

USBメモリ中のhogeフォルダをデスクトップにコピーしておいてください。コピー後のファイルサイズが同じになっているかもチェックしてください。

前回(6/16)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2回

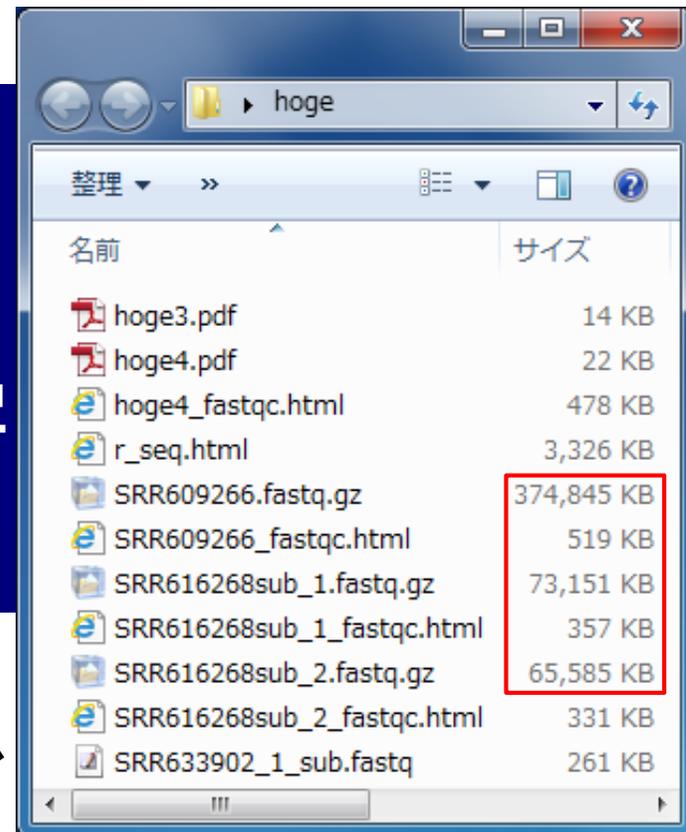
大学院農学生命科学研究科

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

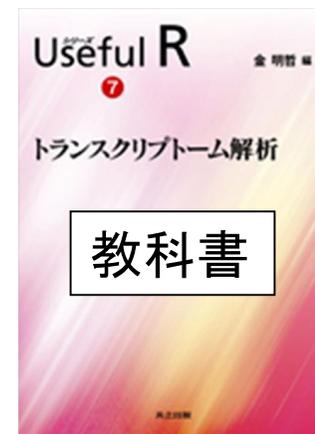
<http://www.iu.a.u-tokyo.ac.jp/~kadota/>



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)
 - アセンブル、マッピング、カウント情報取得、クラスタリング
 - 教科書の2.3節周辺
- 第4回(2015年7月7日)
 - データ正規化、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

■ QC(Quality Control)

- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

FastQC結果を眺める

FastQC Report Fri 8 May 2015
SRR616268sub_1.fastq

Summary

- Basic Statistics
- Per base sequence quality
- Per tile sequence quality
- Per sequence quality scores
- Per base sequence content **②**
- Per sequence GC content
- Per base N content
- Sequence Length Distribution
- Sequence Duplication Levels
- Overrepresented sequences
- Adapter Content
- Kmer Content

Basic Statistics

| Measure | Value |
|-----------------------------------|-------------------------|
| Filename | SRR616268sub_1.fastq |
| File type | Conventional base calls |
| Encoding | Illumina 1.5 |
| Total Sequences | 1000000 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 107 |
| %GC | 50 |

Per base sequence quality

Quality scores across all bases (Illumina 1.5 encoding)

hoge

| 名前 | サイズ |
|-------------------------------------|------------|
| hoge3.pdf | 14 KB |
| hoge4.pdf | 22 KB |
| hoge4_fastqc.html | 478 KB |
| r_seq.html | 3,326 KB |
| SRR609266.fastq.gz | 374,845 KB |
| SRR609266_fastqc.html | 519 KB |
| SRR616268sub_1.fastq.gz | 73,151 KB |
| SRR616268sub_1_fastqc.html ① | 357 KB |
| SRR616268sub_2.fastq.gz | 55,585 KB |
| SRR616268sub_2_fastqc.html | 331 KB |
| SRR633902_1_sub.fastq | 261 KB |

FastQC結果を眺める

①ポジションごとの塩基の出現確率。②赤枠のような塩基ごとのプロファイルがフラットになっていれば基本的にはOK。③-15番目あたりのプロファイルは、実験プロトコル (oligo dT primer or random hexamer; DNase Iを利用した破碎) に由来するのかもしれない。このデータは、原著論文がなく実験プロトコルの詳細が不明なため、なんともいえない。

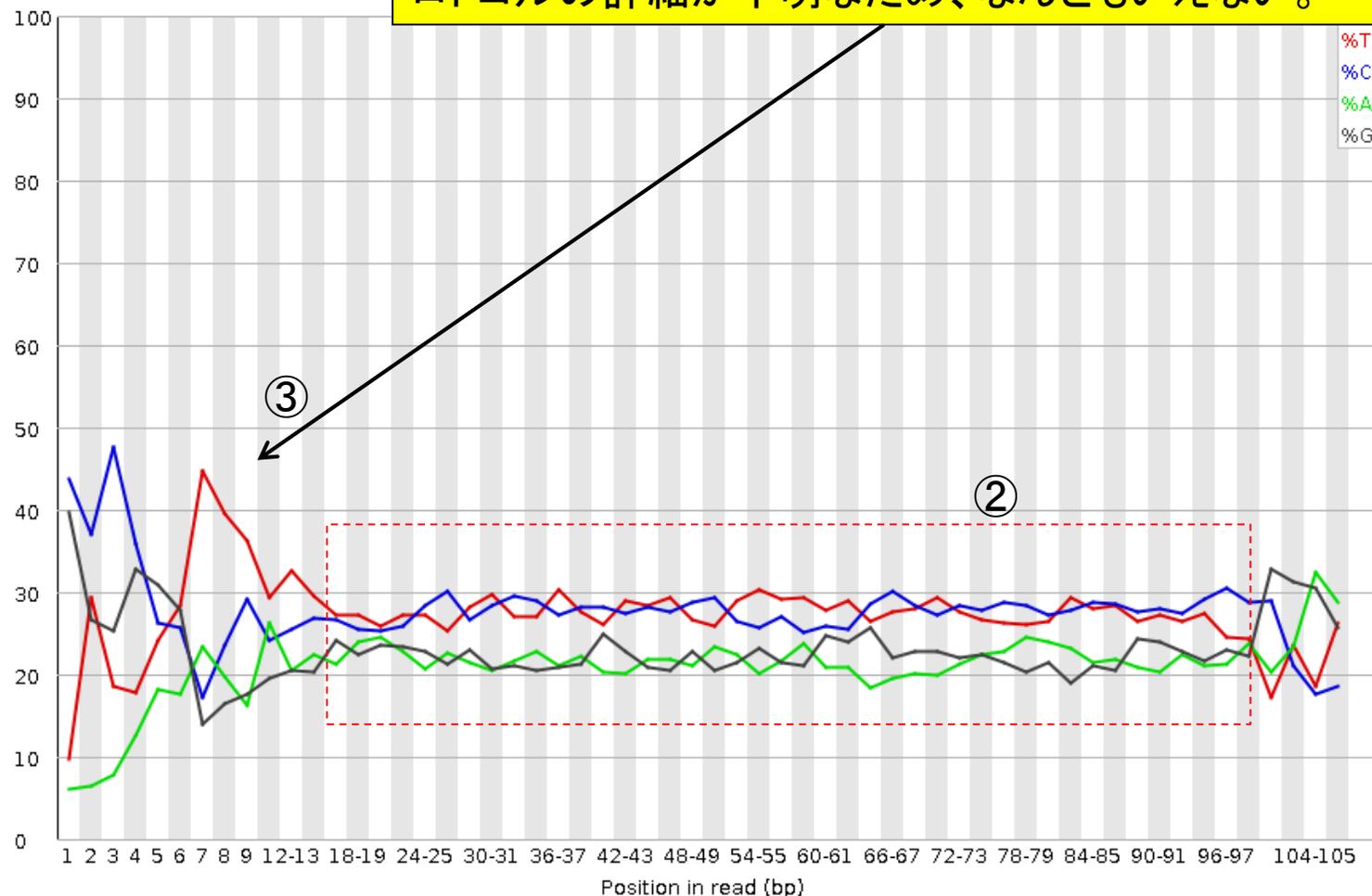
FastQC Report

Summary

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



✗ Per base sequence content



FastQC結果を眺める

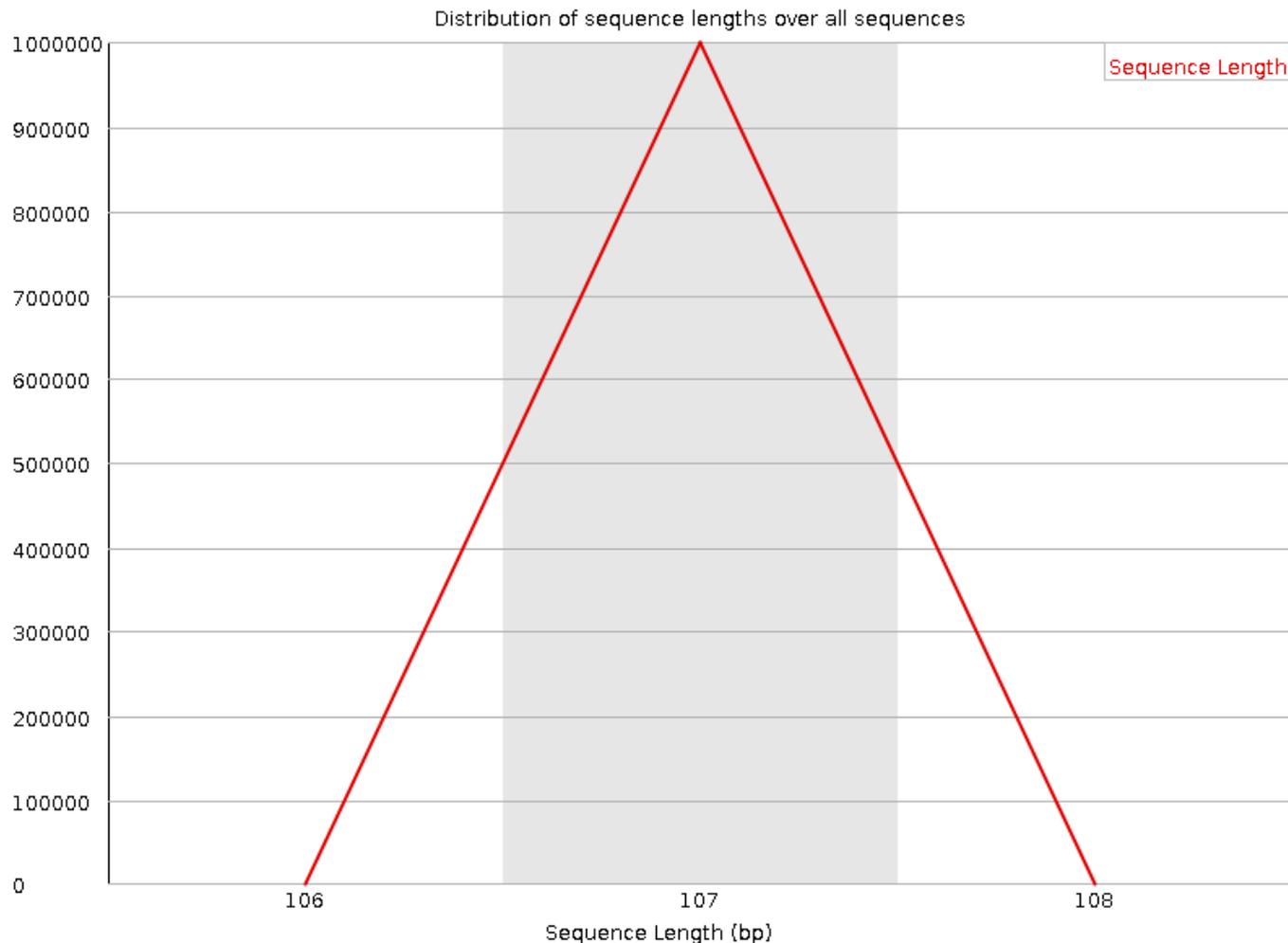
①全部で100万リードのうち、全てが107 bpだということがわかる。PacBioのデータやアダプター配列除去後のデータを入力とすればヒストグラムのようになるだろう

FastQC Report

Summary

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

✔ Sequence Length Distribution



FastQC結果を眺める

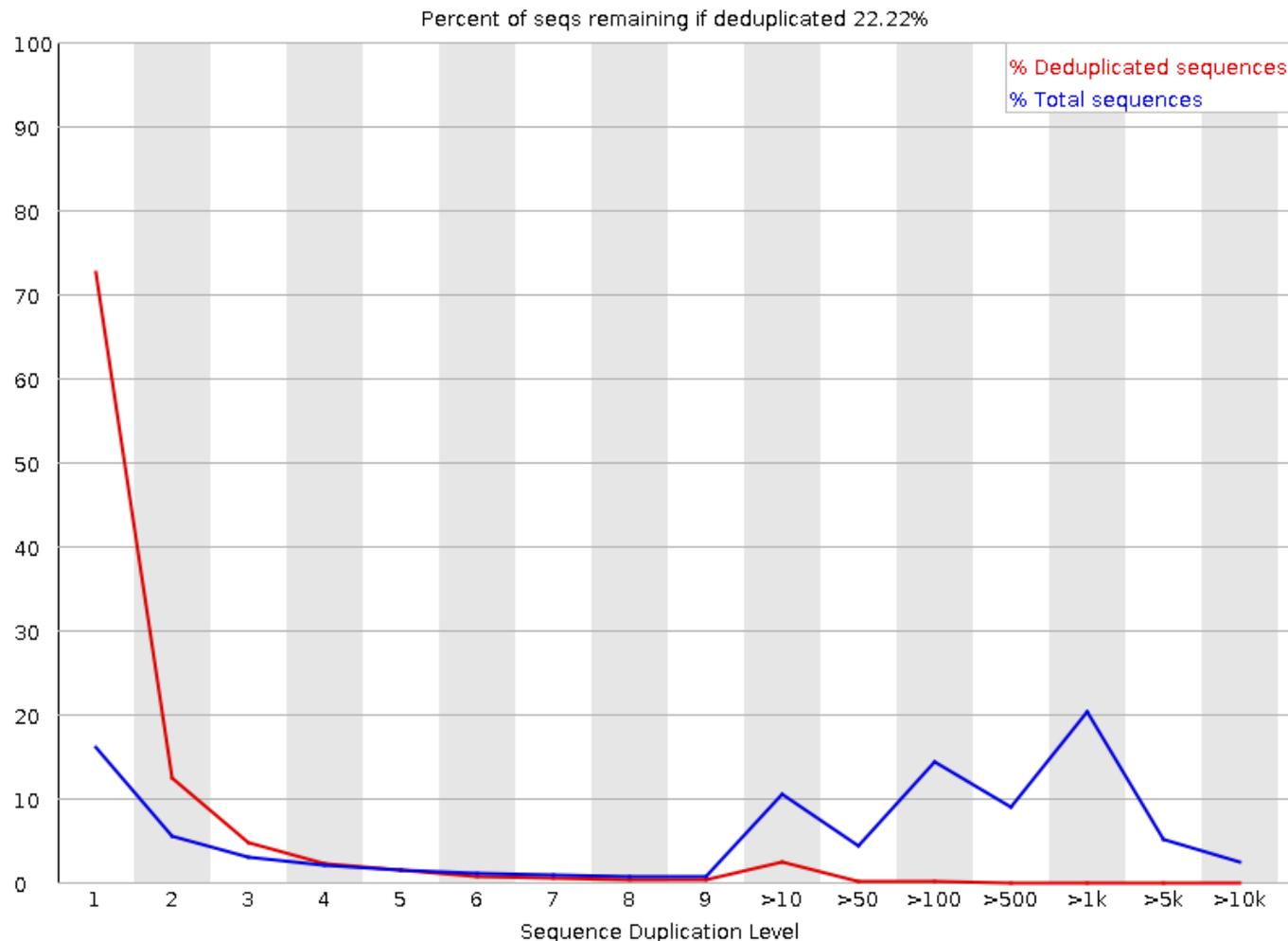
FastQC Report

Fri:
SRR616268

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](#) ①
- ✗ [Overrepresented sequences](#)
- ✓ [Adapter Content](#)
- ✗ [Kmer Content](#)

✗ Sequence Duplication Levels



FastQC結果を眺める

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

FastQC Report

Summary

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



Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|-------|---------------------|---|
| CCCCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCA | 14383 | 1.4383 | No Hit ② |
| GGCCTATTCACTGCGGCTGACCTTGCGGTCAGCACCCCTTCTTCCGAAGT | 11044 | 1.1044 | No Hit |
| GTGCTTTTCACCTTTCCCTCACGGTACTGGTTCACTATCGGTCACTAGGG | 8892 | 0.8892000000000001 | No Hit |
| CCCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCAC | 8474 | 0.8474 | No Hit |
| GTCCTAGGGAGTATTTAGCCTGGGAGATGGTCTCCCGGATCCGACG | 8189 | 0.8189 | No Hit |
| GCCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCACGATCGACCTC | 8132 | 0.8132 | No Hit |
| GTCCAGTCTACAACCCCGAGAAGCAAGCTTCTCGGTTTGGGCTCTTCCC | 6663 | 0.6663 | No Hit |
| GTGCGTTTGGGTACGGGTAGTTTATTTCTCACTAGAAGCTTTTCTTGGC | 6411 | 0.6411 | No Hit |
| GGTCACTTGGTTTCGGGTCTACATCTGCTTACTCATTGCCCCTGTTTCTCAGA | 5502 | 0.5502 | No Hit |
| GCCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCACGATCAACCTTC | 4845 | 0.48450000000000004 | No Hit |
| CCCTCCATCGCTTAAACAAAATAAACTAGTGAGGAATCTCAACCTGCTT | 4395 | 0.43949999999999995 | No Hit |
| CCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACGCACT | 4385 | 0.4385 | No Hit |
| CCCGGCTCTGCCCGCCAGCTATGTATTCCTGACAAGCAATACACTG | 4366 | 0.4366 | No Hit |
| CCACAGTTTCGGTATTATGCTTAGCCCCGGTATATTTTCGGCGCAGTGCC | 4314 | 0.4314 | No Hit |
| CTGGGCTGTTCCCTTTTCGACAATGGACCTATCGCTCACTGTCTGACTC | 4113 | 0.41130000000000005 | No Hit |
| CCGCCGCTACTCAGGATCCTGGACGGAGGGTTCGACGTTTCGCTTACAGGG | 4081 | 0.4081 | No Hit |
| CCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCACGATCGACCTTCA | 3846 | 0.3846 | No Hit |
| GTAGGTCCTTGGTTTCGGGTCTACATCTGCTTACTCATTGCCCCTGTTT | 3823 | 0.3823 | No Hit |

FastQC結果を眺める

①頻出する配列をリストアップ。②ときどき既知のアダプター配列とマッチするものが見つかる。2,415回出現した「GATCGGAA…」という50塩基の配列は、TruSeq Adapter Index 3というものと100%一致していたことを示す。このようなリード中の部分配列はトリムされるべき。

FastQC Report

Summary

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



| | | | |
|--|------|----------------------|--|
| GTTCACTATCGGTCAC TAGGGAGTAT TAGCCT TGGGAGATGGTCCCTCC | | | |
| GCCCCGGTATATTTTCGGCGCAGTGCCACTCGACTAGTGAGCTATTACG | | | |
| GCCCTGTTCAGACTCGCTTTTCGCTACGGCTCCGACTTTTCATCTTAACC | | | |
| GCCGGTTCATTCTACAAAAGGCACGCCATTACCGTTAACGGGCTTTGAT | | | |
| GCCACATCCTTTTCCACTTAGCATAAATTTAGGGACCTTAACCTGGTGATC | 2486 | 0.2486 | No Hit |
| CCTCCATCGCTTAAACAAAATAAACTAGTGCAGGAATCTCAACCTGCTTG | 2451 | 0.2451 | No Hit |
| CTTGGGAGATGGTCCCTCCCGGATTCGACGGAATTTACGTGTTCCGCCG | 2438 | 0.243800000000000002 | No Hit |
| CTTGGTTTCGGGTCTACATCTGCTTACTCATTGCCCCTGTTTCAGACTCGC | 2430 | 0.243 | No Hit |
| GTCATGGGTAGGTCACTTGGTTTCGGGTCTACATCTGCTTACTCATTGCG | 2429 | 0.2429 | No Hit |
| CTAGGGAGATTTAGCCTTGGGAGATGGTCCCTCCCGGATTCGACGGAAT | 2421 | 0.2421 | No Hit |
| GATCGGAAGAGCACAGTCTGAACTCCAGTCACTTAGGCATCTCGTATGC | 2415 | 0.2415 | TruSeq Adapter, Index 3 (100% over 50bp) |
| CCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCAGATCAACCTTCA | 2393 | 0.2393 | No Hit |
| GGGAGTATTTAGCCTTGGGAGATGGTCCCTCCCGGATTCGACGGAATTTTC | 2365 | 0.2365 | No Hit |
| TGGGCCTATTCAGTGGGCTGACCTTGGGTCAGCACCCCTTCTTCGGAA | 2332 | 0.2332 | No Hit |
| CTCCTGCGTCCCTCCATCGCTTAAACAAAATAAACTAGTGCAGGAATCTC | 2331 | 0.233100000000000003 | No Hit |
| CGCCGCTACTATGGGAATCGAGTTTTCTTCTCTTCTCCTGCGGGTACTGAG | 2255 | 0.2255 | No Hit |
| CGCCGGTTCATTCTACAAAAGGCACGCCATTACCGTTAACGGGCTTTGA | 2172 | 0.217199999999999998 | No Hit |
| CCGCCGGCCAGCTATGTATTCAGTACAAGCAATACACTGATGTGACTG | 2081 | 0.2081 | No Hit |
| CCCCTTCGACAATGGACCTTATCGCTCACTGTCTGACTCCCGGAGTAAG | 2070 | 0.207 | No Hit |
| CTTAACCTGGTGATCTGGGCTGTTCCCTTCGACAATGGACCTTATCGCT | 2005 | 0.200499999999999998 | No Hit |

他にもちらほらと既知のアダプター配列と一致するものが見つかる。

FastQC結果を眺める

Fri 8 May 2015

SRR616268sub_1.fastq

FastQC Report

Summary

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



| Sequence | Count | Percentage | Match |
|--|-------|---------------------|--|
| CACGGTTTCAGGAACTGTTTCACTCCCTCCGGGGTGCTTTTACACCTT | 1610 | 0.161 | No Hit |
| GGGCTGTTCCTTTTCGACAATGGACCTTATCGCTCACTGTCTGACTCC | 1560 | 0.156 | No Hit |
| GTCCTCTTCTGCACTCAAGTTTCCAGTTTCCGATGCGCTTCTCGGTTA | 1559 | 0.1559 | No Hit |
| GATCGGAAGAGCACACGTCTGAACTCCAGTACCGATGTATCTCGTATGC | 1539 | 0.1539 | TruSeq Adapter, Index 2 (100% over 50bp) |
| CTCGCCTTAGATCCCGACTAACCCCTGGGAGGACGAGCCTTCCCCAGAAA | 1499 | 0.14989999999999998 | No Hit |
| GGGACCTTAACTGGTGATCTGGGCTGTTCCCTTTTCGACAATGGACCTTA | 1497 | 0.1497 | No Hit |
| ATTAGCCTTGGGAGATGGTCTCCCGATTCCGACGGAATTTACGTGT | 1487 | 0.1487 | No Hit |
| CATGGGTAGGTCACCTGGTTCGGGTCTACATCTGCTTACTCATTCGCC | 1464 | 0.1464 | No Hit |
| CCCTACTGCTGCCTCCCGTAGGAGTTTGGGCCGTGTCTCAGTCCCAATGT | 1455 | 0.1455 | No Hit |
| CTTTCGACAATGGACCTTATCGCTCACTGTCTGACTCCCGGAGTAAGATC | 1451 | 0.1451 | No Hit |
| CGTACTCATGCCGGCATTCTCACTTCTAAGCGCTCCAGCCGCTCCTCACG | 1405 | 0.1405 | No Hit |
| CCGGGGTGCTTTTACCTTTCCCTCACGGTACTGGTTCACATCGGTCAC | 1361 | 0.1361 | No Hit |
| CTGGTGATCTGGGCTGTTCCTTTTCGACAATGGACCTTATCGCTCACTG | 1338 | 0.1338 | No Hit |
| GGGCCTATTCACCTGCGGCTGACCTTGGGTCAGCACCCCTTCTCCGAAG | 1337 | 0.1337 | No Hit |
| GCCGCCGGCCAGCTATGTATTCACTGACAAGCAATACTGATGTGTACT | 1324 | 0.13240000000000002 | No Hit |
| AGATCGGAAGAGCACACGTCTGAACTCCAGTACCGATGTATCTCGTATG | 1318 | 0.1318 | TruSeq Adapter, Index 2 (100% over 49bp) |
| GCCGAGTTCCTTAACGAGAGTTCGCTCGCTCACCTGAGGATACTCTCCTC | 1310 | 0.131 | No Hit |
| CGTCCCTCCATCGCTTAAACAAAATAAAGTGCAGGAATCTCAACCTG | 1276 | 0.1276 | No Hit |
| CACACGGTTTCAGGAACTGTTTCACTCCCTCCGGGGTGCTTTTACCT | 1275 | 0.1275 | No Hit |

