291:6...652:2235 length=107
[4] 45 SRR616268.22 2291:6...505:2236 length=107
> head(id(fastq2), n=4) #descripti$
A BStringSet instance of length 4
width seq
[1] 43 SRR616268.7 2291:6:...1412:2249 length=93
[2] 44 SRR616268.48 2291:6...2243:2242 length=93
[3] 44 SRR616268.61 2291:6...2360:2248 length=93
[4] 44 SRR616268.73 2291:6...2584:2248 length=93
> |
```

フィルタリング：発展形

共通リード抽出用の基本情報として、リードIDの積集合(intersection)を得るのが基本戦略。しかし、description部分の情報をそのまま使うと、length=107とlength=93という部分文字列の違いにより、①同一リードを同一文字列として認識できないという問題がある

1. gzip圧縮FASTQ形式ファイル(SRR616268sub_trim_1.fastq.gzとSRR616268sub_trim_2.fastq.gz)

前処理 | トリミング | アダプター配列除去(基礎) | QuasR(Gaidatzis 2015) の例題7と8の実行
 100万リード弱のpaired-endデータです。SRR616268sub_trim_1.fastq.gzは、998,658リード、71,111,111 bytes (約71MB)です。SRR616268sub_trim_2.fastq.gzは、999,136リード、63,442,904 bytes (約63MB)です。
 のところで長ったらしい作業をしているのは、FASTQファイル中のdescription行でリードのIDのままではできないので、対応付け可能な部分文字列を抽出しているためです。イントロ | NGS | 読み込み | FASTQ形式 | description行の記述を整形のコードをテンプレートにしています。共通リードは998,428個です。
 リード数がほとんど減っていないにもかかわらず、ファイルサイズが(71,111,111 bytes -> 63,442,904 bytes)と大幅に減っているのは、description行のID部分に起因します。

```
in_f1 <- "SRR616268sub_trim_1.fastq.gz" #入力ファイル
in_f2 <- "SRR616268sub_trim_2.fastq.gz" #入力ファイル
out_f1 <- "hoge1_1.fastq.gz" #出力ファイル
out_f2 <- "hoge1_2.fastq.gz" #出力ファイル

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

#入力ファイルの読み込み
fastq1 <- readFastq(in_f1) #in_f1で指定
fastq2 <- readFastq(in_f2) #in_f2で指定
head(id(fastq1), n=4) #description行
head(id(fastq2), n=4) #description行

#前処理(FASTQのdescription行を整形; in_f1)
fastq <- fastq1 #fastqとして
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #分割
description <- BStringSet(sapply(hoge, "[", 1)) #description行
fastq <- ShortReadQ(sread(fastq), quality(fastq)) #ShortReadQ
head(id(fastq), n=4) #description行
fastq1 <- fastq #元に戻す
```

```
R Console
> id(fastq1)[1:4]
A BStringSet instance of length 4
width seq
[1] 44 SRR616268.7 2291:6:1...412:2249 length=107
[2] 45 SRR616268.20 2291:6:...720:2221 length=107
[3] 45 SRR616268.21 2291:6:...652:2235 length=107
[4] 45 SRR616268.22 2291:6:...505:2236 length=107
> id(fastq2)[1:4]
A BStringSet instance of length 4
width seq
[1] 43 SRR616268.7 2291:6:1...1412:2249 length=93
[2] 44 SRR616268.48 2291:6:...2243:2242 length=93
[3] 44 SRR616268.61 2291:6:...2360:2248 length=93
[4] 44 SRR616268.73 2291:6:...2584:2248 length=93
> intersect(c("AA", "B", "C"), c("A", "C"))
[1] "C"
> intersect(id(fastq1), id(fastq2))
A BStringSet instance of length 0
> |
```



フィルタリング：発展形

description部分を眺め、部分文字列で同一リードを認識できる箇所を探す。ここでは、① description部分を区切り文字「スペース(“ ”)で」分断し、②1番目の要素を取り出せば、赤下線部分のみの部分文字列で評価できるという戦略のもとでプログラムを組んでいる。

```
#前処理(FASTQのdescription行を整形; in_f1)
fastq <- fastq1 #fastq1として取り扱う
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE)#id(fastq)をスペースで区切る
description <- BStringSet(sapply(hoge,"[[", 1))#hogeのリスト中の1番目の要素を取り出す
fastq <- ShortReadQ(sread(fastq), quality(fastq), description)#ReadFastqQオブジェクトを作成
head(id(fastq), n=4) #description部分を表示
fastq1 <- fastq #元に戻す
```

```
#前処理(FASTQのdescription行を整形; in_f2)
fastq <- fastq2 #fastq2として取り扱う
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE)#id(fastq)をスペースで区切る
description <- BStringSet(sapply(hoge,"[[", 1))#hogeのリスト中の1番目の要素を取り出す
fastq <- ShortReadQ(sread(fastq), quality(fastq), description)#ReadFastqQオブジェクトを作成
head(id(fastq), n=4) #description部分を表示
fastq2 <- fastq #元に戻す
```

```
#本番
common <- intersect(id(fastq1), id(fastq2))#共通リードを抽出
length(fastq1) #fastq1のリード数
length(fastq2) #fastq2のリード数
length(common) #共通リード数
obj1 <- is.element(as.character(id(fastq1)), as.character(id(fastq2)))
```

```
R Console
> #入力ファイルの読み込み
> fastq1 <- readFastq(in_f1) #in_f1で指$
> fastq2 <- readFastq(in_f2) #in_f2で指$
> head(id(fastq1), n=4) #descripti$
A BStringSet instance of length 4
width seq
[1] 44 SRR616268.7 2291:6:...412:2249 length=107
[2] 45 SRR616268.20 2291:6...720:2221 length=107
[3] 45 SRR616268.21 2291:6...652:2235 length=107
[4] 45 SRR616268.22 2291:6...505:2236 length=107
> head(id(fastq2), n=4) #descripti$
A BStringSet instance of length 4
width seq
[1] 43 SRR616268.7 2291:6:...1412:2249 length=93
[2] 44 SRR616268.48 2291:6...2243:2242 length=93
[3] 44 SRR616268.61 2291:6...2360:2248 length=93
[4] 44 SRR616268.73 2291:6...2584:2248 length=93
> |
```

Forward側。うまくリードID (SRR6161268.7)のみにできていることがわかります

フィルタリング：発展形

#前処理 (FASTQのdescription行を整形; in_f1)

```
fastq <- fastq1 #fastqとして取り扱う
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切って
description <- BStringSet(sapply(hoge,"[[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFastQ
head(id(fastq), n=4) #description部分を表示
fastq1 <- fastq #元に戻す
```

#前処理 (FASTQのdescription行を整形; in_f2)

```
fastq <- fastq2 #fastqとして取り扱う
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切って
description <- BStringSet(sapply(hoge,"[[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFastQ
head(id(fastq), n=4) #description部分を表示
fastq2 <- fastq #元に戻す
```

#本番

```
common <- intersect(id(fastq1), id(fastq2)) #共通リードID
length(fastq1) #fastq1のリード数
length(fastq2) #fastq2のリード数
length(common) #共通リード数
obj1 <- is.element(as.character(id(fastq1)), as.character(id(fastq2)))
```

```
R Console
> head(id(fastq1), n=4)
A BStringSet instance of length 4
width seq
[1] 44 SRR616268.7 2291:6:1...412:2249 length=107
[2] 45 SRR616268.20 2291:6:...720:2221 length=107
[3] 45 SRR616268.21 2291:6:...652:2235 length=107
[4] 45 SRR616268.22 2291:6:...505:2236 length=107

> fastq <- fastq1 #fastqとして取り扱う
> hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切って
> description <- BStringSet(sapply(hoge,"[[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
> fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFastQ
> head(id(fastq), n=4) #description部分を表示
A BStringSet instance of length 4
width seq
[1] 11 SRR616268.7
[2] 12 SRR616268.20
[3] 12 SRR616268.21
[4] 12 SRR616268.22

> fastq1 <- fastq #元に戻す
> |
```

Reverse側。うまくリードID (SRR6161268.7)のみにできていることがわかります

フィルタリング：発展形

```
#前処理(FASTQのdescription行を整形; in_f1)
fastq <- fastq1 #fastqとして取り扱う
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切っ
description <- BStringSet(sapply(hoge,"[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFas
head(id(fastq), n=4) #description部分を表示
fastq1 <- fastq #元に戻す
```

```
#前処理(FASTQのdescription行を整形; in_f2)
fastq <- fastq2 #fastqとして取り扱う
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切っ
description <- BStringSet(sapply(hoge,"[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFas
head(id(fastq), n=4) #description部分を表示
fastq2 <- fastq #元に戻す
```

```
#本番
common <- intersect(id(fastq1), id(fastq2)) #共通リードID
length(fastq1) #fastq1のリード数
length(fastq2) #fastq2のリード数
length(common) #共通リード数
obj1 <- is.element(as.character(id(fastq1)), as.character(id(fastq2))) #共通リードID
```

```
R Console
> head(id(fastq2), n=4)
A BStringSet instance of length 4
width seq
[1] 43 SRR616268.7 2291:6:1...1412:2249 length=93
[2] 44 SRR616268.48 2291:6:...2243:2242 length=93
[3] 44 SRR616268.61 2291:6:...2360:2248 length=93
[4] 44 SRR616268.73 2291:6:...2584:2248 length=93

> fastq <- fastq2 #fastqとして取り扱う
> hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切っ
> description <- BStringSet(sapply(hoge,"[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
> fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFastQ
> head(id(fastq), n=4) #description部分を表示
A BStringSet instance of length 4
width seq
[1] 11 SRR616268.7
[2] 12 SRR616268.48
[3] 12 SRR616268.61
[4] 12 SRR616268.73

> fastq2 <- fastq #元に戻す
> |
```


Tips: コピペで使いまわす

わざわざ別のオブジェクト名に置き換えているので、一見ややこしいですが、黒枠内のコードを使いまわせるので、間違いを減らせるというメリットがあります。

```
#前処理(FASTQのdescription行を整形; in_f1)
```

```
fastq <- fastq1 #fastqとして扱う
```

```
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切
```

```
description <- BStringSet(sapply(hoge, "[[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
```

```
fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFas
```

```
head(id(fastq), n=4) #description部分を表示
```

```
fastq1 <- fastq #元に戻す
```

```
#前処理(FASTQのdescription行を整形; in_f2)
```

```
fastq <- fastq2 #fastqとして扱う
```

```
hoge <- strsplit(as.character(id(fastq)), " ", fixed=TRUE) #id(fastq)中の文字列を" "で区切
```

```
description <- BStringSet(sapply(hoge, "[[", 1)) #hogeのリスト中の1番目の要素のみ取り出してdesc
```

```
fastq <- ShortReadQ(sread(fastq), quality(fastq), description) #ReadFastQ関数を用いてReadFas
```

```
head(id(fastq), n=4) #description部分を表示
```

```
fastq2 <- fastq #元に戻す
```

```
#本番
```

```
common <- intersect(id(fastq1), id(fastq2)) #共通して現れるリードID情報を抽出
```

```
length(fastq1) #fastq1のリード数を表示
```

```
length(fastq2) #fastq2のリード数を表示
```

```
length(common) #共通リード数を表示
```

```
obj1 <- is.element(as.character(id(fastq1)), as.character(common)) #共通リードの位置をTRUE、
```

```
<
```

```
>
```

フィルタリング：発展形

得られた共通リード数(998,428)は、forward側(998,658)およびreverse側(999,136)より少ないので妥当。

```
#本番
common <- intersect(id(fastq1), id(fastq2))#共通して現れるリードID情報を抽出
length(fastq1) #fastq1のリード数を表示
length(fastq2) #fastq2のリード数を表示
length(common) #共通リード数を表示
obj1 <- is.element(as.character(id(fastq1)), as.character(common))#共通リードの位置をTRUE、
obj2 <- is.element(as.character(id(fastq2)), as.character(common))#共通リードの位置をTRUE、

#ファイルに保存
writeFastq(fastq1[obj1], out_f1, compress=T)#fastqの中身を指定したファイル名で保存
writeFastq(fastq2[obj2], out_f2, compress=T)#fastqの中身を指定したファイル名で保存
file.size(in_f1) #入力ファイルのサイズを表示
file.size(in_f2) #入力ファイルのサイズを表示
file.size(out_f1) #出力ファイルのサイズを表示
file.size(out_f2) #出力ファイルのサイズを表示
```

```
R Console
> common <- intersect(id(fastq1), id(fastq2))#共通$
> length(fastq1) #fastq1の$
[1] 998658
> length(fastq2) #fastq2の$
[1] 999136
> length(common) #共通リー-$
[1] 998428
> obj1 <- is.element(as.character(id(fastq1)), as.c$
> obj2 <- is.element(as.character(id(fastq2)), as.c$
>
```

フィルタリング: 発展形

```
#本番
common <- intersect(id(fastq1), id(fastq2))#共通して現れるリードID情報を抽出
length(fastq1) #fastq1のリード数を表示
length(fastq2) #fastq2のリード数を表示
length(common) #共通リード数を表示
obj1 <- is.element(as.character(id(fastq1)), as.character(common))#共通リードの位置をTRUE、
obj2 <- is.element(as.character(id(fastq2)), as.character(common))#共通リードの位置をTRUE、

#ファイルに保存
writeFastq(fastq1[obj1], out_f1, compress=T)#fastq
writeFastq(fastq2[obj2], out_f2, compress=T)#fastq
file.size(in_f1) #入力ファイ
file.size(in_f2) #入力ファイ
file.size(out_f1) #出力ファイ
file.size(out_f2) #出力ファイ
```

```
R Console
> as.character(id(fastq1)) [1:24]
 [1] "SRR616268.7" "SRR616268.20" "SRR616268.21"
 [4] "SRR616268.22" "SRR616268.30" "SRR616268.31"
 [7] "SRR616268.32" "SRR616268.33" "SRR616268.44"
[10] "SRR616268.45" "SRR616268.46" "SRR616268.47"
[13] "SRR616268.48" "SRR616268.58" "SRR616268.59"
[16] "SRR616268.60" "SRR616268.61" "SRR616268.69"
[19] "SRR616268.70" "SRR616268.71" "SRR616268.72"
[22] "SRR616268.73" "SRR616268.86" "SRR616268.87"
> as.character(id(fastq2)) [1:9]
 [1] "SRR616268.7" "SRR616268.48" "SRR616268.61"
 [4] "SRR616268.73" "SRR616268.101" "SRR616268.117"
 [7] "SRR616268.140" "SRR616268.141" "SRR616268.154"
> as.character(common) [1:9]
 [1] "SRR616268.7" "SRR616268.48" "SRR616268.61"
 [4] "SRR616268.73" "SRR616268.101" "SRR616268.117"
 [7] "SRR616268.140" "SRR616268.141" "SRR616268.154"
> |
```

フィルタリング：発展形

is.element関数は、右側の集合の中の要素が存在する位置をTRUE、それ以外をFALSEとして返す。共通であるfastq1オブジェクトの1, 13, 17, 22番目の要素がTRUEになっていることがわかる。

```
#本番
common <- intersect(id(fastq1), id(fastq2))#共通して現れるリードID情報を抽出
length(fastq1) #fastq1のリード数を表示
length(fastq2) #fastq2のリード数を表示
length(common) #共通リード数を表示
```

```
obj1 <- is.element(as.character(id(fastq1)), as.character(common))#共通リードの位置をTRUE、
obj2 <- is.element(as.character(id(fastq2)), as.character(common))#共通リードの位置をTRUE、
```

```
#ファイルに保存
writeFastq(fastq1[obj1], out_f1, compress=T)#fastq1を出力ファイルに保存
writeFastq(fastq2[obj2], out_f2, compress=T)#fastq2を出力ファイルに保存
file.size(in_f1) #入力ファイルのサイズ
file.size(in_f2) #入力ファイルのサイズ
file.size(out_f1) #出力ファイルのサイズ
file.size(out_f2) #出力ファイルのサイズ
```

```
R Console
> is.element(c(3, 8, 2, 4), c(4, 8))
[1] FALSE TRUE FALSE TRUE
> is.element(c(3, 8, 2, 4), c(4, 8, 5))
[1] FALSE TRUE FALSE TRUE
> obj1[1:27]
[1] TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
[9] FALSE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
[17] TRUE FALSE FALSE FALSE FALSE TRUE FALSE FALSE
[25] FALSE FALSE FALSE
> obj2[1:27]
[1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
[10] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
[19] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
> sum(obj1)
[1] 998428
> sum(obj2)
[1] 998428
> |
```

共通リードのみなので、リード数が揃っていることがわかる。

フィルタリング：発展形

```
#本番
common <- intersect(id(fastq1), id(fastq2))#共通して現れるリードID情報を抽出
length(fastq1) #fastq1のリード数を表示
length(fastq2) #fastq2のリード数を表示
length(common) #共通リード数
obj1 <- is.element(as.character(id(fastq1)), as.character(id(fastq2)))
obj2 <- is.element(as.character(id(fastq2)), as.character(id(fastq1)))

