

機能ゲノム学 第4回

¹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

²微生物科学イノベーション連携研究機構

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/>

Contents

■ マッピング (アラインメント) の続き

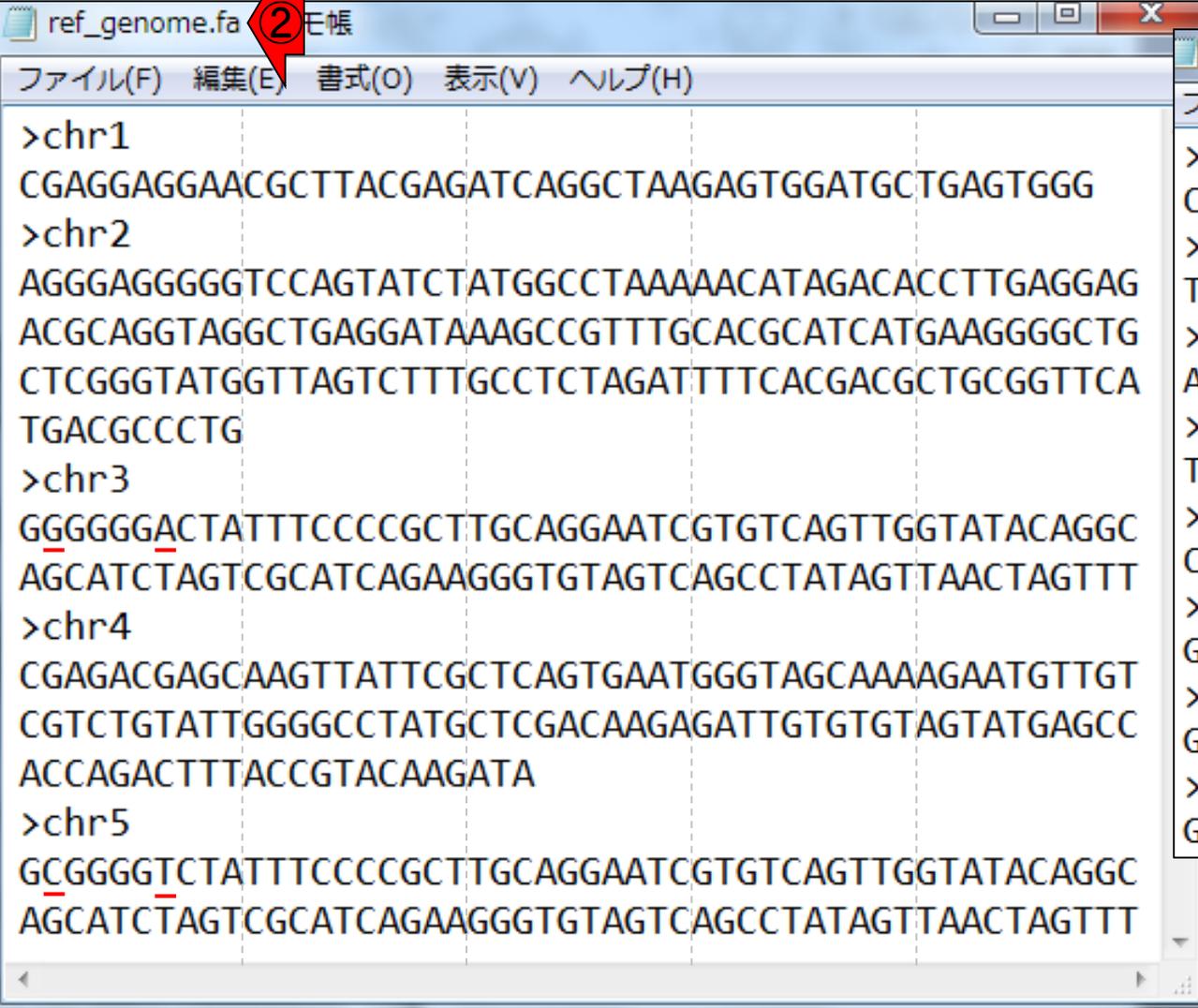
- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