FastQC結果を眺める

7-merの塩基配列の出現頻度解析結果。出現頻度の期待値に比べて実測値が位置特異的に極端に多いk-merの上位がリストアップされている。

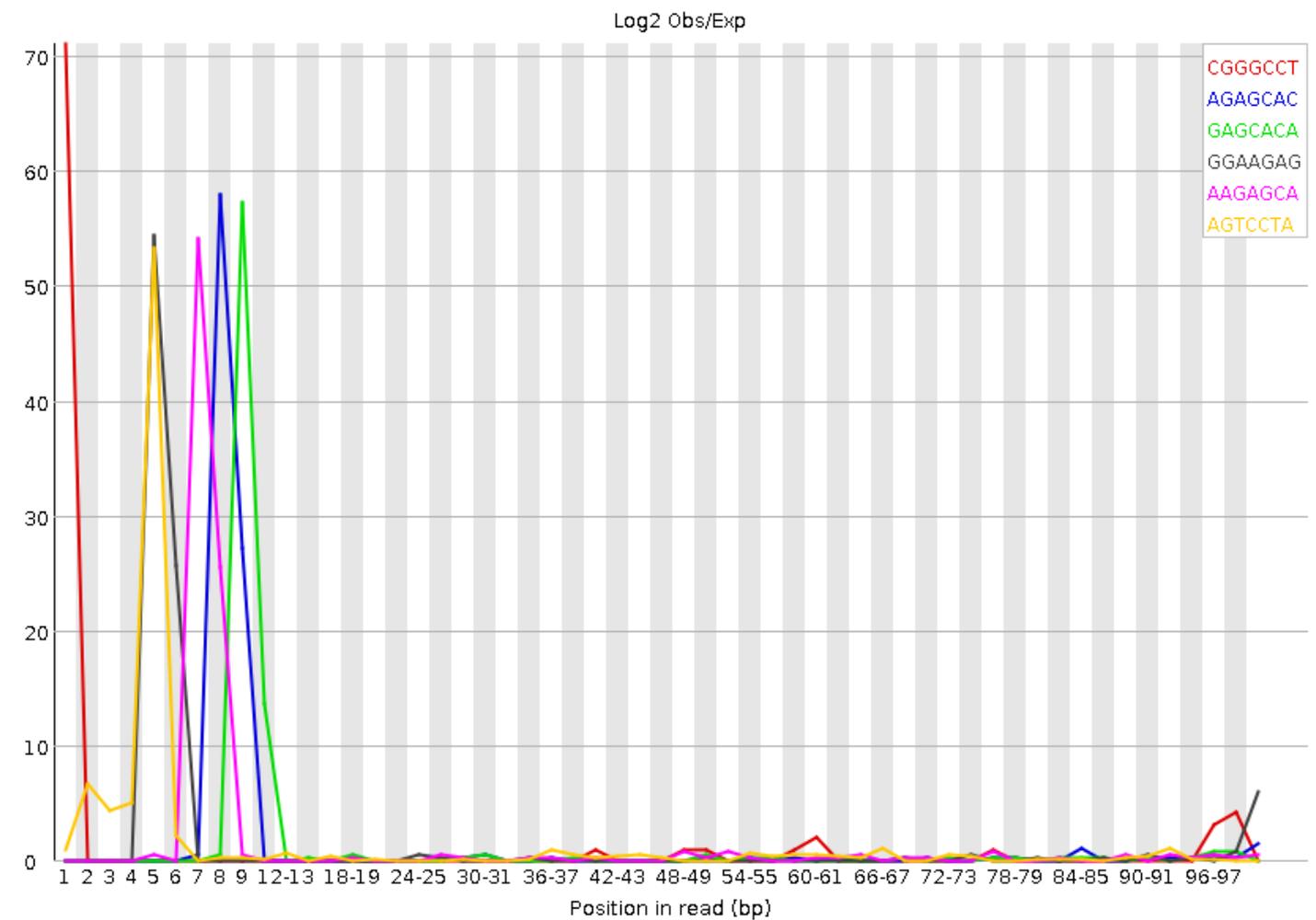
FastQC Report

Summary

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



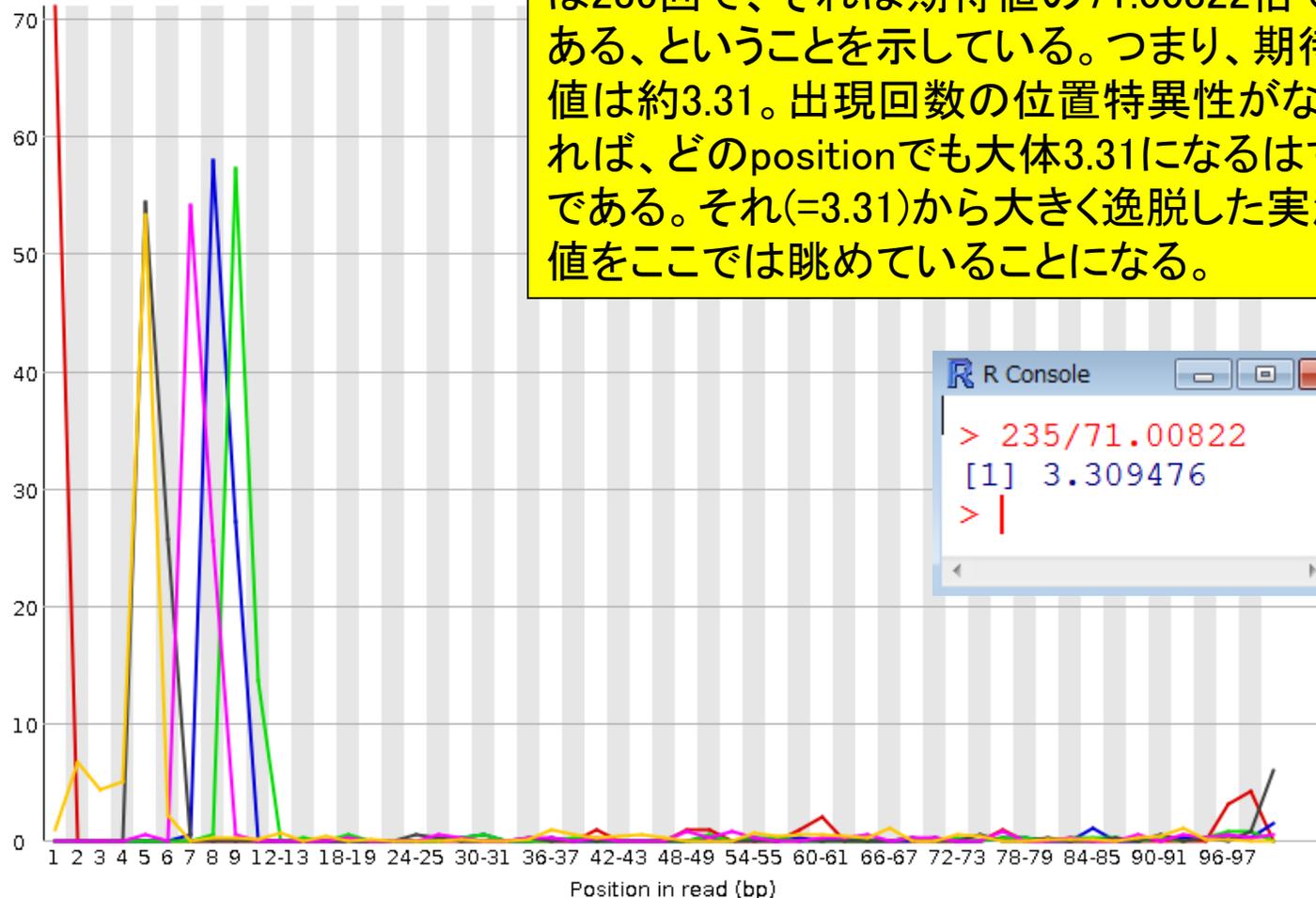
✗ Kmer Content



FastQC結果を眺める

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



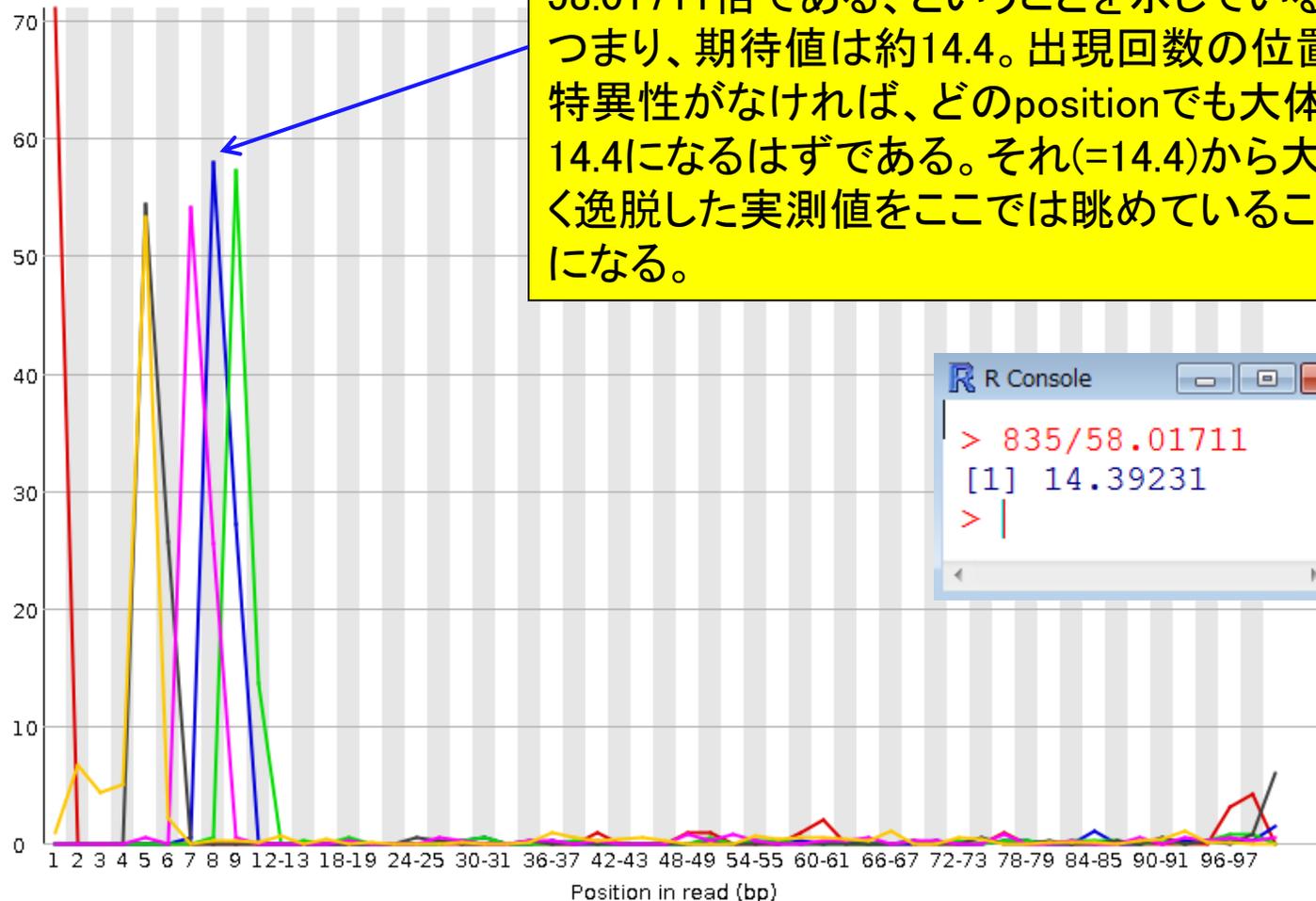
| Sequence | Count | PValue | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|--------|-------------|--|
| CGGGCCT | 235 | 0.0 | 71.00822 | 1 ② |
| AGAGCAC | 835 | 0.0 | 58.01711 | 8 |
| GAGCACA | 845 | 0.0 | 57.330517 | 9 |

FastQC結果を眺める

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



②次に多いのはAGAGCAC。これはposition 8で非常に多く出現する。実際のposition 8での出現回数は835回で、それは期待値の58.01711倍である、ということを示している。つまり、期待値は約14.4。出現回数の位置特異性がなければ、どのpositionでも大体14.4になるはずである。それ(=14.4)から大きく逸脱した実測値をここでは眺めていることになる。

```
R Console
> 835/58.01711
[1] 14.39231
> |
```

| Sequence | Count | PValue | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|--------|-------------|----------------------|
| CGGGCCT | 235 | 0.0 | 71.00822 | 1 |
| AGAGCAC | 835 | 0.0 | 58.01711 | 8 |
| GAGCACA | 845 | 0.0 | 57.330517 | 9 |

②

Contents

■ QC(Quality Control)

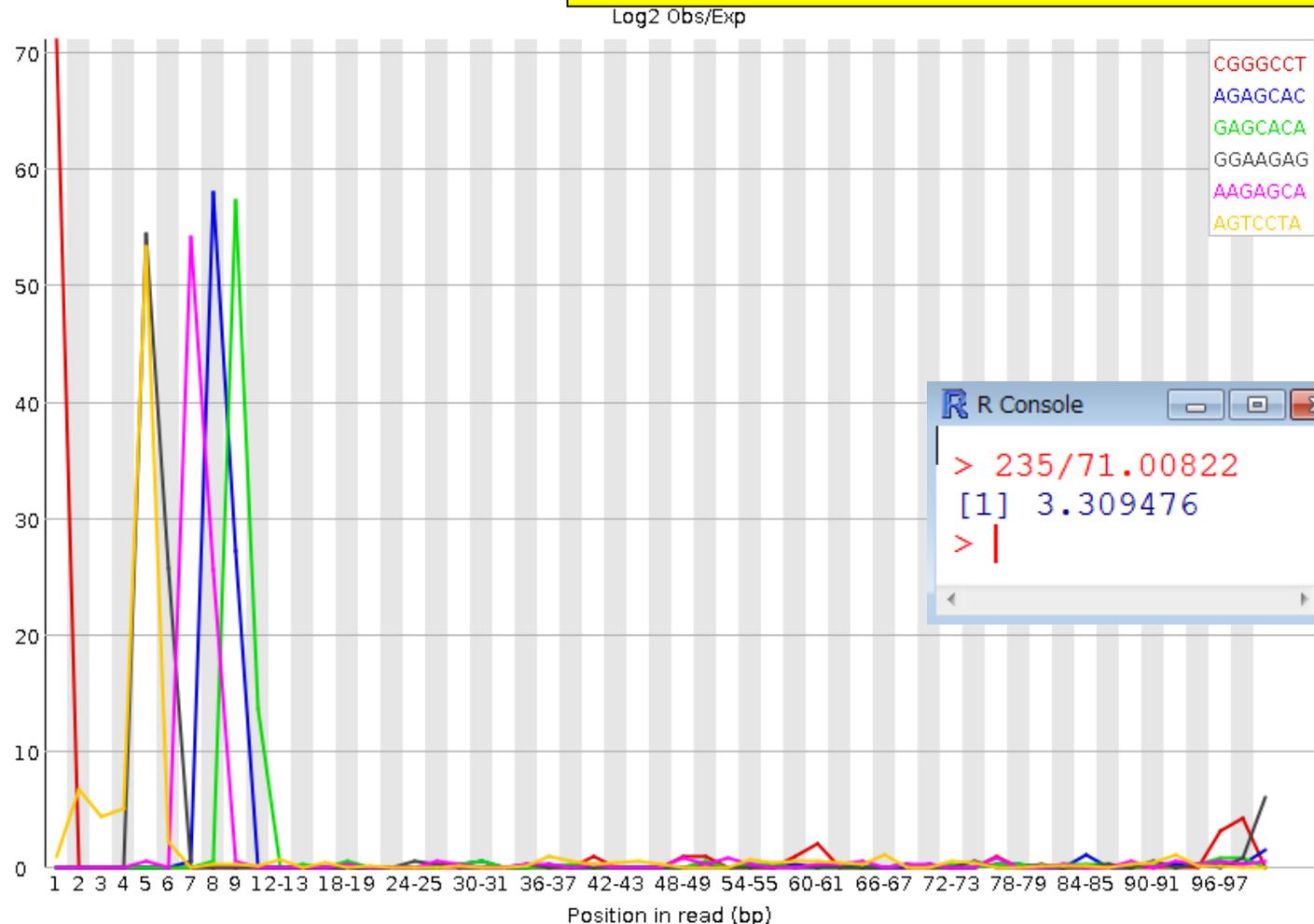
- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去 (QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去 (QuasRパッケージ)

k-mer解析

CGGGCCTがposition 1で235回出現しており、それが期待値(=3.31)の71.00822倍だということを、Rで確認してみよう。

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**



```
R Console  
> 235/71.00822  
[1] 3.309476  
> |
```

| Sequence | Count | PValue | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|--------|-------------|----------------------|
| CGGGCCT | 235 | 0.0 | 71.00822 | 1 ② |
| AGAGCAC | 835 | 0.0 | 58.01711 | 8 |
| GAGCACA | 845 | 0.0 | 57.330517 | 9 |

k-mer解析

CGGGCCTがposition 1で235回出現しており、それが期待値(=3.31)の71.00822倍だということを、Rで確認してみよう。

- ・ [イントロ](#) | [NGS](#) | [アノテーション情報取得](#) | [TranscriptDb](#) | [GFF/GTF形式ファイルから](#)
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [BSgenome](#) | [基本情報を取得](#) (last modified 2015/06/10)
- ・ [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 2015/06/19) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [基礎](#) (last modified 2015/06/18) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [description行の記述を整形](#) (last modified 2015/06/18) **NEW**

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

- ・ [イントロ](#) | [ファイル形式](#) | [Sanger FASTQ形式ファイルを読み込み](#)
- ・ [イントロ](#) | [ファイル形式](#) | [FASTQ形式ファイルを読み込み](#)
- ・ [前処理](#) | [クオリティチェック](#)

1. サンプルデータのFASTQ形式

[SRR037439](#)から得られるFASTQ形式のFASTQ形式ファイルの読み込み
quality情報を除く塩基配列情報の読み込み

```
in_f <- "SRR037439.fastq"
#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
```

2. サンプルデータのFASTQ形式

[SRR037439](#)から得られるFASTQ形式のFASTQ形式ファイルの読み込み

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。

```
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") #in_fで指定したファイルの読み込み
fasta #確認してるだけです
#以下はおまけ(部分配列を抽出し出現頻度上位を概観)
hoge <- subseq(fasta, start=1, end=7) #指定した始点と終点の範囲の配列を抽出
hoge #確認してるだけです
head(table(hoge)) #出現頻度を順番に表示
head(sort(table(hoge), decreasing=T)) #出現頻度上位を表示
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
hoge <- subseq(fasta, start=3, width=7) #指定した始点から一定範囲の配列を抽出
hoge #確認してるだけです
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
hoge <- subseq(fasta, start=5, width=7) #指定した始点から一定範囲の配列を抽出
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

Tips: 正規表現

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。

```
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(Biobase) #パッケージの読み込み
```

```
#入力ファ
```

```
fasta <- readFastq(in_f) #in_fで指定したファイルの読み込み
```

```
#以下はお
```

```
hoge <- sampleNames(fasta) #各点と終点の範囲
```

```
hoge #各点と終点の範囲
```

```
head(table(hoge)) #各点と終点の範囲
```

```
head(sort(table(hoge))) #各点と終点の範囲
```

```
table(hoge) #各点と終点の範囲
```

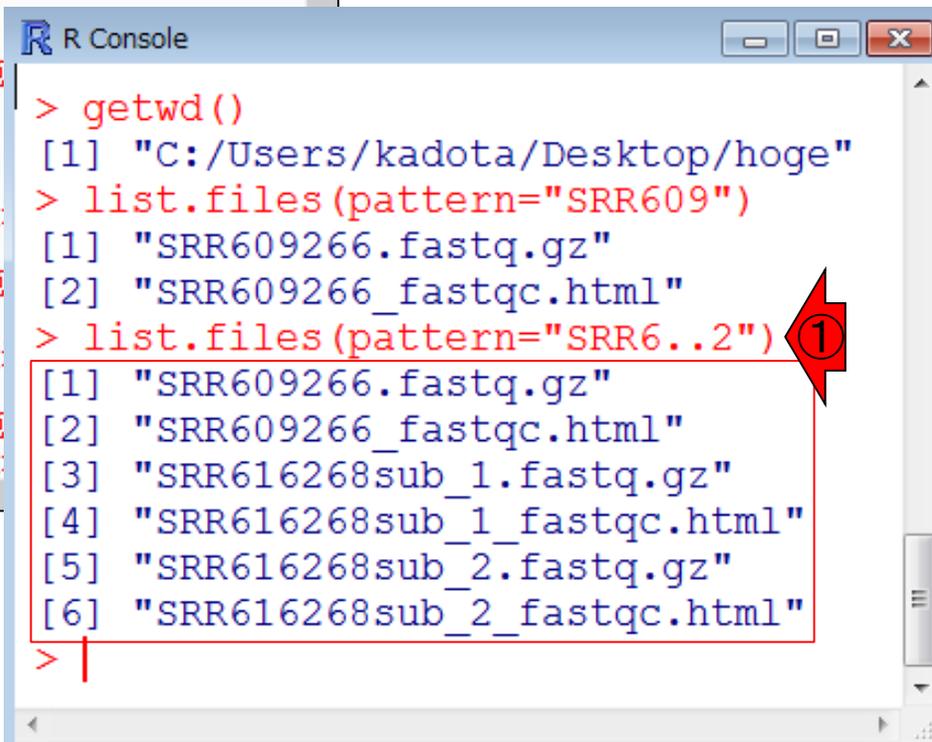
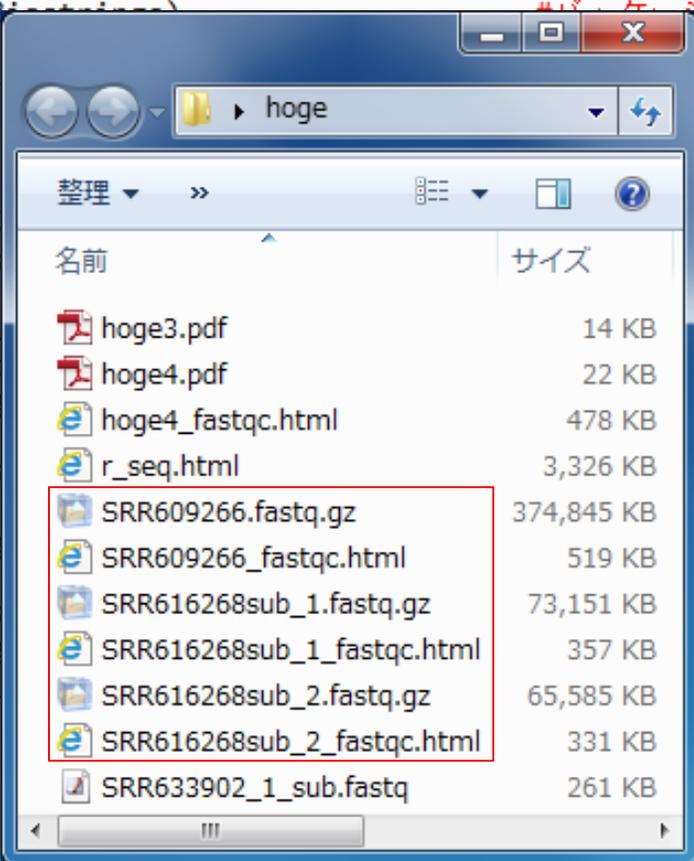
```
hoge <- sampleNames(fasta) #各点から一定範囲
```

```
hoge #各点から一定範囲
```

```
table(hoge) #各点から一定範囲
```

```
hoge <- sampleNames(fasta) #各点から一定範囲
```

```
table(hoge) #各点から一定範囲
```



Tips: 正規表現

list.files関数中のpatternオプションでは、正規表現を用いた検索が可能。「.(どっと)」は任意の1文字を表す。ここでは「..」とが2つ並んでいるので、「SRR6..2」の条件に該当する「SRR6092」と「SRR6162」を含むファイル群が表示される。アダプター配列中のIndex配列部分を曖昧にして検索したいときに便利。

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

```
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(Biobase) #パッケージの読み込み
```

```
#入力ファ
fasta <- readFasta(in_f) #in_fで指定したファイルの読み込み
fasta #ただです
```

```
#以下はお
hoge <- sampleNames(fasta) #各点と終点の範囲
hoge #ただです
head(table(hoge)) #順番に表示
head(sort(table(hoge))) #上位を表示
#出現頻度を表
```

```
hoge <- sampleNames(fasta) #各点から一定範囲
hoge #ただです
table(hoge) #出現頻度を表
```

```
hoge <- sampleNames(fasta) #各点から一定範囲
table(hoge) #出現頻度を表
```

| 名前 | サイズ |
|----------------------------|------------|
| hoge3.pdf | 14 KB |
| hoge4.pdf | 22 KB |
| hoge4_fastqc.html | 478 KB |
| r_seq.html | 3,326 KB |
| SRR609266.fastq.gz | 374,845 KB |
| SRR609266_fastqc.html | 519 KB |
| SRR616268sub_1.fastq.gz | 73,151 KB |
| SRR616268sub_1_fastqc.html | 357 KB |
| SRR616268sub_2.fastq.gz | 65,585 KB |
| SRR616268sub_2_fastqc.html | 331 KB |
| SRR633902_1_sub.fastq | 261 KB |

```
R Console
> getwd()
[1] "C:/Users/kadota/Desktop/hoge"
> list.files(pattern="SRR609")
[1] "SRR609266.fastq.gz"
[2] "SRR609266_fastqc.html"
> list.files(pattern="SRR6..2")
[1] "SRR609266.fastq.gz"
[2] "SRR609266_fastqc.html"
[3] "SRR616268sub_1.fastq.gz"
[4] "SRR616268sub_1_fastqc.html"
[5] "SRR616268sub_2.fastq.gz"
[6] "SRR616268sub_2_fastqc.html"
> |
```

k-mer解析

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

```

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") #in_fで指定したファイルの読み込み
fasta #確認してるだけです

```

```

#以下はおまけ(部分配列を抽出し出現頻度)
hoge <- subseq(fasta, start=1, end=100000)
hoge
head(table(hoge))
head(sort(table(hoge), decreasing=T))
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, start=3, width=100000)
hoge
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, start=5, width=100000)
table(hoge)["CGGGCCT"]

```

Biostringsパッケージ中のreadDNASTringSet関数は、formatオプションを"fastq"にすることで、FASTQ形式ファイルを読み込める。gzip圧縮ファイルを入力として与えても自動判定して読み込んでくれる。ただし、quality score情報は読み飛ばすようだ。ここではリード塩基配列情報しか取り扱わないので特に問題ない。

```

R Console
> fasta #確認してるだけです
A DNASTringSet instance of length 1000000
      width seq
[1] 107 AGCCCGACTTT...CTTCTCTAAC SRR616268.7 2291:...
[2] 107 GATCTGGGCTG...ACCTCCGAAA SRR616268.20 2291:...
[3] 107 CCGGTATATTT...ACGATCTTAC SRR616268.21 2291:...
[4] 107 CTTCGATACCG...CATAAAGACC SRR616268.22 2291:...
[5] 107 CCCCGGTATAT...CAAGGTCCAT SRR616268.30 2291:...
... ..
[999996] 107 CCCCGGTATAT...CAAGGGTTAC SRR616268.1000856...
[999997] 107 TTCGGGTCTAC...CGGGTGAAAT SRR616268.1000857...
[999998] 107 CGTCCATCCCG...AGCCAAGAGT SRR616268.1000858...
[999999] 107 CTAGGGAGTAT...TTCCTGGACA SRR616268.1000859...
[1000000] 107 GCCTTGTC AAT...GACAGCAGCT SRR616268.1000860...
> |

```

k-mer解析

目的: CGGGCCTがposition 1で235回出現しており、それが期待値(=3.31)の71.00822倍だということを、Rで確認。まずはposition 1から7-mer分の部分配列を切り出し、7-merの種類ごとの出現頻度を計測。

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

```
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") #in_fで指定したファイルの読み込み
fasta #確認してるだけです
```

```
#以下はおまけ(部分配列を抽出し出現頻度)
hoge <- subseq(fasta, start=1, end=7)
hoge
head(table(hoge))
head(sort(table(hoge), decreasing=T))
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, start=3, width=7)
hoge
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, start=5, width=7)
table(hoge)["CGGGCCT"]
```

```
R Console
> fasta #確認してるだけです
A DNASTringSet instance of length 1000000
      width seq
[1] 107 AGCCCGACTTT...CTTCTCTAAC SRR616268.7 2291:...
[2] 107 GATCTGGGCTG...ACCTCCGAAA SRR616268.20 2291:...
[3] 107 CCGGTATATTT...ACGATCTTAC SRR616268.21 2291:...
[4] 107 CTTCGATACCG...CATAAAGACC SRR616268.22 2291:...
[5] 107 CCCCGGTATAT...CAAGGTCCAT SRR616268.30 2291:...
... ..
[999996] 107 CCCCGGTATAT...CAAGGGTTAC SRR616268.1000856...
[999997] 107 TTCGGGTCTAC...CGGGTGAAT SRR616268.1000857...
[999998] 107 CGTCCATCCCG...AGCCAAGAGT SRR616268.1000858...
[999999] 107 CTAGGGAGTAT...TTCCTGGACA SRR616268.1000859...
[1000000] 107 GCCTTGTC AAT...GACAGCAGCT SRR616268.1000860...
> |
```

k-mer解析

position 1から7-mer分の部分配列を切り出したところまで。subseq関数については、2015年4月14日の「ゲノム情報解析基礎」スライド72-75あたりを復習せよ。

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。

```
#以下はおまけ(部分配列を抽出し出現頻度上位を概観)
hoge <- subseq(fasta, start=1, end=7) #指定した始点と終点の範囲の配列を抽出
hoge #確認してるだけです
head(table(hoge)) #出現頻度を順番に表示
head(sort(table(hoge), decreasing=T)) #出現頻度上位を表示
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
hoge <- subseq(fasta, start=3, width=7) #指定した始点から一定範囲の配列を抽出
hoge #確認してるだけです
table(hoge)["CGGGCCT"]
```

```
hoge <- subseq(fasta, start=5, width=7)
table(hoge)["CGGGCCT"]
```

```
table(subseq(fasta, start=1, width=7))
table(subseq(fasta, start=2, width=7))
table(subseq(fasta, start=3, width=7))
table(subseq(fasta, start=4, width=7))
table(subseq(fasta, start=5, width=7))
table(subseq(fasta, start=6, width=7))
table(subseq(fasta, start=7, width=7))
table(subseq(fasta, start=8, width=7))
```