#ファイルに保存
writeFastq(fastq1[obj1], out_f1, compress=T)#fastq1を出力
writeFastq(fastq2[obj2], out_f2, compress=T)#fastq2を出力
file.size(in_f1) #入力ファイルサイズ
file.size(in_f2) #入力ファイルサイズ
file.size(out_f1) #出力ファイルサイズ
file.size(out_f2) #出力ファイルサイズ
```

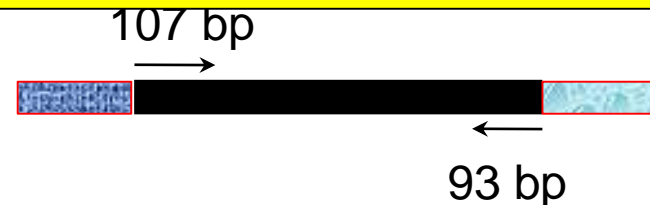
```
R Console
> fastq1
class: ShortReadQ
length: 998658 reads; width: 15..107 cycles
> fastq2
class: ShortReadQ
length: 999136 reads; width: 14..93 cycles
> fastq1[obj1]
class: ShortReadQ
length: 998428 reads; width: 15..107 cycles
> fastq2[obj2]
class: ShortReadQ
length: 998428 reads; width: 14..93 cycles
> writeFastq(fastq1[obj1], out_f1, compress=T)#fastq1を出力
> writeFastq(fastq2[obj2], out_f2, compress=T)#fastq2を出力
> file.size(in_f1) #入力ファイルサイズ
[1] 71227695
> file.size(in_f2) #入力ファイルサイズ
[1] 63442904
> file.size(out_f1) #出力ファイルサイズ
[1] 63446240
> file.size(out_f2) #出力ファイルサイズ
[1] 56019410
> |
```

乳酸菌RNA-seqデータ

①はデフォルトの「許容するN数が2、最短配列長が14」のときの結果。②は若干厳しめの「許容するN数が1、最短配列長が18」のときの結果なので、②のほうがリード数が若干少ない。

100万リードのオリジナル?!ファイル

- SRR616268sub_1.fastq.gz: 107 bp, 74,906,576 bytes
- SRR616268sub_2.fastq.gz: 93 bp, 67,158,462 bytes



100万リードのアダプター除去後のファイル

- SRR616268sub_trim2_1.fastq.gz(hoge4.fastq.gz): 71,343,605 bytes
- SRR616268sub_trim2_2.fastq.gz : 63,524,791 bytes

- hoge2_1.fastq.gz: 998,208リード、71,194,586 bytes
- hoge2_2.fastq.gz: 998,208リード、63,375,473 bytes



アダプター除去およびフィルタリング後のファイル

- SRR616268sub_trim_1.fastq.gz(hoge7.fastq.gz): 998,658リード、71,227,695 bytes
- SRR616268sub_trim_2.fastq.gz(hoge8.fastq.gz): 999,136リード、63,442,904 bytes

- hoge1_1.fastq.gz: 998,428リード、63,446,240 bytes
- hoge1_2.fastq.gz: 998,428リード、56,019,410 bytes



Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

アセンブル:ゲノム用

ゲノム用。②非モデル生物やヘテロ接合度の高い生物種用、③微生物など小～中規模ゲノム配列決定用。アセンブルのアルゴリズム周辺は「2014.06.25」でwebページ内を検索

- 前処理 | フィルタリング | paired-end | 配列長とN数 | [QuasR\(Gaidatzis 2015\)](#) (last mod
- 前処理 | フィルタリング | paired-end | 共通リード抽出 | [ShortRead\(Morgan 2009\)](#) (las

- アセンブル | [ゲノム用](#) ①
- アセンブル | [トランスクリプト](#)
- マッピング | [マッピング](#) | [basic aligner](#) (last m
- マッピング | [splice-aware ali](#)
- マッピング | [Bisulfite sequen](#)
- マッピング | [\(ESTレベルの長](#)
- マッピング | [基礎](#) (last modif
- マッピング | [single-end](#) | ゲノ
- マッピング | [single-end](#) | ゲノ
- マッピング | [single-end](#) | ゲノ

アセンブル | ゲノム用

Rパッケージはおそらくありません。

プログラム:

- Phusion: [Mullikin and Ning, Genome Res., 2003](#)
- PCAP: [Huang et al., Genome Res., 2003](#)
- Newbler: [Margulies et al., Nature, 2005](#)
- SSAKE: [Warren et al., Bioinformatics, 2007](#)
- Minimus: [Sommer et al., BMC Bioinformatics,](#)
- VCAKE: [Jeck et al., Bioinformatics, 2007](#)
- SHARCGS: [Dohm et al., Genome Res., 2007](#)
- Edena: [Hernandez et al., Genome Res., 2008](#)
- Velvet: [Zerbino and Birney, Genome Res., 200](#)
- CABOG: [Miller et al., Bioinformatics, 2008](#)
- EULER-SR: [Chaisson et al., Genome Res., 200](#)
- SHORTY: [Hossain et al., BMC Bioinformatics](#)
- QSRA: [Bryant et al., BMC Bioinformatics, 200](#)
- ABYSS: [Simpson et al., Genome Res., 2009](#)
- Taipan: [Schmidt et al., Bioinformatics, 2009](#)
- Forge: [Diguistini et al., Genome Biol., 2009](#)
- ALLPATHS-LG: [Maccallum et al., Genome B](#)
- SOAPdenovo: [Li et al., Genome Res., 2010](#)
- SGA: [Simpson et al., Genome Res., 2012](#)
- Mapsembler: [Peterlongo and Chikhi, BMC Bio](#)
- IDBA-UD: [Peng et al., Bioinformatics, 2012](#)
- SPAdes: [Bankevich et al., J Comput Biol., 201](#)

- SparseAssembler: [Ye et al., BMC Bioinformatics, 2012](#)
- Readjoinder: [Gonnella et al., BMC Bioinformatics, 2012](#)
- Fermi: [Li H., Bioinformatics, 2012](#)
- MetaVelvet(メタゲノム用): [Namiki et al., Nucleic Acids Res., 2012](#)
- TIGER: [Wu et al., BMC Bioinformatics, 2012](#)
- CloudBrush: [Chang et al., BMC Genomics, 2012](#)
- RACA: [Kim et al., PNAS, 2013](#)
- BayesHammer: [Nikolenko et al., BMC Genomics, 2013](#)
- SPA: [Yang and Yooseph, Nucleic Acids Res., 2013](#)
- SOAPdenovo2: [Luo et al., Gigascience, 2012](#)
- JR-Assembler: [Chu et al., Proc Natl Acad Sci U S A., 2013](#)
- Platanus: [Kajitani et al., Genome Res., 2014](#)
- Omega(メタゲノム用): [Haider et al., Bioinformatics, 2014](#)

Review、ガイドライン、パイプライン系:

- Review: [Miller et al., Genomics, 2010](#)
- CG pipeline(パイプライン; 454用): [Kislyuk et al., Bioinformatics, 2010](#)
- A5 pipeline(パイプライン): [Tritt et al., PLoS One, 2012](#)
- Review: [El-Metwally et al., PLoS Comput Biol., 2013](#)
- iMetAMOS(パイプライン): [Koren et al., BMC Bioinformatics, 2014](#)
- MyPro(パイプライン; Prokaryote用): [Liao et al., J Microbiol Methods., 2015](#)

アセンブル: 転写物用

- [アセンブル | について](#) (last modified 2014/06/20)
- [アセンブル | ゲノム用](#) (last modified 2015/05/11)
- [アセンブル | トランスクリプトーム\(転写物\)用](#) (last modified 2015/04/22)
- [マッピング | について](#) (last modified 2015/01/16)