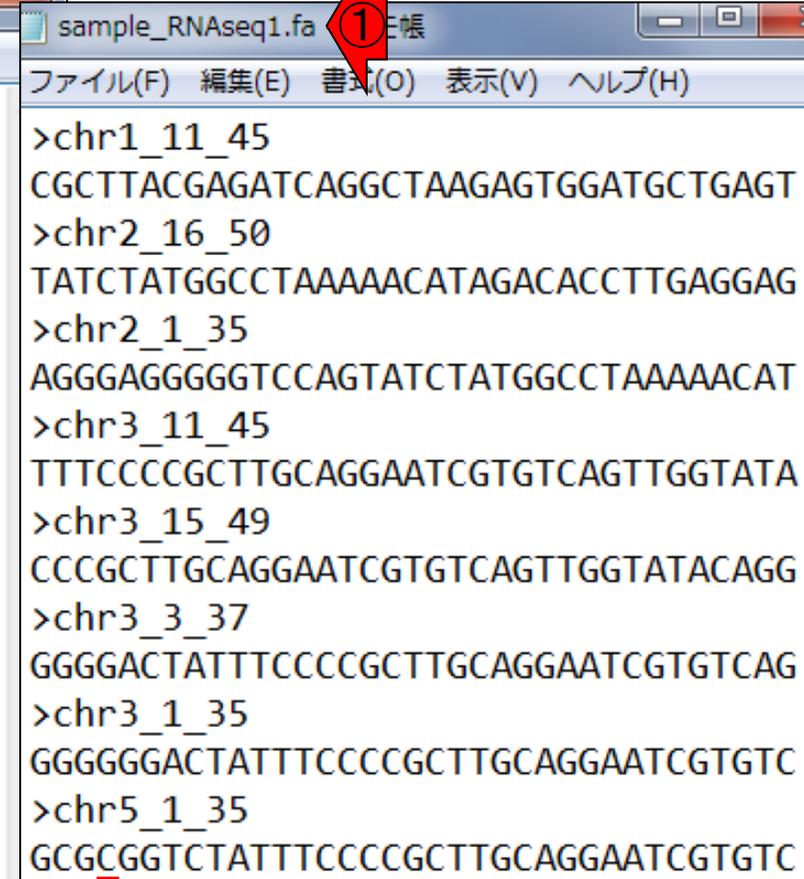
- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`でgenenameの個数を確認

おさらい

マッピングは、マップする側の仮想RNA-seqデータが①sample_RNAseq1.fa、マップされる側のリファレンス配列が② ref_genome.faとして行われた



```
ref_genome.fa ② 帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG
>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACGCCCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
```



```
sample_RNAseq1.fa ① 帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
>chr3_3_37
GGGGACTATTTCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGGGTCTATTTCCCGCTTGCAGGAATCGTGTC
```

マップされる側

```
ref_genome.fa - ノート帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG
>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACCGCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAAGTTTACCGTACAAGATA
>chr5
GGGGGCTCTATTTCCCGCTTGCAGGAATCGTGTCAAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
```

マップされる側

①ref_genome.faのおさらい。②chr3と③chr5の違いは、④2番目と⑤7番目の塩基のみ。従って、8番目の塩基以降は全く同じ

```
ref_genome.fa - ノート帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)

>chr1
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG

>chr2
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGCTG
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA
TGACGCCCTG
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT

>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA

>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
```

①sample_RNAseq1.faのおさらい。全8リードのうち、②最後のリード(chr5_1_35)のみ、③4番目の塩基を変えている

マップする側

```
ref_genome.fa ② 帳  
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)  
>chr1  
CGAGGAGGAACGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGTGGG  
>chr2  
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACATAGACACCTTGAGGAG  
ACGCAGGTAGGCTGAGGATAAAGCCGTTTGCACGCATCATGAAGGGGCTG  
CTCGGGTATGGTTAGTCTTTGCCTCTAGATTTTCACGACGCTGCGGTTCA  
TGACGCCCTG  
>chr3  
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC  
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT  
>chr4  
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT  
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC  
ACCAGACTTTACCGTACAAGATA  
>chr5  
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC  
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
```

```
sample_RNAseq1.fa ① 帳  
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)  
>chr1_11_45  
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT  
>chr2_16_50  
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG  
>chr2_1_35  
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT  
>chr3_11_45  
TTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATA  
>chr3_15_49  
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG  
>chr3_3_37  
GGGGACTATTTCCCGCTTGCAGGAATCGTGTCAG  
>chr3_1_35  
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTC  
>chr5_1_35  
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTC
```