```
R Console
> hoge
A DNAStringSet instance of length 1000000
      width seq
[1]      7 AGCCCGA
[2]      7 GATCTGG
[3]      7 CCGGTAT
[4]      7 CTCGAT
[5]      7 CCCCGGT
...      ...
[999996] 7 CCCCGGT
[999997] 7 TTCGGGT
[999998] 7 CGTCCAT
[999999] 7 CTAGGGA
[1000000] 7 GCCTTGT
```

Tips: 講義資料再取得

ウェブページ内での「subseq」でのキーワード検索でもいいでしょうし、「2015.04.14」でキーワード検索して講義資料PDFを再取得するのもアリでしょう。

(Rで)塩基配列解析

～NGS、RNA-seq、ゲノム、トランスクリプトーム、正規化、発現変動、統計、モデル、バイオインフォマティクス

(last modified

- はじめに (last modified 2015/03/31)
- 参考資料(講義、講習会、本など) **NEW** (last modified 2015/06/16)
- 過去のお知らせ (last modified 2015/05/28) **NEW**
- インストール | について (last modified 2015/04/04)

What's new?

- このウェブページはMacintosh2015とWindows2015を体系的にまとめたものであるという
- 参考資料(講習会、講習会、講習会)
- 「イントロ | NGS」の更新に伴って発生しているエラーを修正
- Rパッケージのリリースに対する
- Rパッケージの

- インストール | R本体 | 最新版 | Win用 (last modified 2015/03/22) 推奨
- インストール | R本体 | 最新版 | Mac用 (last modified 2015/04/22) 推奨
- インストール | R本体
- インストール | R本体
- インストール | Rパッケージ
- (削除予定) Rのインストール
- (削除予定) 個別パッケージ
- 基本的な利用法 (last modified 2015/03/31)
- サンプルデータ (last modified 2015/03/31)
- バイオインフォマティクス
- 書籍 | トランスクリプトーム
- 書籍 | トランスクリプトーム

参考資料(講義、講習会、本など) **NEW**

基本的に私(門田)の個人ページに記載してあるものです。かなり古い講演資料などの情報をもとに勉強されている方もいらっしゃるようですので、ここでは2013年秋以降の情報を載せておくとともに、大まかな内容についても述べておきます。講演予定のものについては、資料のアップは講演当日のファイルがほしい方はお電話や私個人への謝辞

書籍、学会誌

- 孫建強, 三浦文彦. 第3回Linux環境構築. 2015. 内容: 日本乳酪とMacのLinux環境の違い」として資料などは「書籍」の項目をご参照ください。
- 孫建強, 湯敏. 手法: 第2回GU. 174, 2014.

越しくたさい。門田が出勤していれば対応できます。2コマ(2×90 min)分。

- 門田幸二「[講義資料](#)」, [アグリバイオインフォマティクス教育研究プログラム](#)の大学院講義科目:[ゲノム情報解析基礎](#), 東京大学(東京), 2015.04.21
内容: CRANとBioconductor, BSgenomeパッケージを利用してヒトゲノム中のCpG出現確率が低いことを確認。2連続塩基の出現頻度解析。作図(box plot)。1コマ(90 min)分。
- 門田幸二「[講義資料\(Win版とMac版\)](#)」, [アグリバイオインフォマティクス教育研究プログラム](#)の大学院講義科目:[ゲノム情報解析基礎](#), 東京大学(東京), **2015.04.14**
内容: このウェブページの基本的な使い方、ありがちなミスや警告メッセージの読み取り。コード内部の説明、関数の使用法、タブ補完、二重クォーテーション問題などのTips。multi-FASTAファイルの解析。GC含量計算など。2コマ(180 min)分。
- 門田幸二「[ウェブページと講義資料](#)」, [アグリバイオインフォマティクス教育研究プログラム](#)の大学院講義科目:[ゲノム情報解析基礎](#), 東京大学(東京), 2015.04.07
内容: 初心者向けバイオインフォマティクス全般およびゲノム情報解析系のイントロダクションの話。Rのイントロダクションやこのウェブページの簡単な使い方を含む。1コマ(90 min)分。

k-mer解析

```
#以下はおまけ(部分配列を抽出し出現頻度上位を概観)
hoge <- subseq(fasta, start=1, end=7) #指定した始点と終点の範囲の配列を抽出
hoge #確認してるだけです
head(table(hoge)) #出現頻度を順番に表示
head(sort(table(hoge), decreasing=T)) #出現頻度上位を表示
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
hoge <- subseq(fasta, start=3, width=7) #指定した始点から一定範囲の配列を抽出
hoge #確認してるだけです
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
hoge <- subseq(fasta, start=5, width=7) #指定した始点から一定範囲の配列を抽出
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
table(subseq(fasta, start=1, width=7))["CGGGCCT"]
table(subseq(fasta, start=2, width=7))["CGGGCCT"]
table(subseq(fasta, start=3, width=7))["CGGGCCT"]
table(subseq(fasta, start=4, width=7))["CGGGCCT"]
table(subseq(fasta, start=5, width=7))["CGGGCCT"]
table(subseq(fasta, start=6, width=7))["CGGGCCT"]
table(subseq(fasta, start=7, width=7))["CGGGCCT"]
table(subseq(fasta, start=8, width=7))["CGGGCCT"]
```

次に、7-merの種類ごとの出現頻度をtable関数を用いて計測。①デフォルトはalphabet順に作成したk-mer (k=7)の種類ごとに出現回数を返す。②table(hoge)実行結果を降順(decreasing=T)にsortした上位6個を表示。一番出現回数が多いのはCCCCGGTの14,995回。

```
R Console
> head(table(hoge)) #出現頻度
hoge
AAAAAAA AAAAAAC AAAAAAG AAAAAAT AAAAAACA AAAAAACG
      238      26      3      1      1      1
> head(sort(table(hoge), decreasing=T)) #出現頻度
hoge
CCCCGGT GCCGGCA GGCCTAT GTGCTTT CCCGGTA GTCACTA
      14995      13663      11403      9206      8811      8451
> table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度
CGGGCCT
      1804
> length(table(hoge))
[1] 15639
> 4^7
[1] 16384
> |
```

k-mer解析

③ピンポイントで目的のCGGGCCTの出現回数(=1,804)を表示。FastQC内部の計算事情を知らないヒトは、この段階でオカシと思わなければいけない。

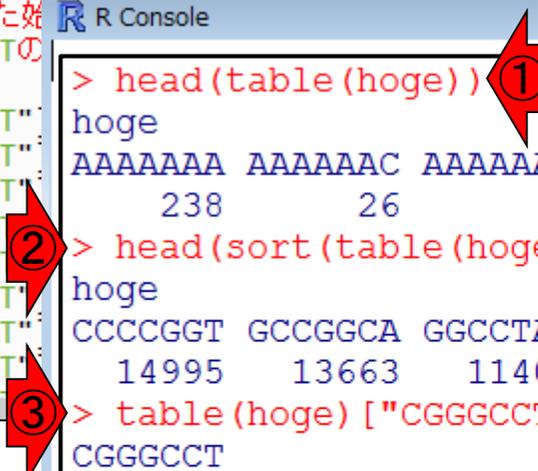
```
#以下はおまけ(部分配列を抽出し出現頻度上位を概観)
hoge <- subseq(fasta, start=1, end=7) #指定した始点と終点の範囲の配列を抽出
hoge #確認してるだけです
head(table(hoge)) #出現頻度を順番に表示
head(sort(table(hoge), decreasing=T)) #出現頻度上位を表示
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
hoge <- subseq(fasta, start=3, width=7) #指定した始点から一定範囲の配列を抽出
hoge #確認してるだけです
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
hoge <- subseq(fasta, start=5, width=7) #指定した始点から一定範囲の配列を抽出
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
table(subseq(fasta, start=1, width=7))["CGGGCCT"]
table(subseq(fasta, start=2, width=7))["CGGGCCT"]
table(subseq(fasta, start=3, width=7))["CGGGCCT"]
table(subseq(fasta, start=4, width=7))["CGGGCCT"]
table(subseq(fasta, start=5, width=7))["CGGGCCT"]
table(subseq(fasta, start=6, width=7))["CGGGCCT"]
table(subseq(fasta, start=7, width=7))["CGGGCCT"]
table(subseq(fasta, start=8, width=7))["CGGGCCT"]
```

```
R Console
> head(table(hoge)) #出現頻度上位
hoge
AAAAAAA AAAAAAC AAAAAAG AAAAAAT AAAAAACA AAAAAACG
      238      26      3      1      1      1
> head(sort(table(hoge), decreasing=T)) #出現頻度上位
hoge
CCCCGGT GCCGGCA GGCCTAT GTGCTTT CCCGGTA GTCACTA
      14995      13663      11403      9206      8811      8451
> table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻度
CGGGCCT
      1804
> length(table(hoge))
[1] 15639
> 4^7
[1] 16384
> |
```

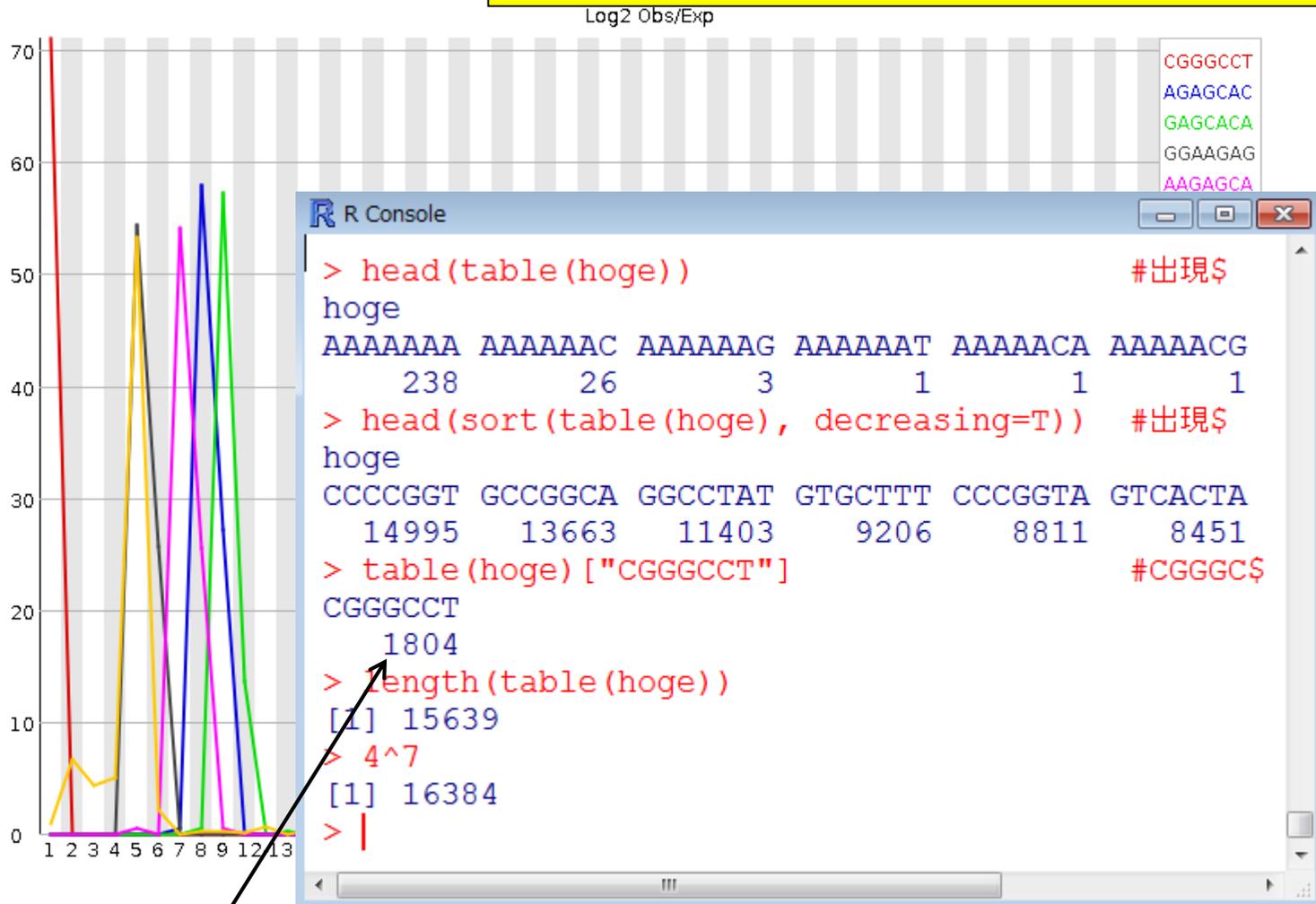


k-mer解析

理由は、FastQC実行結果では、「CGGGCCTがposition 1で235回出現しており、それが期待値(=3.31)の71.00822倍」だったからである

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



| Sequence | Count | PValue | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|--------|-------------|----------------------|
| CGGGCCT | 235 | 0.0 | 71.00822 | 1 |
| AGAGCAC | 835 | 0.0 | 58.01711 | 8 |
| GAGCACA | 845 | 0.0 | 57.330517 | 9 |

FastQCは一部を計算

FastQCはプログラム名通り(半分冗談半分本気)、リーズナブルな計算時間でQC結果を返すことに重きを置いている。全リードの2%のみで解析しているため、Rの結果とは合わない。

 Babraham Bioinformatics

[About](#) | [People](#) | [Services](#) | [Projects](#) | [Training](#) | [Publications](#)

FastQC

Documentation

A copy of the [FastQC](#) documentation is available for you to try before you buy (well download..).

Example Reports

- [Good Illumina](#)
- [Bad Illumina](#)
- [Adapter d](#)
- [Small RN](#)
- [Reduced](#)
- [PacBio](#)
- [454](#)

Index of /projects/fastqc/Help

①

| Name | |
|---|----|
|  Parent Directory | |
|  1 Introduction/ | 25 |
|  2 Basic Operations/ | 25 |
|  3 Analysis Modules/ | 25 |

②

Index of /projects/fastqc/Help/3 Analysis Modules

| Name | Last modified | Size | Description |
|---|-------------------|------|-------------|
|  Parent Directory | | | |
|  1 Basic Statistics.html | 25-Mar-2015 09:40 | 1.8K | |
|  2 Per Base Sequence Quality.html | 25-Mar-2015 09:40 | 3.6K | |
|  3 Per Sequence Quality Scores.html | 25-Mar-2015 09:40 | 1.7K | |
|  4 Per Base Sequence Duplicates.html | | | |
|  5 Per Sequence Duplicates.html | | | |
|  6 Per Base N Content.html | | | |
|  7 Sequence Length Distribution.html | | | |
|  8 Duplicate Sequences.html | | | |
|  9 Overrepresented Sequences.html | 25-Mar-2015 09:40 | 2.4K | |
|  10 Adapter Content.html | 25-Mar-2015 09:40 | 2.4K | |
|  11 Kmer Content.html | 25-Mar-2015 09:40 | 2.5K | |
|  12 Per Tile Sequence Quality.html | 25-Mar-2015 09:40 | 2.2K | |
| duplication_levels.png | 25-Mar-2015 09:40 | 20K | |

③

To allow this module to run in a reasonable time only 2% of the whole library is analysed and the results are extrapolated to the rest of the library. Sequences longer than 500bp are truncated to 500bp for this analysis.

k-mer解析

```
#以下はおまけ(部分配列を抽出し出現頻度上位を概観)
hoge <- subseq(fasta, start=1, end=7) #指定した始点と終点の範囲の
hoge #確認してるだけです
head(table(hoge)) #出現頻度を順番に表示
head(sort(table(hoge), decreasing=T)) #出現頻度上位を表示
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示

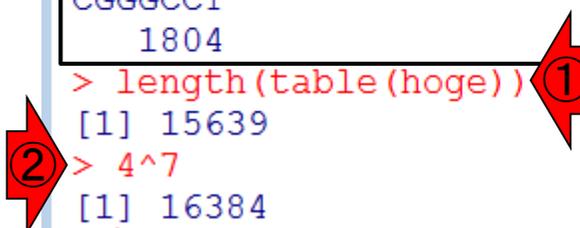
hoge <- subseq(fasta, start=3, width=7) #指定した始点から一定範囲の配列を抽出
hoge #確認してるだけです
table(hoge)["CGGGCCT"] #CGGGCCTの出現頻度を表示

hoge <- subseq(fasta, start=5, width=7) #指定した始点から一定範囲の配列を抽出
table(hoge)["CGGGCCT"] #確認してるだけです

table(subseq(fasta, start=1, width=7))["CGGGCCT"]
table(subseq(fasta, start=2, width=7))["CGGGCCT"]
table(subseq(fasta, start=3, width=7))["CGGGCCT"]
table(subseq(fasta, start=4, width=7))["CGGGCCT"]
table(subseq(fasta, start=5, width=7))["CGGGCCT"]
table(subseq(fasta, start=6, width=7))["CGGGCCT"]
table(subseq(fasta, start=7, width=7))["CGGGCCT"]
table(subseq(fasta, start=8, width=7))["CGGGCCT"]
```

①table(hoge)実行結果の要素数は15,639。それに対し、②理論上の7-merの種類数は $4^7 = 16,384$ 。この結果より、出現回数が0個の7-merが $(16,384 - 15,639) = 745$ 個あったのだろうと推測する。

```
R Console
> head(table(hoge)) #出現$
hoge
AAAAAAA AAAAAAC AAAAAAG AAAAAAT AAAAAACA AAAAAACG
      238      26       3       1       1       1
> head(sort(table(hoge), decreasing=T)) #出現$
hoge
CCCCGGT GCCGGCA GGCCTAT GTGCTTT CCCGGTA GTCACTA
      14995      13663      11403      9206      8811      8451
> table(hoge) ["CGGGCCT"] #CGGGCCT$
CGGGCCT
      1804
> length(table(hoge))
[1] 15639
> 4^7
[1] 16384
> |
```

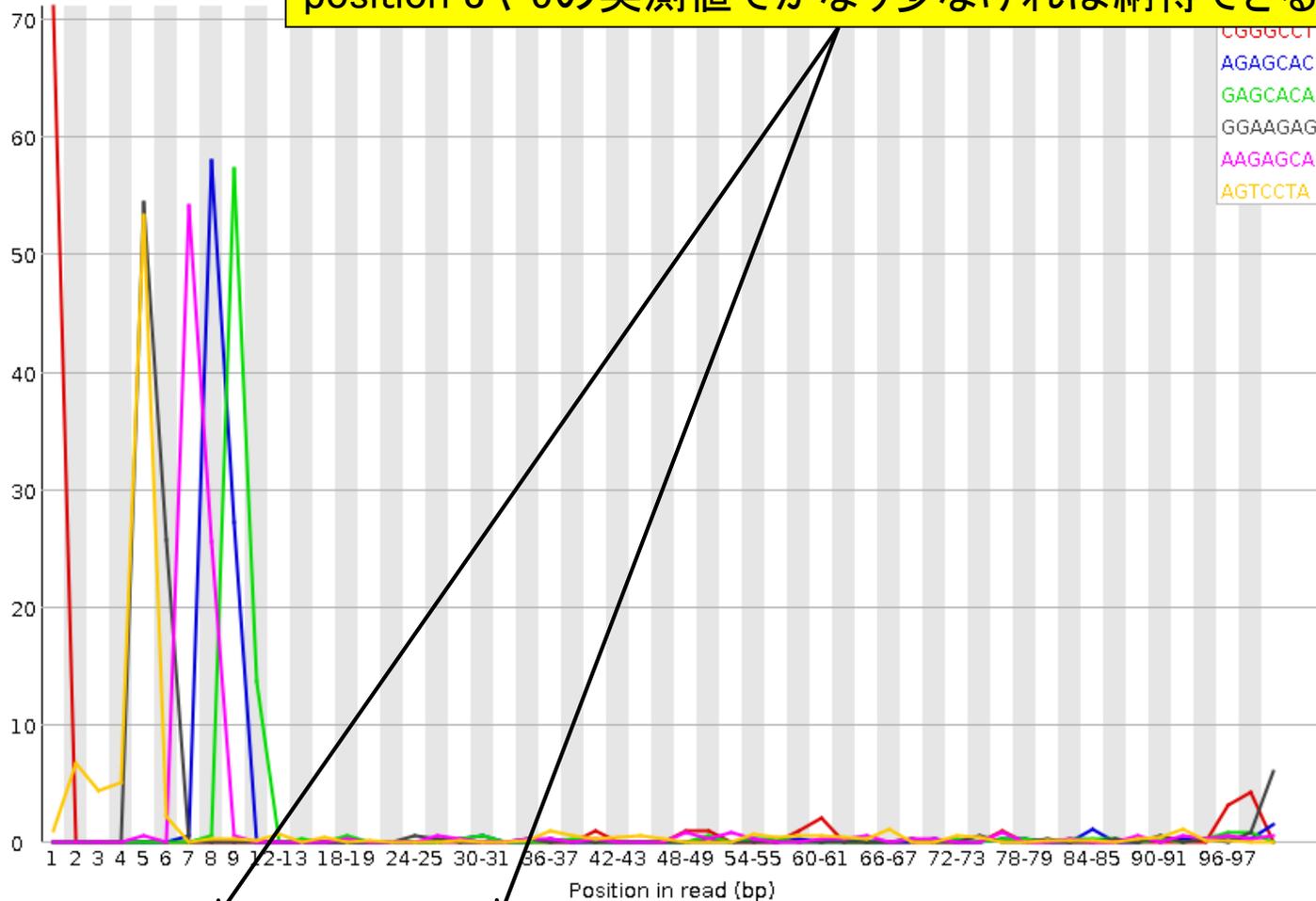


k-mer解析

今ここで調べたいのは、「CGGGCCTがposition 1で235回と、期待値の約71倍多く出現していた」こと。Rの全データを用いた結果はposition 1で1,804回であったことから、position 3や5の実測値でかなり少なければ納得できる。

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



| Sequence | Count | PValue | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|--------|-------------|----------------------|
| CGGGCCT | 235 | 0.0 | 71.00822 | 1 |
| AGAGCAC | 835 | 0.0 | 58.01711 | 8 |
| GAGCACA | 845 | 0.0 | 57.330517 | 9 |

position 3におけるCGGGCCTの実測値を得るためには、赤枠部分の出現頻度を調べればよい。

k-mer解析

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。

```
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") #in_fで指定したファイルの読み込み
fasta #確認してるだけです
```

```
#以下はおまけ(部分配列を抽出し出現頻度)
hoge <- subseq(fasta, start=1, end=100000)
hoge
head(table(hoge))
head(sort(table(hoge), decreasing=T))
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, start=3, width=100000)
hoge
table(hoge)["CGGGCCT"]

hoge <- subseq(fasta, start=5, width=100000)
table(hoge)["CGGGCCT"]
```

```
R Console
> fasta #確認してるだけです
A DNASTringSet instance of length 100000
      width seq names
[1] 107 AGCCCGACTTT...CTTCTCTAAC SRR616268.7 2291:...
[2] 107 GATCTGGGCTG...ACCTCCGAAA SRR616268.20 2291:...
[3] 107 CGGGTATATTT...ACGATCTTAC SRR616268.21 2291:...
[4] 107 CTTCGATACCG...CATAAAGACC SRR616268.22 2291:...
[5] 107 CCGGTATAT...CAAGGTCCAT SRR616268.30 2291:...
... ..
[999996] 107 CCGGTATAT...CAAGGGTTAC SRR616268.1000856...
[999997] 107 TTCGGGTCTAC...CGGGTGAAT SRR616268.1000857...
[999998] 107 CGTCCATCCG...AGCCAAGAGT SRR616268.1000858...
[999999] 107 CTAGGGAGTAT...TTCCTGGACA SRR616268.1000859...
[1000000] 107 GCCTTGTCAT...GACAGCAGCT SRR616268.1000860...
> |
```

k-mer解析

```
hoge <- subseq(fasta, start=3, width=7) #指定した始点から一定範囲の配列を抽出
hoge #確認してるだけです
table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
hoge <- subseq(fasta, start=5, width=7) #指定した始点から一定範囲の配列を抽出
table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
table(subseq(fasta, start=1, width=7))
table(subseq(fasta, start=2, width=7))
table(subseq(fasta, start=3, width=7))
table(subseq(fasta, start=4, width=7))
table(subseq(fasta, start=5, width=7))
table(subseq(fasta, start=6, width=7))
table(subseq(fasta, start=7, width=7))
table(subseq(fasta, start=8, width=7))
table(subseq(fasta, start=9, width=7))
table(subseq(fasta, start=10, width=7))
table(subseq(fasta, start=11, width=7))
table(subseq(fasta, start=12, width=7))
table(subseq(fasta, start=13, width=7))
table(subseq(fasta, start=14, width=7))
table(subseq(fasta, start=15, width=7))
```

k=7と部分配列長がfixされている場合は、subseq関数利用時にwidthオプションを用いるほうがシンプル。
position 3におけるCGGGCCTの出現回数は41回。position 1の1,804回に比べて大幅に少ない。

```
R Console
> hoge <- subseq(fasta, start=3, width=7) #指定した始点から$
> hoge #確認してるだけで$
A DNASTringSet instance of length 1000000
      width seq
[1] 7 CCCGACT SRR616268.7 2291:...
[2] 7 TCTGGGC SRR616268.20 2291...
[3] 7 GGTATAT SRR616268.21 2291...
[4] 7 TCGATAC SRR616268.22 2291...
[5] 7 CCGGTAT SRR616268.30 2291...
... ..
[999996] 7 CCGGTAT SRR616268.1000856...
[999997] 7 CGGGTCT SRR616268.1000857...
[999998] 7 TCCATCC SRR616268.1000858...
[999999] 7 AGGGAGT SRR616268.1000859...
[1000000] 7 CTTGTCA SRR616268.1000860...
> table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻$
CGGGCCT
      41
> |
```

k-mer解析

position 5におけるCGGGCCTの出現回数は3回。position 1の1,804回に比べて大幅に少ない。このように動作確認をしつつ、徐々に必要最小限の記述に進化させていく。

```
hoge <- subseq(fasta, start=5, width=7) #指定した始点から一定範囲の配列を抽出
table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻度を表示
```

```
table(subseq(fasta, start=1, width=7)) ["CGGGCCT"] #一気に計算
table(subseq(fasta, start=2, width=7)) ["CGGGCCT"] #一気に計算
table(subseq(fasta, start=3, width=7)) ["CGGGCCT"] #一気に計算
table(subseq(fasta, start=4, width=7)) ["CGGGCCT"] #一気に計算
```

```
table(subseq(fasta, start=5, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=6, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=7, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=8, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=9, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=10, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=11, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=12, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=13, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=14, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=15, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=16, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=17, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=18, width=7)) ["CGGGCCT"]
table(subseq(fasta, start=19, width=7)) ["CGGGCCT"]
```

```
R Console
[1] 7 CCGGACT SRR616268.7 2291:...
[2] 7 TCTGGGC SRR616268.20 2291...
[3] 7 GGTATAT SRR616268.21 2291...
[4] 7 TCGATAC SRR616268.22 2291...
[5] 7 CCGGTAT SRR616268.30 2291...
... ..
[999996] 7 CCGGTAT SRR616268.1000856...
[999997] 7 CGGGTCT SRR616268.1000857...
[999998] 7 TCCATCC SRR616268.1000858...
[999999] 7 AGGGAGT SRR616268.1000859...
[1000000] 7 CTTGTCA SRR616268.1000860...
> table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻度
CGGGCCT
41
> hoge <- subseq(fasta, start=5, width=7) #指定した始点から$
> table(hoge) ["CGGGCCT"] #CGGGCCTの出現頻度
CGGGCCT
3
> |
```

k-mer解析

一気にコピペで計算。start=1, ..., 35までのCGGGCCTの出現頻度。start=1-3あたりまでは高いがそれ以降は10前後になっていることがわかる。

The image shows a series of overlapping R Console windows. Each window displays the command `> table(subseq(fasta, start=X, width=7))["CGGGCCT"]$` and its output. The windows are arranged from left to right, showing the results for start positions 1 through 35. The frequency of 'CGGGCCT' starts at 1804 for start=1 and decreases to 13 for start=35. Arrows from the yellow text box point to the first few windows.

| start | frequency |
|-------|-----------|
| 1 | 1804 |
| 2 | 3 |
| 3 | 69 |
| 4 | 6 |
| 5 | 41 |
| 6 | 4 |
| 7 | 12 |
| 8 | 9 |
| 9 | 11 |
| 10 | 4 |
| 11 | 3 |
| 12 | 7 |
| 13 | 4 |
| 14 | 1 |
| 15 | 18 |
| 16 | 1 |
| 17 | 8 |
| 18 | 4 |
| 19 | 6 |
| 20 | 3 |
| 21 | 3 |
| 22 | 6 |
| 23 | 3 |
| 24 | 1 |
| 25 | 1 |
| 26 | 10 |
| 27 | 4 |
| 28 | 6 |
| 29 | 8 |
| 30 | 2 |
| 31 | 2 |
| 32 | 3 |
| 33 | 9 |
| 34 | 7 |
| 35 | 13 |

k-mer解析

start=36, ..., 70までのCGGGCCTの出現頻度。このあたりも10前後と低い出現回数。