- [マッピング | basic aligner](#) (last modified 2014/08/08)
- [マッピング | splice-aware aligner](#) (last modified 2015/06/02)
- [マッピング | Bisulfite sequencing用](#) (last modified 2014/07/09)
- [マッピング | \(ESTレベルの長さの\)contig](#) (last modified 2014/06/24)
- [マッピング | 基礎](#) (last modified 2013/06/19)
- [マッピング | single-end | ゲノム | basic aligner\(基礎\)](#) | [QuasR\(Gaidatzis 2011\)](#)
- [マッピング | single-end | ゲノム | basic aligner\(応用\)](#) | [QuasR\(Gaidatzis 2011\)](#)
- [マッピング | single-end | ゲノム | splice-aware aligner](#) | [QuasR\(Gaidatzis 2011\)](#)
- [マップ後 | について](#) (last modified 2013/06/19)
- [マップ後 | 出力ファイル形式について](#) (last modified 2013/11/05)
- [マップ後 | 出力ファイルの読み込み](#) | [BAM形式](#) (last modified 2014/06/21)
- [マップ後 | 出力ファイルの読み込み](#) | [Bowtie形式](#) (last modified 2013/06/18)
- [マップ後 | 出力ファイルの読み込み](#) | [SOAP形式](#) (last modified 2013/06/19)

アセンブル | トランスクリプトーム(転写物)用

Rパッケージはおそらくありません。

プログラム:

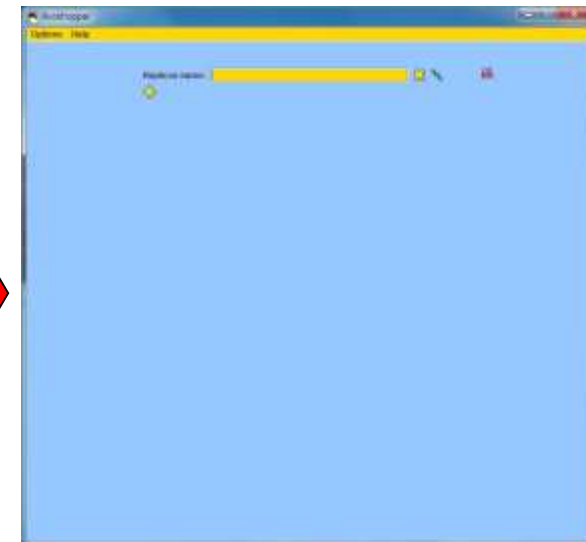
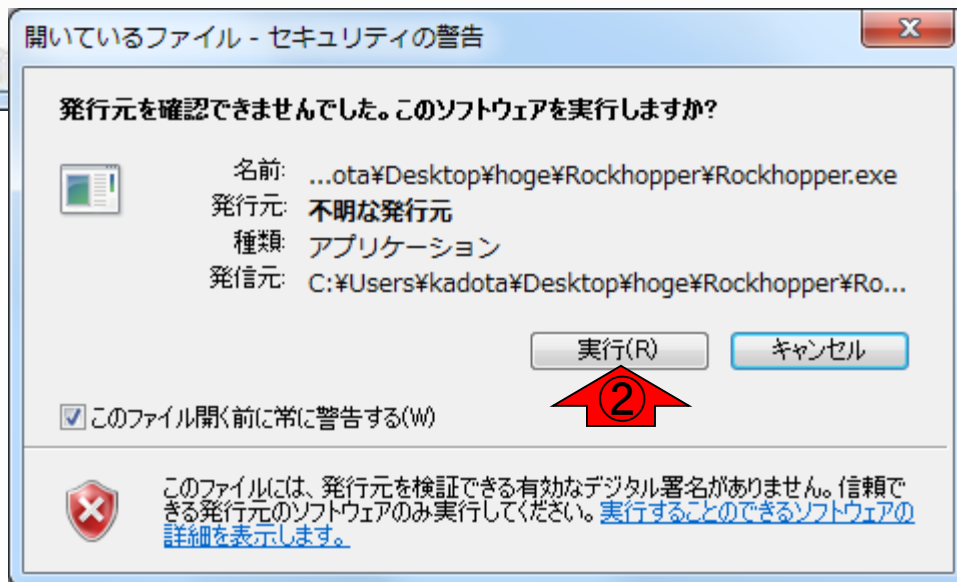
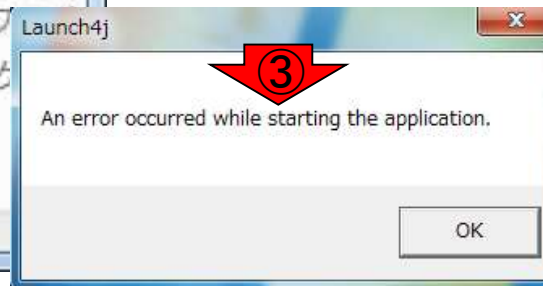
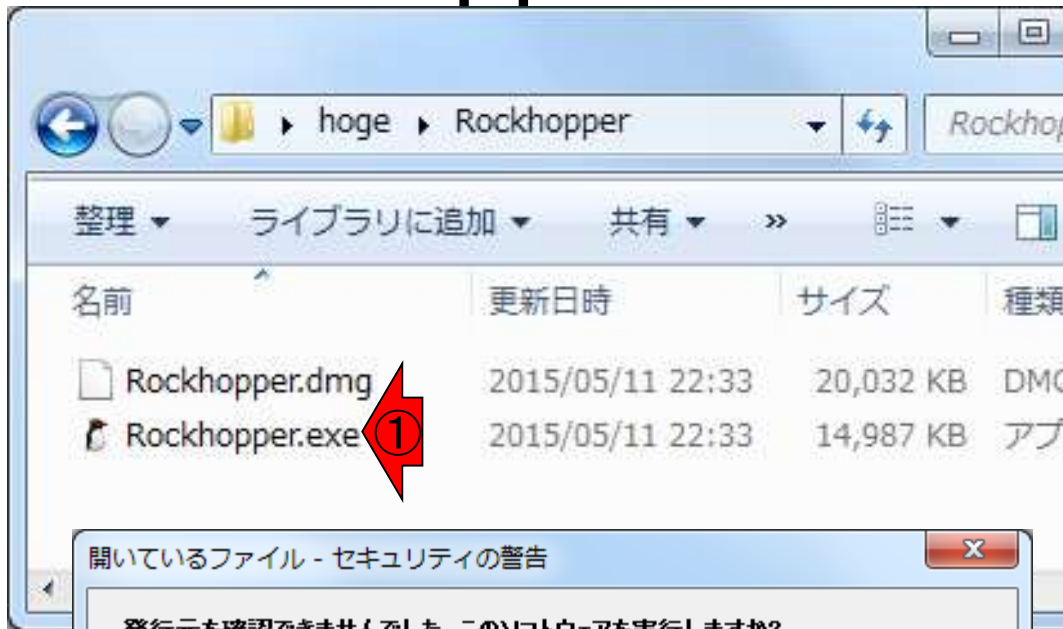
- [Multiple-k: Surget-Groba and Montoya-Burgos, Genome Res., 2010](#)
- [Trans-ABYSS: Robertson et al., Nat Methods, 2010](#)
- [Rnnotator: Martin et al., BMC Genomics, 2010](#)
- [Trinity: Grabherr et al., Nat Biotechnol, 2011](#)
- [SOAPdenovo-trans: Luo et al., Gigascience, 2012](#)
- [Oases: Schulz et al., Bioinformatics, 2012](#)
- [EBARDenovo: Chu et al., Bioinformatics, 2013](#)
- [BRANCH: Bao et al., Bioinformatics, 2013](#)
- [IDBA-tran: Peng et al., Bioinformatics, 2013](#)
- [SOAPdenovo-Trans: Xie et al., Bioinformatics, 2014](#)
- [Rockhopper 2\(バクテリア用\): Tjaden B, Genome Biol., 2015](#)
- [DETONATE\(RSEM-EVAL\): Li et al., Genome Biol., 2014](#)
- [Bridger: Chang et al., Genome Biol., 2015](#)
- [IFRAT: Mbandi et al., BMC Bioinformatics, 2015](#)

Review、ガイドライン、パイプライン系:

- [Review: Martin and Wang, Nat Rev Genet., 2011](#)
- [ガイドライン: Haznedaroglu et al., BMC Bioinformatics, 2012](#)
- [Review: Góngora-Castillo, Nat Prod Rep., 2013](#)
- [ガイドライン: Yang and Smith, BMC Genomics, 2013](#)
- [ガイドライン: O'Neil et al., BMC Genomics, 2013](#)
- [ガイドライン: Feldmesser et al., BMC Genomics, 2014](#)
- [パイプライン\(454用\): Melicher et al., BMC Genomics, 2014](#)
- [パイプライン\(組合せ系; SAMP and CDTA\): He et al., BMC Genomics, 2015](#)

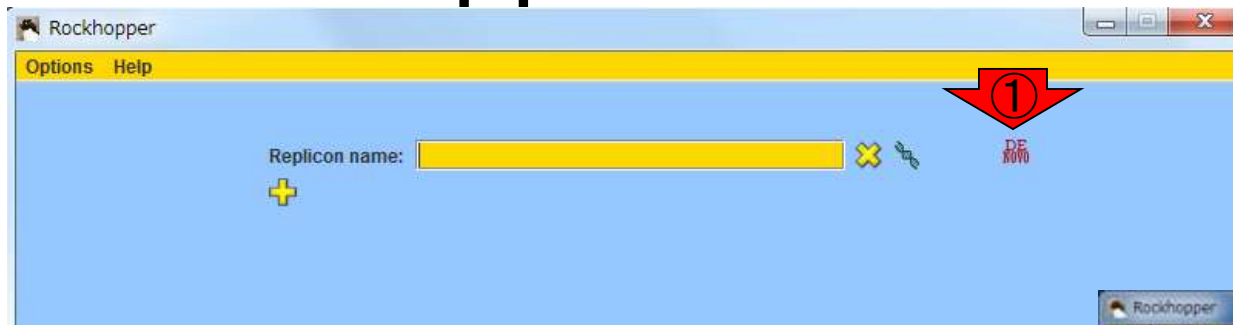
Rockhopper

①Windowsのヒトは.exeファイル、Macintoshのヒトは.dmgファイルでRockhopperを起動。「hoge - Rockhopper」フォルダ中にRockhopper_Resultsというフォルダがないヒトは、ダブルクリックで起動したときに自動的に作成されます。③最初「An error occurred while opening the file」的なポップアップが出現しますが、30秒ほど待つと、④RockhopperのGUIが起動します。



Rockhopper

①(見づらいが)DE NOVOと赤字で書いている部分をクリック。②入力ファイルを指定すべく、BROWSEボタンを押す。

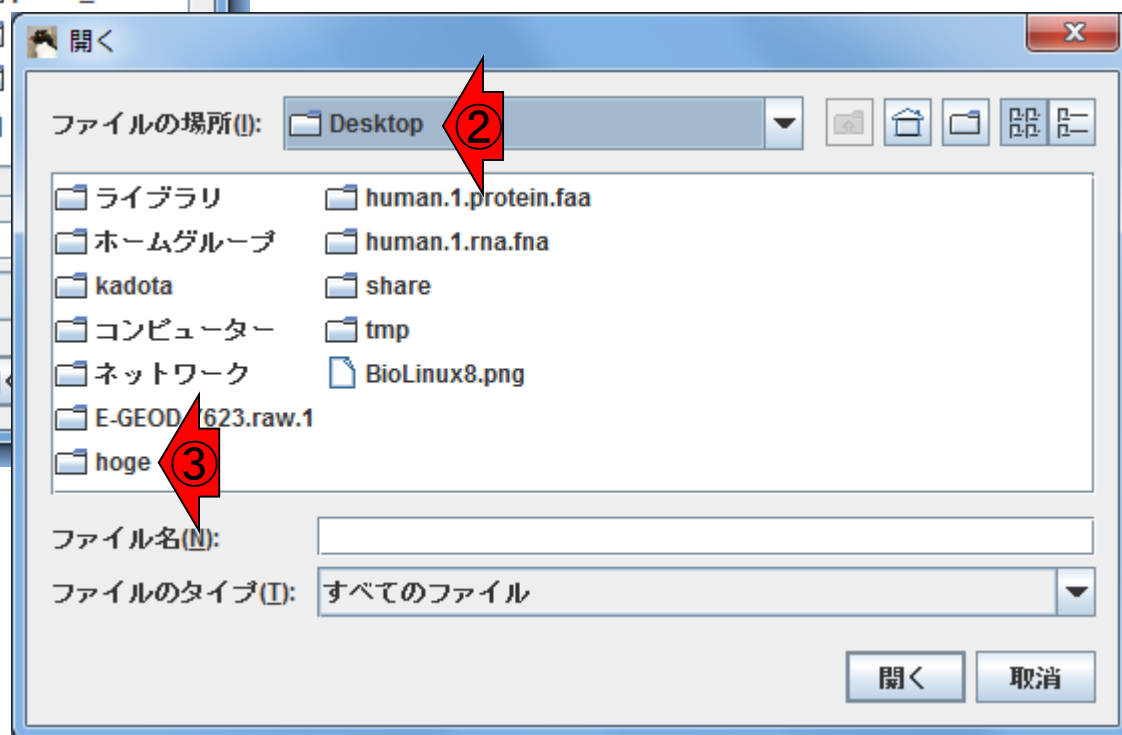
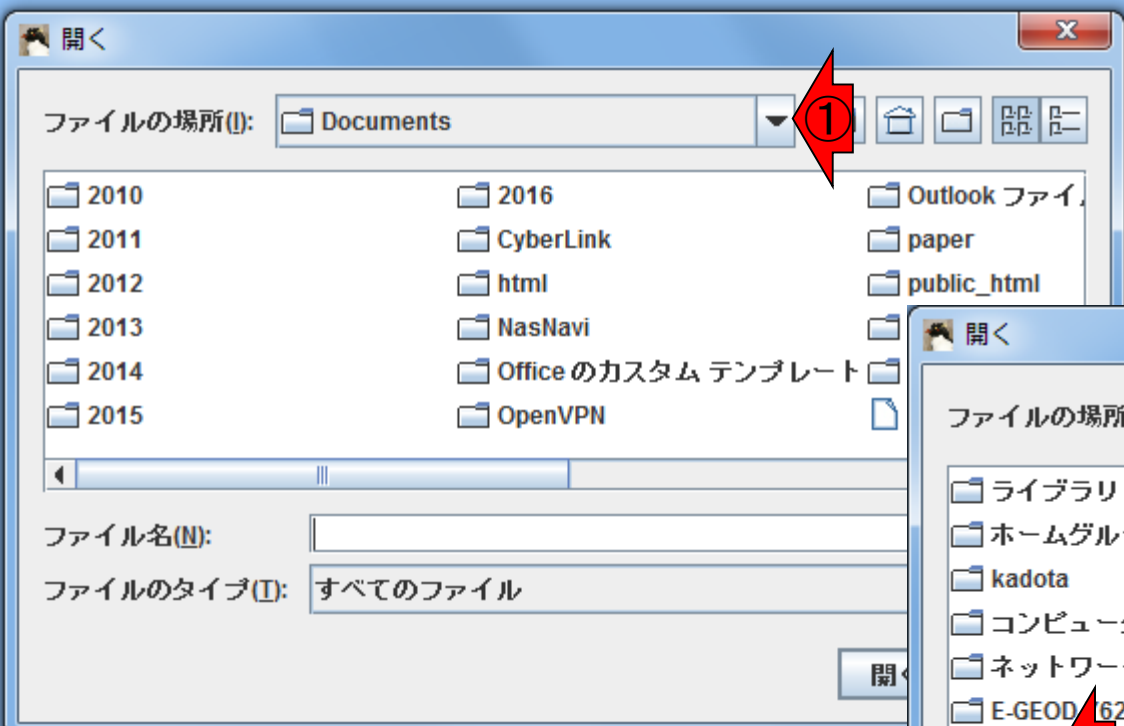


Rockhopper

①ファイルの場所を指定、②Desktop、③hoge。

DE
NOVO

Perform de novo transcript assembly without reference replicons



乳酸菌RNA-seqデータ

100万リードのオリジナル?!ファイル

- ▶ SRR616268sub_1.fastq.gz: 107 bp, 74,906,576 bytes
- SRR616268sub_2.fastq.gz: 93 bp, 67,158,462 bytes



100万リードのアダプター除去後のファイル

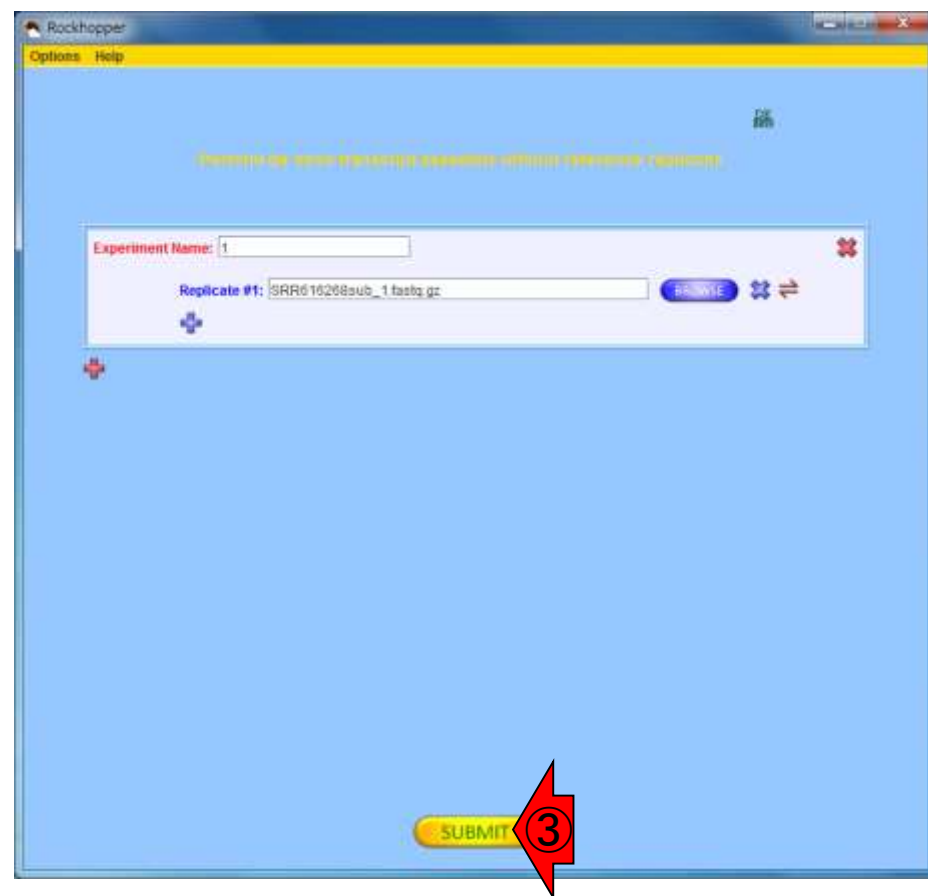
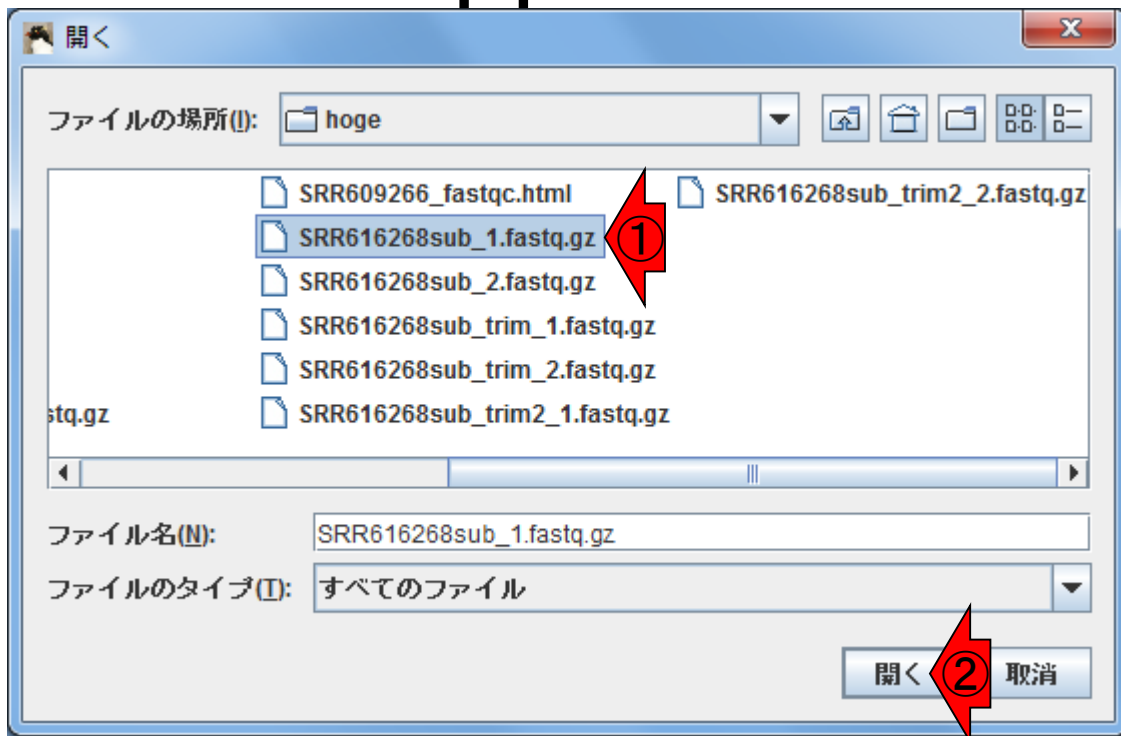
- SRR616268sub_trim2_1.fastq.gz(hoge4.fastq.gz): 71,343,605 bytes
- SRR616268sub_trim2_2.fastq.gz : 63,524,791 bytes
- hoge2_1.fastq.gz: 998,208リード、71,194,586 bytes
- hoge2_2.fastq.gz: 998,208リード、63,375,473 bytes

アダプター除去およびフィルタリング後のファイル

- SRR616268sub_trim_1.fastq.gz(hoge7.fastq.gz): 998,658リード、71,227,695 bytes
- SRR616268sub_trim_2.fastq.gz(hoge8.fastq.gz): 999,136リード、63,442,904 bytes
- hoge1_1.fastq.gz: 998,428リード、63,446,240 bytes
- hoge1_2.fastq.gz: 998,428リード、56,019,410 bytes

Rockhopper: SE

まずはsingle-end (SE)データとして実行。①入力ファイルを選び、②開く。③SUBMIT。約1分。



Rockhopper: SE

①入力ファイルは1,000,000リードだったので、一定の基準で内部的にフィルタリングしているのだろう。②アセンブル後のリファレンス配列に対してマッピングした結果、1,964リードがマップされたということだろう。③アセンブルによって得られた転写物配列数は8個しかなかった。④その平均配列長(121 bp)と⑤中央値(106 bp)。⑥アセンブルされたトータルの塩基数は974。これは③×④の結果(8×121=968)とほぼ同じなので、まあよしとしよう。

Progress

Initializing RNAseq analysis...

Assembling transcripts from reads in file:

C:\Users\kadota\Desktop\hoge\SRR616268s

Aligning reads to assembled transcripts using file:

C:\Users\kadota\Desktop\hoge\SRR616268s

Total reads in file:	996739	①	
Perfectly aligned reads:	1964	②	0%
Total number of assembled transcripts:	8	③	
Average transcript length:	121	④	
Median transcript length:	106	⑤	
Total number of assembled bases:	974	⑥	

Summary of results written to file: Rockhopper_Results/summary.txt

Details of assembled transcripts written to file: Rockhopper_Results/transcripts.txt

FINSIHED.

次にreverse側をsingle-end (SE) データとしてde novo assemble。

乳酸菌RNA-seqデータ

100万リードのオリジナル?!ファイル

- SRR616268sub_1.fastq.gz: 107 bp, 74,906,576 bytes
- ▶ □ SRR616268sub_2.fastq.gz: 93 bp, 67,158,462 bytes