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`でgenenameの個数を確認

QuasRのマッピング結果

①使用したオプション。計8リードのうち、マップされなかったのは赤枠の3リード。Bowtie ver.1に相当するRbowtieというパッケージを利用

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

```
ref_genome.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1
CGAGGAGGA chr1 11 45
>chr2
AGGGAGGGC chr2 1 35 TGGATGCTGAGTGGG
ACGCAGGTA chr2 16 50
CTCGGGTAT chr3 1 35 ATAGACACCTTGAGGAG
TGACGCCCTG chr3 3 37 GCATCATGAAGGGGCTG
>chr3
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
>chr4
CGAGACGAGCAAGTTATTCGCTCAGTGAATGGGTAGCAAAAGAATGTTGT
CGTCTGTATTGGGGCCTATGCTCGACAAGAGATTGTGTGTAGTATGAGCC
ACCAGACTTTACCGTACAAGATA
>chr5
GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGGC
AGCATCTAGTCGCATCAGAAGGGTGTAGTCAGCCTATAGTAACTAGTTT
```

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

完全一致でも複数個所にマップされるために落とされたのは2リード

QuasRのマッピング結果

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

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

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

QuasRのマッピング結果

1か所にのみマップされるがミスマッチのため落とされたのは1リード

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

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

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

GCGGGGTCTATTTCCCGCTTGCAGGAATCGTGTC

Contents

- マッピング (アラインメント) の続き
 - おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
 - マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
 - SAM形式の解説、マッピング結果の違い、課題
 - Linux環境以外でのBowtie2実行手段
- カウント情報取得
 - アノテーション情報がない場合: 単一サンプル、複数サンプル
 - アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - grepでgenenameの個数を確認

bowtie2実行結果

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1] 16:18
iu@bielinux[mapping_kiso1] bowtie2 --version [ 4:18午後 ]
/usr/bin/../../lib/bowtie2/bin/bowtie2-align-s version 2.2.4
64-bit
Built on lgw01-04
Fri Dec 12 17:13:13 UTC 2014
Compiler: gcc version 4.8.2 (Ubuntu 4.8.2-19ubuntu1)
Options: -O3 -m64 -msse2 -funroll-loops -g3 -DPOPCNT_CAPABILIT
Y
Sizeof {int, long, long long, void*, size_t, off_t}: {4, 8, 8,
8, 8, 8}
iu@bielinux[mapping_kiso1] [ 4:18午後 ]
```

bowtie2実行結果

①1回だけマップされたリードは3個 (37.50%)。これは1か所へのみマップされたリードと解釈すればよい。②2回以上(複数個所に)マップされたリードは5個 (62.50%)。③マップ率 (alignment rate) は100%

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
/home/iu/Desktop/mac_share/mapping_kiso1
iu@bielinux[mapping_kiso1] ls -l
total 8197
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu 53 5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu 133 5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu 590 9月 29 2013 ref_genome.fa
-rwxrwxrwx 1 iu iu 396 10月 1 2013 sample_RNAseq1.fa
iu@bielinux[mapping_kiso1] bowtie2 -x pigya -f sample_RNAseq1.f
a -S sample_RNAseq1.sam
8 reads; of these:
  8 (100.00%) were unpaired; of these:
    0 (0.00%) aligned 0 times
    ① 3 (37.50%) aligned exactly 1 time
    ② 5 (62.50%) aligned >1 times
③ 100.00% overall alignment rate
iu@bielinux[mapping_kiso1] █
```