```

R Console
> table(fasta, start=36, width=7)
CGGGCCT
12
> table(fasta, start=37, width=7)
CGGGCCT
4
> table(fasta, start=38, width=7)
CGGGCCT
9
> table(fasta, start=39, width=7)
CGGGCCT
13
> table(fasta, start=40, width=7)
CGGGCCT
12
> table(fasta, start=41, width=7)
CGGGCCT
4
> table(fasta, start=42, width=7)
CGGGCCT
6
> table(fasta, start=43, width=7)
CGGGCCT
5
> table(fasta, start=44, width=7)
CGGGCCT
8
> table(fasta, start=45, width=7)
CGGGCCT
2
> table(fasta, start=46, width=7)
CGGGCCT
17
> table(fasta, start=47, width=7)
CGGGCCT
7
> table(fasta, start=48, width=7)
CGGGCCT
3
> |

R Console
> table(fasta, start=36, width=7)
CGGGCCT
4
> table(fasta, start=37, width=7)
CGGGCCT
9
> table(fasta, start=38, width=7)
CGGGCCT
7
> table(fasta, start=39, width=7)
CGGGCCT
14
> table(fasta, start=40, width=7)
CGGGCCT
14
> table(fasta, start=41, width=7)
CGGGCCT
3
> table(fasta, start=42, width=7)
CGGGCCT
2
> table(fasta, start=43, width=7)
CGGGCCT
45
> table(fasta, start=44, width=7)
CGGGCCT
6
> table(fasta, start=45, width=7)
CGGGCCT
22
> table(fasta, start=46, width=7)
CGGGCCT
8
> table(fasta, start=47, width=7)
CGGGCCT
6
> table(fasta, start=48, width=7)
CGGGCCT
16
> table(fasta, start=49, width=7)
CGGGCCT
8
> table(fasta, start=50, width=7)
CGGGCCT
16
> table(fasta, start=51, width=7)
CGGGCCT
9
> table(fasta, start=52, width=7)
CGGGCCT
16
> |

R Console
> table(subseq(fasta, start=64, width=7) ["CGGGCCT"])$
CGGGCCT
14
> table(subseq(fasta, start=65, width=7) ["CGGGCCT"])$
CGGGCCT
12
> table(subseq(fasta, start=66, width=7) ["CGGGCCT"])$
CGGGCCT
6
> table(subseq(fasta, start=67, width=7) ["CGGGCCT"])$
CGGGCCT
45
> table(subseq(fasta, start=68, width=7) ["CGGGCCT"])$
CGGGCCT
6
> table(subseq(fasta, start=69, width=7) ["CGGGCCT"])$
CGGGCCT
8
> table(subseq(fasta, start=70, width=7) ["CGGGCCT"])$
CGGGCCT
8
> |
    
```

k-mer解析

start=71, ..., 101までのCGGGCCTの出現頻度。start=96-98あたりで出現頻度が増加していることがわかる。

```

> table(fasta, start=71, width=7)
CGGGCCT
23
> table(fasta, start=72, width=7)
CGGGCCT
2
> table(fasta, start=73, width=7)
CGGGCCT
4
> table(fasta, start=74, width=7)
CGGGCCT
1
> table(fasta, start=75, width=7)
CGGGCCT
7
> table(fasta, start=76, width=7)
CGGGCCT
2
> table(fasta, start=77, width=7)
CGGGCCT
9
> table(fasta, start=78, width=7)
CGGGCCT
4
> table(fasta, start=79, width=7)
CGGGCCT
5
> table(fasta, start=80, width=7)
CGGGCCT
6
> table(fasta, start=81, width=7)
CGGGCCT
3
> table(fasta, start=82, width=7)
CGGGCCT
5
> table(fasta, start=83, width=7)
CGGGCCT
2
> |
> table(fasta, start=84, width=7)
CGGGCCT
3
> table(fasta, start=85, width=7)
CGGGCCT
8
> table(fasta, start=86, width=7)
CGGGCCT
3
> table(fasta, start=87, width=7)
CGGGCCT
12
> table(fasta, start=88, width=7)
CGGGCCT
7
> table(fasta, start=89, width=7)
CGGGCCT
5
> table(fasta, start=90, width=7)
CGGGCCT
4
> table(fasta, start=91, width=7)
CGGGCCT
4
> table(fasta, start=92, width=7)
CGGGCCT
5
> table(fasta, start=93, width=7)
CGGGCCT
1
> table(fasta, start=94, width=7)
CGGGCCT
50
> table(fasta, start=95, width=7)
CGGGCCT
13
> table(fasta, start=96, width=7)
CGGGCCT
68
> table(fasta, start=97, width=7)
CGGGCCT
124
> table(fasta, start=98, width=7)
CGGGCCT
124
> table(fasta, start=99, width=7)
CGGGCCT
22
> table(fasta, start=100, width=7)
CGGGCCT
5
> table(fasta, start=101, width=7)
CGGGCCT
4
>

```

k-mer解析

position 1, 2, ..., 101までのCGGGCCTの実測した出現頻度をベクトルとしてObsというオブジェクトに格納。①head関数でベクトルObsの最初の6個の要素を表示。②length関数でベクトルObsの要素数(=101)を表示。

```

R Console
> table(CGGGCCT)
23
> table(CGGGCCT)
2
> table(CGGGCCT)
4
> table(CGGGCCT)
1
> table(CGGGCCT)
7
> table(CGGGCCT)
2
> table(CGGGCCT)
9
> table(CGGGCCT)
4
> table(CGGGCCT)
5
> table(CGGGCCT)
6
> table(CGGGCCT)
3
> table(CGGGCCT)
2
> |

R Console
> table(CGGGCCT)
8
> table(CGGGCCT)
3
> table(CGGGCCT)
7
> table(CGGGCCT)
5
> table(CGGGCCT)
4
> table(CGGGCCT)
4
> table(CGGGCCT)
5
> table(CGGGCCT)
1
> table(CGGGCCT)
5
> table(CGGGCCT)
13
> table(CGGGCCT)
2
> |

R Console
> table(CGGGCCT)
8
> table(CGGGCCT)
3
> table(CGGGCCT)
7
> table(CGGGCCT)
5
> table(CGGGCCT)
4
> table(CGGGCCT)
5
> table(CGGGCCT)
1
> table(CGGGCCT)
50
> table(CGGGCCT)
13
> table(CGGGCCT)
68
> table(CGGGCCT)
8
> |

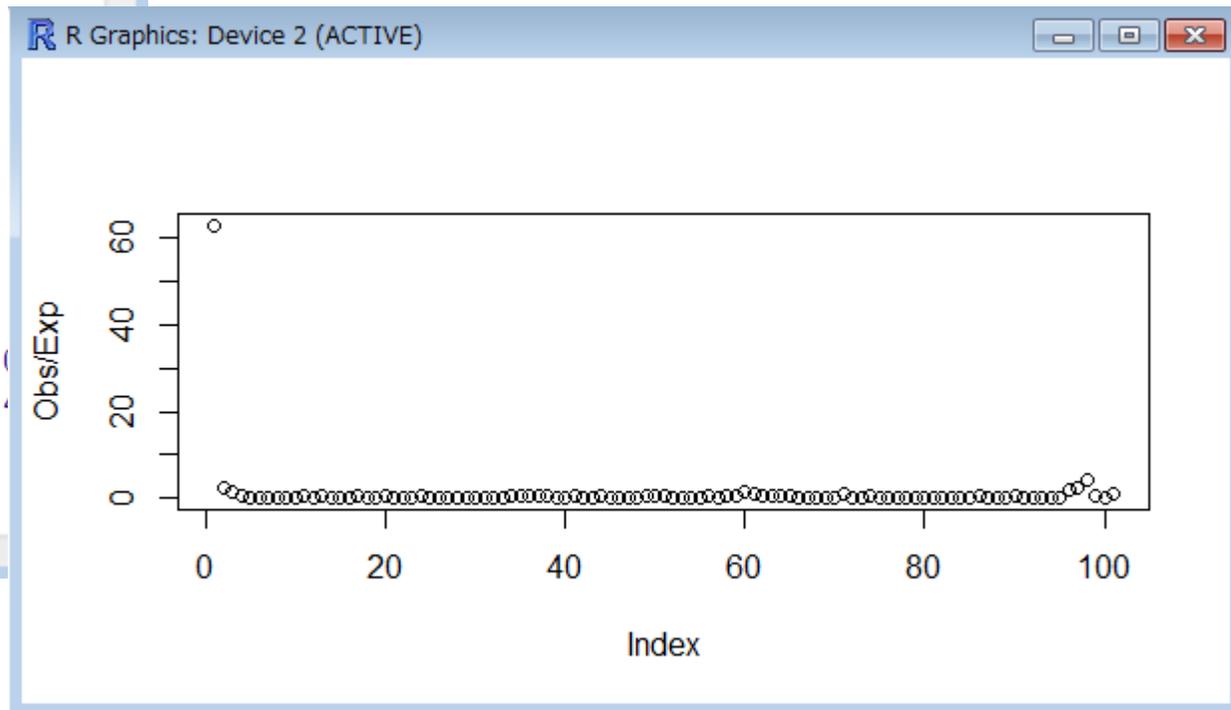
R Console
> table(CGGGCCT)
12
> table(CGGGCCT)
22
> |

R Console
> Obs <- c(
+ 1804, 69, 41, 9, 3, 1, 1, # position1-7
+ 3, 6, 4, 11, 7, 18, 8, # position8-14
+ 6, 7, 12, 4, 4, 15, 4, # position15-21
+ 8, 2, 9, 3, 3, 1, 6, # position22-28
+ 2, 3, 7, 6, 1, 10, 13, # position29-35
+ 12, 9, 12, 6, 8, 17, 3, # position36-42
+ 4, 13, 4, 5, 2, 7, 10, # position43-49
+ 9, 10, 3, 2, 6, 6, 9, # position50-56
+ 7, 14, 14, 45, 22, 16, 16, # position57-63
+ 14, 12, 6, 6, 8, 8, 7, # position64-70
+ 23, 4, 7, 9, 5, 3, 2, # position71-77
+ 2, 1, 2, 4, 6, 5, 3, # position78-84
+ 8, 9, 5, 4, 1, 13, 8, # position85-91
+ 3, 7, 4, 5, 50, 68, 124, # position92-98
+ 12, 5, 22) # position99-101
> head(Obs) ①
[1] 1804 69 41 9 3 1
> length(Obs) ②
[1] 101
> |
    
```

k-mer解析

①平均出現頻度(=期待値; ②Exp)は25.58、
 ③position 1の実測値(=1,804)をExpで割った
 ものがFastQCのKmer_Contentという項目。
 ここでは63.11という値が得られている。④R
 上ではObs/Expとやるだけでよい。⑤plotの
 基本形で全体像を把握。

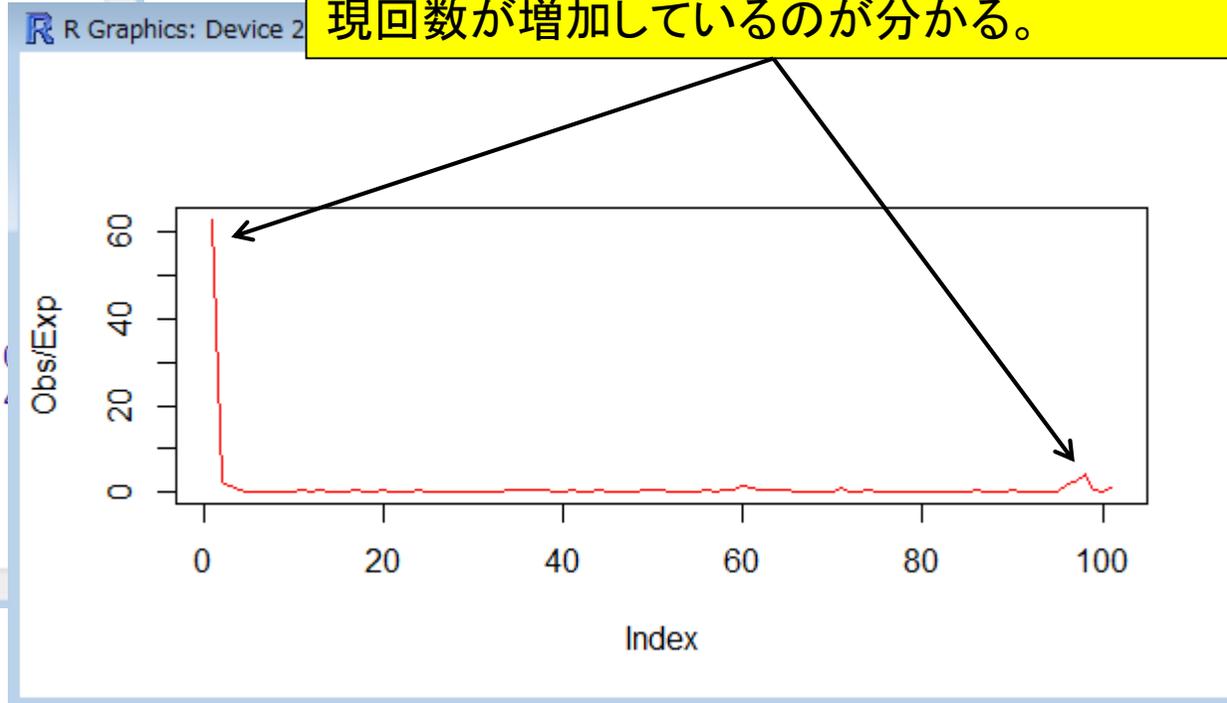
```
R Console
> head(Obs)
[1] 1804 69 41 9 3 1
> length(Obs)
[1] 101
> mean(Obs) ①
[1] 28.58416
> Exp <- mean(Obs) ②
> 1804/Exp ③
[1] 63.11188
> head(Obs)
[1] 1804 69 41 9 3 1
> head(Obs/Exp) ④
[1] 63.11188085 2.41392449 1.43436
[4] 0.31485972 0.10495324 0.03498
> plot(Obs/Exp) ⑤
> |
```



k-mer解析

①平均出現頻度(=期待値; ②Exp)は25.58、
③position 1の実測値(=1,804)をExpで割った
ものがFastQCのKmer_Contentという項目。
ここでは63.11という値が得られている。④R
上ではObs/Expとやるだけでよい。⑤plotの
基本形で全体像を把握。⑥赤の折れ線グラ
フで再度プロット。position 1と96-98周
辺で出現回数が増加しているのが分かる。

```
R Console
> head(Obs)
[1] 1804 69 41 9 3 1
> length(Obs)
[1] 101
> mean(Obs) ①
[1] 28.58416
> Exp <- mean(Obs) ②
> 1804/Exp ③
[1] 63.11188
> head(Obs)
[1] 1804 69 41 9 3 1
> head(Obs/Exp) ④
[1] 63.11188085 2.41392449 1.43436
[4] 0.31485972 0.10495324 0.03498
> plot(Obs/Exp) ⑤
⑥ plot(Obs/Exp, type="l", col="red")
>
```

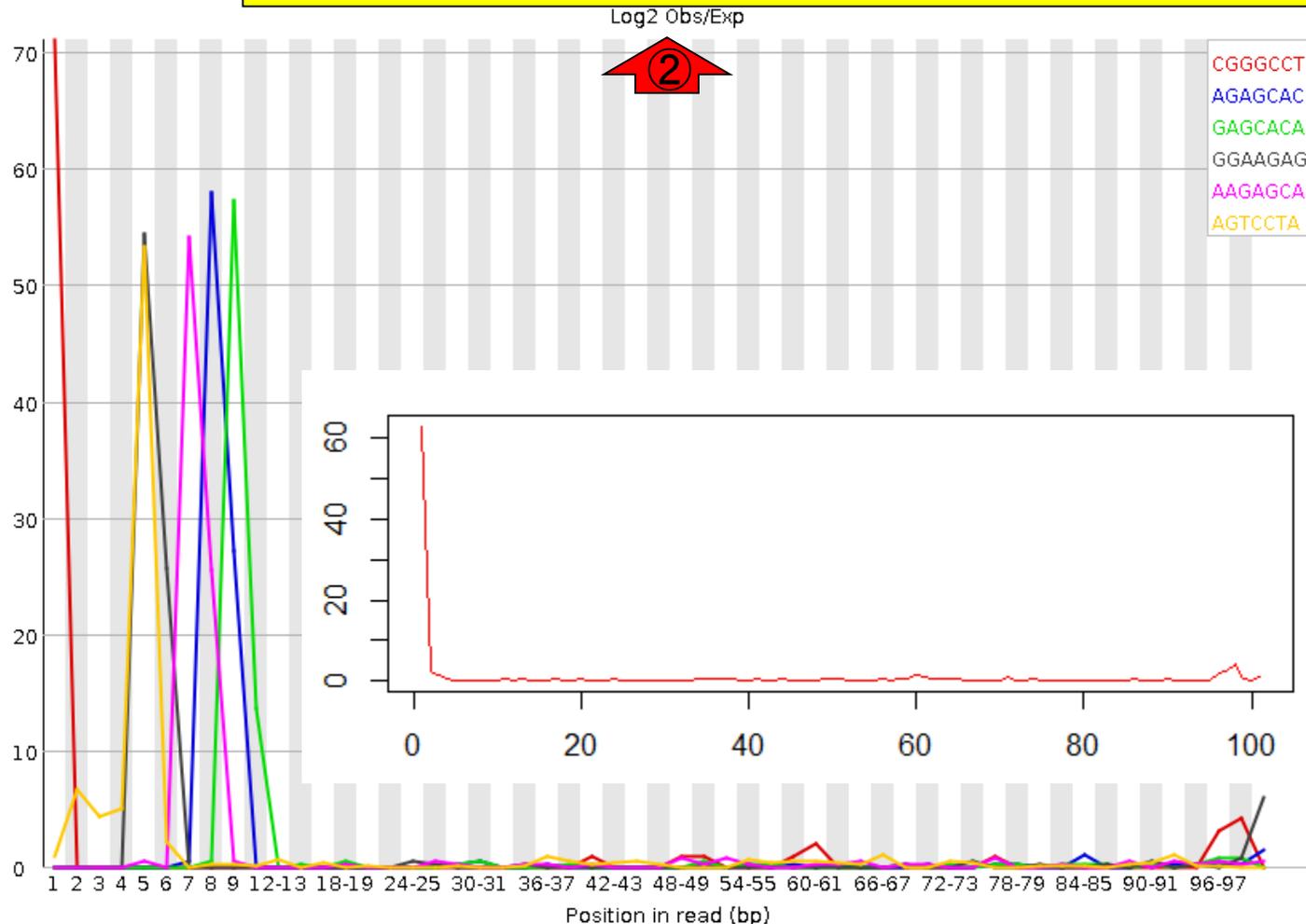


k-mer解析

①FastQCのKmer_Content中のCGGGCCTの結果とよく似ていることが分かります。②FastQC中では「Log2 Obs/Exp」となっていますが、おそらくこれは「Obs/Exp」が正解です。

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**



| Sequence | Count | PValue | Obs/Exp Max | Max Obs/Exp Position |
|----------|-------|--------|-------------|----------------------|
| CGGGCCT | 235 | 0.0 | 71.00822 | 1 |
| AGAGCAC | 835 | 0.0 | 58.01711 | 8 |
| GAGCACA | 845 | 0.0 | 57.330517 | 9 |

Tips: オブジェクトの消去

①オブジェクトの消去はrm関数を用いる。rm(Obs)実行後はObsが存在しないので、②head(Obs)をやっても正しくErrorが出る。③このR Console画面上で利用している全オブジェクト(in_f, Exp, fasta, hogeなど)を一度に消去する場合は「rm(list = ls())」。ここでは後の作業に支障があるのでやらないで。

```
R Console
> head(Obs)
[1] 1804 69 41 9 3 1
> length(Obs)
[1] 101
> mean(Obs)
[1] 28.58416
> Exp <- mean(Obs)
> 1804/Exp
[1] 63.11188
> head(Obs)
[1] 1804 69 41 9 3 1
> head(Obs/Exp)
[1] 63.11188085 2.41392449 1.43436093
[4] 0.31485972 0.10495324 0.03498441
> plot(Obs/Exp)
> plot(Obs/Exp, type="l", col="red")
> rm(Obs) ①
> head(Obs) ②
Error in head(Obs) :
  引数 'x' の評価中にエラーが起きました
> rm(list = ls()) ③
```

Tips: forループ

forループで一気にCGGGCCTのposition 1-101の出現回数を計算することもできる。これはforループの基本形(初級)

- ・ [イントロ](#) | [NGS](#) | [アノテーション情報取得](#) | [TranscriptDb](#) | [GFF/GTF形式ファイルから](#)
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [BSgenome](#) | [基本情報を取得](#) (last modified 2015/06/10)
- ・ [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 2015/06/19) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [基礎](#) **①** (last modified 2015/06/19) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [応用](#) (last modified 2015/06/18) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [description行の記述を整形](#) (last modified 2015/06/19) **NEW**

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

② 8. FASTQ形式ファイル(SRR616268sub_1.fastq.gz)の場合:

Sanger FASTQ形式ファイルを読み込みます。入力は「ファイル」メニューから行います。

1. サンプルデータ SRR0374397 (SRR616268sub_1.fastq.gz)の最初の100万リード分(約73MB)です。長さは全て107 bpです。

```

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

#入力ファイル名を指定してfastaに格納
fasta <- readFastq(in_f)

#以下はおまけ(forループを用いて美しく...初級)
param_kmer <- 7 #k-merのkの値を指定
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
for(i in 1:101){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=param_kmer))[param_obj])
}
head(Obs)
mean(Obs)
Exp <- mean(Obs)
head(Obs/Exp)
plot(Obs/Exp, type="l", col="red")
    
```

Tips: forループ

「手計算?!で得られた結果」と「forループを用いて得られた結果」の値が同じであることを頼りにこのような記述形式であっていることを確認。tail(Obs)もやったほうがいいかもしれない。

```
table(subseq(fasta, start=99, width=7))["CGGGCCT"]#一気に計算
table(subseq(fasta, start=100, width=7))["CGGGCCT"]#一気に計算
table(subseq(fasta, start=101, width=7))["CGGGCCT"]#一気に計算
```

#以下はおまけ(forループを用いて美しく...初級)

```
param_kmer <- 7 #k-merのkの値を指定
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
for(i in 1:101){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=param_kmer)))
}
head(Obs)
mean(Obs)
Exp <- mean(Obs)
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...中級)

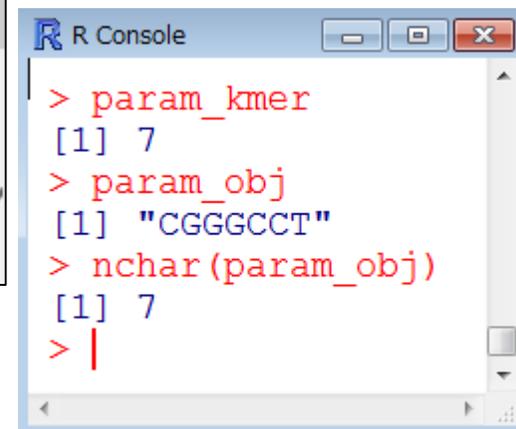
```
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
for(i in 1:101){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=nchar(param_obj))))
}
```

```
R Console
> #以下はおまけ (forループを用いて美しく$
> param_kmer <- 7 $
> param_obj <- "CGGGCCT" $
> Obs <- NULL $
> for(i in 1:101){ $
+   Obs <- c(Obs, table(subseq(fasta,$
+ })
> head(Obs)
CGGGCCT CGGGCCT CGGGCCT CGGGCCT CGGGCCT
 1804      69      41      9      3
CGGGCCT
 1
> mean(Obs)
[1] 28.58416
> Exp <- mean(Obs)
> head(Obs/Exp)
      CGGGCCT      CGGGCCT      CGGGCCT
63.11188085  2.41392449  1.43436093
      CGGGCCT      CGGGCCT      CGGGCCT
 0.31485972  0.10495324  0.03498441
> plot(Obs/Exp, type="l", col="red")
>
```

Tips: スキルアップ

```
#以下はおまけ(forループを用いて美しく...初級)
param_kmer <- 7 #k-merのkの値を指定
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
for(i in 1:101){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=param_kmer))[param_obj])
}
head(Obs)
mean(Obs)
Exp <- mean(Obs)
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...中級)
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
for(i in 1:101){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=nchar(param_obj)))[param_obj])
}
head(Obs)
mean(Obs)
Exp <- mean(Obs)
```



```
R Console
> param_kmer
[1] 7
> param_obj
[1] "CGGGCCT"
> nchar(param_obj)
[1] 7
> |
```

Tips: スキルアップ

```
#以下はおまけ(forループを用いて美しく...中級)
param_obj <- "CGGGCCT" #調べたいk-merを指定
Obs <- NULL #おまじない
for(i in 1:101){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=nchar(par
})
head(Obs)
mean(Obs)
Exp <- mean(Obs)
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hogelに
for(i in 1:hoge){ #ループを回す
  Obs <- c(Obs, table(subseq(fasta, start=i, width=nchar(param_obj)))[param
})
head(Obs)
mean(Obs)
```

中級 → 上級1。なるべくプログラム本体を気にすることなく、パラメータ部分をいじるだけで実行させることを心掛ける。中級は、k=7以外の際に101の数値を変更しなければいけないという問題点がある。この101という数値は、「107 - 7 + 1」として計算できる。つまり、リード長Lとkの値が決まればおのずと「L - k + 1」で決まるということ。一見ややこしいが、このほうがパラメータ部分のみしか気になくていいので後々心穏やかになれる。

```
R Console
> hoge <- param_len_ngs - nchar(param_obj) + 1
> hoge
[1] 101
> |
```

Tips: スキルアップ

上級1 ⇔ 上級2。107 bpというリード長はfastaオブジェクトのwidth列、つまりwidth(fasta)部分で取得可能。全リードの長さは一定という前提条件のもと、ここではwidth(fasta)[1]として、最初のリード長を代表値として採用している。この前提条件を満たさない場合にエラーが出る。それゆえ、上級1程度にしておくほうが門田好み(蓼食う虫も好き好き; some prefer nettles)。

```
#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)
for(i in 1:hoge){
  Obs <- c(Obs, table(subseq(fasta, st
})
head(Obs)
mean(Obs)
Exp <- mean(Obs)
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
param_obj <- "CGGGCCT"
Obs <- NULL
hoge <- width(fasta)[1] - nchar(param_ob
for(i in 1:hoge){
  Obs <- c(Obs, table(subseq(fasta, st
})
head(Obs)
```