100万リードのアダプター除去後のファイル

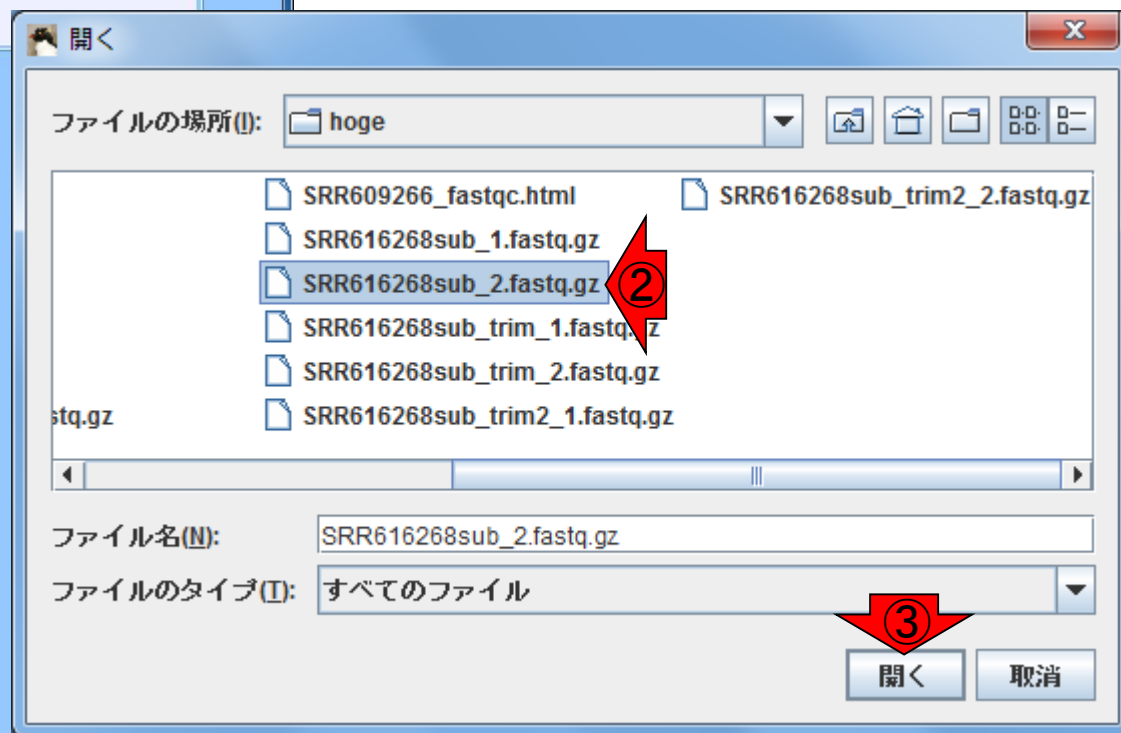
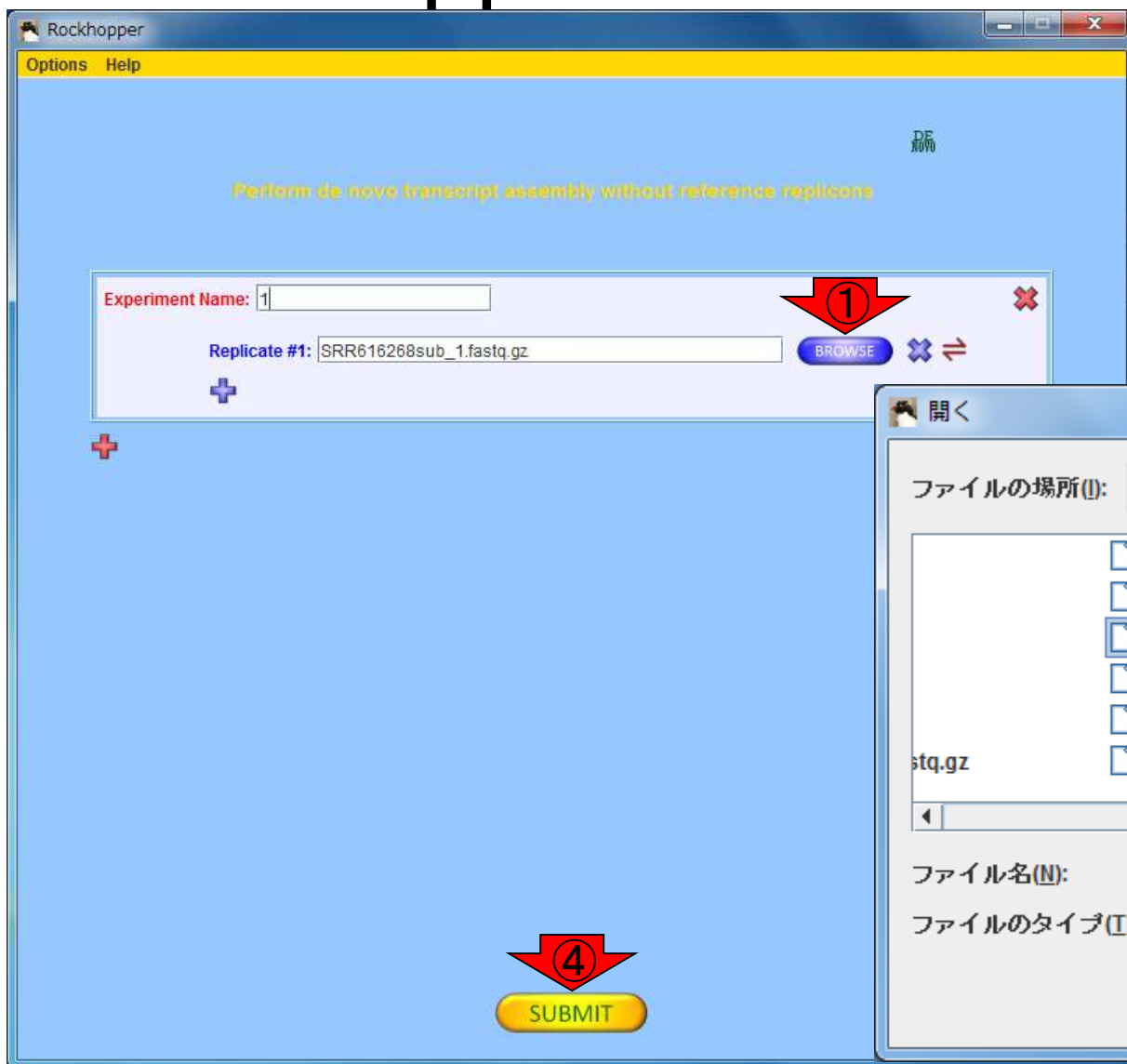
- SRR616268sub_trim2_1.fastq.gz(hoge4.fastq.gz): 71,343,605 bytes
- SRR616268sub_trim2_2.fastq.gz : 63,524,791 bytes
- hoge2_1.fastq.gz: 998,208リード、71,194,586 bytes
- hoge2_2.fastq.gz: 998,208リード、63,375,473 bytes

アダプター除去およびフィルタリング後のファイル

- SRR616268sub_trim_1.fastq.gz(hoge7.fastq.gz): 998,658リード、71,227,695 bytes
- SRR616268sub_trim_2.fastq.gz(hoge8.fastq.gz): 999,136リード、63,442,904 bytes
- hoge1_1.fastq.gz: 998,428リード、63,446,240 bytes
- hoge1_2.fastq.gz: 998,428リード、56,019,410 bytes

Rockhopper: SE

まずはsingle-end (SE)データとして実行。①入力ファイルを選び、②開く。③SUBMIT。約1分。



Rockhopper: SE

①入力ファイルは1,000,000リード。内部的にフィルタリングした結果、983,854リードになったのだろう。②アセンブル後のリファレンス配列に対してマッピングした結果、710,393リードがマップされたということだろう。③アセンブルによって得られた転写物配列数は424個。④その平均配列長(436 bp)と⑤中央値(228 bp)。⑥アセンブルされたトータルの塩基数は185,233。これは③×④の結果(424×436=184,864)とほぼ同じなので、まあよしとしよう

Progress

Initializing RNAseq analysis...

Assembling transcripts from reads in file:

C:\Users\kadota\Desktop\hoge\SRR616268s

Aligning reads to assembled transcripts using file:

C:\Users\kadota\Desktop\hoge\SRR616268s

Total reads in file:	983854	①	
Perfectly aligned reads:	710393	②	72%
Total number of assembled transcripts:	424	③	
Average transcript length:	436	④	
Median transcript length:	228	⑤	
Total number of assembled bases:	185233	⑥	

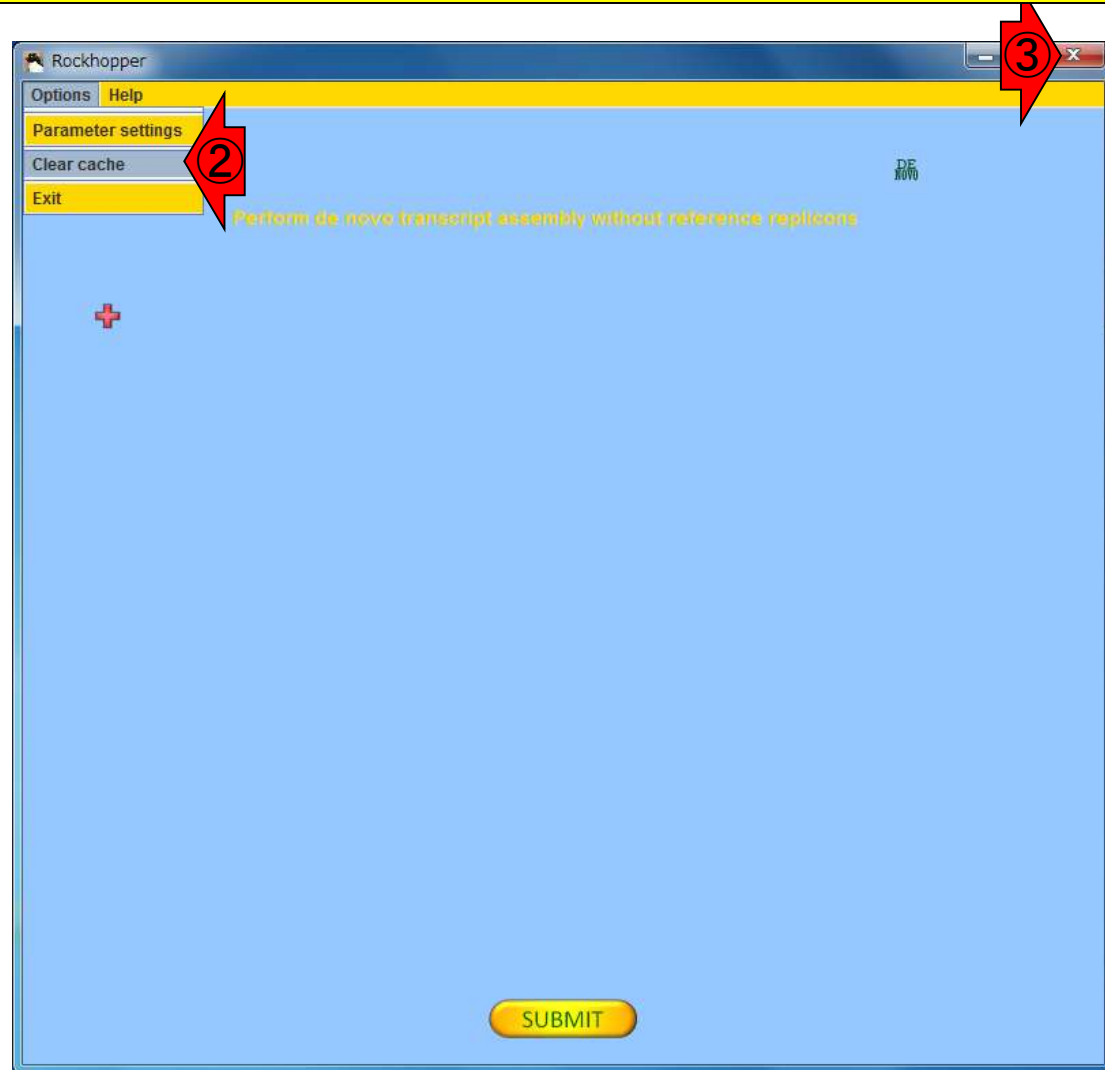
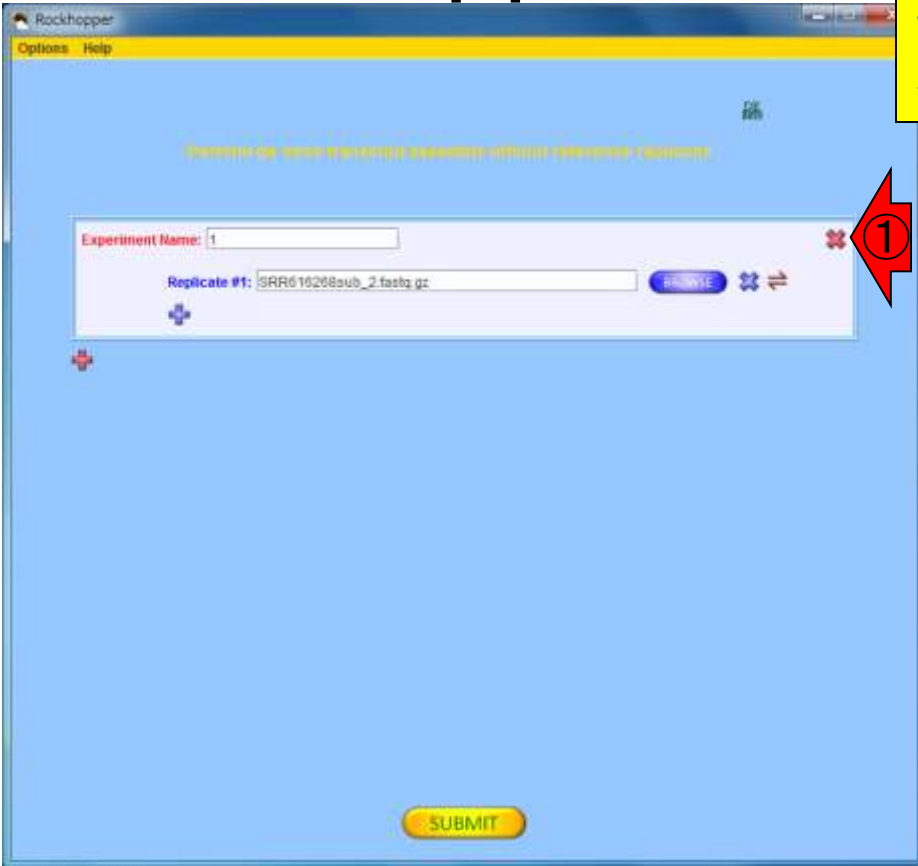
Summary of results written to file: Rockhopper_Results/summary.txt

Details of assembled transcripts written to file: Rockhopper_Results/transcripts.txt

FINISHED.

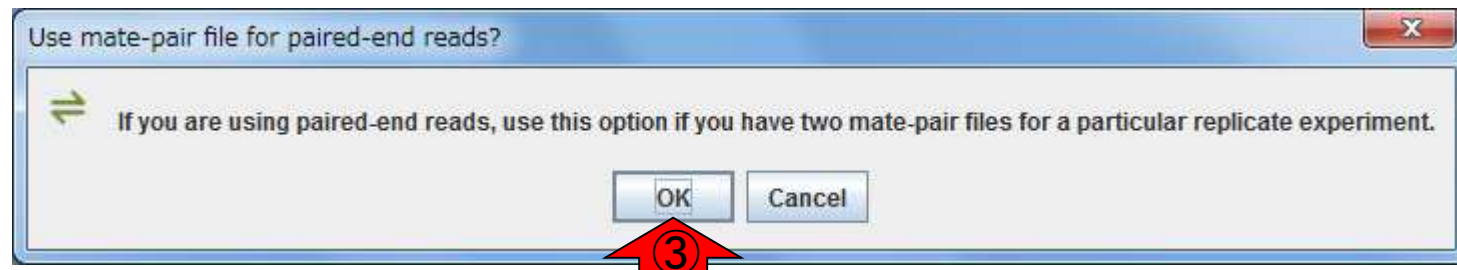
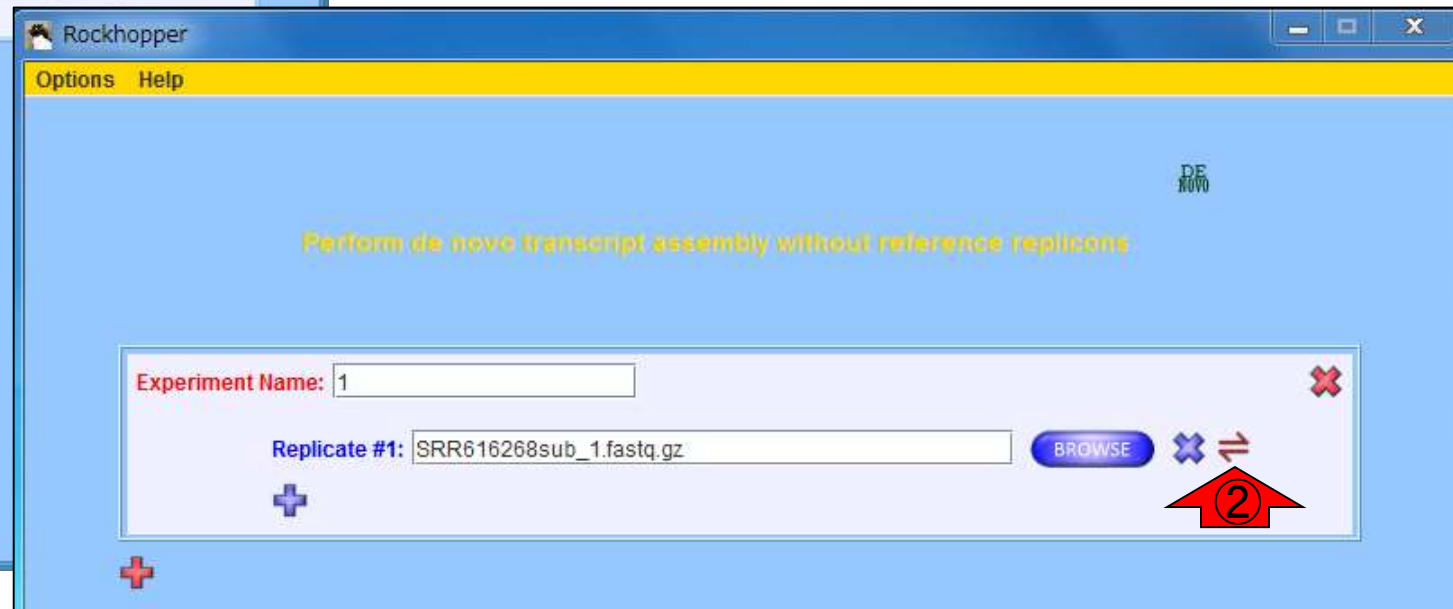
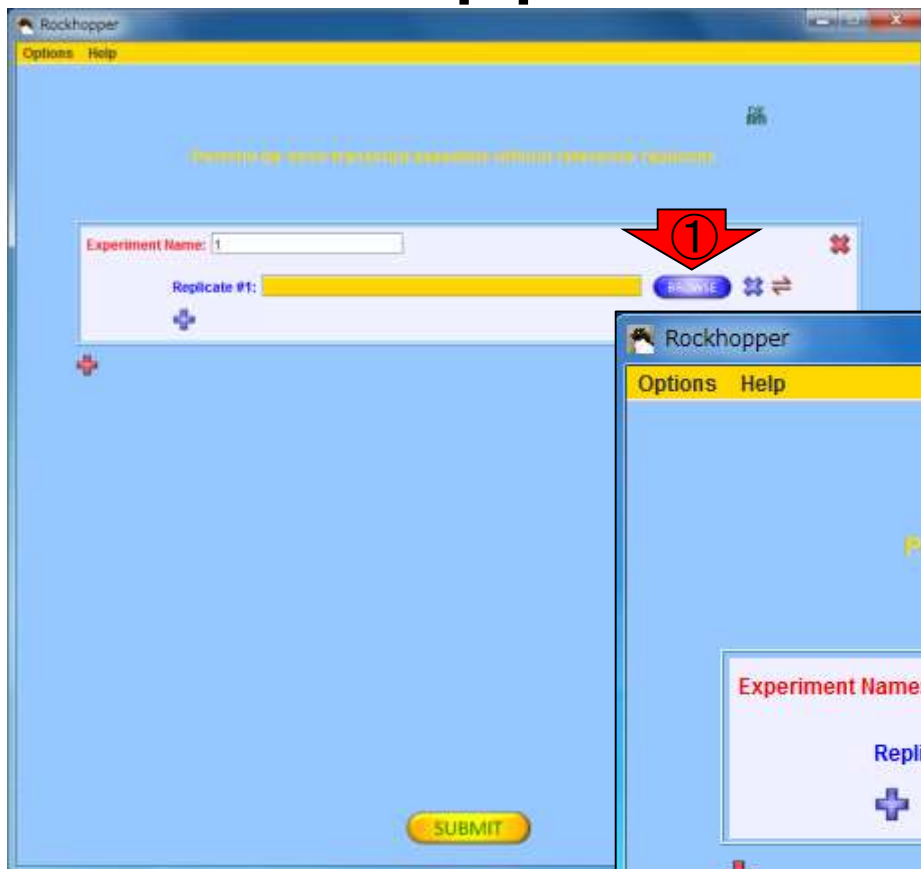
Rockhopper: PE

Paired-endデータとしてアセンブル。①念のため一旦消す。②「Options - Clear cache」というのがあったので、こちらも念のためやってみる。数分フリーズ状態になりました。③Rockhopperアプリごと再起動しても、キャッシュに残っているようなので、Clear cacheをやるのが無難かも



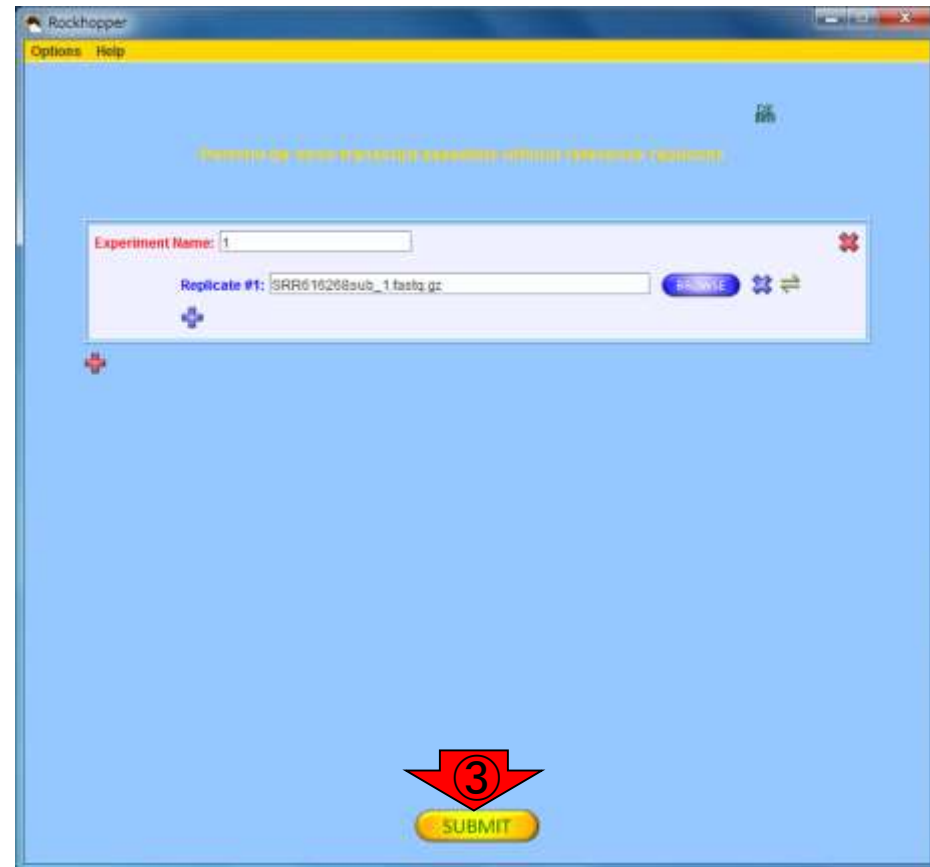
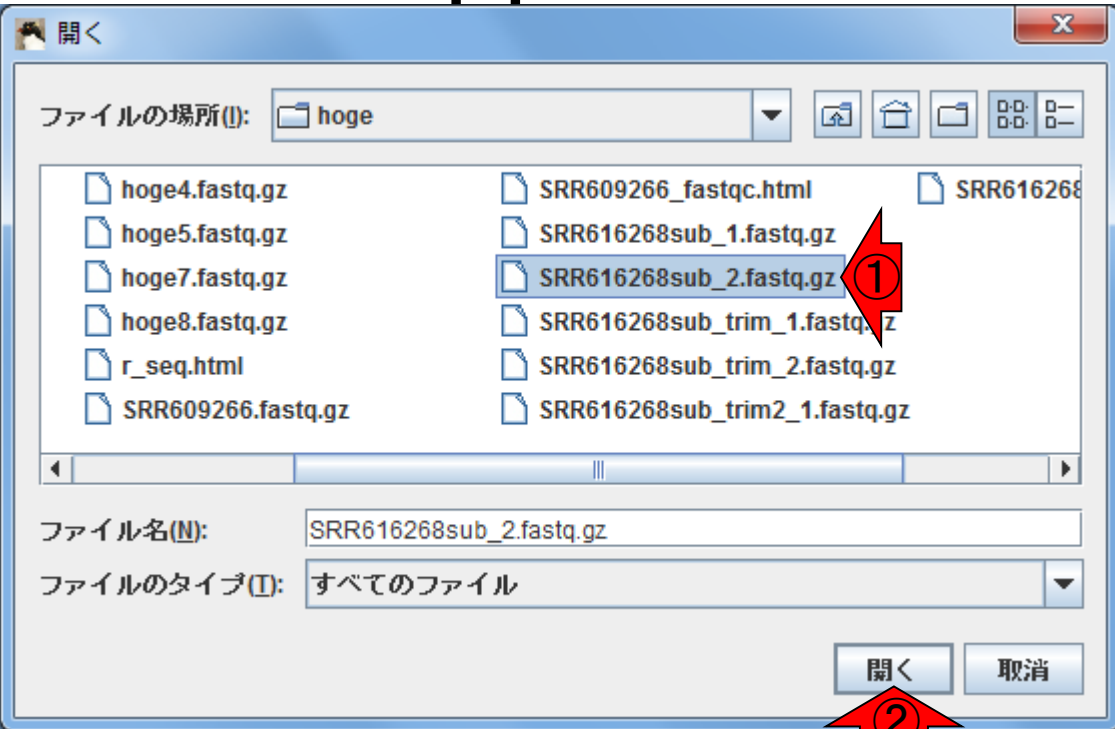
Rockhopper: PE

しばらく経つと左のような画面になるので、① paired-endの1つ目のforward側のファイルを指定。②「双方向矢印」のところにカーソルをもっていくと、「Using paired-end reads; ...」と書いてあるので、ここをクリック。③OKを押す。



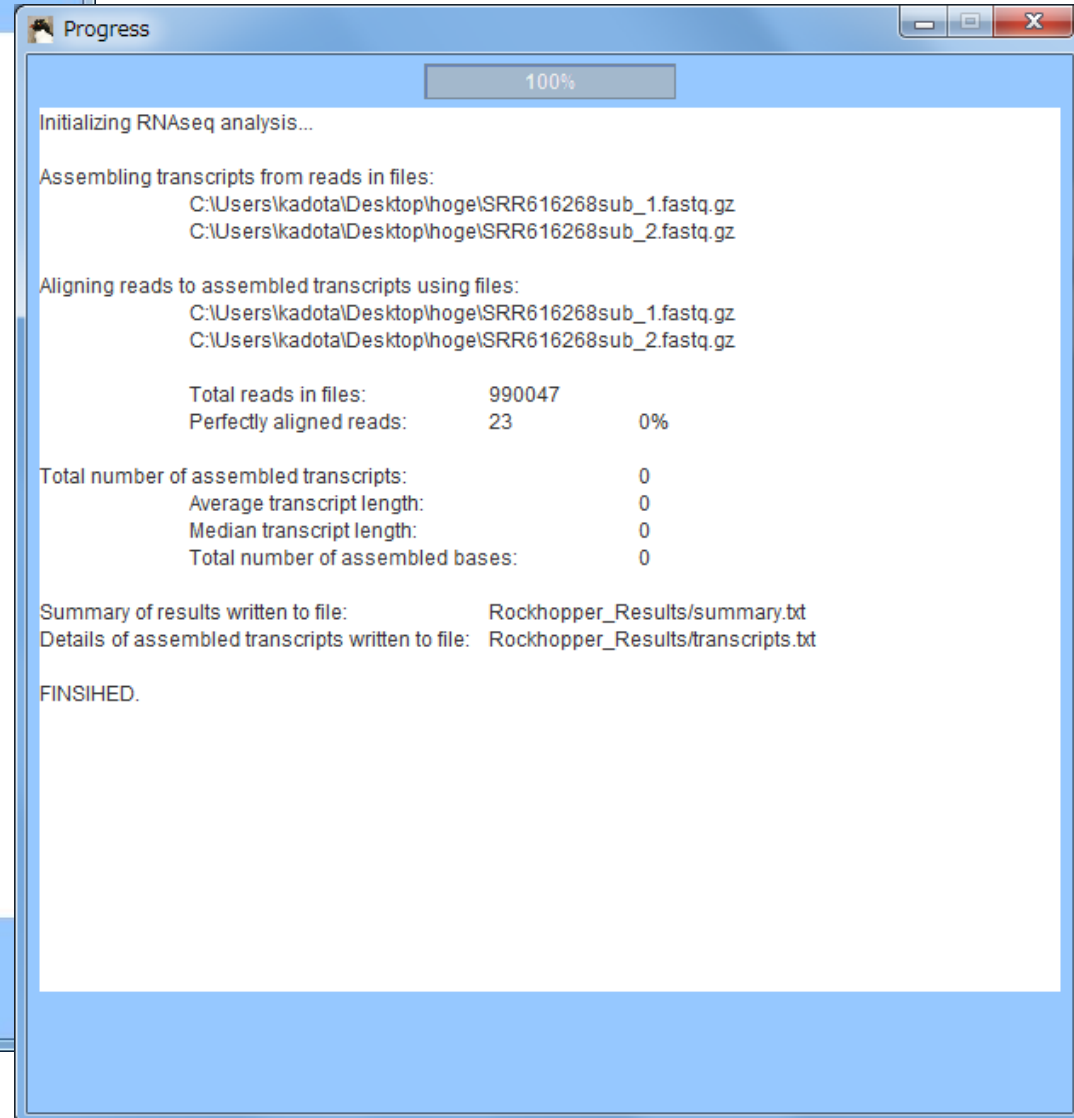
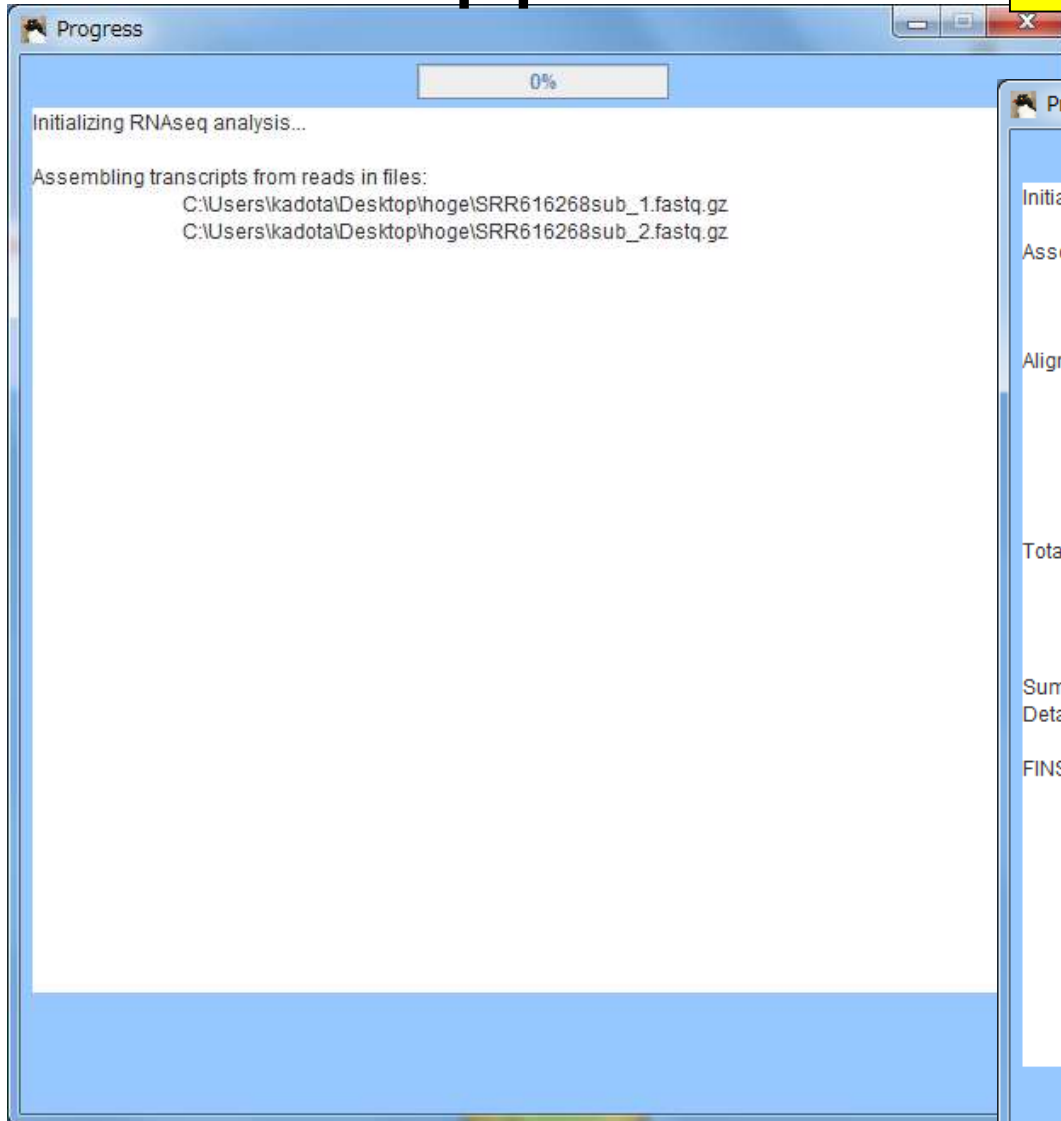
Rockhopper: PE

①reverse側のファイルを選択して、②開く。Paired-end の2つのファイルが見られるわけではないが、③SUBMIT



Rockhopper: PE

約2分で終わるが、残念ながらうまくアセンブルできていないようだ。後にこの原因は、forwardファイル中の3'側のadapter/primer配列が主原因と判明



Rockhopper: PE

それゆえ、①「Options - Parameter settings」で、
②Strand specific やReverse complement reads
のあたりをいじっても状況に変化はない。

Rockhopper

Options Help

Parameter settings

Clear cache

Exit

Perform de novo transcript assembly without reference

Experiment Name: 1

Replicate #1: SRR616268sub_trim_1.fastq.gz

+

+

SUBMIT

Parameter Settings

General parameters

Strand specific

Reverse complement reads

Test for differential expression

Orientation of mate-pair reads fr

Verbose output

Output SAM file

Max bases between paired reads 500

Number of processors 0

Output file location: Rockhopper_Results/ BROWSE

Parameters for reference based assembly

Allowed mismatches 0.15

Minimum seed length 0.33

Identify transcript boundaries

Predict operons

Minimum expression of UTRs and ncRNAs 0.0 0.5 1.0

Parameters for de novo assembly

Min reads mapping to a transcript 20

Minimum transcript length 50

Min count to seed a transcript 50

Min count to extend a transcript 5

DEFAULTS SAVE

アセンブル結果まとめ

ファイル名	Total number of assembled transcripts	Average length	Median length	Total number of assembled bases	Total reads in file	Perfectly aligned reads
・100万リードのオリジナル?!ファイル						
SRR616268sub_1.fastq.gz	8	121	106	974	996739	1964
SRR616268sub_2.fastq.gz	424	436	228	185233	983854	710393
paired-end	0	0	0	0	990047	23
・100万リードのアダプター除去後のファイル						
SRR616268sub_trim2_1.fastq.gz	8	121	104	974	993311	1231
SRR616268sub_trim2_2.fastq.gz	425	433	228	184437	981649	708951
hoge2_1.fastq.gz	8	121	104	974	993001	1231
hoge2_2.fastq.gz	424	434	229	184210	980997	708395
paired-end	0	0	0	0	989869	25
・アダプター除去およびフィルタリング後のファイル						
SRR616268sub_trim_1.fastq.gz	8	121	104	974	993311	1231
SRR616268sub_trim_2.fastq.gz	425	433	228	184437	981649	708951
hoge1_1.fastq.gz	8	121	104	974	993081	1231
hoge1_2.fastq.gz	423	435	229	184211	981125	708482
paired-end	0	0	0	0	989869	25

Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

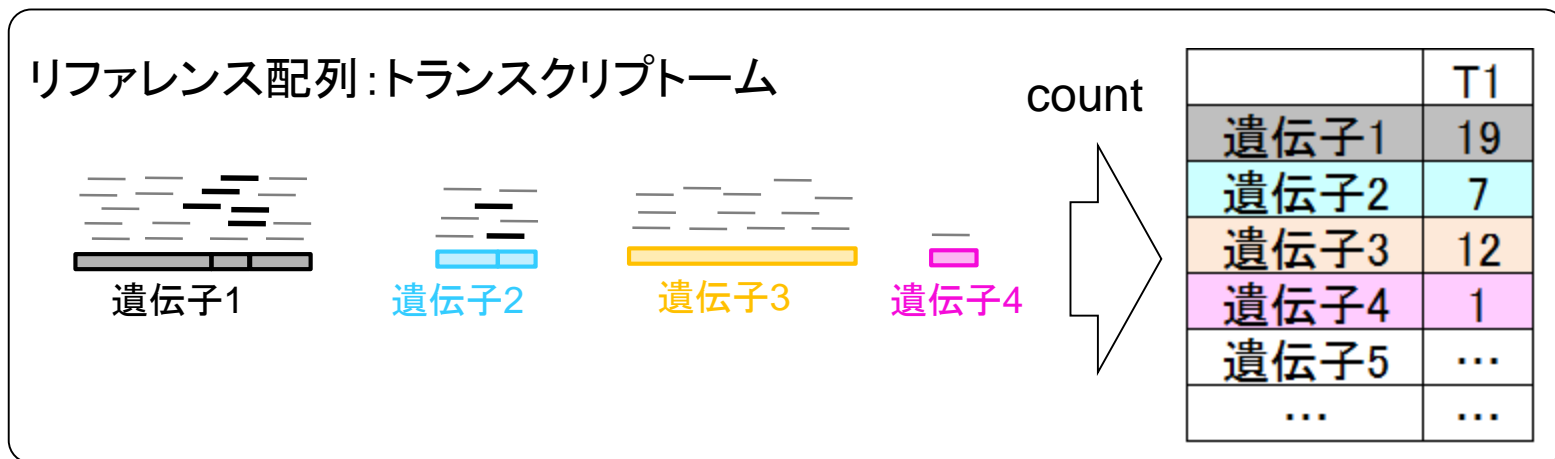
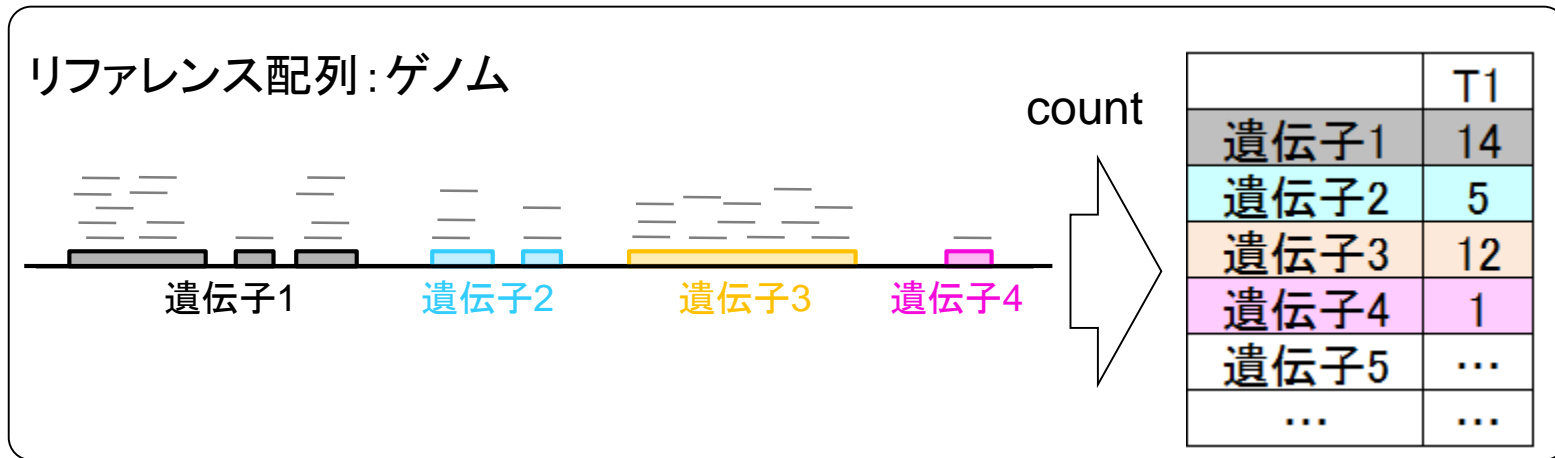
基本イメージ。「マップされたリード数 = 発現量」ではないが、マップされたリード数のカウント情報は、発現量推定の基本情報

マッピング基礎

■ 基本的なマッピングプログラム (bowtieなど) を用いた場合

T1サンプルの
RNA-Seqデータ

mapping



マッピング基礎

- マップされる側のリファレンス配列: `hoge4.fa`
- マップする側のRNA-seqデータ(リードと呼ばれる): "AGG"

「マッピング = 大量高速文字列検索」という捉え方でよい。マッピングプログラムの出力は、(どのリードが)リファレンス配列上のどの位置から転写されたものかという座標情報