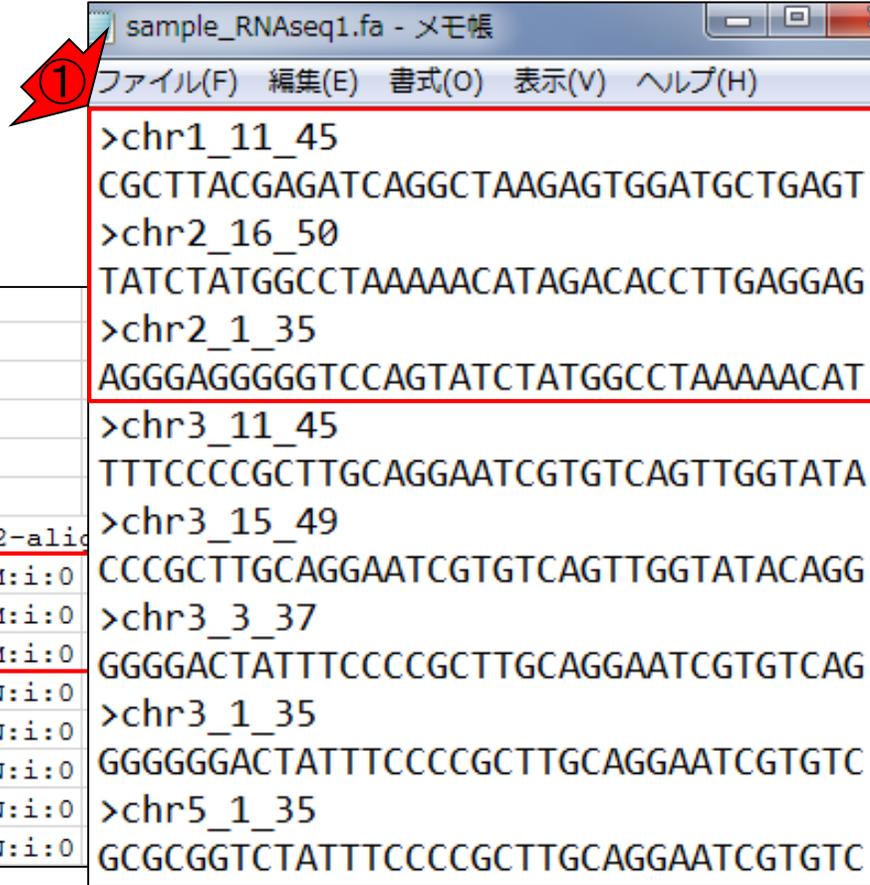
[3:08午後]

[3:08午後]

Bowtie2実行結果

①1か所にのみマップされた
3リードの、②マッピング結果

@HD	VN:1.0	SO:unsorted																		
@SQ	SN:chr1	LN:48																		
@SQ	SN:chr2	LN:160																		
@SQ	SN:chr3	LN:100																		
@SQ	SN:chr4	LN:123																		
@SQ	SN:chr5	LN:100																		
@PG	ID:bowt	PN:bow	VN:2	CL:"/usr/bin/../../lib/bowtie2/bin/bowtie2-align																
chr1_11_45	0	chr1	11	42	35M	*	0	0	CGIIAS:i:0	XN:i:0	XM:i:0									
chr2_16_50	0	chr2	16	42	35M	*	0	0	TIIAS:i:0	XN:i:0	XM:i:0									
chr2_1_35	0	chr2	1	42	35M	*	0	0	ACIIAS:i:0	XN:i:0	XM:i:0									
chr3_11_45	0	chr5	11	1	35M	*	0	0	TIIAS:i:0	XS:i:0	XN:i:0									
chr3_15_49	0	chr3	15	1	35M	*	0	0	CCIIAS:i:0	XS:i:0	XN:i:0									
chr3_3_37	0	chr3	3	31	35M	*	0	0	GIIAS:i:0	XS:i:-6	XN:i:0									
chr3_1_35	0	chr3	1	35	35M	*	0	0	GIIAS:i:0	XS:i:-12	XN:i:0									
chr5_1_35	0	chr5	1	16	35M	*	0	0	GIIAS:i:-6	XS:i:-18	XN:i:0									



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

Bowtie2実行結果

①複数個所にマップされた5リードの、②マッピング結果

@HD	VN:1.0	SO:unsorted																		
@SQ	SN:chr1	LN:48																		
@SQ	SN:chr2	LN:160																		
@SQ	SN:chr3	LN:100																		
@SQ	SN:chr4	LN:123																		
@SQ	SN:chr5	LN:100																		
@PG	ID:bowt2	PN:bowt2	VN:2	CL:"/usr/bin/../../lib/bowtie2/bin/bowtie2-align																
chr1_11_45	0	chr1	11	42	35M	*	0	0	CGIIAS:i:0	XN:i:0	XM:i:0									
chr2_16_50	0	chr2	16	42	35M	*	0	0	TIIAS:i:0	XN:i:0	XM:i:0									
chr2_1_35	0	chr2	1	42	35M	*	0	0	ACIIAS:i:0	XN:i:0	XM:i:0									
chr3_11_45	0	chr5	11	1	35M	*	0	0	TIIAS:i:0	XS:i:0	XN:i:0									
chr3_15_49	0	chr3	15	1	35M	*	0	0	CCIIAS:i:0	XS:i:0	XN:i:0									
chr3_3_37	0	chr3	3	31	35M	*	0	0	GIIAS:i:0	XS:i:-6	XN:i:0									
chr3_1_35	0	chr3	1	35	35M	*	0	0	GIIAS:i:0	XS:i:-12	XN:i:0									
chr5_1_35	0	chr5	1	16	35M	*	0	0	GIIAS:i:-6	XS:i:-18	XN:i:0									