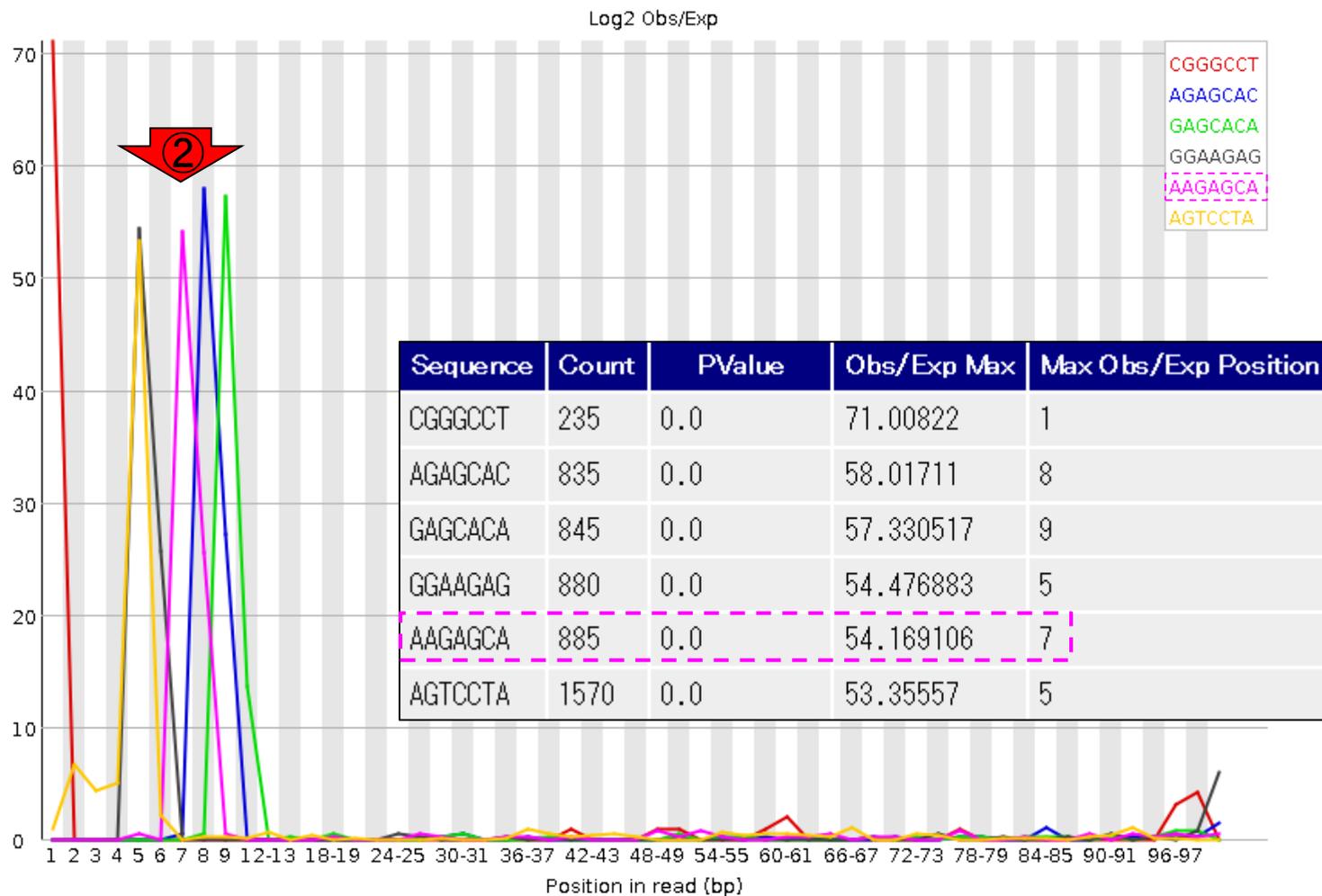
```
> fasta
A DNASTringSet instance of length 1000000
      width seq          names
[1] 107 AGCCCGA...CTCTAAC SRR616268.7 2291:...
[2] 107 GATCTGG...TCCGAAA SRR616268.20 2291...
[3] 107 CCGGTAT...ATCTTAC SRR616268.21 2291...
[4] 107 CTTCGAT...AAAGACC SRR616268.22 2291...
[5] 107 CCCCGGT...GGTCCAT SRR616268.30 2291...
...
[999996] 107 CCCCGGT...GGGTAC SRR616268.1000856...
[999997] 107 TTCGGGT...GTGAAAT SRR616268.1000857...
[999998] 107 CGTCCAT...CAAGAGT SRR616268.1000858...
[999999] 107 CTAGGGA...CTGGACA SRR616268.1000859...
[1000000] 107 GCCTTGT...AGCAGCT SRR616268.1000860...
> width(fasta)[1]
[1] 107
> head(width(fasta))
[1] 107 107 107 107 107 107
> median(width(fasta))
[1] 107
> |
```

「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。

k-mer解析

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



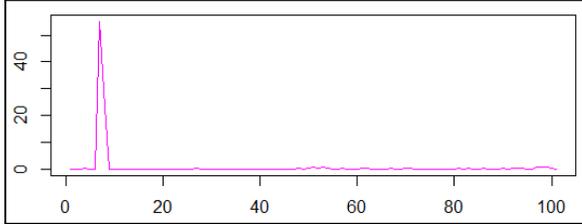
k-mer解析

「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 #リード長を$
> param_obj <- "AAGAGCA" #調べたいk-m$
> Obs <- NULL #おまじない
> hoge <- param_len_ngs - nchar(param_obj) + 1 #positi$
> for(i in 1:hoge){ #ループを回$
+   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")
> |
```

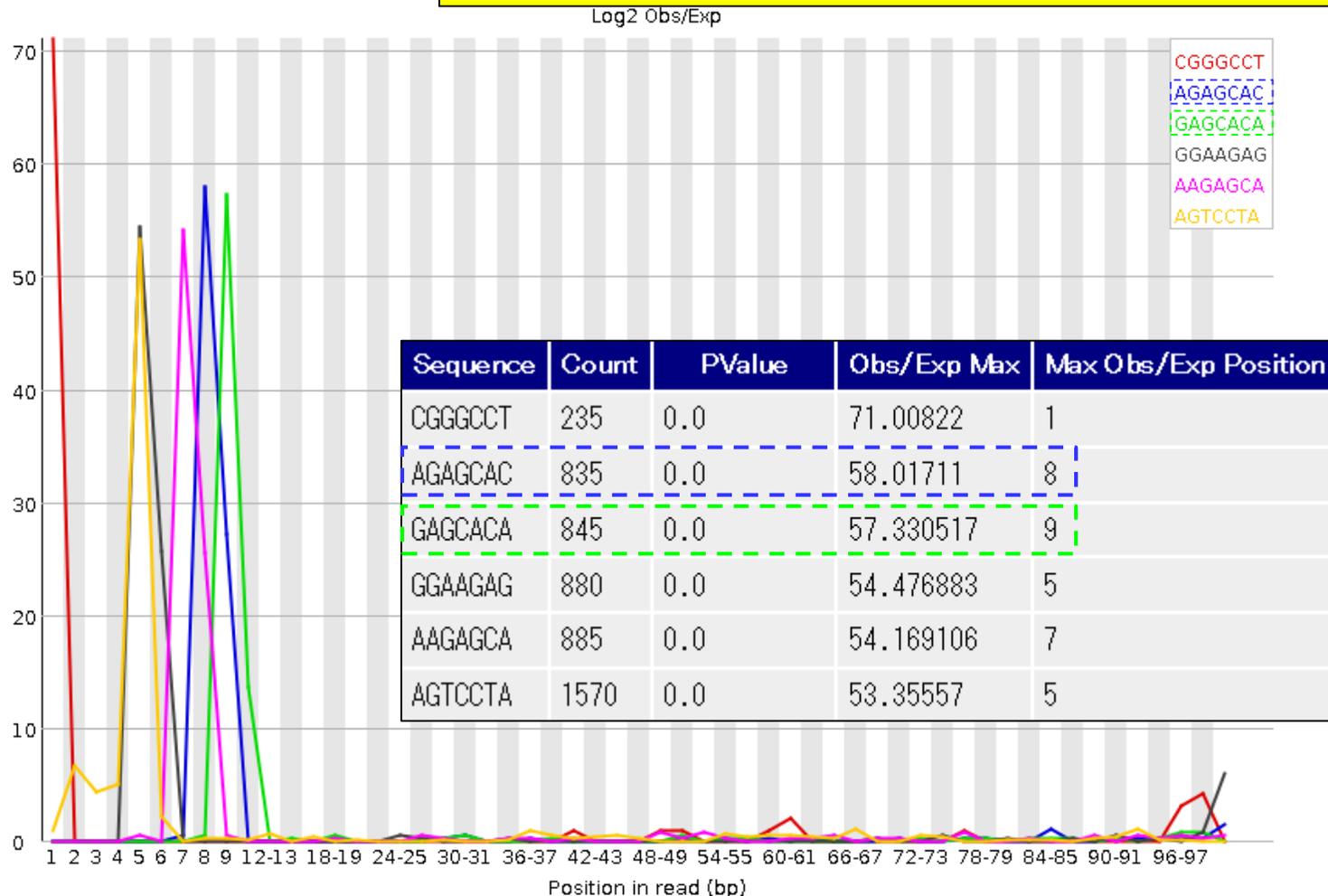


課題1

「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**



Contents

■ QC(Quality Control)

- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

FastQC結果を眺める

カイコsRNA-seqデータのFastQC結果。①総リード数は11,928,428(約1,200万)。②リードの長さは49 bp。③GC含量は52%。④赤枠の、「×や!」になっているものが要チェック項目。信号の黄色や赤色をイメージすればよい

FastQC Report

Summary

- Basic Statistics
- Per base sequence quality
- Per tile sequence quality
- Per sequence quality scores
- Per base sequence content** (4)
- Per sequence GC content** (4)
- Per base N content
- Sequence Length Distribution
- Sequence Duplication Levels** (4)
- Overrepresented sequences** (4)
- Adapter Content
- Kmer Content

Basic Statistics

| Measure | Value |
|-----------------------------------|-------------------------|
| Filename | SRR609266.fastq.gz |
| File type | Conventional base calls |
| Encoding | Illumina 1 |
| Total Sequences | 11928428 (1) |
| Sequences flagged as poor quality | 0 |
| Sequence length | 49 (2) |
| %GC | 52 (3) |

Per base sequence quality

Quality scores across all bases (Illumina 1.5 encoding)

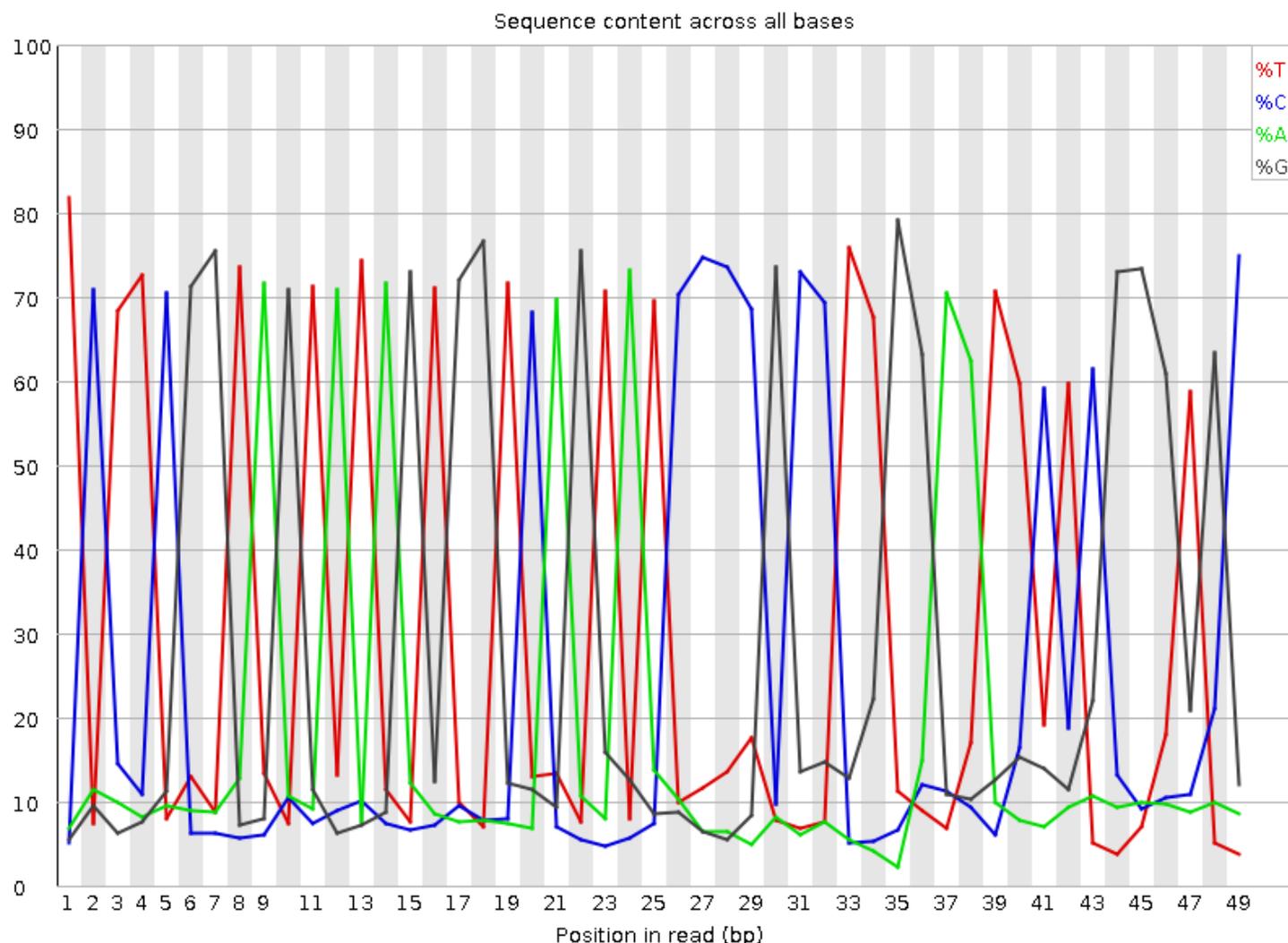
FastQC結果を眺める

①ポジションごとの塩基の出現確率。明らかに変。あたかも全く同じ塩基配列のリード(TCTTCGGTAG...)ばかりが読まれているようなパターンに見える。

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

✗ Per base sequence content



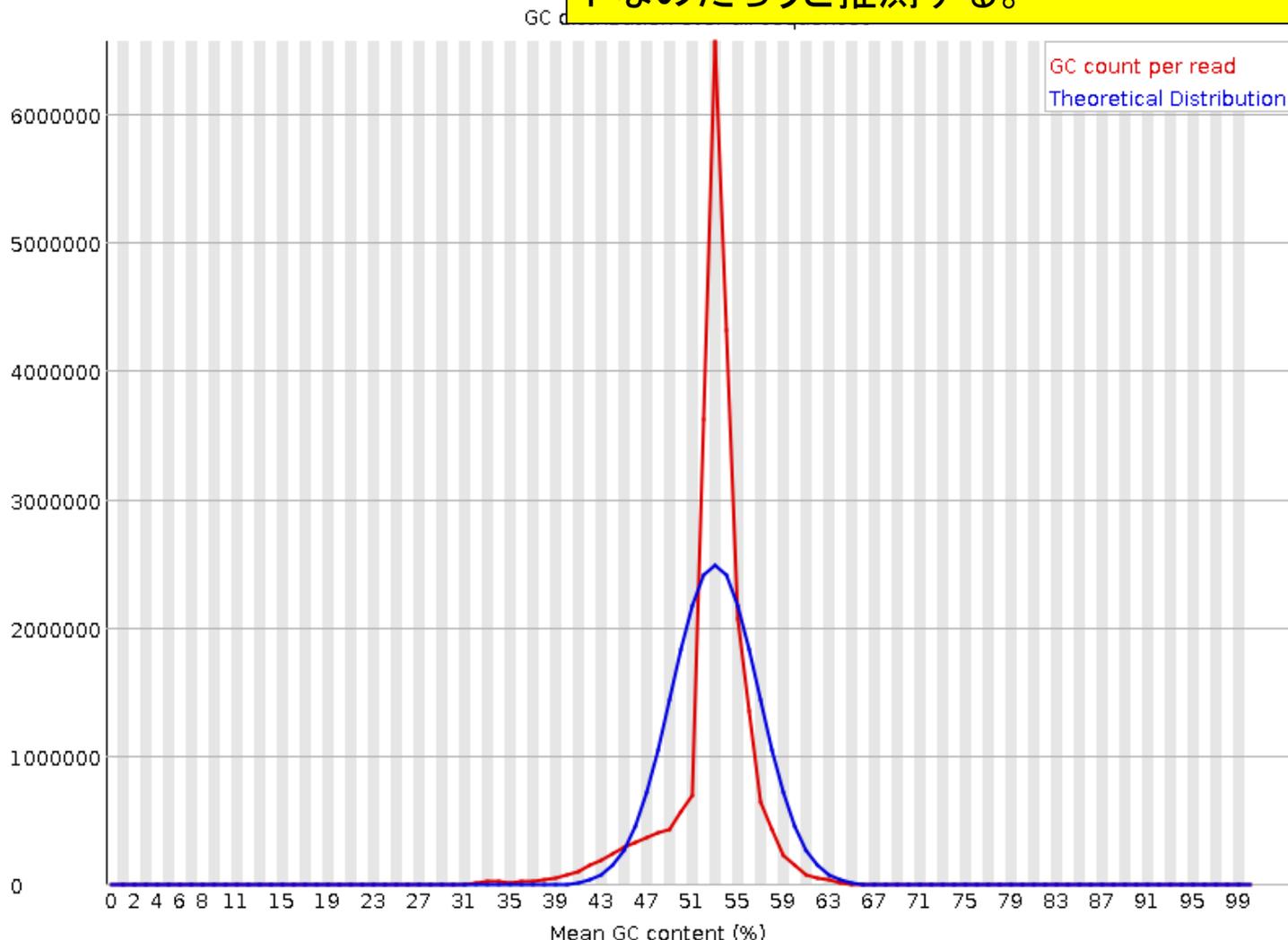
FastQC結果を眺める

おそらく赤色の実測値のピークは52%のところ。TCTTCGGTAG...ばかりが読まれているであろうこと、全部で49 bpであることから25 bp分がG or Cのリードなのだろうと推測する。

Summary

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

✗ Per sequence GC content

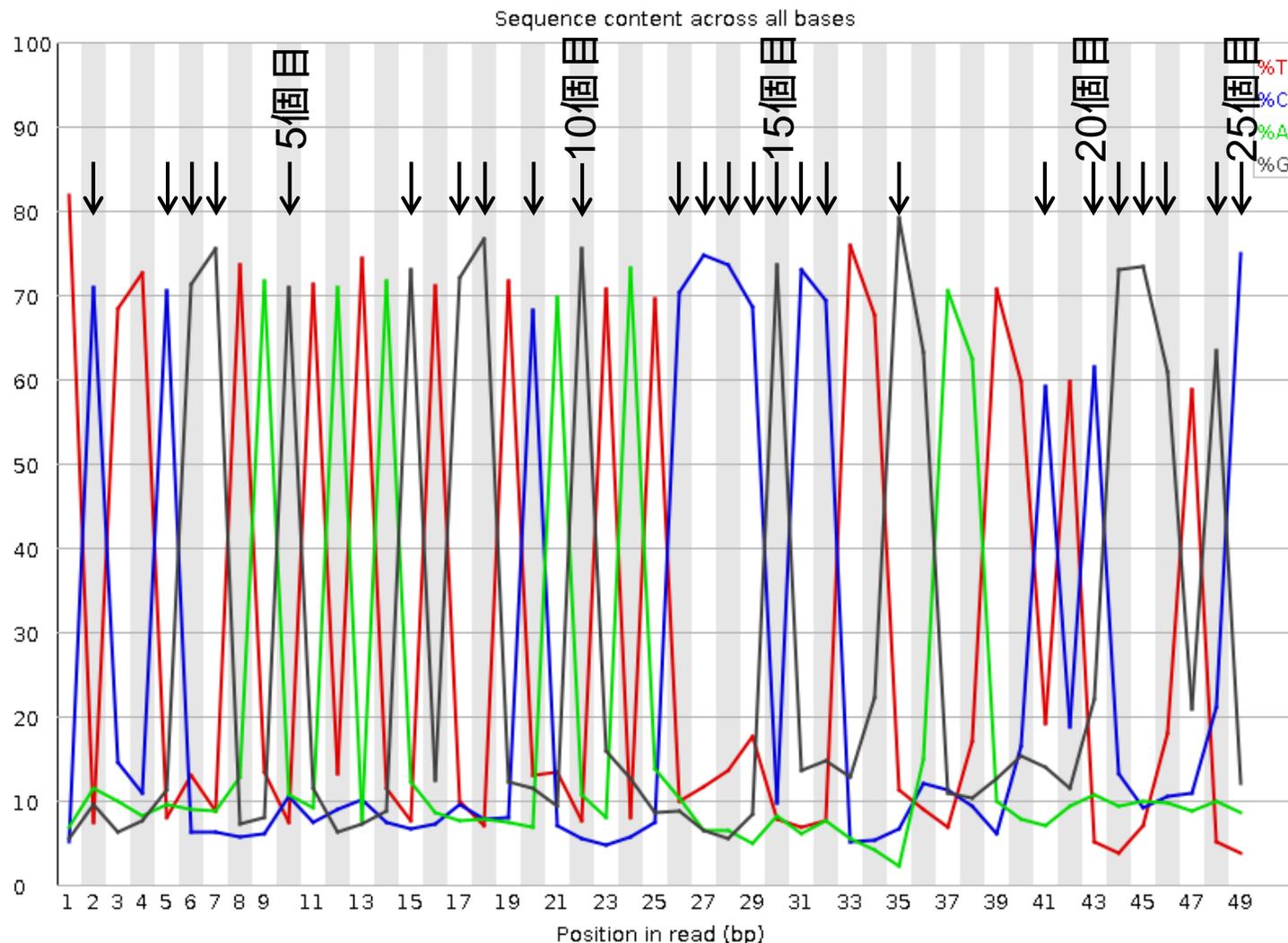


FastQC結果を眺める

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

✗ Per base sequence content



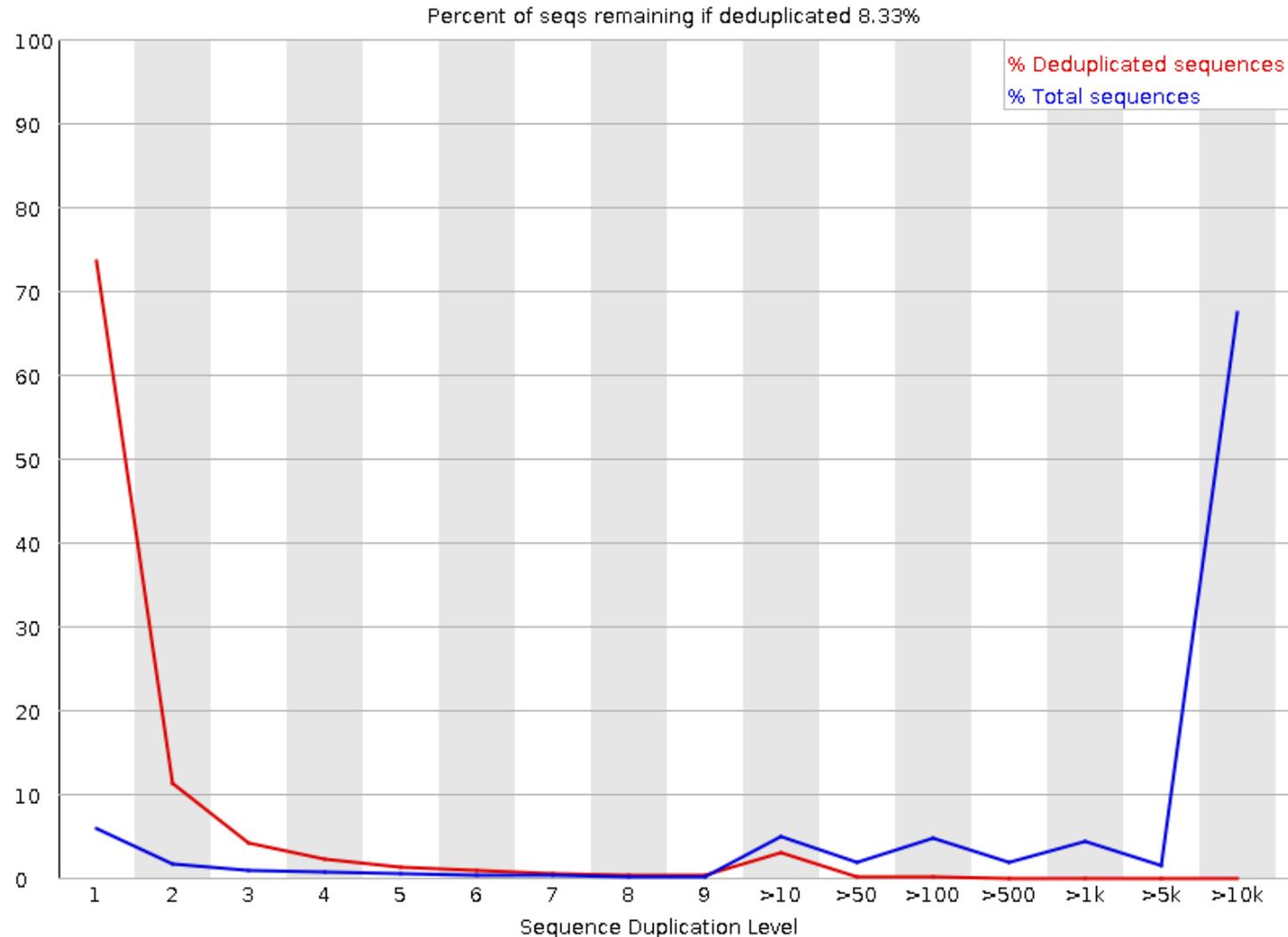
FastQC結果を眺める

私はこの段階でもまだこの解釈が困難です(のちに% Total sequencesは納得した)。
発展課題:% Deduplicated sequencesの意味を説明せよ。

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

✗ Sequence Duplication Levels



1

FastQC結果を眺める

① Overrepresented sequences。総リード数 11,928,428のうち、②赤枠の塩基配列からなるリードが5,646,094個(47.3%)を占める。

Summary

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

Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|--|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164982 | 1.3800000000000001 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.9232000000000001 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6895000000000001 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGCC | 75176 | 0.6305000000000001 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
| TGGAATTCTCGGGTGCCAAGGAACTCCAGTCACATCACGATCTCGTATG | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
| AAATCATTACCCTGGACGGTGGATCACTGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
| TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