```

hoge4.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>contig_1
CGGACAGCTCCTCGGCATCCGGAT
>contig_2
GTCTGCCTCAAGCGCCCCAAGTGGGTTTGGAGGCCTAACATCGCAAGTCG
ACACTCAGTCCGGCCGTCTGGTTGGCAGGGGCAGAGACCCAGCACACCCT
GTC
>contig_3
TGTAGGAGAAGGGCGGTATCAGCGTCCACTTACACGATCCGTTACTAATT
GTATGAGGTCGGGCA
>contig_4
CGTGCTGATTCCACACAGCAGTAAACGCGGACCTCTACCTATGAACATG

```

出力ファイル

	start	end
contig_2	31	33
contig_2	77	79
contig_3	4	6
contig_3	10	12
contig_3	56	58

マッピング

Rパッケージとして提供されているものは圧倒的に少ないですが、QuasRパッケージ (Gaidatzis et al., 2015)のおかげでWindows OS上でもマッピングを気軽に利用できるようになりました。これは内部的にBowtie (Langmead et al., 2009)プログラムを利用しています。

- 前処理 | トリミング | [指定した末端塩基数だけ除去](#) (last modified 2013/06/15)
- [アセンブル | について](#) (last modified 2014/06/20)
- アセンブル | [ゲノム用](#) (last modified 2014/07/08)
- アセンブル | [トランスクリプトーム\(転写物\)用](#) (last modified 2014/07/08)
- [マッピング | について](#) (last modified 2015/01/16)
- マッピング | [basic aligner](#) (last modified 2014/08/08)
- マッピング | [splice-aware aligner](#) (last modified 2014/07/09)
- マッピング | [Bisulfite sequencing用](#) (last modified 2014/07/09)
- マッピング | [\(ESTレベルの長さの\)contig](#) (last modified 2014/07/09)
- マッピング | [基礎](#) (last modified 2013/06/19)
- マッピング | [single-end | ゲノム | basic aligner\(基礎\)](#) | [QuasR](#)
- マッピング | [single-end | ゲノム | basic aligner\(応用\)](#) | [QuasR](#)
- マッピング | [single-end | ゲノム | splice-aware aligner](#) | [QuasR](#)
- [マップ後 | について](#) (last modified 2013/06/19)
- [マップ後 | 出力ファイル形式 | について](#) (last modified 2013/11/11)

マッピング | について

リファレンス配列にマッピングを行うプログラム達です。basic aligner (unspliced aligner)はsplice-aware aligner (spliced aligner)内部で使われていたりします。

R用:

- [Rsubread](#)(Windows版なし): [Liao et al., Nucleic Acids Res., 2013](#)
- [QuasR](#)(Windows版あり): [Gaidatzis et al., Bioinformatics, 2014](#)
- [HTSeqGenie](#)(Windows版なし): [原著論文はまだみたいです](#)

R以外(basic aligner; unspliced aligner):

- [SSAHA2](#): [Ning et al., Genome Res., 2001](#)
- [RMAP](#): [Smith et al., BMC Bioinformatics, 2008](#)
- [MAQ](#): [Li et al., Genome Res., 2008](#)
- [PASS](#): [Campagna et al., Bioinformatics, 2009](#)
- [MOM](#): [Zaves and Gao, Bioinformatics, 2009](#)
- [Bowtie](#): [Langmead et al., Genome Biol., 2009](#)
- [BWA](#): [Li and Durbin, Bioinformatics, 2009](#)(BWA-shortの論文)
- [SHRiMP](#): [Rumble et al., PLoS Comput Biol., 2009](#)

オプションは、デフォルトである程度よきに計らってくれるが...実際の挙動を完全に把握できる状況で様々なオプションを試したい

オプション

■ QuasRパッケージは内部的にBowtieを利用

- マッピング時に多くのオプションを指定可能
- “-v”: 許容するミスマッチ(mismatch)数を指定するオプション。“-v 0”は、リードがリファレンスに完全一致するもののみレポート。“-v 2”は、2塩基ミスマッチまで許容してマップされうる場所を探索。
- “-m”: 出力するリード条件を指定するオプション。“-m 1”は、複数個所にマップされるリードを除外して、1か所にのみマップされたリードをレポート。“-m 3”は、合計3か所にマップされるリードまでをレポート。
- “--best --strata”: 最も少ないミスマッチ数でマップされるもののみ出力する、という意味表示。これをつけずに“-v 2 -m 1”などと指定すると、たとえ完全一致(ミスマッチ数0)で1か所にのみマップされるリードがあったとしても、どこか別の場所で1塩基ミスマッチでマップされる個所があれば、マップされうる場所が2か所ということを意味し、そのリードは出力されなくなる。それを防ぐのが主な目的
- ...

仮想リファレンス配列

- マップされる側のリファレンス配列: `ref_genome.fa`

chr3とchr5の違いは、2番目と7番目の塩基のみ。主に"-m"オプション(-m 1:1ヶ所にマップされるリードのみ出力)の違いの把握が可能

```

ref_genome.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG
>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACGCCCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
  
```

仮想リファレンス配列

■ マップされる側のリファレンス配列: `ref_genome.fa`

- 個別パッケージのインストール (last modified 2015/02/20) **NEW**
- 基本的な利用法 (last modified 2015/01/16)
- サンプルデータ (last modified 2015/02/15) **NEW**
- 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](#) (last modified 2015/02/15)

- 書籍
- 書籍
- 書籍

サンプルデータ **NEW**

1. Illumina/36bp/single-end/human (SRA000299) data (Marioni et al., Genome Res., 2008)

「Kidney 7 samples vs Liver 7 samples」のRNA-seqの遺伝子発現行列データ(SupplementaryTable2.txt)です。サンプルは二つの濃度(1.5 pM and 3 pM)でシーケンスされており、「3 pMのものが5 samples vs. 5 samples」の構成です。

18. ランダムな塩基配列から生成したリファレンスゲノム配列データ(`ref_genome.fa`)。48, 160, 100, 123, 100 bpの配列長をもつ、計5つの塩基配列を生成しています。description行は"contig"という記述を基本としています。塩基の存在比はAが28%, Cが22%, Gが26%, Tが24%にしています。set.seed関数を利用し、chr3の配列と同じものをchr5としてコピーして作成したのち、2番目と7番目の塩基置換を行っています。そのため、実際に指定するのは最初の4つ分の配列長のみです。

```

out_f <- "ref_genome.fa" #出力ファイル名を指定してout_fに格納
param_len_ref <- c(48, 160, 100, 123) #配列長を指定
narabi <- c("A", "C", "G", "T") #以下の数値指定時にACGTの並びを間違えないよう
param_composition <- c(28, 22, 26, 24) # (A, C, G, Tの並びで)各塩基の存在比率を指定
param_desc <- "chr" #FASTA形式ファイルのdescription行に記述する
param4 <- 3 #コピーを作成したい配列番号を指定
param5 <- c(2, 7) #コピー先配列の塩基置換したい位置を指定

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

#塩基置換関数の作成
genkichikan <- function(fa, n) { #関数名や引数の作成

```

許容するミスマッチ数による違いや、マップされるべき場所が完全に把握できるように、リードのdescription行に情報を記載。

仮想RNA-seqリード

- マップする側のRNA-seqデータ: `sample_RNAseq1.fa`

```
ref_genome.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG
>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACGCCCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
```

```
sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1 11 45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
>chr3_3_37
GGGGACTATTTCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTC
```


仮想RNA-seqリード

- マップする側のRNA-seqデータ: `sample_RNAseq1.fa`

- [個別パッケージのインストール](#) (last modified 2015/02/20) **NEW**
- [基本的な利用法](#) (last modified 2015/01/16)
- [サンプルデータ](#) (last modified 2015/02/15) **NEW**
- [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](#) (last modified 2015/02/15)

- 書籍
- 書籍
- 書籍

サンプルデータ **NEW**

1. Illumina/36bp/single-end/human (SRA000299) data (Marioni et al., Genome Res., 2008)
 「Kidney 7 samples vs Liver 7 samples」のRNA-seqの遺伝子発現行列データ(Supplement 1)です。サンプルは二つの濃度(1.5 pM and 3 pM)でシーケンスされており、「3 pMのものが5

19. 上記リファレンスゲノム配列データ([ref_genome.fa](#))に対してbasic alignerでマッピングする際の動作確認用RNA-seqデータ ([sample_RNAseq1.fa](#))とそのリファレンス配列を読み込んで、[list_sub3.txt](#)で抽出した部分配列を抽出したものです。どこに置換を入れているかがわかっています。DNAStrngSetオブジェクトを入力として塩基置換を行うDNAStrngSetオブジェクトを用いて、最後のリードのみ4番目の塩基にミスマッチを入れています。

```

in_f1 <- "ref_genome.fa"      #入力ファイル
in_f2 <- "list_sub3.txt"     #入力ファイル
out_f <- "sample_RNAseq1.fa" #出力ファイル
param <- 4                  #塩基置換した回数

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

#塩基置換関数の作成
DNAStrng_chartr <- function(fa, p) { #関数名や引数の作成
  #文字列に変更

```

```

sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
>chr3_3_37
GGGGACTATTTCCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGGACTATTTCCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGCGGTCTATTTCCCCGCTTGCAGGAATCGTGTC

```

Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

マッピング本番

- マッピング | (ESTレベルの長さの)contig (last modified 2014/06/24)
- マッピング | 基礎 (last modified 2013/06/19)
- マッピング | single-end | ゲノム | basic aligner(基礎) | QuasR(Gaidatzis 2014) (last modified 2014/06/21)
- マッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis 2014) (last modified 2015/02/24)
- マッピング | single-end | ゲノム | splice-aware aligner | QuasR(Gaidatzis 2014) (last modified 2014/06/21)
- マップ後 | について (last modified 2013/06/10)
- マップ後 | 出力ファイル
- マップ後 | 出力ファイル

マッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2014) NEW

QuasRパッケージを用いて single-end RNA-seqデータのリファレンスゲノム配列へのマッピングを行うやり方を示します。basic alignerの一つであるBowtie (Langmead et al., Genome Biol., 2009)を実装した Rbowtieパッケージを内部的に使っています。Bowtie自体は、複数個所にマップされるリードの取り扱い(uniqely mapped reads or multi-mapped reads)を"-m"オプションで指定したり、許容するミスマッチ数を指定する"-v"などの様々なオプションを利用可能ですが、「基礎」のところではやり方を示しませんでした。ここでは、マッピングのオプションをいくつか変更して挙動を確認したり、複数のRNA-seqファイルを一度にマッピングするやり方を示します。尚、出力ファイルは、"*bam", "*_QC.pdf", "*.bed"の3つです。それ以外のファイルは基本無視で大丈夫です。

「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。



2. サンプルデータ18,19のRNA-seqデータ(sample RNAseq1.fa)のref genome.faへのマッピングの場合(mapping_single_genome1.txt):

オプションを"-m 1 --best --strata -v 0"とした例です。sample_RNAseq1.faでマップされないのは計3リードです。2リード("chr3_11_45"と"chr3_15_49")はchr5にもマップされるので、"-m 1"オプションで落とされます。1リード("chr5_1_35")は該当箇所と完全一致ではない(4番目の塩基にミスマッチをいれている)ので落とされます。

```

in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納(リファレンス配列)
param_mapping <- "-m 1 --best --strata -v 0"#マッピング時のオプションを指定

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

#本番(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping)#マッピングを行うqAlign関数を実行した結果をout
time_e <- proc.time() #計算時間を計測するため
time_e - time_s #計算時間を表示(一番右側の数字。単位はsecond)
out #マッピングに用いたパラメータや入力ファイルの情報などを表示
alignmentStats(out) #マッピング結果(alignment_statistics)の表示。seqlength:リファレンス配列
    
```

マッピング本番

- マッピング | (ESTレベルの長さの)contig (last modified 2014/06/24)
- マッピング | 基礎 (last modified 2013/06/19)
- マッピング | single-end | ゲノム | basic aligner(基礎) | QuasR(Gaidatzis 2014) (last modified 2014/06/21)
- マッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis 2014) (last modified 2015/02/24)
- マッピング | single-end | ゲノム | splice-aware aligner | QuasR(Gaidatzis 2014) (last modified 2014/06/21)
- マップ後 | 出力ファイル (last modified 2013/06/10)
- マップ後 | 出力ファイル
- マップ後 | 出力ファイル

マッピング | single-end | ゲノム | basic aligner(応用) | QuasR

QuasRパッケージを用いて single-end RNA-seqデータのリファレンスゲノム配列へのマッピングを行う。あるBowtie (Langmead et al., Genome Biol., 2009)を実装した Rbowtieパッケージを内部的に使用しています。Bowtie自体は、複数個所にマップされるリードの取り扱い(uniqely mapped reads or multi-mapped reads)を"-m"オプションで指定したり、許容するミスマッチ数を指定する"-v"などの様々なオプションを利用可能ですが、「基礎」のところではやり方を示しませんでした。ここでは、マッピングのオプションをいくつか変更して挙動を確認したり、複数のRNA-seqファイルを一度にマッピングするやり方を示します。尚、出力ファイルは、"*.bam", "*_QC.pdf", "*.bed"の3つです。それ以外のファイルは基本無視で大丈夫です。

「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

	A	B
1	FileName	SampleName
2	sample RNAseq1.fa	naamae

複数のRNA-seqサンプルを実行できるようにリストファイルとして与える

1. サンプルデータ18,19のRNA-seqデータ(sample RNAseq1.fa)のref_genome.faへのマッピング

オプションを"-m 1 --best --strata -v 0"とした例です。sample_RNAseq1.faでマップされないのは計3リードです。2リード("chr3_11_45"と"chr3_15_49")はchr5にもマップされるので、"-m 1"オプションで落とされます。1リード("chr5_1_35")は該当箇所と完全一致ではない(4番目の塩基にミスマッチをいれている)ので落とされます。

```

in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納(リファレンス配列)
param_mapping <- "-m 1 --best --strata -v 0" #マッピング時のオプションを指定

#必要なパッケージをロード
library(QuasR) #パッケージの読
library(GenomicAlignments) #パッケージの読

#本番(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッピングを行うqAlign関数を実行した結果をout
time_e <- proc.time() #計算時間を計測するため
time_e - time_s #計算時間を表示(一番右側の数字。単位はsecond)
out #マッピングに用いたパラメータや入力ファイルの情報などを表示
alignmentStats(out) #マッピング結果(alignment_statistics)の表示 seqlength:リファレンス
    
```

許容するミスマッチ数は0個("-v 0")、1か所にマップされるリードのみ出力("-m 1")

マッピング本番

「hoge - mapping_kiso1」に最低限必要な3つの入力ファイルが揃っていることを確認してコピペ。ファイヤーフォールか何かでアクセス許可関連のポップアップが出たらOKを押す

1. サンプルデータ18,190のRNA-seqデータ(sample_RNAseq1.fa)のref_genome.faへのマッピングの場合(m

オプションを"-m 1 --best --strata -v 0"とした例です。sample_RNAseq1.faでマップされないのは計3リードで("chr3_11_45"と"chr3_15_49")はchr5にもマップされるので、"-m 1"オプションで落とされます。1リード("chr一致ではない(4番目の塩基にミスマッチをいれている)ので落とされます。

```

in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定し
in_f2 <- "ref_genome.fa" #入力ファイル名を指定し
param_mapping <- "-m 1 --best --strata -v 0" #マッピング時のオ

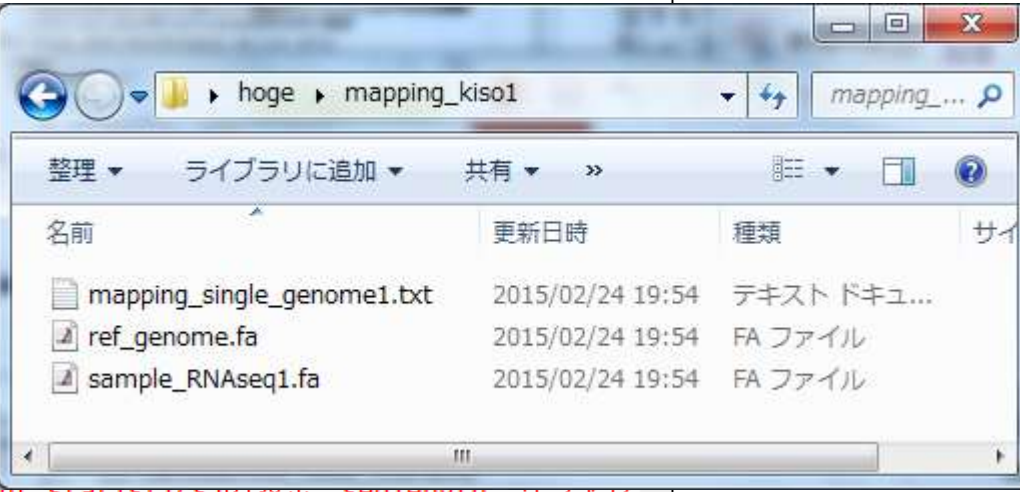
#必要なパッケージをロード
library(QuasR) #パッケージの読み込み
library(GenomicAlignments) #パッケージの読み込み

#本番(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #計算時間を計測するため
time_e <- proc.time() #計算時間を表示(一番右側)
time_e - time_s #マッピングに用いたパラ
out #マッピング結果(alignment_statistics)の表示。seqlength: ファイル

#ファイルに保存(QCレポート用のpdfファイル作成)
out_f <- sub(".bam", "_QC.pdf", out@alignments[,1]) #QCレポート結果
qQCReport(out, pdfFilename=out_f) #ファイル名を表
out_f

#ファイルに保存(BED形式ファイル)

```



```

R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "ref_genome.fa"
[3] "sample_RNAseq1.fa"
> |

```

出力ファイル形式

出力ファイルとして実際に取り扱うのは①BAM形式ファイルです。②赤枠部分はランダムに発生させる文字列部分なので、ヒトによって異なります

1. サンプルデータ18,19のRNA-seqデータ(sample RNAseq1.fa)のref_genome.faへのマッピングの場合

オプションを"-m 1 --best --strata -v 0"とした例です。sample_RNAseq1.faでマップされないのは計3リードです。2リード("chr3_11_45"と"chr3_15_49")はchr5にもマップされるので、"-m 1"オプションで落とされます。1リード("chr5_1_35")は該当箇所と完全一致ではない(4番目の塩基にミスマッチをいれている)ので落とされます。

```
in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "ref_genome.fa"
param_mapping <- "-m 1 --best --st
```

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

#本番(マッピング)
time_s <- proc.time()
out <- qAlign(in_f1, in_f2, alignm
time_e <- proc.time()
time_e - time_s
out
alignmentStats(out)

#ファイルに保存(QCレポート用のpdfフ
out_f <- sub("%.bam", "_QC.pdf", ou
qQCReport(out, pdfFilename=out_f)
out_f

#ファイルに保存(BED形式ファイル)
```

```
R Console
+ out_f <- sub("%.bam", ".bed", tmpfname[i])#BED形式ファイル名$
+ out_f #ファイル名を表示して$
+ write.table(tmp, out_f, sep="\t", append=F, quote=F, row.na$
+ }
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_21744443273a.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_217479832d2f.bam"
[9] "sample_RNAseq1_217479832d2f.bam.bai"
[10] "sample_RNAseq1_217479832d2f.bam.txt"
[11] "sample_RNAseq1_217479832d2f.bed"
[12] "sample_RNAseq1_217479832d2f_QC.pdf"
> |
```

出力ファイル形式

- ゲノム上のどの位置にどのリードがマッピングされたか(トランスクリプトームの場合どの転写物配列上のどの位置にどのリードがマッピングされたか)を表す出力ファイル形式は複数あります。
 - SAM (Sequence Alignment/Map) format
 - SAMtools (Li et al., *Bioinformatics*, **25**: 2078-2079, 2009)
 - **BAM** (Binary Alignment/Map) format
 - SAMtools (Li et al., *Bioinformatics*, **25**: 2078-2079, 2009)
 - **BED** (Browser Extensible Data) format
 - BEDtools (Quinlan et al., *Bioinformatics*, **26**: 841-842, 2010)
 - ...

出力ファイル形式

BAM形式ファイル

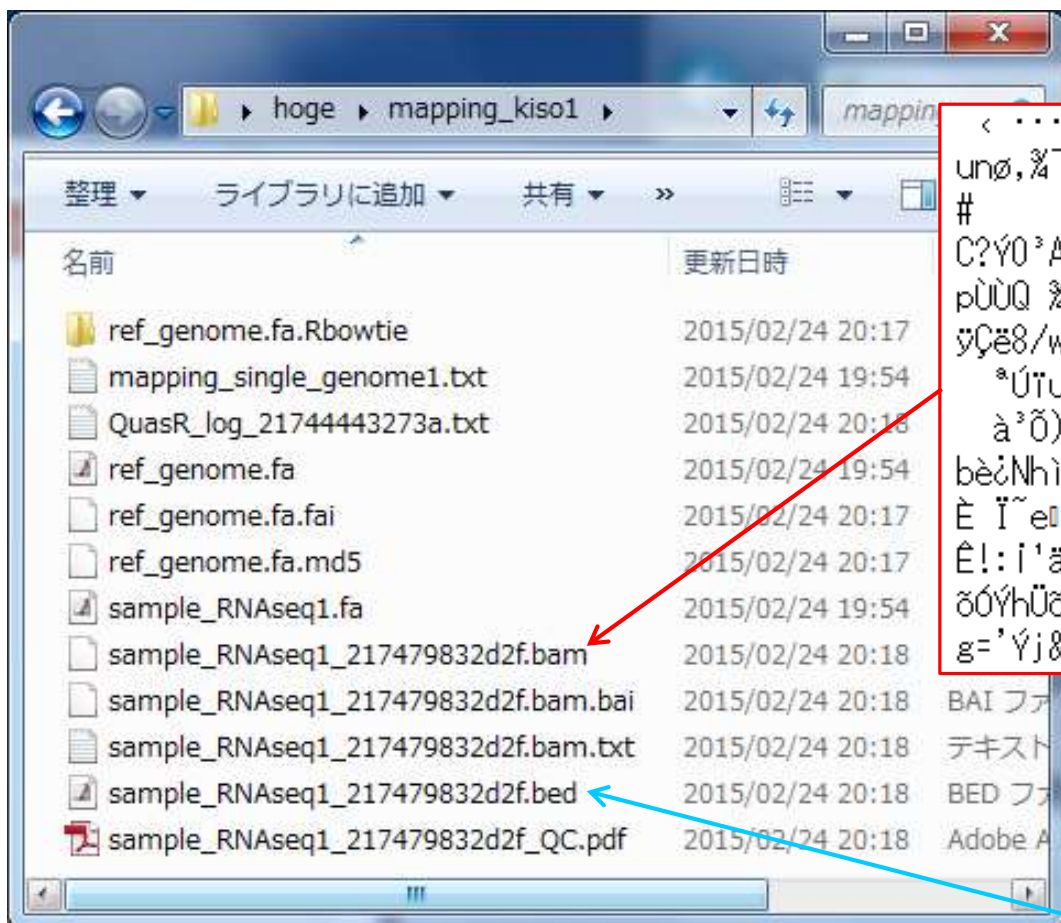
```

< .....ÿ ·BC ·P ]NNAO É #á î·µë0!WAFt:nø7MEØ: °
unø,%_g $CM%hpbB9 [GÓIW á<= Hf4 çJ7E,yw =hQ@u5? øUÿ Y /Ñ9
#
C?Y0³A0wø -y. 7'á § )ô Oú eóW50Æ@lsgÍ|:ùÀ,dÍ^ úÚ ±
pÙÚQ %çðif 5QóáøVò¥<7+-V% 97²Ø8¥í&xEE÷
ÿÇè8/w| ðì`u ôsiø/y-ô²¥¯ A=z øRÖü ðftøéc%Hà {äyÉBânSÜ²?à ù
*Úïu/··uø ZciÖ ·¶Æ ÿ- 4oË a ·· < .....ÿ ·BC · ò=oÃ
à³ô) R cI<0ÆuHÖ- ¶o ;W] òiký
bè¿Nhì( c;‡ ò=z9jhj ö÷F&p.Âbµøü ·ýøSif¥` $ò& $×
È Ì~e }ä%Tø)x¯Jø ]&ü->ôd
É!:i'ä8·ñhÝç³OQ °6jyP-³%Pî(a»øÆýQ³STDøöíí é §
øÓYhÜøðø! øë !· I _éBó^óÇ!bUFÄ eV¯p¥P6(Yp ¶
g='Ýj&W èÆ0-<0 ù d,·· < .....ÿ ·BC · · .....