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



Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`でgenenameの個数を確認

bowtie2実行結果

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
/home/iu/Desktop/mac_share/mapping_kiso1
iu@bielinux[mapping_kiso1] ls -l [ 3:08午後 ]
total 8197
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu 53 5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu 133 5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu 590 9月 29 2013 ref_genome.fa
-rwxrwxrwx 1 iu iu 396 10月 1 2013 sample_RNAseq1.fa
iu@bielinux[mapping_kiso1] bowtie2 -x pigya -f sample_RNAseq1.f
a -S sample_RNAseq1.sam
8 reads; of these:
  8 (100.00%) were unpaired; of these:
    0 (0.00%) aligned 0 times
    3 (37.50%) aligned exactly 1 time
    5 (62.50%) aligned >1 times
① 100.00% overall alignment rate
iu@bielinux[mapping_kiso1] [ 3:08午後 ]
```


- ①basenameをagriにして、
- ②bowtie-buildを実行

リファレンス配列の前処理

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
iu@bielinux[mapping_kiso1] pwd [ 9:51午後 ]
/home/iu/Desktop/mac_share/mapping_kiso1
iu@bielinux[mapping_kiso1] ls -l [ 9:51午後 ]
total 8199
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu 53 5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu 133 5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu 590 9月 29 2013 ref_genome.fa
-rwxrwxrwx 1 iu iu 396 10月 1 2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu 1616 5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie-build ref_genome.fa agri
```



bowtie-build完了

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
eftabSz: 80
ftabLen: 1048577
ftabSz: 4194308
offsLen: 17
offsSz: 68
isaLen: 0
isaSz: 0
lineSz: 64
sideSz: 64
sideBwtSz: 56
sideBwtLen: 224
numSidePairs: 2
numSides: 4
numLines: 4
ebwtTotLen: 256
ebwtTotSz: 256
reverse: 0
Total time for backward call to driver() for mirror index: 00:00:00
iu@bielinux[mapping_kiso1] clear [ 9:54午後 ]
```

lsで確認

①lsで確認。作成されたインデックスファイルの拡張子部分が
②bowtie1ではebwtとなっており、③bowtie2のbt2とは異なるこ
とが分かります。このことから、bowtie2で作成したインデックス
ファイルはbowtie1では利用できないのだろうと思ったりします

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
iu@bielinux[mapping_kiso1] pwd
/home/iu/Desktop/mac_share/mapping_kiso1
iu@bielinux[mapping_kiso1] ls -l
total 16394
-rwxrwxrwx 1 iu iu 4194810  5月 17 21:54 agri.1.ebwt
-rwxrwxrwx 1 iu iu      72  5月 17 21:54 agri.2.ebwt
-rwxrwxrwx 1 iu iu      53  5月 17 21:54 agri.3.ebwt
-rwxrwxrwx 1 iu iu     133  5月 17 21:54 agri.4.ebwt
-rwxrwxrwx 1 iu iu 4194810  5月 17 21:54 agri.rev.1.ebwt
-rwxrwxrwx 1 iu iu      72  5月 17 21:54 agri.rev.2.ebwt
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu      53  5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu     133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu     590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu     396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu    1616  5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1]
```



Bowtie実行

①bowtie実行コマンド。Bowtie2の時との違いは、②マ
ップされる側のbasename情報(agri)を指定する際に-x
オプションをつけていない点。Bowtie (ver. 1)のときは-x
をつけてはいけません(つけると動きません)。尚、出
力SAMファイル名は③bowtie1_default.samとしています