```
R Console
> 5646094/11928428
[1] 0.4733309
> |
```



FastQC結果を眺める

もう少し下のほうまで眺める。10,000回以上出現しているリードが結構あることに気づく。そして「10,000は10kと表現」することを思い出し、「Sequence Duplication Levels」の項目の意味についてひらめく。

| | | | |
|---|-------|---------------------|---|
| TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.1676415368 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCCCCCTTGAATTCTCGGGTGC | 19932 | 0.1670966199 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGTGGAATTCTCG | 19650 | 0.16473251965808067 | No Hit |
| GTCAGTCGATCCTAAGCTCAAGGAGAGATGGAATTCTCGGGTGCCAAGG | 19014 | 0.15940071902181913 | Illumina Small RNA Adapter 2 (100% over 21bp) |
| GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTATGGAATTCTCGGGTG | 17862 | 0.14974311786934538 | No Hit |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCGAATTTAAGCTGGAATTC | 16808 | 0.14090708348157863 | No Hit |
| TGGACGGAGAACTGATAAGGGCTGGAATTCTCGGGTGCCAAGGAACTCC | 16488 | 0.13822441649478037 | RNA PCR Primer, Index 1 (100% over 27bp) |
| GGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 15993 | 0.13407466599957682 | No Hit |
| CATTTGGATCGCGGAGATCTGGAATTCTCGGGTGCCAAGGAACTCCAGT | 15322 | 0.1284494486616342 | RNA PCR Primer, Index 1 (100% over 30bp) |
| CTTCGGTAGTATAGTGGTCAGTATCCCCGCTTGAATTCTCGGGTGCC | 15147 | 0.12698236515322892 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 14993 | 0.12569133166583224 | No Hit |
| AGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 14895 | 0.12486976490112528 | No Hit |
| TGAGATCATTGTGAAAGCTAATTGGAATTCTCGGGTGCCAAGGAACTCC | 14165 | 0.11874993083749175 | RNA PCR Primer, Index 1 (100% over 27bp) |
| TATTATGATGACAAACAACTAAGGAACCACTGATTGCATTGGAATTCT | 13934 | 0.11681338060639676 | No Hit |
| TCCTCGGTAGTATAGTGGTCAGTATCCCCGCTTGAATTCTCGGGTGC | 13352 | 0.11193428002415741 | No Hit |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCGAATTTAAGTGAATTCT | 13152 | 0.11025761315740851 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGGTGAATTCTC | 12713 | 0.10657732938489464 | No Hit |

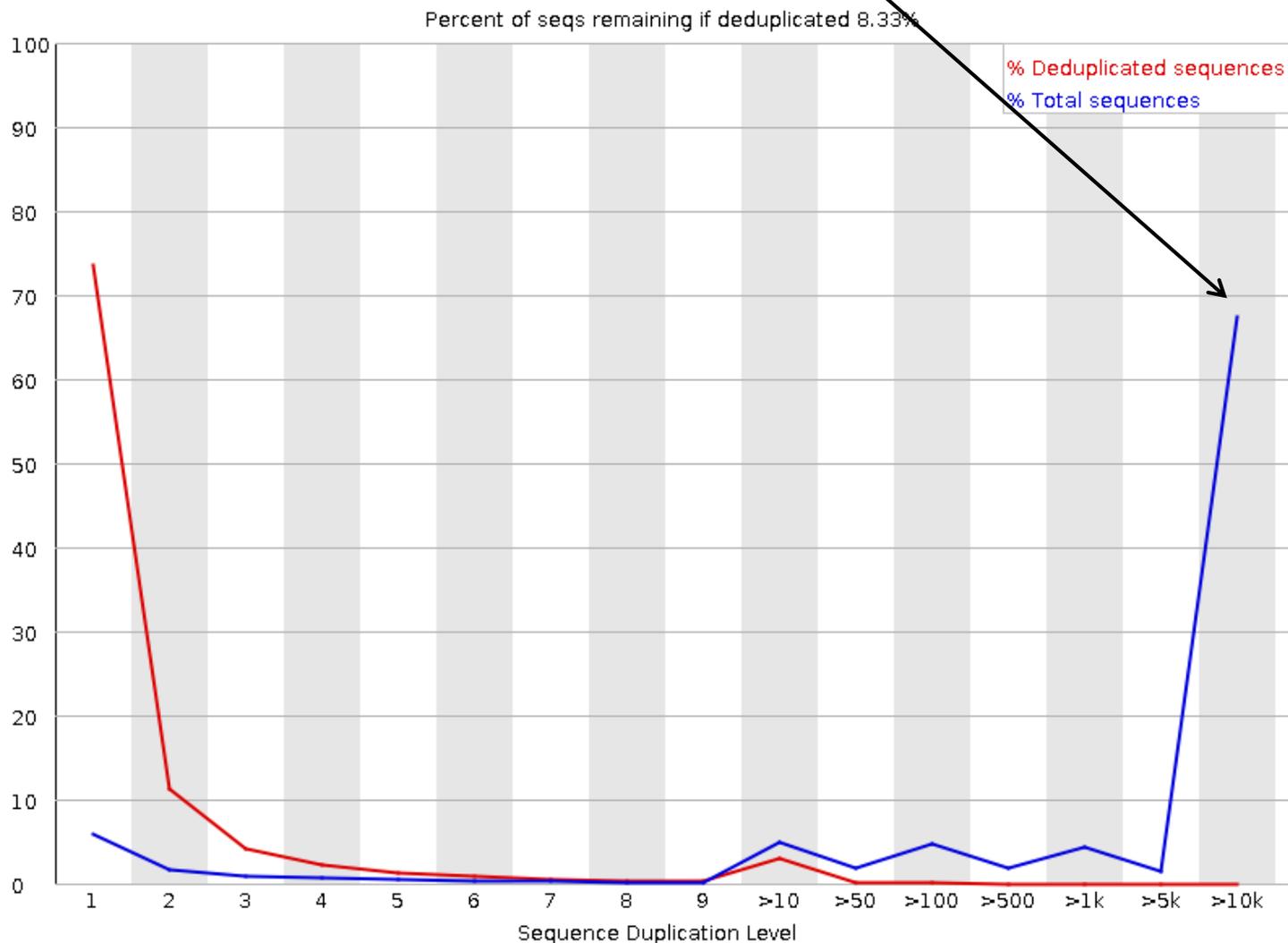
FastQC結果を眺める

つまり、「Sequence Duplication Levelが10kより大きいものが全体(Total sequences; 11,928,428)の70%弱を占めるということを行っているのだろうと予想を立てる。

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 **①**
- ✗ Overrepresented sequences
- ! Adapter Content
- ✗ Kmer Content

✗ Sequence Duplication Level



FastQC結果を眺める

FastQCのhtmlレポートで見られるものはCount列が12713個までだが、それを全部足すと8,033,706個(67.35%)。残りの10,000回出現するリードをまで足したとしても、たかが知れている。つまり、70%は超えることはなく、せいぜい68%程度だと考えればリーズナブル。

| | | | |
|---|-------|-----------|---|
| TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.1676415 | |
| TCCTCGGTAGTATAGTGGTGAGTATGCCCCCTTGAATTCTCGGGTGC | 19932 | 0.1670966 | |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGTGGAATTCTCG | 19650 | 0.1647325 | No Hit |
| GTCAGTCGATCCTAAGCTCAAGGAGAGATGGAATTCTCGGGTGCCAAGG | 19014 | 0.1594007 | Illumina Small RNA Adapter 2 (100% over 21bp) |
| GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTATGGAATTCTCGGGTG | 17862 | 0.1497431 | No Hit |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCGAATTTAAGCTGGAATTC | 16808 | 0.1409070 | No Hit |
| TGGACGGGAACTGATAAGGGCTGGAATTCTCGGGTGCCAAGGAACTCC | 16488 | 0.1382241 | RNA PCR Primer, Index 1 |
| GGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 15993 | 0.1340 | |
| CATTTGGATCGCGGAGATCTGGAATTCTCGGGTGCCAAGGAACTCCAGT | 15322 | 0.1284 | |
| CTTCGGTAGTATAGTGGTCAGTATCCCCGCTTGAATTCTCGGGTGCC | 15147 | 0.1269 | |
| TGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 14993 | 0.1256 | |
| AGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 14895 | 0.1248 | |
| TGAGATCATTGTGAAAGCTAATTGGAATTCTCGGGTGCCAAGGAACTCC | 14165 | 0.1187 | |
| TATTATGATGACAAACAACTAAGGAACCACTGATTGCATTGGAATTCT | 13934 | 0.1168 | |
| TCCTCGGTAGTATAGTGGTCAGTATCCCCGCTTGAATTCTCGGGTGC | 13352 | 0.1119 | |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCGAATTTAAGTGAATTCT | 13152 | 0.1102 | |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGGTGAATTCTC | 12713 | 0.1065773 | No Hit |

```

R Console
> 5646094 + 1201120 + 172150 + 164982 +
+ 110129 + 82246 + 75176 + 68069 +
+ 58409 + 38416 + 35680 + 28678 +
+ 28671 + 27757 + 22712 + 19997 +
+ 19932 + 19650 + 19014 + 17862 +
+ 16808 + 16488 + 15993 + 15322 +
+ 15147 + 14993 + 14895 + 14165 +
+ 13934 + 13352 + 13152 + 12713
[1] 8033706
> 8033706/11928428
[1] 0.6734924
> |
    
```

Contents

■ QC(Quality Control)

- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

トリミング (Trimming)

原著論文(Nie et al., BMC Genomics, 2013)では、Illuminaのadaptersやsmall RNA primer setを使ったと記載。これとPossible Sourceを比較

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](#) ①
- [Adapter Content](#)
- [Kmer Content](#)

Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|--|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164082 | 1.3830992650498455 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 170129 | 0.923248226840955 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 75176 | 0.6302255418735813 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
| TGGAATTCTCGGGTGCCAAGGAAGTCCAGTCACATCAGGATCTCGTATG | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
| AAATCATTACCCTGGACGGTGGATCACTGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCT | 8006 | | No Hit |
| TCTT | 31266 | | No Hit |
| TCCT | 279037 | | No Hit |
| AAAT | 174787 | | No Hit |
| CGGG | 800652 | | No Hit |
| TCCC | 189002 | | RNA PCR Primer, Index 1 (100% over 31bp) |

RNA-seq

A small RNA library was constructed using the total RNAs extracted above. Briefly, the RNAs were size-fractionated on a 15% polyacrylamide gel, and the 18-50nt fraction was collected. The collected small RNAs were ligated with 5' and 3' Illumina adaptors and subsequently used as a template to synthesize first-strand cDNA. The cDNA was amplified by PCR with the Illumina small RNA primer set and sequenced on the

トリミング (Trimming)

①赤下線の配列は「100% over 49bp」という記載から、RNA PCR Primer Index 1という配列の一部なのだろうと判断。

Summary

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

Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|--|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164982 | 1.3830992650498455 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.923248226840955 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGCC | 75176 | 0.6302255418735813 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
| <u>TGGAATTCTCGGGTGCCAAGGAACTCCAGTCACATCACGATCTCGTATG</u> | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
| AAATCATTACCCTGGACGGTGGATCACTGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
| TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |



①

トリミング (Trimming)

②のリードは「100% over 22bp」という記載から、リード中の22 bp分がRNA PCR Primer Index 1という配列の一部と100%一致なのだろうと判断。つまり、②のリードの黒下線部分以外がsmall RNA配列だろうと判断。

Summary

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

Overrepresented sequences

| Sequence | | | |
|--|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164982 | 1.3830992650498455 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.923248226840955 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGCC | 75176 | 0.6302255418735813 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
|  <u>TGGAATTCTCGGGTGCCAAGGAACTCCAGTCACATCAGTCTCGTATG</u> | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
|  AAATCATTACCCTGGACGGTGGATCACTGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
| TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

トリミング (Trimming)

③のリードについても同じ流れで黒下線部分の22 bpの配列を同定可能。「100% over 31bp」と書いてあるので、黒下線右側を①のリードと見比べれば理解しやすい。

Summary

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

Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|--|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164982 | 1.3830992650498455 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.923248226840955 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGCC | 75176 | 0.6302255418735813 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
|  <u>TGGAATTCTCGGGTGCCAAGGAACTCCAGTCACATCAGTCTCGTATG</u> | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
|  AAATCATTACCCTGGACGGTGGATCACTGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
|  TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

トリミング (Trimming)

Summary

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

Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164982 | 1.3830992650498455 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.923248226840955 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGCC | 75176 | 0.6302255418735813 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
|  TGGAATTCTCGGGTGCCAAGGA ACTCCAGTCACATCAGGATCTCGTATG | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
|  AAATCATTACCCTGGACGGTGGATCAC TGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
|  TCCTGGTTGATCCTGCC TGGAATTCTCGGGTGCCAAGGA ACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

RNA PCR Primer Index 1の最初の22 bpをコピーして文字列検索。④⑤⑥の結果ともに妥当。

トリミング (Trimming)

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](#)

| | | | | |
|---|--|-------|---------------------|---|
|  | TCCCTGAT TGA TCC TGC TGGAA TCTCGGGTGC | 19937 | 0.10704133007163002 | Index 1 (100% over 31bp) |
| | TCCTCGGTAGTATAGTGGTGAGTATGCCCGCCTTGAATTCTCGGGTGC | 19932 | 0.16709661994019664 | No Hit |
| | AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGTGGAATTCTCG | 19650 | 0.16473251965808067 | No Hit |
| | GTCAGTCGATCCTAAGCTCAAGGAGAGATGGAATTCTCGGGTGCCAAGG | 19014 | 0.15940071902181913 | Illumina Small RNA Adapter 2 (100% over 21bp) |
| | GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTATGGAATTCTCGGGTG | 17862 | 0.14974311786934538 | No Hit |
| | GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGCTGGAATTC | 16808 | 0.14090708348157863 | No Hit |
|  | TGGACGGAGA ACTGATAAGGGC TGGAA TCTCGGGTGCCAAGGA ACTCC | 16488 | 0.13822441649478037 | RNA PCR Primer, Index 1 (100% over 27bp) |
| | GGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 15993 | 0.13407466599957682 | No Hit |
|  | CATTTGGATCGCGGAGATC TGGAA TCTCGGGTGCCAAGGA ACTCCAGT | 15322 | 0.1284494486616342 | RNA PCR Primer, Index 1 (100% over 30bp) |
| | CTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGAATTCTCGGGTGCC | 15147 | 0.12698236515322892 | No Hit |
| | TGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 14993 | 0.12569133166583224 | No Hit |
| | AGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 14895 | 0.12486976490112528 | No Hit |
|  | TGAGATCATTGTGAAAGCTAAT TGGAA TCTCGGGTGCCAAGGA ACTCC | 14165 | 0.11874993083749175 | RNA PCR Primer, Index 1 (100% over 27bp) |
| | TATTATGATGACAAACAACTAAGGAACCACTGATTGCATTGGAATTCT | 13934 | 0.11681338060639676 | No Hit |
| | TCCTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGAATTCTCGGGTGC | 13352 | 0.11193428002415741 | No Hit |
| | GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGTGAATTCT | 13152 | 0.11025761315740851 | No Hit |
| | AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGTGGAATTCTC | 12713 | 0.10657732938489464 | No Hit |

トリミング (Trimmin

RNA PCR Primer Index 1の最初の22 bpをコピーして文字列検索。⑦はこれまでと違ってIllumina Small RNA Adapter 2 (100% over 21bp)という記述。しかし、⑦のリードの右側21 bpは検索文字列の左側21 bpと全く同じことに気づく。Illumina Small RNA Adapter 2というのがどこにあるかは(ちゃんと調べてないので)不明。

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](#)



| | | | | |
|--|-------|---------------------|--|---|
| TCCCTGAT TGA TCC TACC TT | | | | |
| TCCTCGGTAGTATAGTGGT | | | | |
| AAATCATTACCCTGGACGGT | | | | |
| ⑦ GTCAGTCGATCCTAAGCTCAAGGAGAGATGGAATTCTCGGGTGCCAAGG | 19014 | 0.15940071902181913 | | Illumina Small RNA Adapter 2 (100% over 21bp) |
| GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTATGGAATTCTCGGGTG | 17862 | 0.14974311786934538 | | No Hit |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGCTGGAATTC | 16808 | 0.14090708348157863 | | No Hit |
| TGGACGGAGA ACTGATAAGGGCTGGAATTCTCGGGTGCCAAGGA ACTCC | 16488 | 0.13822441649478037 | | RNA PCR Primer, Index 1 (100% over 27bp) |
| GGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 15993 | 0.13407466599957682 | | No Hit |
| CATTTGGATCGCGGAGATCTGGAATTCTCGGGTGCCAAGGA ACTCCAGT | 15322 | 0.1284494486616342 | | RNA PCR Primer, Index 1 (100% over 30bp) |
| CTTCGGTAGTATAGTGGTCAGTATCCCCGCC TTGGAATTCTCGGGTGCC | 15147 | 0.12698236515322892 | | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGCCCTACGTTGGAATTCTCGGGTGC | 14993 | 0.12569133166583224 | | No Hit |
| AGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 14895 | 0.12486976490112528 | | No Hit |
| TGAGATCATTGTGAAAGCTAAT TGGAAATTCTCGGGTGCCAAGGA ACTCC | 14165 | 0.11874993083749175 | | RNA PCR Primer, Index 1 (100% over 27bp) |
| TATTATGATGACAAACAACTAAGGAACCACTGATTGCATTGGAATTCT | 13934 | 0.11681338060639676 | | No Hit |
| TCCTCGGTAGTATAGTGGTCAGTATCCCCGCC TTGGAATTCTCGGGTGC | 13352 | 0.11193428002415741 | | No Hit |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGTGGAAATTC | 13152 | 0.11025761315740851 | | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGGGGTGGAATTCTC | 12713 | 0.10657732938489464 | | No Hit |

トリミング (Trimmomatic)

⑧よく見ると、22 bp分の長さの黒下線より短い「RNA PCR Primer Index 1の左側の部分配列」がOverrepresented sequencesの上位の右側に存在することがわかる。これらはおそらく短いためにNo Hitとなったのであろう。

Summary

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

Overrepresented sequences

| | Sequence | Count | Percentage | Possible Source |
|---|---|---------|---------------------|--|
|  | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| | GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| | TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAAATTCTCGGGTGC | 164982 | 1.3830992650498455 | No Hit |
| | TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.923248226840955 | No Hit |
| | TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| | GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAAATTCTCGGGTGCC | 75176 | 0.6302255418735813 | No Hit |
| | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
|  | <u>TGGAATTCTCGGGTGCCAAGGA</u> ACTCCAGTCACATCAGGATCTCGTATG | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
|  | AAATCATTACCCTGGACGGTGGATCACT <u>GGAATTCTCGGGTGCCAAGGA</u> | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| | TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| | TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| | AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| | CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
|  | TCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

トリミング (Trimming)

これらの部分文字列をハイライトさせるべく、「黒よりも短く、赤枠をハイライトさせられる程度の長さの部分文字列を用いて調べる。①の左側の一部をコピーして検索。黄色でハイライトされている部分より右側でもきれいに一致していることがわかる。

Summary

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

Overrepresented

| | | | | | |
|---|---|---|---------|---------------------|--|
| ⑧ | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCT | <u>TGGAATTCTCGGGTGC</u> | 5646094 | 47.333093681749176 | No Hit |
| | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCC | <u>TGGAATTCTCGGGTGCC</u> | 1201120 | 10.06939053494727 | No Hit |
| | GGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | <u>TGGAATTCTCGGGTGC</u> | 172150 | 1.4431910055541266 | No Hit |
| | TGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | <u>TGGAATTCTCGGGTGC</u> | 164982 | 1.3830992650498455 | No Hit |
| | TCCTCGGTAGTATAGTGGTGAGTATGCACGCCT | <u>TGGAATTCTCGGGTGC</u> | 110129 | 0.923248226840955 | No Hit |
| | TCCTCGGTAGTATAGTGGTGAGTATGCTCGCC | <u>TGGAATTCTCGGGTGCC</u> | 82246 | 0.6894957156131554 | No Hit |
| | GCCGTGATCGTCTAGTGGTTAGGACCCTACGT | <u>TGGAATTCTCGGGTGCC</u> | 75176 | 0.6302255418735813 | No Hit |
| | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCG | <u>TGGAATTCTCGGGTGC</u> | 68069 | 0.5706451847636587 | No Hit |
| ① | <u>TGGAATTCTCGGGTGC</u> CCAAGGAACCTCCAGTCACATCAGTATCTCGTATG | | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
| | AAATCATTACCCTGGACGGTGGATCAC | <u>TGGAATTCTCGGGTGCCAAGGA</u> | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| | TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCT | <u>TGGAATTCTCGGGTGC</u> | 35680 | 0.299117369028006 | No Hit |
| | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTG | <u>TGGAATTCTCGGGTG</u> | 28678 | 0.2404172620231266 | No Hit |
| | TCCTCGGTAGTATAGTGGTGAGTATGCACGCC | <u>TGGAATTCTCGGGTGCC</u> | 28671 | 0.24035857868279037 | No Hit |
| | AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAATTCTCGG | | 27757 | 0.23269621110174787 | No Hit |
| | CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | | 22712 | 0.19040228938800652 | No Hit |
| | TCCTGGTTGATCCTGCC | <u>TGGAATTCTCGGGTGC</u> CCAAGGAACCTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

トリミング (Trimming)

出現回数下位を眺めている。黒下線で示すように、ハイライトされていないNo Hitのリード中にも、さらに短い「RNA PCR Primer Index 1の左側の塩基配列」を含んでいることがわかる。

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](#)

| | | | | | |
|---|-----------------------|-------|---------------------|---|-------|
| TCCCTGGTTGATCCTGCC | TGGAATTCTCGG | | | | 31bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCCCGCCT | TGGAATTCTCGGGTGC | 19932 | 0.16709661994019664 | No Hit | |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGT | GGAATTCTCG | 19650 | 0.16473251965808067 | No Hit | |
| GTCAGTCGATCCTAAGCTCAAGGAGAGA | TGGAATTCTCGGGTCCAAGG | 19014 | 0.15940071902181913 | ILLUMINA SMALL RNA ADAPTER 2 (100% OVER 21BP) | |
| GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTA | TGGAATTCTCGGGTG | 17862 | 0.14974311786934538 | No Hit | |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGCT | GGAATTC | 16808 | 0.14090708348157863 | No Hit | |
| TGGACGGAGAACTGATAAGGGC | TGGAATTCTCGGGTCCAAGGA | 16488 | 0.13822441649478037 | RNA PCR PRIMER, INDEX 1 (100% OVER 27BP) | |
| GGCCGTGATCGTCTAGTGGTTAGGCCCTACGT | TGGAATTCTCGGGTGC | 15993 | 0.13407466599957682 | No Hit | |
| CATTTGGATCGCGGAGATC | TGGAATTCTCGGGTCCAAGGA | 15322 | 0.1284494486616342 | RNA PCR PRIMER, INDEX 1 (100% OVER 30BP) | |
| CTTCGGTAGTATAGTGGTCAGTATCCCGCCT | TGGAATTCTCGGGTGCC | 15147 | 0.12698236515322892 | No Hit | |
| TGCCGTGATCGTCTAGTGGTTAGGCCCTACGT | TGGAATTCTCGGGTGC | 14993 | 0.12569133166583224 | No Hit | |
| AGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | TGGAATTCTCGGGTGC | 14895 | 0.12486976490112528 | No Hit | |
| TGAGATCATTGTGAAAGCTAAT | TGGAATTCTCGGGTCCAAGGA | 14165 | 0.11874993083749175 | RNA PCR PRIMER, INDEX 1 (100% OVER 27BP) | |
| TATTATGATGACAAACAACTAAGGA | ACTGATTGCATTGGAATTCT | 13934 | 0.11681338060639676 | No Hit | |
| TCCTCGGTAGTATAGTGGTCAGTATCCCGCCT | TGGAATTCTCGGGTGC | 13352 | 0.11193428002415741 | No Hit | |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGT | GGAATTCT | 13152 | 0.11025761315740851 | No Hit | |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGT | GGAATTCTC | 12713 | 0.10657732938489464 | No Hit | |

トリミング (Trimming)

これらのハイライトされた部分配列や黒下線部分のアダプター配列 (adapter sequences) は、予め取り除かれるべき (このデータはカイコの small RNA なのでカイコ中に本来存在せず、マップもされない)。

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](#)

| | | | | | |
|---|-----------------------|-------|---------------------|---|-------|
| TCCCTGGTTGATCCTGCC | TGGAATTC | | | | 31bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCCCGCCT | TGGAATTCTCGGGTGC | 19932 | 0.16709661994019664 | No Hit | |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGT | GGAATTCTCG | 19650 | 0.16473251965808067 | No Hit | |
| GTCAGTCGATCCTAAGCTCAAGGAGAGA | TGGAATTCTCGGGTCCAAGG | 19014 | 0.15940071902181913 | ILLUMINA SMALL RNA ADAPTER 2 (100% OVER 21BP) | |
| GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTA | TGGAATTCTCGGGTG | 17862 | 0.14974311786934538 | No Hit | |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGCT | GGAATTC | 16808 | 0.14090708348157863 | No Hit | |
| TGGACGGAGAACTGATAAGGGC | TGGAATTCTCGGGTCCAAGGA | 16488 | 0.13822441649478037 | RNA PCR PRIMER, INDEX 1 (100% OVER 27BP) | |
| GGCCGTGATCGTCTAGTGGTTAGGCCCTACGT | TGGAATTCTCGGGTGC | 15993 | 0.13407466599957682 | No Hit | |
| CATTTGGATCGCGGAGATC | TGGAATTCTCGGGTCCAAGGA | 15322 | 0.1284494486616342 | RNA PCR PRIMER, INDEX 1 (100% OVER 30BP) | |
| CTTCGGTAGTATAGTGGTCAGTATCCCGCCT | TGGAATTCTCGGGTGCC | 15147 | 0.12698236515322892 | No Hit | |
| TGCCGTGATCGTCTAGTGGTTAGGCCCTACGT | TGGAATTCTCGGGTGC | 14993 | 0.12569133166583224 | No Hit | |
| AGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | TGGAATTCTCGGGTGC | 14895 | 0.12486976490112528 | No Hit | |
| TGAGATCATTGTGAAAGCTAAT | TGGAATTCTCGGGTCCAAGGA | 14165 | 0.11874993083749175 | RNA PCR PRIMER, INDEX 1 (100% OVER 27BP) | |
| TATTATGATGACAAACAACTAAGGA | ACTGATTGCATTGGAATTCT | 13934 | 0.11681338060639676 | No Hit | |
| TCCTCGGTAGTATAGTGGTCAGTATCCCGCCT | TGGAATTCTCGGGTGC | 13352 | 0.11193428002415741 | No Hit | |
| GTTGGCCTCAGATCAGGGAGGATCACCCGCCAATTTAAGT | GGAATTCT | 13152 | 0.11025761315740851 | No Hit | |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGGT | GGAATTCTC | 12713 | 0.10657732938489464 | No Hit | |

トリミング (Trimming)

このデータ中のアダプター配列除去を行う場合は、RNA PCR Primer, Index 1の実際の配列長が49 bp以上であろうがなかろうが、①の塩基配列をそのままアダプター配列として与えてもうまくトリムできるだろうと判断。

Summary

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

Overrepresented sequences

| Sequence | | | |
|--|---------|---------------------|--|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTTGGAAATTCTCGGGTGC | 5646094 | 47.333093681749176 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGGAATTCTCGGGTGCC | 1201120 | 10.06939053494727 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 172150 | 1.4431910055541266 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGC | 164982 | 1.3830992650498455 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTTGGAAATTCTCGGGTGC | 110129 | 0.923248226840955 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTGGAATTCTCGGGTGCC | 82246 | 0.6894957156131554 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGTTGGAATTCTCGGGTGCC | 75176 | 0.6302255418735813 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCGTGGAATTCTCGGGTGC | 68069 | 0.5706451847636587 | No Hit |
|  <u>TGGAATTCTCGGGTGCCAAGGAACTCCAGTCACATCAGGATCTCGTATG</u> | 58409 | 0.48966217509968624 | RNA PCR Primer, Index 1 (100% over 49bp) |
| AAATCATTACCCTGGACGGTGGATCACTGGAATTCTCGGGTGCCAAGGA | 38416 | 0.3220541717651312 | RNA PCR Primer, Index 1 (100% over 22bp) |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCTTGGAAATTCTCGGGTGC | 35680 | 0.299117369028006 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTGTGGAATTCTCGGGTG | 28678 | 0.2404172620231266 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCCTGGAATTCTCGGGTGCC | 28671 | 0.24035857868279037 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGTGGAAATTCTCGG | 27757 | 0.23269621110174787 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTGTGGAATTCTCGG | 22712 | 0.19040228938800652 | No Hit |
| TCCCTGGTTGATCCTGCCTGGAATTCTCGGGTGCCAAGGAACTCCAGTC | 19997 | 0.16764153667189002 | RNA PCR Primer, Index 1 (100% over 31bp) |

Contents

■ QC(Quality Control)

- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

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A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

- 前処理 | トリミング | アダプター配列除去(基礎) | QuasR(Gaidatzis_2015) NEW
- アセンブル | [ゲノム用](#)
- アセンブル | [トランススクリプト](#)
- マッピング | [basic aligner](#)
- マッピング | [splice-aware aligner](#)
- マッピング | [Bisulfite](#)
- マッピング | [\(ESTレベル\)](#)
- マッピング | [基礎 \(last modified 2015/01/01\)](#)
- マッピング | [single-end](#)

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"#アダプター配列
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)

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

```

```

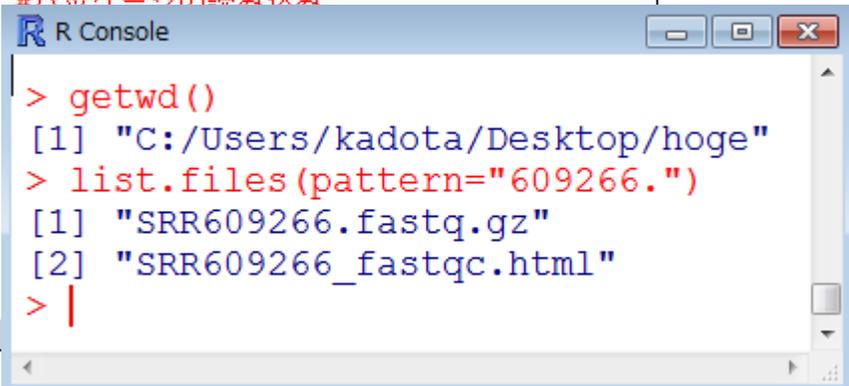
in_f <- "SRR609266.fastq.gz"
out_f <- "hoge4.fastq.gz"
param_adapter <- "TGGAATTCTCGGGTCCAAGGAAGTCCAGTCACATCACGATCTCGTATG"#アダプター配列
param_nrec <- 500000

#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)

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

```



アダプター配列除去

この状態がしばらく続きます。
「応答なし」的な状態になっても、しばらく放置すると復活します

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ムしないやり方です。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 <- "TGGAATTCTCGGGTGCCAAGGAAGTCCAGTCACATCACGATCTCGTATG"#アダプター配列
param_nrec <- 500000 #一度に処理するリード数を指定
```

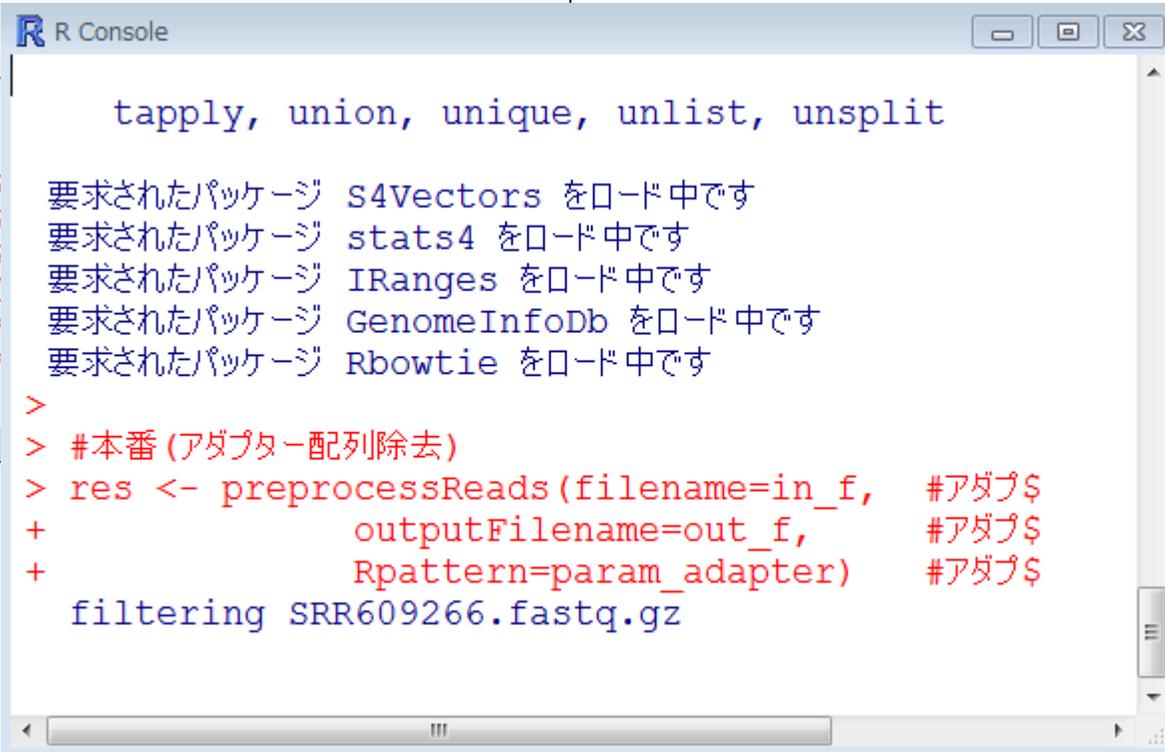
#必要なパッケージをロード