```

BED形式ファイル

chr1	11	45
chr2	1	35
chr2	16	50
chr3	1	35
chr3	3	37



マップされなかったのは、
赤枠の8リード中3リード

オプションと結果の関係

- “-m 1 --best --strata -v 0”: 0 mismatches with 1 read in 1 place only. Output reads mapped to 1 place only with 0 mismatches.

```
ref_genome.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGA chr1 11 45 TGGG
>chr2
AGGGAGGGGGTCCAGTATC chr2 1 35
ACGCAGGTAGGCTGAGGAT chr2 16 50 GAGGAG
CTCGGGTATGGTTAGTCTT chr3 1 35 GGGCTG
TGACGCCCTG chr3 3 37 GGTTC
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTTAACTAGTTT
```

```
sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
>chr3_3_37
GGGACTATTTCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGCGGTCTATTTCCCGCTTGCAGGAATCGTGTC
```

オプションと結果の関係

完全一致でも複数個所にマップされるために落とされたのは2リード。**-m 1**: 1ヶ所にマップされるリードのみ出力

- “**-m 1 --best --strata -v 0**”: 0ミスマッチで**1か所にのみ**マップされるリードを出力

```
ref_genome.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG
>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACGCCCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATAACAGGC
AGCATCTAGTCGCATCAGAAGGGGTGTAGTCAGCCTATAGTTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATAACAGGC
AGCATCTAGTCGCATCAGAAGGGGTGTAGTCAGCCTATAGTTAACTAGTTT
```

```
sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATAACAGG
>chr3_3_37
GGGGACTATTTCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGCGGTCTATTTCCCGCTTGCAGGAATCGTGTC
```

オプションと結果の関係

1か所にのみマップされるが mismatches のため落とされたのは1リード。-v 0: 許容する mismatches 数は0

- “-m 1 --best --strata -v 0”: 0 mismatches で1か所にのみマップされるリードを出力

```
ref_genome.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG
>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACGCCCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTTAACTAGTTT
```

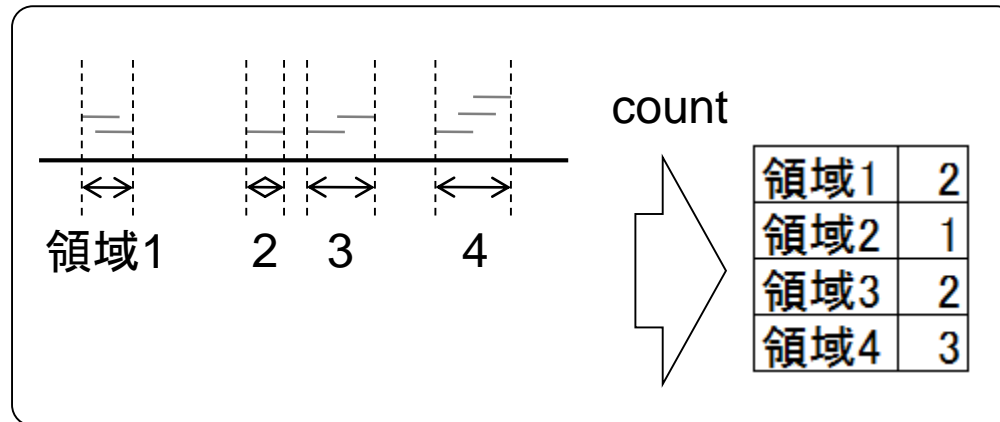
```
sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
>chr3_3_37
GGGACTATTTCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTC
```

Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

カウント情報取得

- アノテーション情報を利用する場合
 - UCSC known Genes, Ensembl Genesなど様々なテーブル名を指定可能
 - gene, exon, promoter, junctionなど様々なレベルを指定可能
- アノテーション情報がない場合
 - マップされたリードの和集合領域を同定したのち、領域ごとのリード数をカウント
 - BEDtools (Quinlan et al., 2010)中のmergeBedプログラムを実行して和集合領域同定後、intersectBedプログラムを実行してリード数をカウントする作業に相当



カウント情報取得

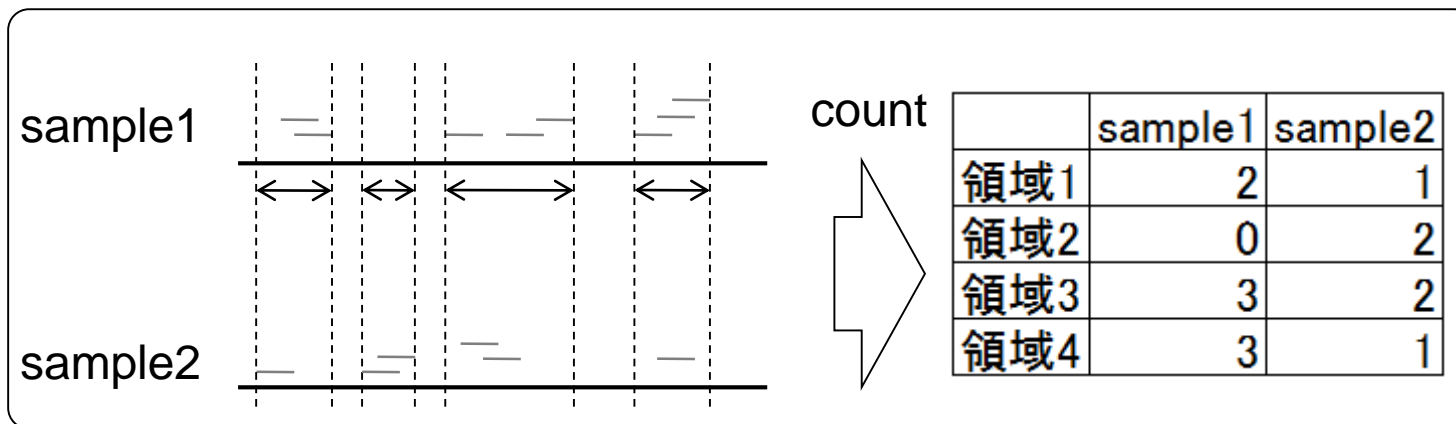
アノテーション情報がない場合の戦略は、複数サンプルの場合には領域が変わりうる。Cufflinksを知っているヒトはcuffmergeと同じイメージだと思えばよい

■ アノテーション情報を利用する場合

- UCSC known Genes, Ensembl Genesなど様々なテーブル名を指定可能
- gene, exon, promoter, junctionなど様々なレベルを指定可能

■ アノテーション情報がない場合

- マップされたリードの和集合領域を同定したのち、領域ごとのリード数をカウント
- BEDtools (Quinlan et al., 2010)中のmergeBedプログラムを実行して和集合領域同定後、intersectBedプログラムを実行してリード数をカウントする作業に相当



カウント情報取得1

- マッピング後 | 出力ファイルの読み込み | [Bowtie形式](#) (last modified 2013/06/18)
- マッピング後 | 出力ファイルの読み込み | [SOAP形式](#) (last modified 2013/06/19)
- マッピング後 | 出力ファイルの読み込み | [htSeqTools\(Planet 2012\)](#) (last modified 2013/06/19)
- マッピング後 | [カウント情報取得](#) | [|](#)について (last modified 2014/12/17)
- マッピング後 | カウント情報取得 | ゲノム | アノテーション有 | [QuasR\(Gaidatzis 2014\)](#) (last modified 2015/02/24)
- マッピング後 | カウント情報取得 | ゲノム | アノテーション無 | [QuasR\(Gaidatzis 2014\)](#) (last modified 2014/06/22)
- マッピング後 | カウント情報取得 | トランスクリプトーム | [BEDファイルから](#) (last modified 2014/06/21)
- マッピング後 | [配列長とカウント数の関係](#) (last modified 2014/07/03)
- [正規化](#) | について



マッピング後 | カウント情報取得 | ゲノム | アノテーション無 | QuasR(Gaidatzis_2014)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マッピングされたリードの和集合領域(union range)を得たのち、領域ごとにマッピングされたリード数をカウントしています。

「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. サンプルデータ18,19のRNA-seqデータ(sample RNAseq1.fa)のref_genome.faへのマッピングの場合(mapping_single_genome1.txt):

オプションを"-m 1 --best --strata -v 0"とした例です。

```

in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "ref_genome.fa"             #入力ファイル名を指定してin_f2に格納(リファレンス配列)
param_mapping <- "-m 1 --best --strata -v 0" #マッピング時のオプションを指定

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

#前処理(マッピング)
time_s <- proc.time()                 #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッピングを行うqAlign関数を実行した結果をoutに
time_e <- proc.time()                 #計算時間を計測するため
time_e - time_s                       #計算時間を表示(一番右側の数字。単位はsecond)
out                                    #マッピングに用いたパラメータや入力ファイルの情報などを表示
    
```

カウント情報取得1

1. サンプルデータ18,19のRNA-seqデータ(sample_RNAseq1.fa)のref_genome.faへのマッピングのオプションを"-m 1 --best --strata -v 0"とした例です。

```

in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納
param_mapping <- "-m 1 --best --strata -v 0" #マッピング時のオプションを指定

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

#前処理(マッピング)
time_s <- proc.time() #計算時間を取得
out <- qAlign(in_f1, in_f2, alignmentParameters=pa) #計算時間を取得
time_e <- proc.time() #計算時間を取得
time_e - time_s #計算時間を取得
out #マッピング結果
alignmentStats(out) #マッピング結果

#本番(マップされたリードの和集合領域同定)
tmpfname <- out@alignments[,1] #ファイル名
tmpsname <- out@alignments[,2] #サンプル名
for(i in 1:length(tmpfname)){ #サンプル数
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式フ
  }
}
    
```

「hoge - mapping_kiso1」に入力ファイルは揃っています。マッピング結果ファイルが既に存在しますが、logファイルを見て同じオプションで実行した結果があるかどうかを自動判定します。再度マッピングを実行する必要がない場合は、以前に得られていたbamファイルを入力として、速やかに①のところの作業に移行します。

```

R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_21744443273a.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_217479832d2f.bam"
[9] "sample_RNAseq1_217479832d2f.bam.bai"
[10] "sample_RNAseq1_217479832d2f.bam.txt"
[11] "sample_RNAseq1_217479832d2f.bed"
[12] "sample_RNAseq1_217479832d2f_QC.pdf"
> |
    
```

①

カウント情報取得1

1. サンプルデータ18,19のRNA-seqデータ(sample_RNAseq1.fa)のref_genome.faへのマッピングの場
(mapping_single_genome1.txt):

オプションを"-m 1 -best --strata -v 0"とした例です。

```
for(i in 1:length(tmpfname)){ #サンプル数(ファイル数)分だけループを
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式フ
  } else{
    k <- c(k, readGAlignments(tmpfname[i]))#BAM形
  }
}
m <- reduce(granges(k)) #GRangesオ

#本番(カウント情報取得)
tmp <- as.data.frame(m) #出力ファイ
for(i in 1:length(tmpfname)){ #サンプル数
  tmpcount <- summarizeOverlaps(m, tmpfname[i])#C
  count <- assays(tmpcount)$counts #Summarize
  colnames(count) <- tmpfname[i] #行列count
  tmp <- cbind(tmp, count) #保存した!
}

#ファイルに保存
out_f <- sub(".bam", "_range.txt", tmpfname[i])#
write.table(tmp, out_f, sep="\t", append=F, quote
```

出力ファイルは何も指定していませんが、*_range.txtという名前のファイルが作成されます。これは、.bamという名前のファイルを内部的に入力として読み込み、その文字列中の.bamを_range.txtに置換したものを出力ファイル名として自動作成しているからそうなります。

```
R Console
> #ファイルに保存
> out_f <- sub(".bam", "_range.txt", tmpfname[i])#$
> write.table(tmp, out_f, sep="\t", append=F, quote$
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_kis01"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_21744443273a.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_217479832d2f.bam"
[9] "sample_RNAseq1_217479832d2f.bam.bai"
[10] "sample_RNAseq1_217479832d2f.bam.txt"
[11] "sample_RNAseq1_217479832d2f.bed"
[12] "sample_RNAseq1_217479832d2f_QC.pdf"
[13] "sample_RNAseq1_217479832d2f_range.txt"
> |
```

カウント情報取得1

.bedファイルと*_range.txtファイルを見比べると理解が深まるでしょう。*_range.txtファイルの一番右側の列がカウント情報です。

1. サンプルデータ18,19のRNA-seqデータ(sample_RNAseq1.fa)のref_genome.faへのマッピング

オプションを"-m 1 --best --strata -v 0"とした例です。

```

in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に
param_mapping <- "-m 1 --best --strata -v 0" #マッピング時のオプションを

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

#前処理(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッピング
time_e <- proc.time() #計算時間を計測するため
time_e - time_s #計算時間を表示(一番右側の数字。)
out #マッピングに用いたパラメータや入
    
```

*.bed

chr1	11	45
chr2	1	35
chr2	16	50
chr3	1	35
chr3	3	37

	A	B
1	FileName	SampleName
2	sample_RNAseq1.fa	naeae

*_range.txt

seqnames	start	end	width	strand	naeae
chr1	11	45	35	+	1
chr2	1	50	50	+	2
chr3	1	37	37	+	2

複数のRNA-seqリードファイルのマッピングからカウントデータ取得までを一気に行う例です。

カウント情報取得2

- マップ後 | 出力ファイルの読み込み | htSeqTools(Planet 2012) (last modified 2013/06/19)
- マップ後 | カウント情報取得 | について (last modified 2014/12/17)
- マップ後 | カウント情報取得 | ゲノム | アノテーション有 | QuasR(Gaidatzis 2014) (last modified 2015/02/24)
- マップ後 | カウント情報取得 | ゲノム | アノテーション無 | QuasR(Gaidatzis 2014) (last modified 2014/06/22)
- マップ後 | カウント情報取得 | トランスクリプトーム | BEDファイルから (last modified 2014/06/21)



マップ後 | カウント情報取得 | ゲノム | アノテーション無 | QuasR(Gaidatzis 2014)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カ
連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用し
合領域(union)領域を得たものを、領域ごとにもマッピングされたリード数をカウントしています



5. サンプルデータ18-20の複数のRNA-seqデータ(sample RNAseq1.faとsample RNAseq2.fa)を ref_genome.faにマッピングする場合(mapping_single_genome4.txt):

1. サンプルデータ

オプションを"-m

全部のマッピング結果をまとめて和集合領域を定め、カウント情報を得るやり方です。一般的なカウントデー
タ行列の形式(2列目以降がカウント情報)にし、配列長情報と別々のファイルにして保存するやり方です。

```
in_f1 <- "ma
in_f2 <- "re
param_mappin

#必要なパッケ
library(Quas
library(Genc

#前処理(マッ
time_s <- pr
out <- qAlign
time_e <- pr
time_e - tin
out
```

```
in_f1 <- "mapping_single_genome4.txt" #入力ファイル名を指定してin_f1に格納(RNA-seq
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納(リファ
out_f1 <- "hoge4_count.txt" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge4_genelength.txt" #出力ファイル名を指定してout_f2に格納
param_mapping <- "-m 1 --best --strata -v 1" #マッピング時のオプションを指定

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

#前処理(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッピングを行うqAl
time_e <- proc.time() #計算時間を計測するため
time_e - time_s #計算時間を表示(一番右側の数字。単位はsecond
out #マッピングに用いたパラメータや入力ファイルの
alignmentStats(out) #マッピング結果(alignment statistics)の表示
```

「hoge - mapping_kiso2」に入カファイルは揃っています。

カウント情報取得2

5. サンプルデータ18-20の複数のRNA-seqデータ(sample_RNAseq1.faとsample_RNAseq2.fa)をref_genome.faにマッピングする場合(mapping_single_genome4.txt):

全部のマッピング結果をまとめて和集合領域を定め、カウント情報を得るやり方です。一般的なカウントデータ行列の形式(2列目以降がカウント情報)にし、配列長情報と別々のファイルにして保存するやり方です。

```

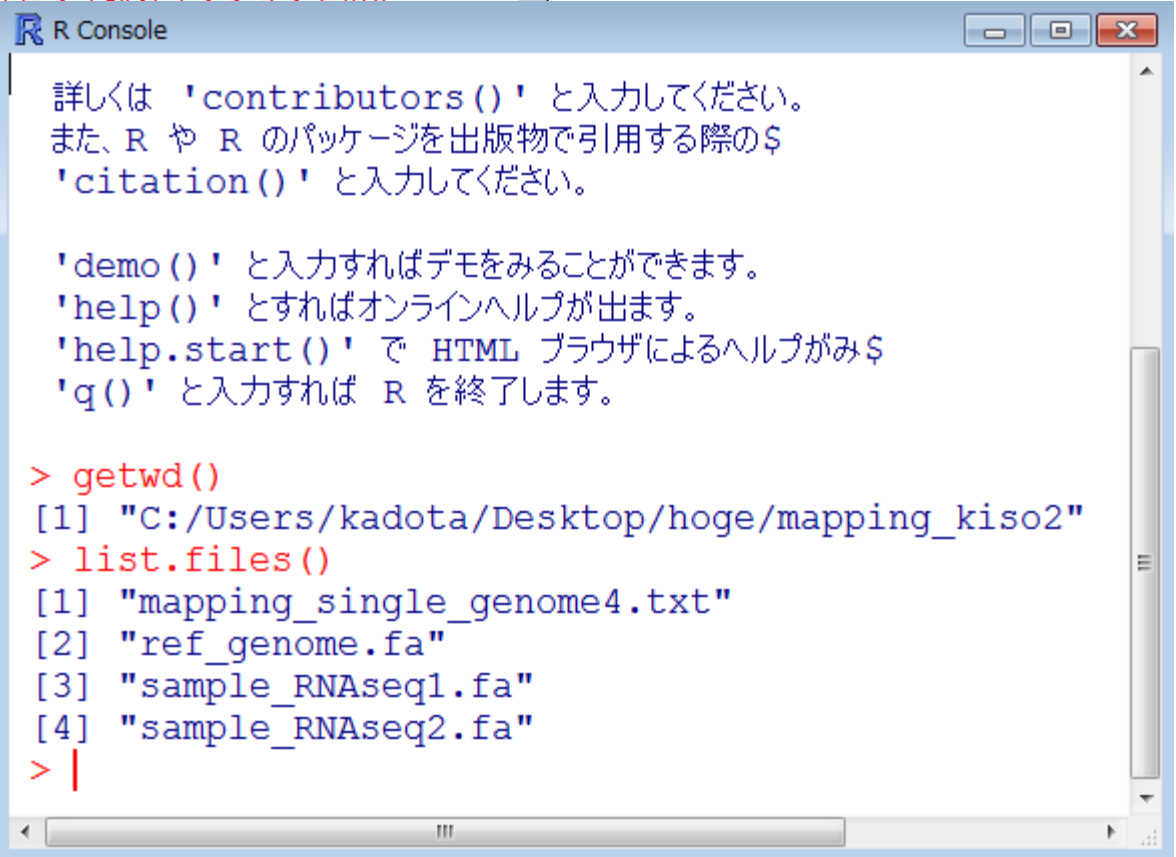
in_f1 <- "mapping_single_genome4.txt" #入力ファイル名を指定してin_f1に格納(RNA-seq)
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納(リファレンス)
out_f1 <- "hoge4_count.txt" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge4_geneLength.txt" #出力ファイル名を指定してout_f2に格納
param_mapping <- "-m 1 --best --strata -v 1" #マッピング時のオプションを指定

#必要なパッケージをロード
library(QuasR) #パッケージ
library(GenomicAlignments) #パッケージ

#前処理(マッピング)
time_s <- proc.time() #計算時間
out <- qAlign(in_f1, in_f2, alignmentParameter=...) #計算時間
time_e <- proc.time() #計算時間
time_e - time_s #計算時間
out #マッピング結果
alignmentStats(out) #マッピング結果

#本番(マップされたリードの和集合領域同定)
tmpfname <- out@alignments[,1] #ファイル名
tmpsname <- out@alignments[,2] #サンプル名
for(i in 1:length(tmpfname)){ #サンプル名

```



カウント情報取得2

コピー。出力ファイル群のうち、主に扱うのは、カウントデータを含むファイル(hoge4_count.txt)。

5. サンプルデータ18-20の複数のRNA-seqデータ(sample_RNAseq1.faとsample_RNAseq2.fa)をref_genome.faにマッピングする場合(mapping_single_genome4.txt):

全部のマッピング結果をまとめて和集合領域を定め、カウント情報を得るやり方です。一般的なカウントデータ行列の形式(2列目以降がカウント情報)にし、配列長情報と別々のファイルにして保存するやり方です。

```

in_f1 <- "mapping_single_genome4.txt" #入力ファイル
in_f2 <- "ref_genome.fa" #入力ファイル
out_f1 <- "hoge4_count.txt" #出力ファイル
out_f2 <- "hoge4_genelength.txt" #出力ファイル
param_mapping <- "-m 1 --best --strata -v 1" #マッピングパラメータ

#必要なパッケージをロード
library(QuasR) #パッケージ
library(GenomicAlignments) #パッケージ

#前処理(マッピング)
time_s <- proc.time() #計算時間
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping)
time_e <- proc.time() #計算時間
time_e - time_s #計算時間
out #マッピング結果
alignmentStats(out) #マッピング結果

#本番(マップされたリードの和集合領域同定)
tmpfname <- out@alignments[,1] #ファイル名
tmpsname <- out@alignments[,2] #サンプル名
for(i in 1:length(tmpfname)){ #サンプルごとに処理

```

```

R Console
> tmp <- cbind(tmp, h$width) #和集合$
> write.table(tmp, out_f2, sep="\t", append=F, q$)
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_kiso2"
> list.files()
[1] "hoge4_count.txt"
[2] "hoge4_genelength.txt"
[3] "mapping_single_genome4.txt"
[4] "QuasR_log_201814983c7.txt"
[5] "ref_genome.fa"
[6] "ref_genome.fa.fai"
[7] "ref_genome.fa.md5"
[8] "ref_genome.fa.Rbowtie"
[9] "sample_RNAseq1.fa"
[10] "sample_RNAseq1_2018acaf46.bam"
[11] "sample_RNAseq1_2018acaf46.bam.bai"
[12] "sample_RNAseq1_2018acaf46.bam.txt"
[13] "sample_RNAseq2.fa"
[14] "sample_RNAseq2_201843d33cbe.bam"
[15] "sample_RNAseq2_201843d33cbe.bam.bai"
[16] "sample_RNAseq2_201843d33cbe.bam.txt"
> |

```

リストファイル中で指定したサンプル名(sample1とsample2)がカウントデータ行列の列名となります

カウント情報取得2

5. サンプルデータ18-20の複数のRNA-seqデータ(sample_RNAseq1.faとsample_RNAseq2.ref_genome.fa)にマッピングする場合(mapping_single_genome4.txt):

全部のマッピング結果をまとめて和集合領域を定め、カウント情報を得るやり方です。一タ行列の形式(2列目以降がカウント情報)にし、配列長情報と別々のファイルにして保存

```

in_f1 <- "mapping_single_genome4.txt" #入力ファイル名を指定してin_f
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f
out_f1 <- "hoge4_count.txt" #出力ファイル名を指定してout_
out_f2 <- "hoge4_genelength.txt" #出力ファイル名を指定してout_
param_mapping <- "-m 1 --best --strata -v 1" #マッピング時のオプション

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

#前処理(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッ
    
```

FileName	SampleName
sample_RNAseq1.fa	sample1
sample_RNAseq2.fa	sample2

tmp	sample1	sample2
chr1_11_45_35_+	1	0
chr2_1_60_60_+	2	1
chr3_1_37_37_+	2	0
chr4_6_65_60_+	0	1
chr5_1_35_35_+	1	0

Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

アノテーション情報がない場合のやり方です

実データ解析

- マッピング | [基礎](#) (last modified 2013/06/19)
- マッピング | [single-end](#) | ゲノム | basic aligner(基礎) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- マッピング | [single-end](#) | ゲノム | basic aligner(応用) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/28)
- マッピング | [single-end](#) | ゲノム | splice-aware aligner | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- マッピング | [paired-end](#) | ゲノム | basic aligner(応用) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/28)
- マップ後 | [|について](#) (last modified 2013/06/19)
- マップ後 | [出力ファイル形式について](#) (last modified 2013/06/19)
- マップ後 | 出力ファイルの読み込み | [BAM形式](#)
- マップ後 | 出力ファイルの読み込み | [Bowtie形式](#)
- マップ後 | 出力ファイルの読み込み | [SOAP形式](#)
- マップ後 | 出力ファイルの読み込み | [htSeqTools](#)
- マップ後 | [カウント情報取得|について](#) (last modified 2013/06/19)
- マップ後 | [カウント情報取得](#) | ゲノム | [アノテーション](#)
- マップ後 | [カウント情報取得](#) | ゲノム | [アノテーション](#)
- マップ後 | [カウント情報取得](#) | [トランスクリプトーム](#)



マッピング | paired-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis 2015) NEW

QuasRパッケージを用いてpaired-end RNA-seqデータのリファレンスゲノム配列へのマッピングを行うやり方を示します。basic alignerの1つであるBowtie (Langmead et al., Genome Biol., 2009)を実装したRbowtieパッケージを内部的に使っています。mapping paired genome1.txtのような2行目の1列目と2列目に「マッピングしたいRNA-seqファイル名1と2」(例: sample_RNAseq_1.faとsample_RNAseq_2.fa)、そして2行目の3列目に「任意のサンプル名」(例: namae)を記載したタブ区切りテキストファイルを用意した上で行います。1行目の文字列は変えてはいけません(つまり"FileName1", "FileName2"および"SampleName"のままにしておくということです)

「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. mapping paired genome1.txt中のFASTQ形式ファイルを乳酸菌ゲノムにマッピングする場合:

乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。SRR616268sub_1.fastq.gzは、約75MB、全リード107 bpです。SRR616268sub_2.fastq.gzは、約67MB、全リード93 bpです。Ensembl (Flicek et al., 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)がリファレンス配列です。マッピングオプションはデフォルトです。

```
in_f1 <- "mapping paired genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus casei 12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa" #入出力ファイル名を指定してin_f2に格納(ゲノム配列ファイル)
#必要なパッケージをロード
```

	A	B	C
1	FileName1	FileName2	SampleName
2	SRR616268sub_1.fastq.gz	SRR616268sub_2.fastq.gz	namae Paired

```
time_s <- proc.time()
out <- qAlign(in_f1, in_f2)
time_e <- proc.time()
#計算時間を計測するための関数
#マッピングを行うqAlign関数を実行した結果をoutに格納
#計算時間を計測するための関数
```


実データ解析

「hoge - mapping_paired1」に解析したいfastq.gzファイルを移動させて実行。約8分。

1. mapping paired genomel.txt中のFASTQ形式ファイルを乳酸菌ゲノムにマッピングする場合:

乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。SRR616268sub_1.fastq.gzは、約75MB、全リード107 bpです。SRR616268sub_2.fastq.gzは、約67MB、全リード93 bpです。Ensembl (Flicek et al., 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)がリファレンス配列です。マッピングオプションはデフォルトです。

```
in_f1 <- "mapping_paired_genomel.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa" #入:

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

#本番(マッピング)
time_s <- proc.time() #計算時間を計測するため
out <- qAlign(in_f1, in_f2) #マッピングを行うqAlign関数を実行した結果をoutに格納
time_e <- proc.time() #計算時間を計測するため
time_e - time_s #計算時間を表示(一番右側の数字。単位はsecond)
out #マッピングに用いたパラメータや入力ファイルの情報など
alignmentStats(out) #マッピング結果(alignment statistics)の表示。seq
```

#ファイルに保存
out_f <- sub("out", "out_f")
qQCReport(out, out_f)
#ファイルに保存
tmpfname <- out

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_paired1"
> list.files()
[1] "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"
[2] "mapping_paired_genomel.txt"
[3] "SRR616268sub_1.fastq.gz"
[4] "SRR616268sub_2.fastq.gz"
> |
```

実データ解析

①このあたりで約1分。リファレンス配列のインデックス作成部分。約3MBの乳酸菌ゲノムなので短時間で終わる。ヒトゲノムだと数十分から数時間といったところか。②のところで約7分。オリジナルの総リード数は1.4億。1/100以下の100万リードだから10分足らずで終わる。

1. [mapping paired genome1.txt](#)中のFASTQ形式ファイルを乳酸菌ゲノムにマッピング
 乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。[SRR616268sub_1.fastq.gz](#)
 107 bpです。[SRR616268sub_2.fastq.gz](#)は、約67MB、全リード93 bpです。[Ensembl \(Flic](#)
 ている [Lactobacillus casei 12A](#)の multi-FASTA形式ゲノム配列ファイル
[\(Lactobacillus casei 12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa\)](#)がリフ
 グオプションはデフォルトです。

```
in_f1 <- "mapping_paired_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイ
in_f2 <- "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"#入:
```

#必要なパッケージ
 library(QuasR)
 library(Genomi
 #本番(マッピング)
 time_s <- proc
 out <- qAlign(
 time_e <- proc
 time_e - time_
 out
 alignmentStats
 #ファイルに保存
 out_f <- sub(""
 qQCReport(out,
 out_f
 #ファイルに保存
 tmpfname <- ou

```
R Console
> #本番(マッピング)
> time_s <- proc.time() #計算時間を計測するため
> out <- qAlign(in_f1, in_f2) #マッピングを行うqAlign関数を実行し$
Creating .fai file for: C:/Users/kadota/Desktop/hoge/mapping_paired1/Lactobac$
alignment files missing - need to:
  create alignment index for the genome
  create 1 genomic alignment(s)
will start in ..9s..8s..7s..6s..5s..4s..3s..2s..1s
Creating an Rbowtie index for C:/Users/kadota/Desktop/hoge/mapping_paired1/La$
Finished creating index
Testing the compute nodes...OK
Loading QuasR on the compute nodes...OK
Available cores:
nodeNames
KADOTA-PC
      1
Performing genomic alignments for 1 samples. See progress in the log file:
```



実データ解析

①マッピングの実体であるqAlign関数
実行前後の時間を計測し、423.14/60
= 7.05分という結果を得ている。

```
#本番(マッピング)
time_s <- proc.time()
out <- qAlign(in_f1, in_f2)
time_e <- proc.time()
time_e - time_s
out
alignmentStats(out)
```

#計算時間を計測するため
#マッピングを行うqAlign関数を実行した結果をoutに格
#計算時間を計測するため
#計算時間を表示(一番右側の数字。単位はsecond)
#マッピングに用いたパラメータや入力ファイルの情報な
#マッピング結果(alignment statistics)の表示。seq

```
#ファイルに保存(QCレポート用のpdfファイル作成)
out_f <- sub(".bam" ". QC.pdf" out@alignments[ 1])#Quality Controlレポートのpdfファイル
qQCReport(out, p
out_f
```

```
#ファイルに保存(B
tmpfname <- out@
for(i in 1:length
hoge <- readGA
hoge <- as.dat
tmp <- hoge[,
out_f <- sub("
out_f
write.table(tm
```

```
R Console
Performing genomic alignments for 1 samples. See progress in the log file:
C:/Users/kadota/Desktop/hoge/mapping_paired1/QuasR_log_d0470eb4872.txt
Genomic alignments have been created successfully

> time_e <- proc.time()
> time_e - time_s
 ユーザ システム 経過
          0.92          0.76          423.14
> out
Project: qProject
Options  : maxHits      : 1
          paired       : fr
          splicedAlignment: FALSE
          bisulfite    : no
          snpFile      : none
Aligner  : Rbowtie v1.8.0 (parameters: -m 1 --best --strata --maxins 500)
Genome   : .../Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Ch$
```

time_e <- proc.time()
time_e - time_s

#計算時間を計測するため
#計算時間を表示(一番右側の数字。単位\$

#マッピングに用いたパラメータや入力\$



実データ解析

①マッピングに用いられたオプション(パラメータ)情報。②これがメインのマッピング結果ファイルの名前(いわゆるBAMファイル)。③マッピング結果の概観。forward側とreverse側合わせて200万リード中8,204リードしかマップされてない!

```
#本番(マッピング)
time_s <- proc.time()
out <- qAlign(in_f1, in_f2)
time_e <- proc.time()
time_s - time_e

#ファイルに保存(Q)
out_f <- sub(".", "bam", out_f)
qQCReport(out, p, out_f)

#ファイルに保存(B)
tmpfname <- out@tmpfname
for(i in 1:length(tmpfname))
  hoge <- readGATC(tmpfname[i], out_f)
  hoge <- as.data.frame(hoge)
  tmp <- hoge[, 1:4]
  out_f <- sub("bam", "bam", out_f)
  write.table(tmp, out_f, append=TRUE, row.names=FALSE, col.names=FALSE)
```

#計算時間を計測するため
#マッピングを行うqAlign関数
#計算時間を計測するため

```
R Console
> out
Project: qProject
Options  : maxHits      : 1
          paired       : fr
          splicedAlignment: FALSE
          bisulfite     : no
          snpFile       : none
Aligner  : Rbowtie v1.8.0 (parameters: -m 1 --best --strata --maxins 500)
Genome   : .../Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.C1

Reads    : 1 pair of files, 1 sample (fastq format):
          1. SRR616268sub_1.fastq.gz  SRR616268sub_2.fastq.gz  namae_paired (phred33)

Genome alignments: directory: same as reads
                  1. SRR616268sub_1_d0466162946.bam

Aux. alignments: none

> alignmentStats(out)
          seqlength mapped unmapped
namae_paired:genome 2907892  8204 1991796
```

#マッピングに用いたパラメータや入力\$

#マッピング結果 (alignment statistics\$)

実データ解析

①list.filesを用いて、ここまでの作業で作成されたファイルを概観。元は黒枠の4つのファイルのみだったが、マッピングが終わった時点で数多くのファイルが作成されていることが分かる。②マッピング結果を取扱うときの基本はBAMファイル。③マッピングに用いられたオプション(パラメータ)情報は、このlogファイルからも閲覧可能。

```

#本番(マッピング)
time_s <- proc.time()
out <- qAlign(in_f1, in_f2)
time_e <- proc.time()
time_e - time_s
out
alignmentStats(out)

#計算時間を計測するため
#マッピングを行うqAlign関数を実行し
#計算時間を計測するため
#計算時間を表示(一番右側の数字。単位は秒)
#マッピングに用いたパラメータや入力ファイル名
#マッピング結果(alignment statistics)

#ファイルに保存(QCレポート用のpdfファイル作成)
out_f <- sub(".bam", "_QC.pdf", out@alignments[,1])#Quqlity Controlレポートのpdfファイル
qQCReport(out, pdfFilename=out_f) #QCレポート結果をファイルに保存
out_f #ファイル名を表示してるだけです

#ファイルに保存(BED形式ファイル)
tmpfname <- out@alignments[,1] #ファイル名(in_f1の1列目に相当)をtmpfnameとして取
for(i in 1:length(tmpfname)) #サンプル数(ファイル数)分だけループを回す

```

```

R Console
> list.files()
[1] "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa"
[2] "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa.fai"
[3] "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa.md5"
[4] "Lactobacillus_casei_12a.GCA_000309565.2.25.dna.chromosome.Chromosome.fa.Rbowtie"
[5] "mapping paired genomel.txt"
[6] "QuasR_log_d0470eb4872.txt"
[7] "SRR616268sub_1.fastq.gz"
[8] "SRR616268sub_1_d0466162946.bam"
[9] "SRR616268sub_1_d0466162946.bam.bai"
[10] "SRR616268sub_1_d0466162946.bam.txt"
[11] "SRR616268sub_2.fastq.gz"
> |

```

QCレポート作成

①sub関数を用いて黒下線で見られるファイル名中の拡張子部分.bamを_QC.pdfに置換したファイル名を作成してout_flに格納している。拡張子部分周辺の文字列だけを変更する程度 of 出力ファイル名であれば、このように自動的に作成することができる。数十秒程度。

```

#ファイルに保存(QCレポート用のpdfファイル作成)
out_f <- sub(".bam", "_QC.pdf", out@alignments[,1])#Quqlity Controlレポート
qQCReport(out, pdfFilename=out_f) #QCレポート結果をファイルに保存
out_f #ファイル名を表示してるだけです

#ファイルに保存(BED形式ファイル)
tmpfname <- out@alignments[,1] #ファイル名(in_flの1列目に相当)をtmpf
for(i in 1:length(tmpfname)){ #サンプル数(ファイル数)分だけループを回す
  hoge <- readGAlignments(tmpfname[i]) #BAM形式ファイルを読み込んだ結果をhogeに格納(これは
  hoge <- as.data.frame(hoge) #データフレーム形式に変換
  tmp <- hoge[, c("seqnames","start","end")]#必要な列の情報のみ抽出した結果をtmpに格納
  out_f <- sub(".bam", ".bed", tmpfname[i])#BED形式ファイル名を作成した結果をout_flに格納
  out_f #ファイル名を表示してるだけです
  write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=F)#tmpの
}

```

```

R Console
> out@alignments[,1]
[1] "C:/Users/kadota/Desktop/hoge/mapping_paired1\\SRR616268sub_1_d0466162946.bam"
① > out_f <- sub(".bam", "_QC.pdf", out@alignments[,1])#Quqlity Controlレポートのpdfフ$
> qQCReport(out, pdfFilename=out_f) #QCレポート結果をファイルに保存
collecting quality control data
creating QC plots
> out_f #ファイル名を表示してるだけです
[1] "C:/Users/kadota/Desktop/hoge/mapping_paired1\\SRR616268sub_1_d0466162946_QC.pdf"
> |

```

BEDファイル作成

```
#ファイルに保存(QCレポート用のpdfファイル作成)
out_f <- sub(".bam", "_QC.pdf", out@alignments[,1])#Quqlity Cont
qQCReport(out, pdfFilename=out_f) #QCレポート結果をファイルに
out_f #ファイル名を表示してるだけ
```