```
library(QuasR) #パッケージをロード
```

#本番(アダプター配列除去)

```
res <- preprocessReads(filename=in_f, #アダプター配列を指定
                       outputFilename=out_f, #アダプター配列を指定
                       Rpattern=param_adapter) #アダプター配列を指定
```

```
res #確認して
file.size(in_f) #入力ファイルサイズ
file.size(out_f) #出力ファイルサイズ
```



```
R Console
tapply, union, unique, unlist, unsplit

要求されたパッケージ S4Vectors をロード中です
要求されたパッケージ stats4 をロード中です
要求されたパッケージ IRanges をロード中です
要求されたパッケージ GenomeInfoDb をロード中です
要求されたパッケージ Rbowtie をロード中です

>
> #本番(アダプター配列除去)
> res <- preprocessReads(filename=in_f, #アダプ$
+                       outputFilename=out_f, #アダプ$
+                       Rpattern=param_adapter) #アダプ$
filtering SRR609266.fastq.gz
```

アダプター配列除去

これが正常終了時の画面。アダプター配列除去後のリードが短いもの (tooShort) が90,242個、Nを多く含むもの (tooManyN) が875個あったことがわかります。①これらを除いた残りの11,837,311リードが出力されているようです。②hoge4.fastq.gzのファイルサイズは281MB。

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。このファイルを実行結果として「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、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のみが除去されないやり方です。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 <- "TGGAATTCTCGGGTGCCAAGGAAGTCCAG" #一度に処理
param_nrec <- 500000 #一度に処理

#必要なパッケージをロード
library(QuasR) #パッケージ

#本番(アダプター配列除去)
res <- preprocessReads(filename=in_f, #アダプター配列を指定してin_fに格納
                       outputFilename=out_f, #アダプター配列を指定してout_fに格納
                       Rpattern=param_adapter) #アダプター配列を指定してin_fに格納

res #確認して
file.size(in_f) #入力ファイルサイズ
file.size(out_f) #出力ファイルサイズ
```

```
R Console
+ outputFilename=out_f, #アダプ$
+ Rpattern=param_adapter) #アダプ$
filtering SRR609266.fastq.gz
> res #確認し$

SRR609266.fastq.gz
totalSequences 11928428
matchTo5pAdapter 0
matchTo3pAdapter 11839903
tooShort 90242
tooManyN 875
lowComplexity 0
totalPassed 11837311 ①
> file.size(in_f) #入力フ$
[1] 383840833
> file.size(out_f) #出力フ$
[1] 281324281
> |
```



トリム後のFastQC結果

hoge4.fastq.gzに対するFastQC実行結果ファイル(hoge4_fastqc.html)を眺める。hogeフォルダ中にあります。

①リード数はQuasR実行結果と同じ11,837,311個。

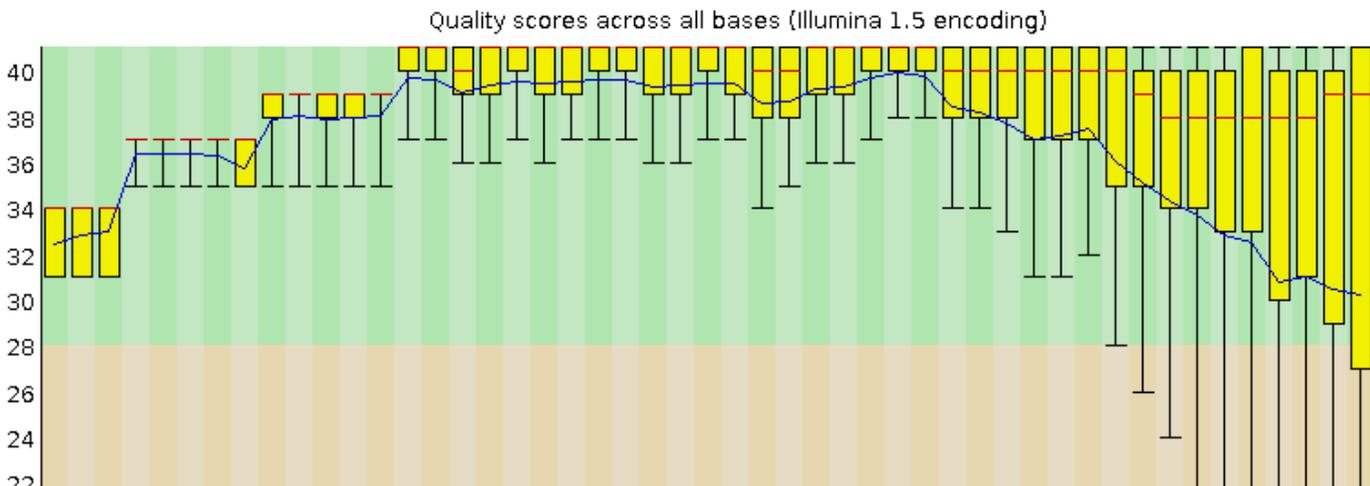
Summary

- ✔ 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

✔ Basic Statistics

| Measure | Value |
|-----------------------------------|-------------------------|
| Filename | hoge4.fastq.gz |
| File type | Conventional base calls |
| Encoding | Illumina 1.5 |
| Total Sequences | 11837311 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 14-49 |
| %GC | 50 |

✔ Per base sequence quality



トリム後のFastQC結果

①配列長の範囲は14-49。このことから、②QuasR実行結果のtooShortで落とされた90,242リードは、アダプター配列除去後に14塩基未満になったものたちだろうと推測。

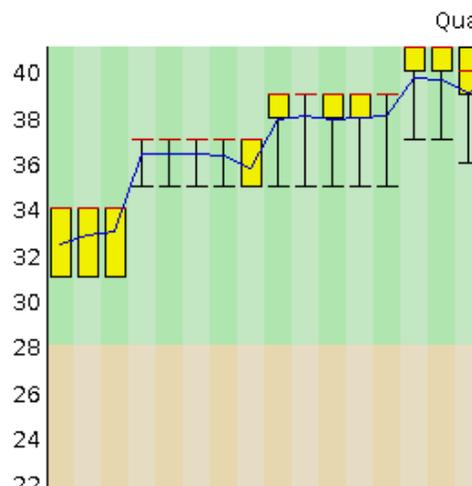
Summary

- ✓ 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

Basic Statistics

| Measure | Value |
|-----------------------------------|-------------------------|
| Filename | hoge4.fastq.gz |
| File type | Conventional base calls |
| Encoding | Illumina 1.5 |
| Total Sequences | 11837311 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 14-49 |
| %GC | 50 |

Per base sequence quality



```
R Console
> res
SRR609266.fastq.gz
totalSequences      11928428
matchTo5pAdapter    0
matchTo3pAdapter    11839903
tooShort            90242
tooManyN            875
lowComplexity        0
totalPassed         11837311
> |
```

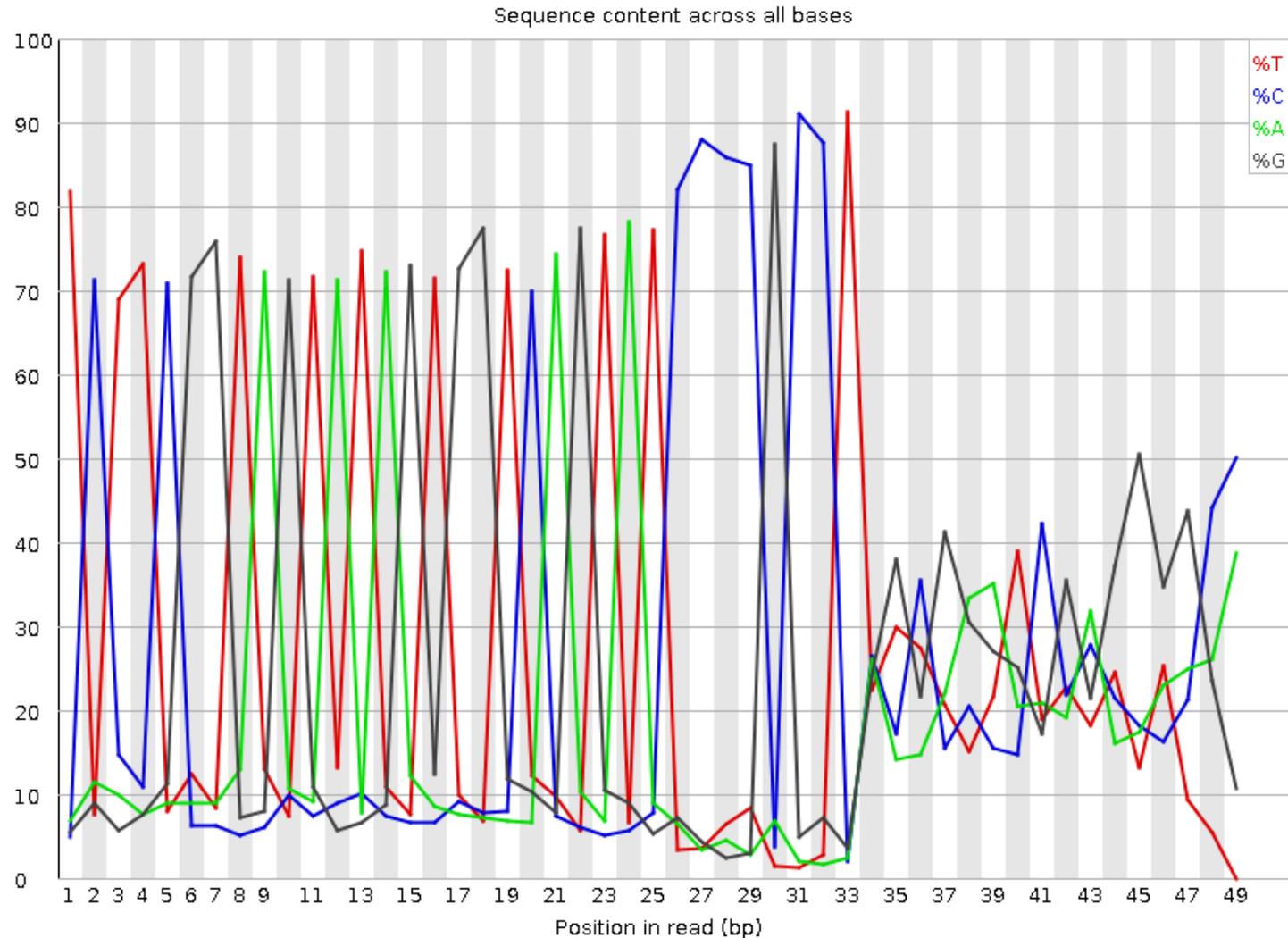
トリム後のFastQC結果

課題2:なぜアダプター配列除去前後で34 bp以降のプロファイルが劇的に変わったのか述べよ。

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

✗ Per base sequence content



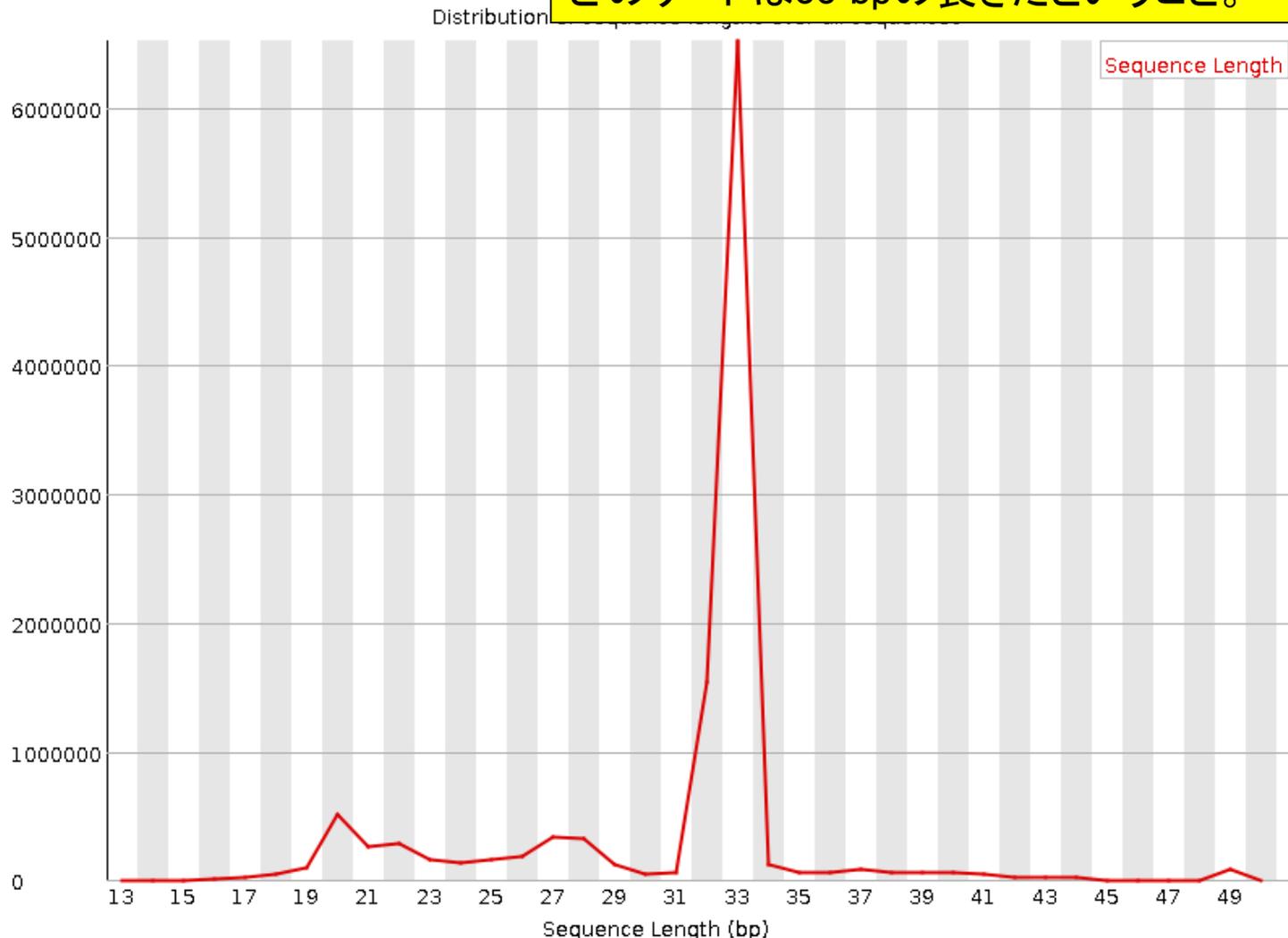
トリム後のFastQC結果

アダプター配列除去前は49 bpのところ
に一本のピークが見えていた(つまり全
リードの長さが49 bp)。除去後は33 bp
のところにピークがある。つまり、ほとん
どのリードは33 bpの長さだということ。

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

! Sequence Length Distribution



トリム後のFastQC結果

Possible Sourceのところ全てNo Hitになっていることが分かる。QuasRを用いたアダプター配列除去手順の正しさは、この結果からもわかる。

Summary

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

✗ Overrepresented sequences

| Sequence | Count | Percentage | Possible Source |
|---------------------------------------|---------|---------------------|-----------------|
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCT | 5711563 | 48.25051061005325 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCC | 1224329 | 10.34296556033714 | No Hit |
| GGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | 174048 | 1.4703339297244111 | No Hit |
| TGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | 166820 | 1.4092727647351666 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCT | 111334 | 0.9405345521461759 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCC | 83415 | 0.7046786216903485 | No Hit |
| GCCGTGATCGTCTAGTGGTTAGGACCCTACGT | 76203 | 0.6437526225339522 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCG | 68748 | 0.5807737922911715 | No Hit |
| AAATCATTACCCTGGACGGTGGATCAC | 39237 | 0.33146886146693283 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCT | 36078 | 0.30478205734393565 | No Hit |
| TCTTCGGTAGTATAGTGGTCAGTATCCCCGCTG | 29252 | 0.24711693390500597 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCACGCC | 29165 | 0.246381969688893 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTGGCTCGCG | 28441 | 0.24026571575250494 | No Hit |
| CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTG | 23202 | 0.19600735335922154 | No Hit |
| TCCCTGGTTGATCCTGCC | 20710 | 0.17495527489309015 | No Hit |
| TCCTCGGTAGTATAGTGGTGAGTATGCCCGCCT | 20141 | 0.17014843996242052 | No Hit |
| AAATCATTACCCTGGACGGTGGATCACTGGCTCGCGG | 19867 | 0.16783372507489241 | No Hit |
| GTCAGTCGATCCTAAGCTCAAGGAGAGA | 19382 | 0.16373651076667667 | No Hit |
| GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTA | 18124 | 0.1531090971589747 | No Hit |
| TGGACGGAGAACTGATAAGGGC | 17985 | 0.15193484398610463 | No Hit |

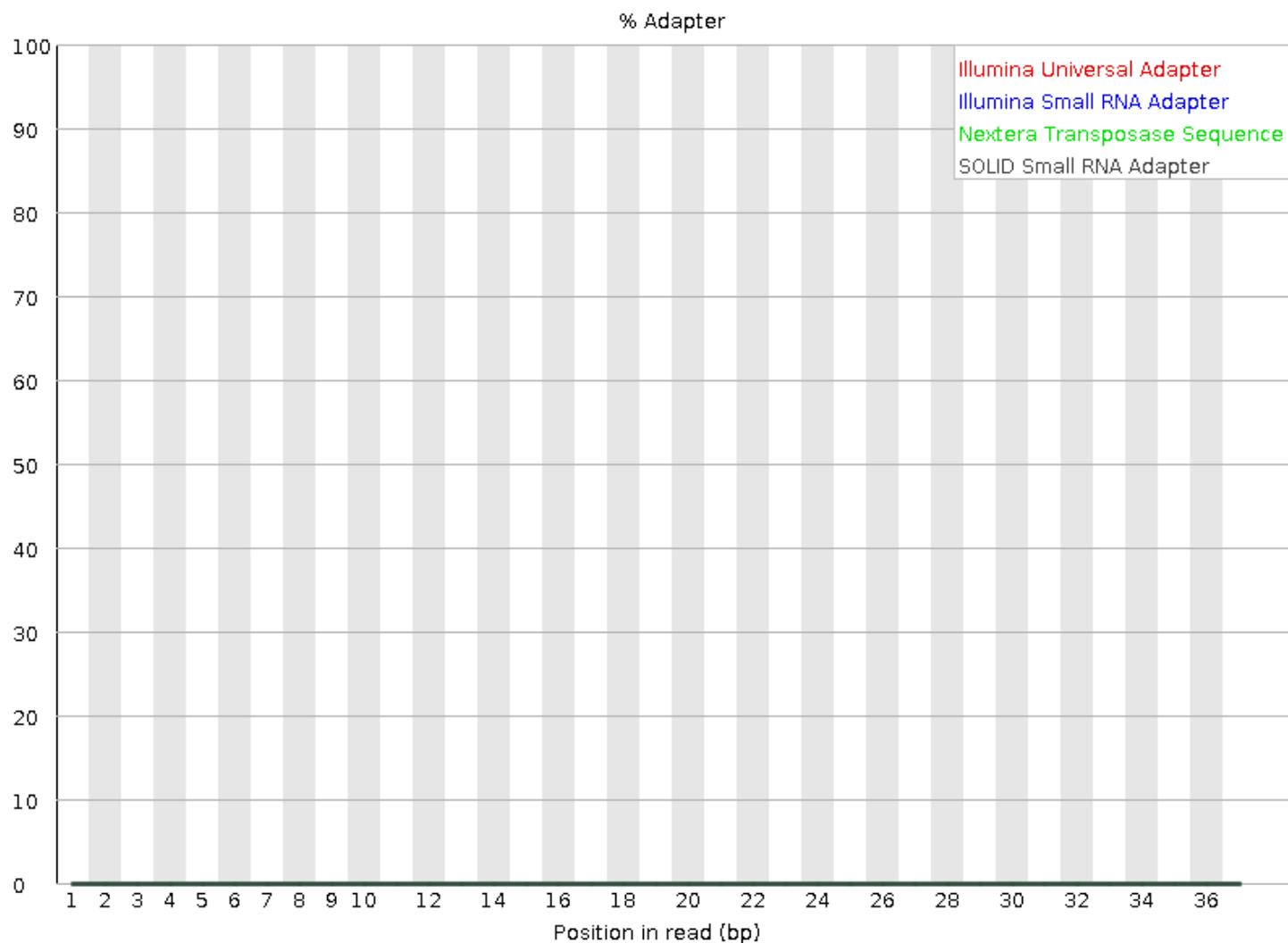
トリム後のFastQC結果

アダプター配列除去前はIllumina Small RNA Adapterの割合が20 bp以降で上昇していたが、除去後は消えていることがわかる。

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

✓ Adapter Content



Contents

■ QC(Quality Control)

- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

FastQCの項目「Per base sequence content」と似たような結果は、sequence logosでも得ることができます。

Sequence logos

- 解析 | 一般 | [アラインメント\(ペアワイズ; 応用\)](#) (last modified 2010/6/8)
- 解析 | 一般 | [パターンマッチング](#) (last modified 2013/06/19)
- 解析 | 一般 | [GC含量\(GC contents\)](#) (last modified 2015/04/18)
- 解析 | 一般 | [Sequence logos\(Schneider 1990\)](#) (last modified 2014/07/23)
- 解析 | 一般 | 上流配列解析 | [LDSS\(Yamamoto 2007\)](#) (last modified 2015/02/19)
- 解析 | 一般 | 上流配列解析 | [Relative Appearance Ratio\(Yamamoto 2011\)](#) (last modified 2015/02/19)

解析 | 一般 | Sequence logos(Schneider_1990)

seqLogoパッケージを用いてsequence logosを作成します。ここでは、multi-FASTA形式のファイルを使用します。上流-35 bpにTATA boxが「ファイル」-「ディレクトリの変更」で解析

9. FASTQ形式ファイル(hoge4.fastq.gz)の場合:

small RNA-seqデータ(280Mb弱、11,928,428リード)です。原著論文([Nie et al., BMC Genomics, 2013](#))中の記述から [GSE41841](#)を頼りに、[SRP016842](#)にたどりつき、[前処理 | トリミング | アダプター配列除去\(応用\) | ShortRead\(Morgan 2009\)](#)の4を実行して得られたものが入力ファイルです。アダプター配列除去後のデータなので、リードごとに配列長が異なる場合でも読み込めるShortReadパッケージ中の readFastq関数を用いています。

1. 入力ファイルがmulti-FASTA形式の場合:

```
in_f <- "test1.fasta"

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

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f)

#本番(sequence logoを実行)
hoge <- consensusMatrix(fasta)
out <- makePWM(hoge[1:4,])
seqLogo(out)
```

```
in_f <- "hoge4.fastq.gz"
out_f <- "hoge9.png"
param_fig <- c(787, 370)

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納
#ファイル出力時の横幅と縦幅を指定(単位はポイント)

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

#パッケージの読み込み
#パッケージの読み込み

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

#in_fで指定したファイルの読み込み
#リード塩基配列情報をfastaに格納

#本番(sequence logoを実行)
hoge <- consensusMatrix(fasta, as.prob=T, baseOnly=T)#各ポジションの塩基組成
out <- makePWM(hoge[1:4,]) #hogeはACGT以外の塩基(例えばN)のprob

#ファイルに保存
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2])#出力ファイル名を指定してout_fに格納
seqLogo(out)
dev.off()
```

Sequence logos

塩基ごとの色の使い分けもG以外は同じ (A: green, C: blue, and T: red)なので分かりやすいと思います。

9. FASTQ形式ファイル(hoge4.fastq.gz)の場合:

small RNA-seqデータ(280Mb弱、11,928,428リード)です。原著論文([Nie et al., BMC Genomics, 2013](#))中の記述から [GSE41841](#)を頼りに、[SRP016842](#)にたどりつき、[前処理 | トリミング | アダプター配列除去\(応用\) | ShortRead\(Morgan 2009\)](#)の4を実行して得られたものが入力ファイルです。アダプター配列除去後のデータなので、リードごとに配列長が異なる場合でも読み込める [ShortRead](#)パッケージ中の `readFastq`関数を用いています。

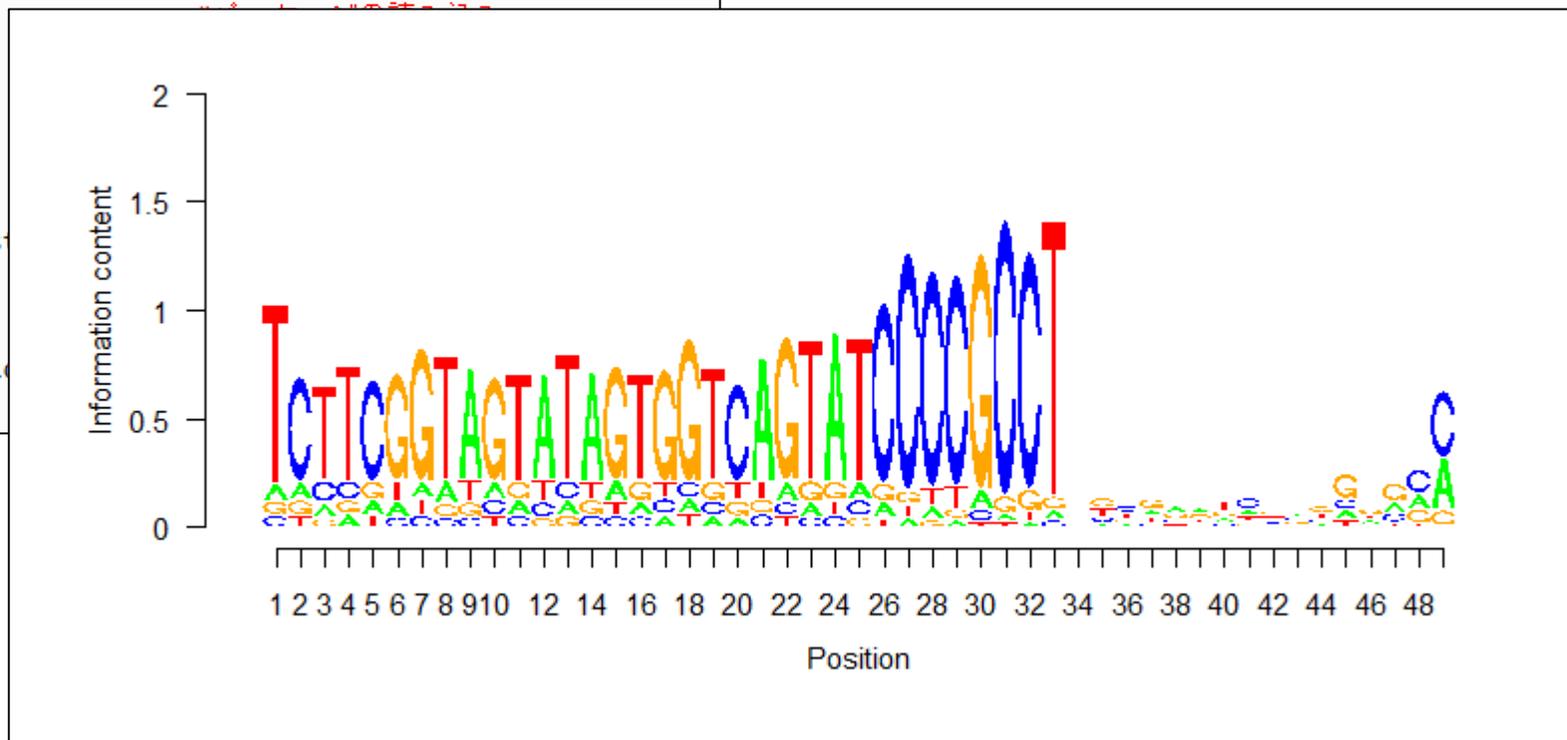
```
in_f <- "hoge4.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge9.png" #出力ファイル名を指定してout_fに格納
param_fig <- c(787, 370) #ファイル出力時の横幅と縦幅を指定(単位はポイント)

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

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

#本番(sequence logoを実行)
hoge <- consensusMatrix(fasta)
out <- makePWM(hoge[1:4,])

#ファイルに保存
png(out_f, pointsize=13, width=param_fig[1], height=param_fig[2])
seqLogo(out)
```