```
#ファイルに保存(BED形式ファイル)
tmpfname <- out@alignments[,1] #ファイル名(in_f1の1列目に格納)
for(i in 1:length(tmpfname)){ #サンプル数(ファイル数)分だけループを回す
  hoge <- readGAlignments(tmpfname[i]) #BAM形式ファイルを読み込んだ結果をhogeに格納(これは)
  hoge <- as.data.frame(hoge) #データフレーム形式に変換
  tmp <- hoge[, c("seqnames", "start", "end")]#必要な列の情報のみ抽出した結果をtmpに格納
  out_f <- sub(".bam", ".bed", tmpfname[i])#BED形式ファイル名を作成した結果をout_fに格納
  out_f #ファイル名を表示してるだけです
  write.table(tmp, out_f, sep="\t", append=F, quote=F, row.names=F, col.names=F)#tmpの
```

①これはBAMファイルがバイナリ形式でそんなものだと慣れてしまえば不必要。単純に視覚的に見やすいようにbamファイルを読み込んでBED形式ファイルを作成しているところ(つまりファイル形式の変換)。②得られたBED形式ファイルの中身。マップされたリード数が8,204個であったことから、BEDファイルの行数が8,204行なのは妥当

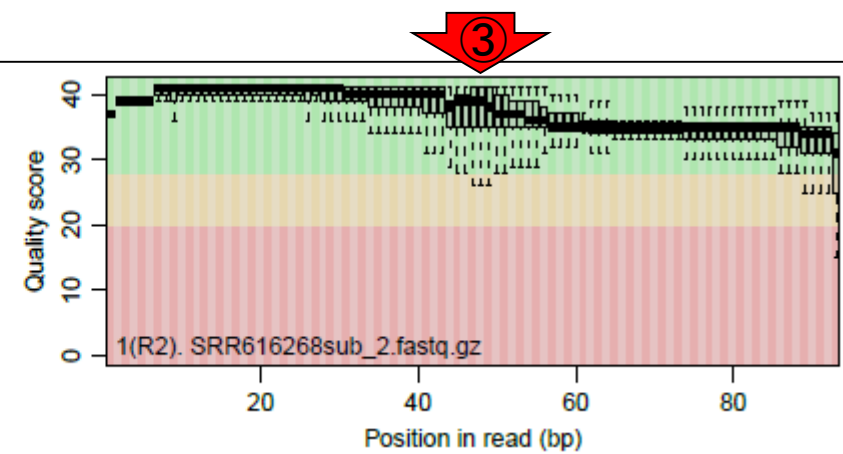
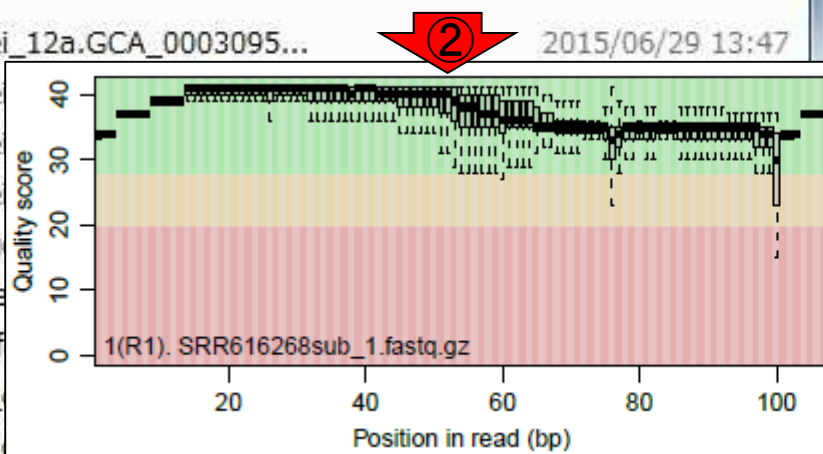


	A	B	C
1	Chromosome	592	684
2	Chromosome	634	740
3	Chromosome	971	1063
4	Chromosome	1032	1138
5	Chromosome	1105	1197
6	Chromosome	1289	1395
...
8198	Chromosome	2906800	2906912
8199	Chromosome	2906881	2906973
8200	Chromosome	2906897	2906989
8201	Chromosome	2906904	2906996
8202	Chromosome	2906919	2907011
8203	Chromosome	2907095	2907187
8204	Chromosome	2907108	2907214

QCLレポート概観

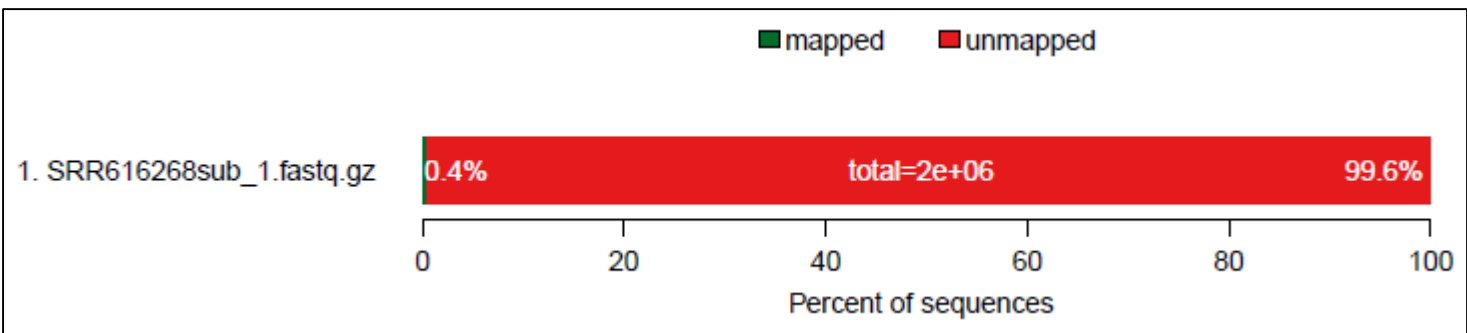
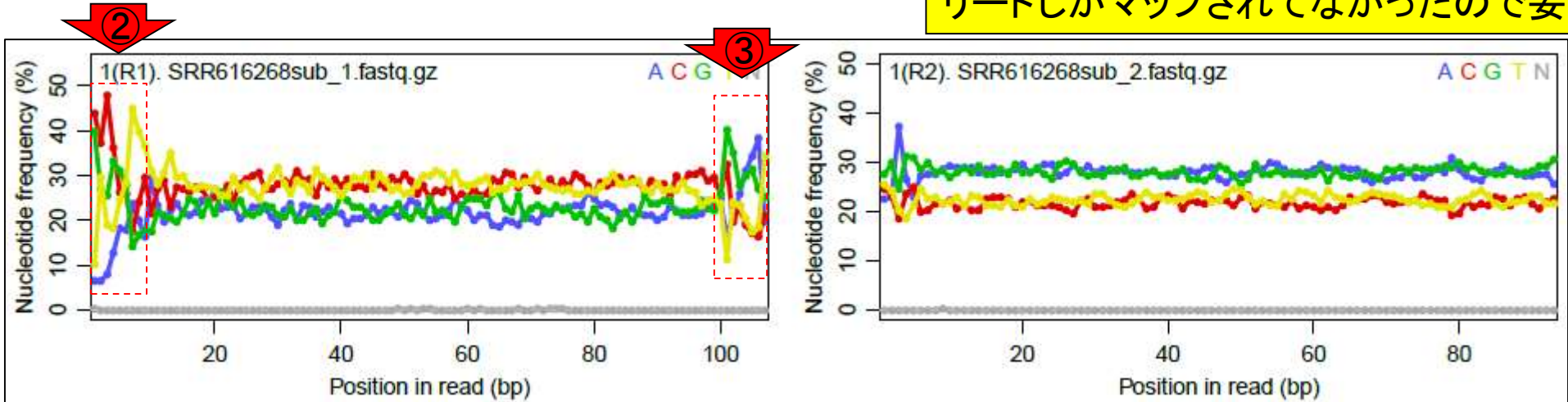
①QCLレポートを眺める。最初のページは quality score分布。②forward側と③ reverse側で左右に表示されている。カイク small RNA-seqデータ (single-end)で左側のみにしか表示されていなかったときは違和感を覚えたが、paired-endデータで眺めると至極妥当な配置であったことが分かる

名前	サイズ	更新日時
Lactobacillus_casei_12a.GCA_0003095...		2015/06/29 13:47
Lactobacillus_casei_12a.GCA_0003095...		
Lactobacillus_casei_12a.GCA_0003095...		
Lactobacillus_casei_12a.GCA_0003095...		
mapping_paired_g...		
QuasR_log_d0470e...		
SRR616268sub_1.f...		
SRR616268sub_1_...		
SRR616268sub_1_...		
SRR616268sub_1_...		
SRR616268sub_1_d0466162946.bam.txt	1 KB	2015/06/29 13:54
SRR616268sub_1_d0466162946.bed	220 KB	2015/06/29 15:02
SRR616268sub_1_d0466162946_QC.pdf	44 KB	2015/06/29 14:51
SRR616268sub_2.fastq.gz	65,585 KB	2015/05/11 16:53



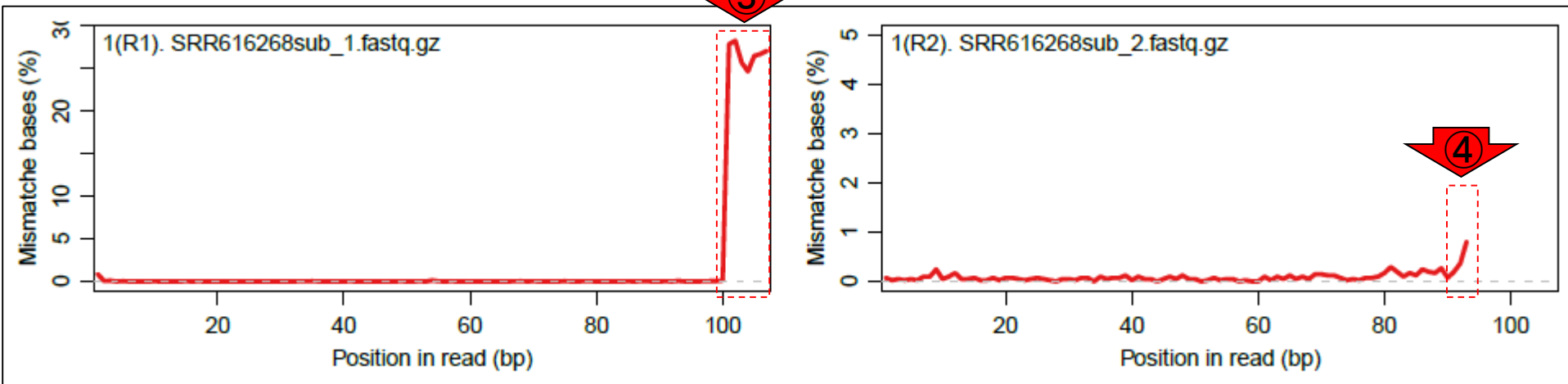
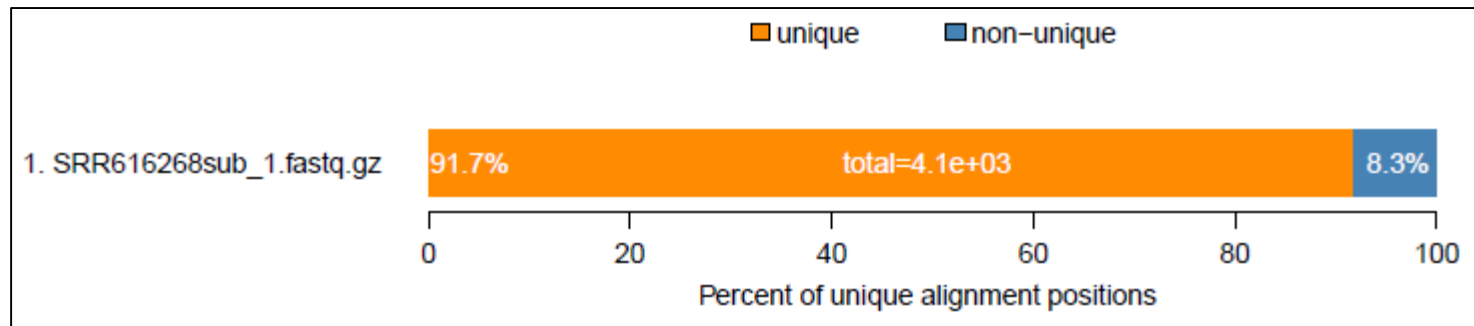
QCレポート概観

①QCレポートの2page目。②Forwardのリードの左側は確かにガタガタしていたので妥当だが、③右側もこんなになってたっけ?!と思いつつ…(後にこれが犯人と判明)。④4page目。リード全体のmapped/unmappedの割合を表示。確かに200万リード中8,204リードしかマップされてなかったのが妥当。



QCレポート概観

①5page目。200万リード中8,204リードマップされていたが、paired-endなので、ここでは分母を $8,204/2 = 4,102 \approx 4.1e+03$ と表現しているようだ。1ヶ所にのみマップされたのは全体の91.7%。②6page目。これはおそらくマップされたリードの中でどのポジションにミスマッチがあったかを示す部分。③と④を眺めて、確かにFastQCで眺めていたのは最初の50 bp分で、adapters/primersの除去は5'側のみだった!



このぶざまな状況を打破する
糸口が完璧に見つかった！

アセンブル結果再考

ファイル名	Total number of assembled transcripts	Average length	Median length	Total number of assembled bases	Total reads in file	Perfectly aligned reads
・100万リードのオリジナル?!ファイル						
SRR616268sub_1.fastq.gz	8	121	106	974	996739	1964
SRR616268sub_2.fastq.gz	424	436	228	185233	983854	710393
paired-end	0	0	0	0	990047	23
・100万リードのアダプター除去後のファイル						
SRR616268sub_trim2_1.fastq.gz	8	121	104	974	993311	1231
SRR616268sub_trim2_2.fastq.gz	425	433	228	184437	981649	708951
hoge2_1.fastq.gz	8	121	104	974	993001	1231
hoge2_2.fastq.gz	424	434	229	184210	980997	708395
paired-end	0	0	0	0	989869	25
・アダプター除去およびフィルタリング後のファイル						
SRR616268sub_trim_1.fastq.gz	8	121	104	974	993311	1231
SRR616268sub_trim_2.fastq.gz	425	433	228	184437	981649	708951
hoge1_1.fastq.gz	8	121	104	974	993081	1231
hoge1_2.fastq.gz	423	435	229	184211	981125	708482
paired-end	0	0	0	0	989869	25

Contents

- 先週の課題やエラーについて
 - 課題1の解説
 - アダプター配列除去(QuasR)実行時のエラーは門田のミス
 - Paired-endデータの取り扱い
- フィルタリング(filtering)
 - 基本形:リード数が同じpaired-endの場合(QuasRパッケージ)
 - 発展形:リード数が異なる場合(ShortReadパッケージ)
- アセンブル(assembly)
 - ゲノム用とトランスクリプトーム用
 - Rockhopper2(バクテリアのトランスクリプトーム用)
- マッピング(mapping)
 - 基礎、オプション、仮想データ説明
 - マッピング本番、出力ファイル形式、オプションと結果の関係
 - カウント情報取得
- 実データ解析
 - QuasRでマッピング → トリミング → 高精度なアセンブルとマッピングへ

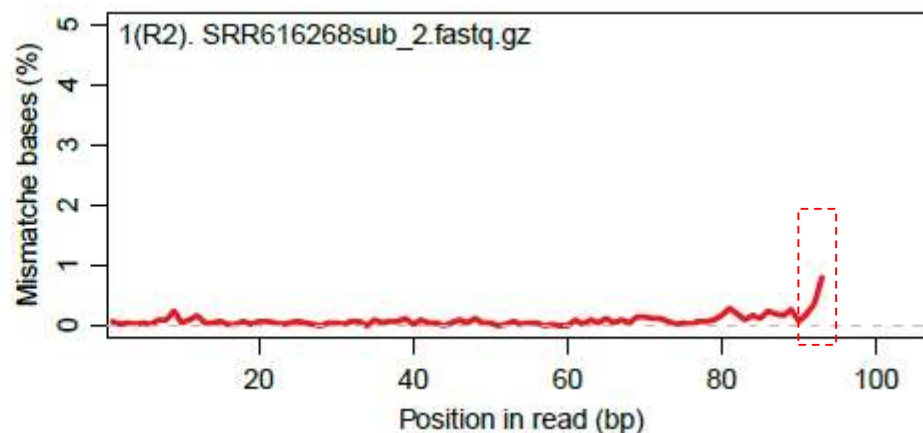
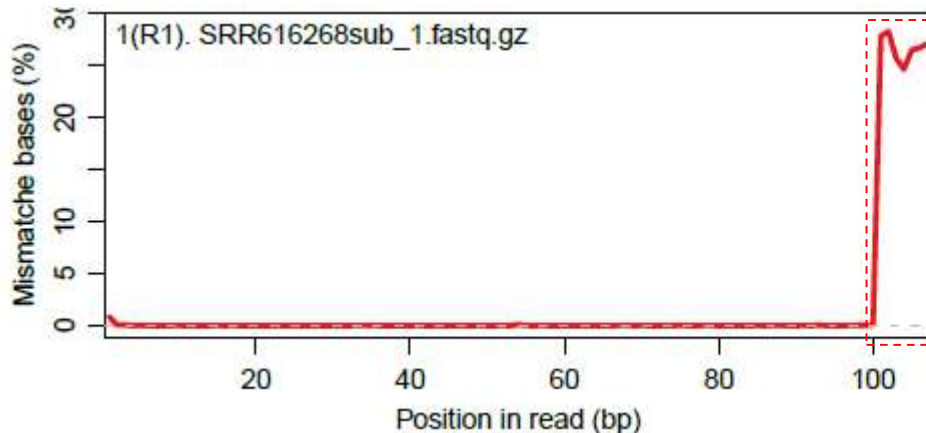
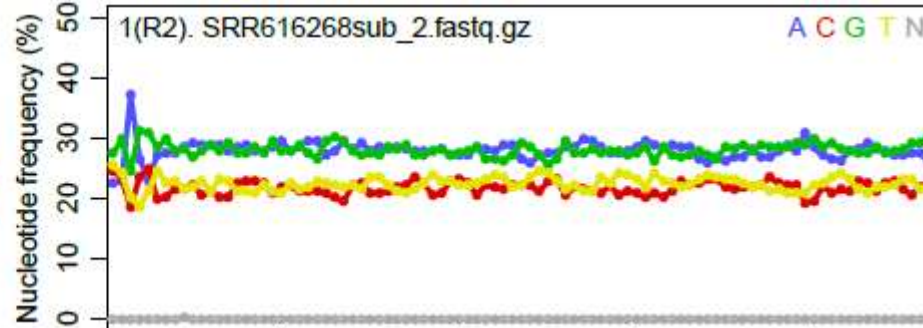
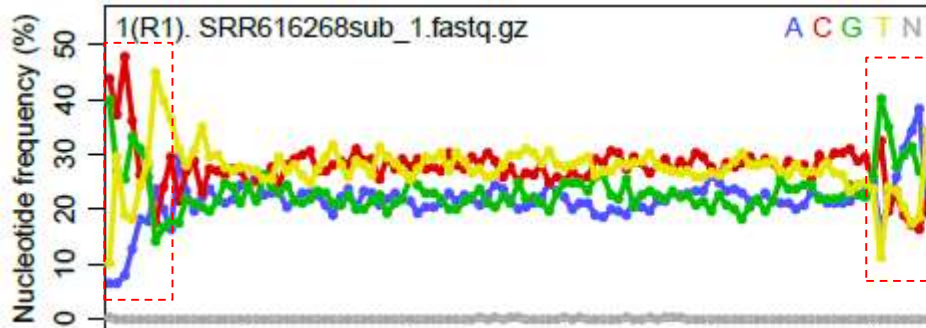
課題1, 2, 3

課題1: 「Rockhopperのアセンブル結果」を「QuasRマッピング結果のQCレポート」と絡めて説明せよ。課題2: どうすればマッピング結果を改善できるのか(マップ率を上げることができるのか)、考えられる戦略を述べよ。課題3: 今回の結果を踏まえ、リファレンスゲノムがない(マッピング結果がない)場合の高精度なアセンブル戦略について述べよ。

ファイル名 Total number of assembled transcripts

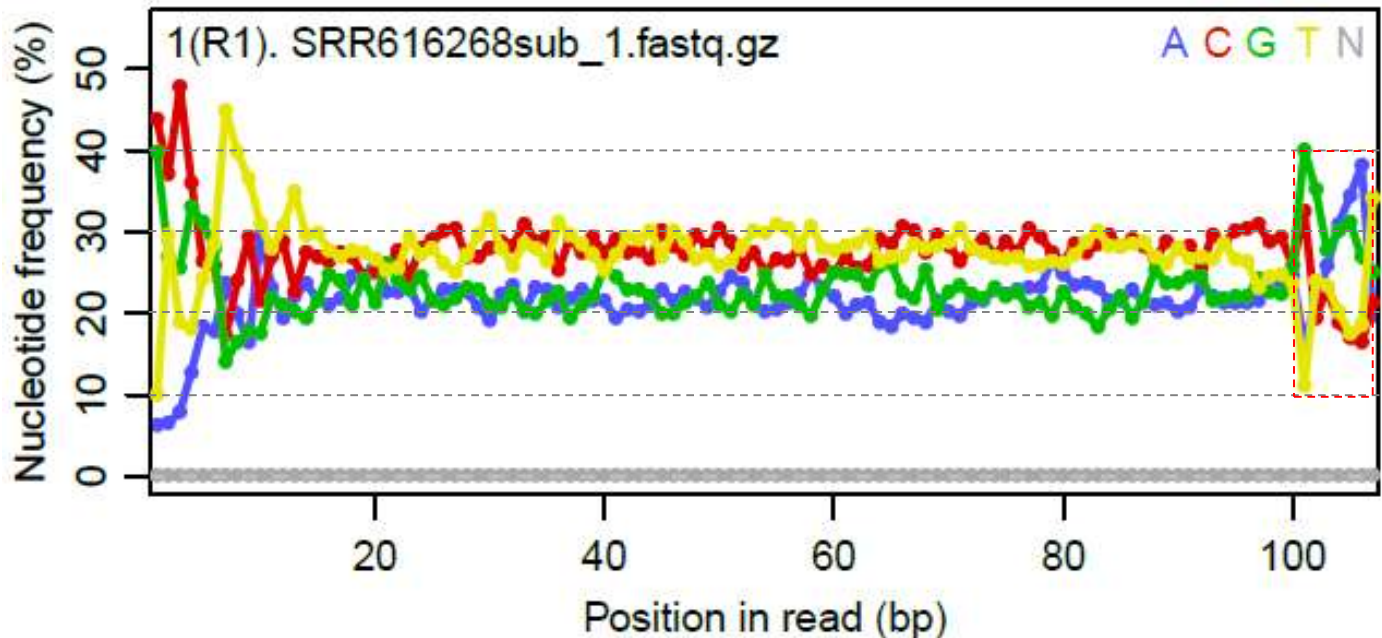
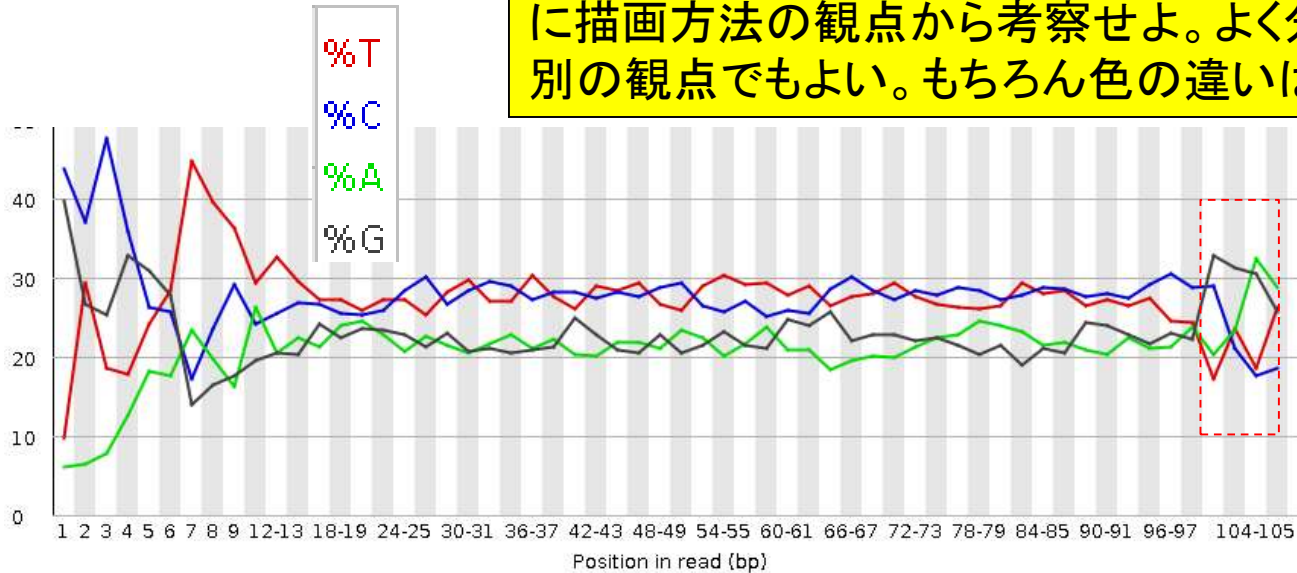
・100万リードのオリジナル?!ファイル

SRR616268sub_1.fastq.gz	8	121	106	974	996739	1964
SRR616268sub_2.fastq.gz	424	436	228	185233	983854	710393
paired-end	0	0	0	0	990047	23



課題4

課題4: ①はFastQCの「Per base sequence content」、②はQuasRのQCレポートファイル中の同様な結果である。赤点線枠部分に違いが見られるが、この理由について主に描画方法の観点から考察せよ。よく分からないヒトは別の観点でもよい。もちろん色の違いは本質的ではない



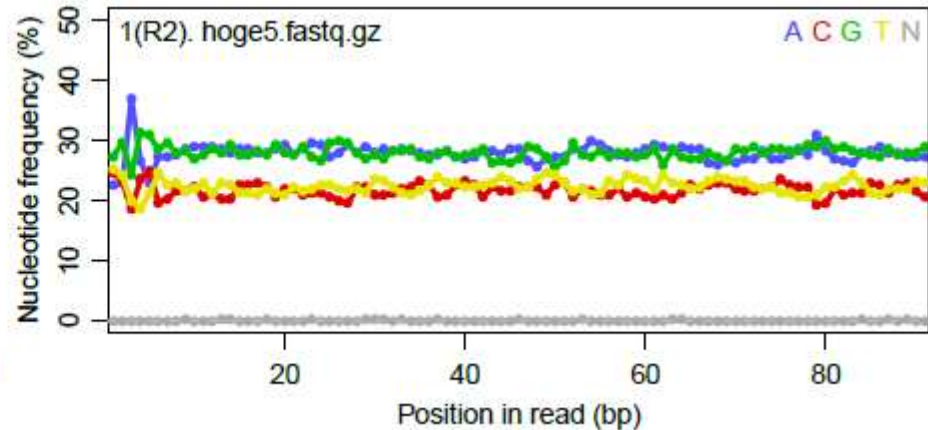
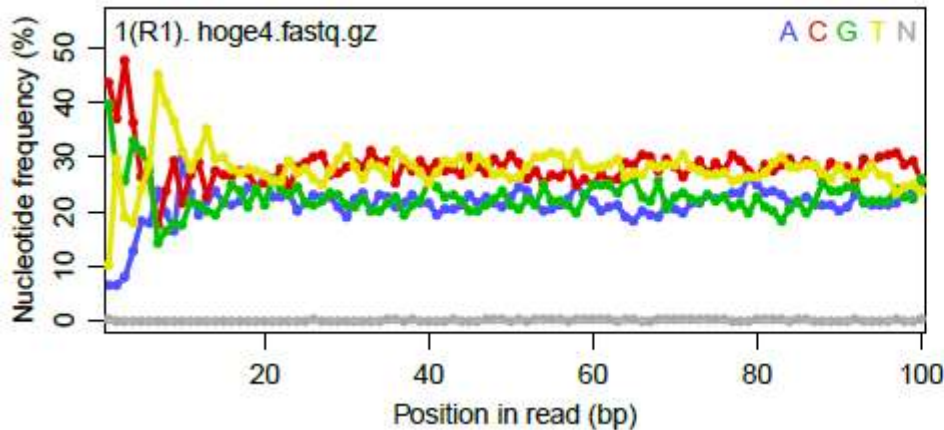
これだけ結果が変わります。
かなりイケてます。

Rockhopperリトライ

ファイル名	Total number of assembled transcripts	Average length	Median length	Total number of assembled bases	Total reads in file	Perfectly aligned reads
・100万リードのオリジナル?!ファイル(トリム後)						
SRR616268sub_1.fastq.gz	738	250	154	184550	994562	635565
SRR616268sub_2.fastq.gz	401	442	228	177642	983854	709473
paired-end	776	577	311	447950	987886	584728
・100万リードのアダプター除去後のファイル(トリム後)						
SRR616268sub_trim2_1.fastq.gz						
SRR616268sub_trim2_2.fastq.gz	388	449	229	174454	981649	706952
hoge2_1.fastq.gz						
hoge2_2.fastq.gz						
paired-end						
・アダプター除去およびフィルタリング後のファイル						
SRR616268sub_trim_1.fastq.gz						
SRR616268sub_trim_2.fastq.gz						
hoge1_1.fastq.gz						
hoge1_2.fastq.gz						
paired-end						

これだけ結果が変わります。
かなりイケてます。

QuasRリトライ



■ mapped ■ unmapped