FastQCの項目「Sequence Length Distribution」と似たような結果も得ることができます。

配列長分布

- 前処理 | クオリティチェック | [qrrc](#) (last modified 2014/07/17)
- 前処理 | クオリティチェック | [PHREDスコアに変換](#) (last modified 2013/06/18)
- 前処理 | クオリティチェック | [配列長分布を調べる](#) **NEW** (last modified 2015/06/22)
- 前処理 | クオリティチェック | [Overrepresented sequences](#) | [ShortRead\(Morgan 2009\)](#) (last modified 2015/06/22)
- 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 | クオリティチェック | 配列長分布を調べる **NEW**

FASTAまたはFASTQ形式ファイルを読み込んで配列長分布を得るやり方を示します。

「ファイル」-「ディレクトリ」

1. FASTA形式ファイル

```
in_f <- "sample2.fasta"
out_f <- "hoge1.txt"
```

```
#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(
  fasta
```

```
#本番
out <- table(width(fasta),
  out
```

```
#ファイルに保存
write.table(out, out_f,
```

5. gzip圧縮状態のFASTQ形式ファイル(hoge4.fastq.gz)の場合:

配列長の異なる架空のファイルです。param_nbinsで50と指定すると、1 bpおきに、1-50 bpまでのリード数が格納されます。QuasRパッケージ(Gaidatzis et al., 2015) マニュアル中の barplotを用いた作図形式に似せています。

```
in_f <- "hoge4.fastq.gz"
out_f1 <- "hoge5.txt"
out_f2 <- "hoge5.png"
param_nbins <- 50
param_fig <- c(700, 400)
```

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

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

```
#本番
out <- tabulate(width(fastq), param_nbins) #出現頻度情報を得た結果をoutに格納
out #確認してるだけです
```

```
#ファイルに保存(テキストファイル)
names(out) <- 1:param_nbins #x軸の数値情報取得
tmp <- cbind(names(out), out) #保存したい情報をtmpに格納
write.table(tmp, out_f1, sep="\t", append=F, quote=F, row.names=F) #tmpの中身を格納
```

配列長分布

5. gzip圧縮状態のFASTQ形式ファイル(hoge4.fastq.gz)の場合:

配列長の異なる架空のファイルです。param_nbinsで50と指定すると、1 bpおきに、1-50 bpまでのリード数が格納されます。QuasRパッケージ(Gaidatzis et al., 2015) マニュアル中の barplot を用いた作図形式に似せています。

```

in_f <- "hoge4.fastq.gz"
out_f1 <- "hoge5.txt"
out_f2 <- "hoge5.png"
param_nbins <- 50
param_fig <- c(700, 400)

#必要なパッケージをロード
library(ShortRead)

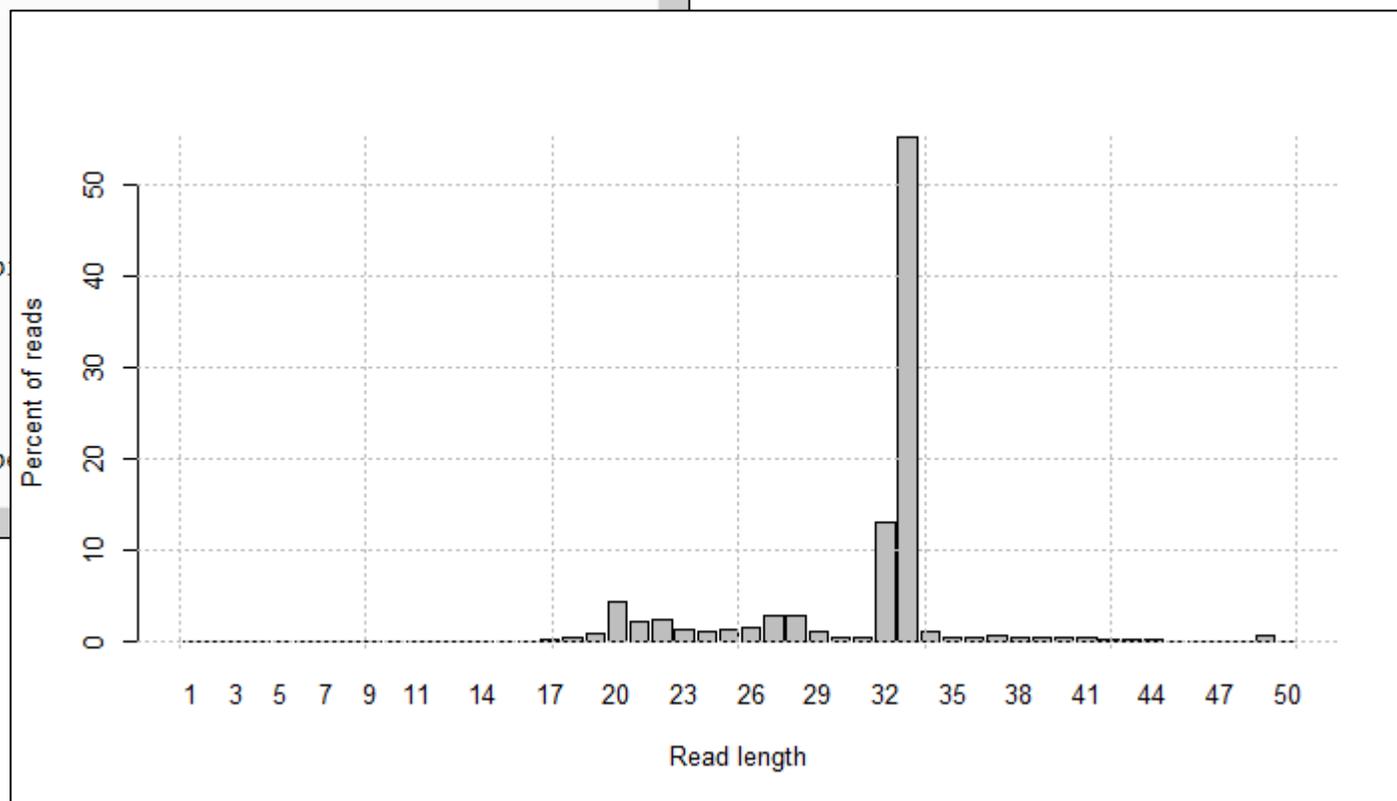
#入力ファイルの読み込み
fastq <- readFastq(in_f)
sread(fastq)

#本番
out <- tabulate(width(fastq), param_nb)
out

#ファイルに保存(テキストファイル)
names(out) <- 1:param_nbins
tmp <- cbind(names(out), out)
write.table(tmp, out_f1, sep="\t", app

```

#入力ファイル名を指定してin_fに格納
 #出力ファイル名を指定してout_f1に格納
 #出力ファイル名を指定してout_f2に格納
 #分割数(binの数)を指定
 #ファイル出力時の横幅と縦幅を指定(単位はピク



FastQCの項目「Overrepresented sequences」と似たような結果も得ることができます。

リードごとの出現回数

2. gzip圧縮状態のFASTQ形式ファイル(hoge4.fastq.gz)の場合:

配列長の異なる架空のファイルです。

```
in_f <- "hoge4.fastq.gz"
out_f <- "hoge2.txt"
```

#入力ファイル名を指定してin_fに格納
#出力ファイル名を指定してout_fに格納

```
#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
fastq <- readFastq(in_f)
sread(fastq)
```

#in
#配列

```
#本番
out <- table(sread(fastq))
out <- sort(out, decreasing=T)
head(out)
```

#リード
#出現
#確認

```
#ファイルに保存
tmp <- cbind(names(out), out)
write.table(tmp, out_f, sep="\t", append=F,
```

#保存

| | A | B |
|----|--|---------|
| 1 | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCT | 5711563 |
| 2 | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCC | 1224329 |
| 3 | GGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | 174048 |
| 4 | TGCCGTGATCGTCTAGTGGTTAGGACCCTACGT | 166820 |
| 5 | TCCTCGGTAGTATAGTGGTGAGTATGCACGCCT | 111334 |
| 6 | TCCTCGGTAGTATAGTGGTGAGTATGCTCGCC | 83415 |
| 7 | GCCGTGATCGTCTAGTGGTTAGGACCCTACGT | 76203 |
| 8 | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCG | 68748 |
| 9 | AAATCATTACCCTGGACGGTGGATCAC | 39237 |
| 10 | TCCTCGGTAGTATAGTGGTGAGTATGCTCGCCT | 36078 |
| 11 | TCTTCGGTAGTATAGTGGTCAGTATCCCCGCCTG | 29252 |
| 12 | TCCTCGGTAGTATAGTGGTGAGTATGCACGCC | 29165 |
| 13 | AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCG | 28441 |
| 14 | CGGGCAACCCGCTGAAACTCCTTCGTGCTGGGGATTG | 23202 |
| 15 | TCCCTGGTTGATCCTGCC | 20710 |
| 16 | TCCTCGGTAGTATAGTGGTGAGTATGCCCGCCT | 20141 |
| 17 | AAATCATTACCCTGGACGGTGGATCACTTGGCTCGCGG | 19867 |
| 18 | GTCAGTCGATCCTAAGCTCAAGGAGAGA | 19382 |
| 19 | GACTGGGAGCGTGGCGTCTCCTGTAATTCGGCTA | 18124 |

Contents

■ QC(Quality Control)

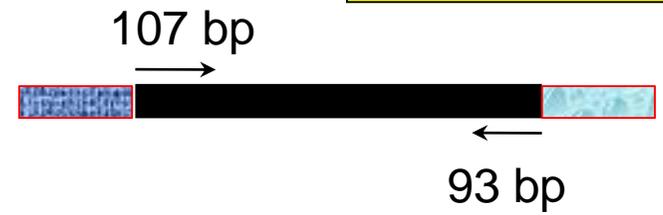
- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

一番下はLinuxコマンド。
Rではどうやるのか？

乳酸菌RNA-seqデータ

paired-endのオリジナルデータ

- 134,755,996リード
- SRR616268_1.fastq.bz2: 107 bp, 7,662,128,101 bytes (約7.7GB)
- SRR616268_2.fastq.bz2: 93 bp, 7,017,031,734 bytes (約7.0GB)



Linuxコマンドを用いて最初の1,000,000リード分からのサブセットを抽出

- 1,000,000リード(one million reads)
- SRR616268sub_1.fastq.gz: 107 bp, 74,906,576 bytes (約75MB)
- SRR616268sub_2.fastq.gz: 93 bp, 67,158,462 bytes (約67MB)

```
C:\Users\kadota\Desktop\share\JSLAB4_1.sh - EmEditor
ファイル(F) 編集(E) 検索(S) 表示(V) ツール(T) ウィンドウ(W) ヘルプ(H)
JSLAB4_1.sh x
#wget -c ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061483/SPX204226/SRR616268_1.fastq.bz2↓
#wget -c ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061483/SPX204226/SRR616268_2.fastq.bz2↓
bzip2 -dc SRR616268_1.fastq.bz2 | head -n 4000000 > SRR616268sub_1.fastq↓
bzip2 -dc SRR616268_2.fastq.bz2 | head -n 4000000 > SRR616268sub_2.fastq↓
```

サブセットの抽出

項目としては、感覚的には黒矢印部分だが、例題数が多すぎると思ったので、とりあえず①のところに書いてます。

- ・ [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 2015/06/20) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [応用](#) **①** (last modified 2015/06/18) **NEW**
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [FASTQ形式](#) | [description行の記述を整形](#) (last modified 2014/08/21)
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [Illuminaの* seq.txt](#) (last modified 2013/06/13)
- ・ [イントロ](#) | [NGS](#) | [読み込み](#) | [Illuminaの* qseq.txt](#) (last modified 2013/06/17)
- ・ [イントロ](#) | [ファイル形式の変換](#) | [|について](#) (last modified 2014/06/09)
- ・ [イントロ](#) | [ファイル形式の変換](#) | [BAM --> BED](#) (last modified 2014/06/21)
- ・ [イントロ](#) | [ファイル形式の変換](#) | [FASTQ --> FASTA](#) (last modified 2015/06/15) **NEW**
- ・ [イントロ](#) | [ファイル形式の変換](#) | [Genbank --> FASTA](#) (last modified 2014/03/10)
- ・ [イントロ](#) | [ファイル形式の変換](#) | [qseq --> FASTA](#) (last modified 2013/06/17)
- ・ [イントロ](#) | [ファイル形式の変換](#) | [qseq --> Illumina FASTQ](#) (last modified 2013/06/17)
- ・ [イントロ](#) | [ファイル形式の変換](#) | [qseq --> Sanger FASTQ](#) (last modified 2013/08/19)
- ・ [前処理](#) | [クオリティコントロール](#) | [|について](#) (last modified 2015/05/04)
- ・ [前処理](#) | [クオリティチェック](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/15) **NEW**
- ・ [前処理](#) | [クオリティチェック](#) | [qrc](#) (last modified 2014/07/17)
- ・ [前処理](#) | [クオリティチェック](#) | [PHREDスコアに変換](#) (last modified 2013/06/18)
- ・ [前処理](#) | [クオリティチェック](#) | [配列長分布を調べる](#) (last modified 2015/06/22) **NEW**
- ・ [前処理](#) | [クオリティチェック](#) | [Overrepresented sequences](#) | [ShortRead\(Morgan 2009\)](#) (last modified 2015/06/22)
- ・ [前処理](#) | [フィルタリング](#) | [PHREDスコアが低い塩基をNに置換](#) (last modified 2014/03/03)
- ・ [前処理](#) | [フィルタリング](#) | [PHREDスコアが低い配列\(リード\)を除去](#) (last modified 2014/08/27)
- ・ [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](#) (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)

サブセットの抽出

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

FASTQ形式ファイルを読み込むやり方を示します。ファイルサイズが数GBというレベルになってきていますので、圧縮ファイルでの読み込みが基本となりつつあります。しかし、一度に圧縮ファイル中の中身を一旦全て読み込むことも難しい、あるいははやできない状態になってきています。この背景を踏まえ、「応用」では、圧縮FASTQファイルを入力として、メモリーオーバーフロー（スタックオーバーフロー/stack overflow）にならないように一部だけを読み込む手順を示します。また途中段階ですが、野間口達洋氏提供情報をもとにいくつか試しています。

「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. small RNA-seqのgzip圧縮FASTQ形式ファイル(SRR609266.fastq.gz)の場合:

イントロ | NGS | 配列取得 | FASTQ or SRA | SRADB(Zhu, 2013)の7を実行して得られたカイコsmall RNA-seqデータ(Nie et al., BMC Genomics, 2013)です。入力ファイルサイズは400Mb弱、11,928,428リードです。この中から最初の100000リード分をgzip圧縮ファイルとして出力しています。

```
in_f <- "SRR609266.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge1.fasta.gz" #出力ファイル名を指定してout_fに格納
param <- 100000 #抽出したいリード数を指定

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

#入力ファイルの読み込み
fasta <- readDNASTringSet(in_f, format="fastq", #in_fで指定したファイルの読み込み
                          skip=0, nrec=param) #in_fで指定したファイルの読み込み
                                              #確認してるだけです

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



2. small RNA-seqのgzip圧縮FASTQ形式ファイル(SRR609266.fastq.gz)の場合:

イントロ | NGS | 配列取得 | FASTQ or SRA | SRADB(Zhu, 2013)の7を実行して得られた 2015年6月18日現在 (R ver. 3.2.0; Biostrings ver. 2.36.1)、エラーを吐かずにFASTQ形式のgzip圧縮ファイルは得られるのですが、クオリティスコア部分が変わります。つまり、このやり方はダメです。おそらくreadDNASTringSet関数で読み込んだ情報の中にクオリティスコアが含まれていないため、ダミーの「J」というスコアがそのまま表示されているだけだと思います。readDNASTringSet関数でクオリティスコアを読み込むやり方(あるいは取り出し方)が分かった方は教えて下さいm(_ _)m

```
in_f <- "SRR609266.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge2.fastq.gz" #出力ファイル名を指定してout_fに格納
param <- 100000 #抽出したいリード数を指定

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

#入力ファイルの読み込み
fastq <- readDNASTringSet(in_f, format="fastq", #in_fで指定したファイルの読み込み
                          skip=0, nrec=param) #in_fで指定したファイルの読み込み
                                              #確認してるだけです

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



3. small RNA-seqのgzip圧縮FASTQ形式ファイル(SRR609266.fastq.gz)の場合:

原始的なやり方ですが、とりあえずうまく動きます。

```
in_f <- "SRR609266.fastq.gz" #入力ファイル名を指定してin_fに格納
out_f <- "hoge3.fastq" #出力ファイル名を指定してout_fに格納
param <- 100000 #抽出したいリード数を指定

#入力ファイルの読み込み
fastq <- readLines(in_f, n=param*4) #in_fで指定したファイルの読み込み

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



①は出力がFASTA形式、②はFASTQ形式で出力しようとして失敗中(爆)。現状では、最もシンプル且つ無難な③がおススメ。但し、出力がFASTQ形式だがgzip圧縮ファイルにはなっていません。発展課題: 例題4として掲載できる、gzip圧縮FASTQファイルとして5000番目から10000番目までのリードを出力するようなコードをレポートせよ。おまけ: QuasRパッケージはbzip2圧縮ファイル(.bz2)も入力として受け付けてくれるようです。

Contents

■ QC(Quality Control)

- Quality Check (FastQC): 乳酸菌RNA-seqデータ
 - k-mer解析、課題1
- Quality Check (FastQC): カイコスRNA-seqデータ
- トリミング: カイコスRNA-seqデータ
 - 全体像を把握
 - アダプター配列除去(QuasRパッケージ)、課題2
 - Rを駆使して結果を確認
- トリミング: 乳酸菌RNA-seqデータ
 - サブセットの抽出
 - paired-endデータのアダプター配列除去(QuasRパッケージ)

Paired-endの取扱い

Paired-endデータに対応できるパッケージは、まだそれほど多くないという例を示します。

- 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 | 指定した末端塩基数だけ除去 | [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

QuasRパッケージをインストールしてR Guiを終了し、2. gzip圧縮FASTQ形式ファイル([SRR616268sub 1.fastq.gz](#)と[SRR616268sub 2.fastq.gz](#))の場合:

6月21日にparam_nrecを構成しません。デフォルトです。処理時間は長くなり「ファイル」-「ディレクトリ」

1. gzip圧縮状態のFASTQ形式ファイルRNA-seqデータ([SRR609211](#))中の記述から [Girafe\(SRAdb\(Zhu 2013\)\)](#)の低いリードを除去しませんが、Table (例: "GCAGTCGTGGC" 列ということになります)は "TGGAATTCTCGG" まで許容して(推定)が18nt以上のものをブあるという前提であり、Rpatternとmax.Rmismatches

```
in_f <- "SRR609211_1.fastq.gz"
out_f <- "hoge1_1.fastq.gz"
param_adapter <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCAC"
param_nrec <- 500000
```

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 "GATCGGAAGAGCACACGTCTGAACTCCAGTCACCTTAGGCATCTCGTATGCCGTCTTCTGCTTG" 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。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 <- "GATCGGAAGAGCACACGTCTGAACTCCAGTCAC"
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,
                        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です。

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

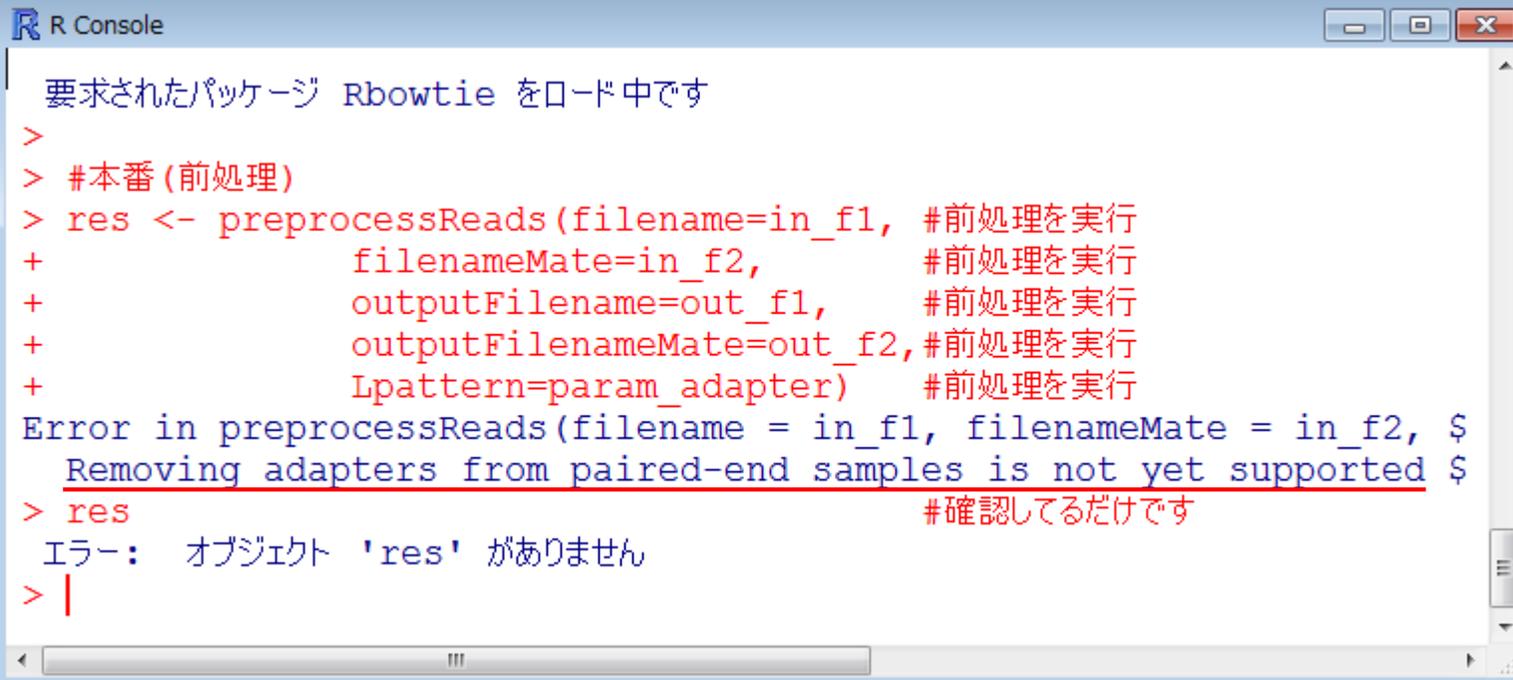
乳酸菌RNA-seqデータSRR616268の最初の100万リード分です。paired-endデータです。SRR616268sub_1.fastq.gzは、約75MB、全リード107 bpです。SRR616268sub_2.fastq.gzは、約67MB、全リード93 bpです。FastQC実行結果として「TruSeq Adapter, Index 3」が含まれているとレポートされました。これでググると塩基配列情報は「GATCGGAAGAGCACACGTCTGAACTCCAGTCACTTAGGCATCTCGTATGCCGTCTTCTGCTTG」と書いてあったので、これを入力として与えます。アダ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。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, nrec=param_nrec)
res
```



```
R Console
要求されたパッケージ Rbowtie をロード中です
>
> #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
+                       Lpattern=param_adapter) #前処理を実行
Error in preprocessReads(filename = in_f1, filenameMate = in_f2, $
  Removing adapters from paired-end samples is not yet supported $
> res #確認してるだけです
エラー: オブジェクト 'res' がありません
> |
```

Tips: sessionInfo

①パッケージ名を指定してバージョン情報を得るやり方。②R本体のバージョンをちゃんと上げていかないと、インストールできるパッケージのバージョンも新しくできないので注意しましょう。

R Console

```
> packageVersion("QuasR")
[1] '1.8.2'
> sessionInfo()
R version 3.2.0 (2015-04-16)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows 7 x64 (build 7601) Service Pack 1

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

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

other attached packages:
[1] QuasR_1.8.2      Rbowtie_1.8.0      GenomicRanges_1.20.3
[4] GenomeInfoDb_1.4.0 IRanges_2.2.1      S4Vectors_0.6.0
[7] BiocGenerics_0.14.0

loaded via a namespace (and not attached):
[1] AnnotationDbi_1.30.1  XVector_0.8.0
[3] zlibbioc_1.14.0      GenomicAlignments_1.4.1
```

原著論文の引用

原著論文はきっちり引用しましょう。
①Bioconductorから提供されている
QuasRパッケージのサイトへのリンク
。②原著論文。③最新バージョン。

2. gzip圧縮FASTQ形式

乳酸菌RNA-seqデータ
は、約75MB、全リード1
して「TruSeq Adapter, In
"GATCGGAAGAGCA
あったので、これを入力
プター配列しかトリムしな
す。残念ながら、QuasR
す。2015年6月23日に開

```
in_f1 <- "SRR6162  
in_f2 <- "SRR6162  
out_f1 <- "hoge1_  
out_f2 <- "hoge1_  
param_adapter <-  
param_nrec <- 500
```

```
#必要なパッケージを  
library(QuasR)
```

```
#本番(前処理)
```

```
res <- preprocess
```

```
fi  
ou  
outp  
Lpat
```

```
res
```

```
<
```

SRP016842: Nie et al

QuasR: Gaidatzis et al., Bioinformatics, 2015



QuasR

platforms some downloads top 20% posts
in Bioc 2 years build ok comm

Quantify and Annotate Short Reads

Bioconductor version: Release (3.1)

This package provides a framework for the quanti
complete workflow starting from raw sequence re
plots, to the quantification of genomic regions of

Author: Anita Lerch, Dimos Gaidatzis and Michael

Maintainer: Michael Stadler <michael.stadler at fr

Citation (from within R, enter `citation("QuasR")`)

Gaidatzis D, Lerch A, Hahne F and Stadler MB (2015). "QuasR: Quantification and annotation of short reads in R." *Bioinformatics*, **31**(7), pp. 1130–1132. <http://dx.doi.org/10.1093/bioinformatics/btu781>, PMID:25417205.

Langmead B, Trapnell C, Pop M and Salzberg SL (2009). "Ultrafast and memory-efficient alignment of short DNA sequences to the human genome." *Genome Biology*, **10**(3), pp. R25. <http://dx.doi.org/10.1186/qb-2009-10-3-r25>, PMID:19261174.

Au KF, Jiang H, Lin L, Xing Y and Wong WH (2010). "Detection of splice junctions from paired-end RNA-seq data by SpliceMap." *Nucleic Acids Research*, **38**(14), pp. 4570–4578. <http://dx.doi.org/10.1093/nar/gkq211>, PMID:20371516.

Details

| | |
|-----------------------|---|
| biocViews | Alignment , ChIPSeq , Coverage , Genetics , MethylSeq , Preprocessing , Quality Control , RNASeq , Sequencing , Software |
| Version | 1.8.3 |
| In Bioconductor since | BioC 2.11 (R-3.0) (2 years) |
| License | GPL-2 |
| Depends | parallel, GenomicRanges(>= 1.13.3) , Rbowtie |
| Imports | methods , zlibbioc , BiocGenerics , S4Vectors , IRanges , BiocInstaller , Biobase , Biostrings , BSgenome , Rsamtools(>= 1.19.38) , GenomicFeatures(>= 1.17.13) , ShortRead(>= 1.19.1) , GenomicAlignments , BiocParallel , GenomeInfoDb |
| LinkingTo | Rsamtools |
| Suggests | rtracklayer , Gviz , RUnit , BiocStyle |
| SystemRequirements | |
| Enhances | |
| URL | |
| Depends On Me | |
| Imports Me | |
| Suggests Me | |