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
/home/iu/Desktop/mac_share/mapping_k
iu@bielinux[mapping_kiso1] ls -l
total 16394
-rwxrwxrwx 1 iu iu 4194810  5月 17 21:54 agri.1.ebwt
-rwxrwxrwx 1 iu iu      72  5月 17 21:54 agri.2.ebwt
-rwxrwxrwx 1 iu iu      53  5月 17 21:54 agri.3.ebwt
-rwxrwxrwx 1 iu iu     133  5月 17 21:54 agri.4.ebwt
-rwxrwxrwx 1 iu iu 4194810  5月 17 21:54 agri.rev.1.ebwt
-rwxrwxrwx 1 iu iu      72  5月 17 21:54 agri.rev.2.ebwt
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu      53  5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu     133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu     590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu     396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu    1616  5月 15 15:08 sample_RNAseq1.sam
① iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
   bowtie1_default.sam
```

①

②

③

①bowtie実行コマンドが無事終了しました。Bowtie2のときとは、出力のされかたが若干異なりますね

Bowtie実行

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu      53  5月 17 21:54 agri.3.ebwt
-rwxrwxrwx 1 iu iu     133  5月 17 21:54 agri.4.ebwt
-rwxrwxrwx 1 iu iu 4194810  5月 17 21:54 agri.rev.1.ebwt
-rwxrwxrwx 1 iu iu      72  5月 17 21:54 agri.rev.2.ebwt
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu      53  5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu     133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu     590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu     396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu    1616  5月 15 15:08 sample_RNAseq1.sam
① iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] [10:30午後]
```

①bowtieもデフォルトオプションで実行すると、全リードがマップされてしまいました…

Bowtie実行

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu      53  5月 17 21:54 agri.3.ebwt
-rwxrwxrwx 1 iu iu     133  5月 17 21:54 agri.4.ebwt
-rwxrwxrwx 1 iu iu 4194810  5月 17 21:54 agri.rev.1.ebwt
-rwxrwxrwx 1 iu iu      72  5月 17 21:54 agri.rev.2.ebwt
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu      53  5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu     133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu     140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu     590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu     396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu    1616  5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] █
```



[10:30午後]

sample_RNAseq1.sam

おさらい。①Bowtie2をデフォルトオプションで実行した結果ファイル(sample_RNAseq1.sam)。①と②が対応する箇所。それに加えて、③のマップされた配列名も異なっていることがわかる。しかしこのこと自体はどちらにマップされててもいいので大した問題ではない

Header	Field 1	Field 2	Field 3	Field 4	Field 5	Field 6	Field 7	Field 8	Field 9	Field 10	Field 11	Field 12	Field 13	Field 14	Field 15	Field 16	Field 17		
@HD	VN:1.0	SO:unsorted																	
@SQ	SN:chr1	LN:48																	
@SQ	SN:chr2	LN:160																	
@SQ	SN:chr3	LN:100																	
@SQ	SN:chr4	LN:123																	
@SQ	SN:chr5	LN:100																	
@PG	ID:bowt	PN:bow	VN:2	CM:"/usr/bin/../../lib/bowtie2/bin/bowtie2-align-s --wrapper basic-0 -x pigya -f sample_RNAseq1.sam															
chr1_11_45	0	chr1	11	42	35M	*	0	0	CGIIAS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU			
chr2_16_50	0	chr2	16	42	35M	*	0	0	TATIAS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU			
chr2_1_35	0	chr2	1	42	35M	*	0	0	ACIIAS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU			
chr3_11_45	0	chr5	11	1	35M	*	0	0	TTIIAS:i:0	XS:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU			
chr3_15_49	0	chr3	15	1	35M	*	0	0	CCIIAS:i:0	XS:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU			
chr3_3_37	0	chr3	3	31	35M	*	0	0	GCTIAS:i:0	XS:i:-6	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU			
chr3_1_35	0	chr3	1	35	35M	*	0	0	GCTIAS:i:0	XS:i:-12	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU		
chr5_1_35	0	chr5	1	16	35M	*	0	0	GCTIAS:i:-6	XS:i:-18	XN:i:0	XM:i:1	XO:i:0	XG:i:0	NM:i:1	MD:Z:3G31	YT:Z:UU		

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`で`genename`の個数を確認

QuasR

QuasRは、内部的にBowtie (ver.1)を用いている。①で示すオプションを指定してマッピングを行った結果として、8リード中3リードがマップされなかった。その結果を再現すべく、①のオプションを付けて、Bio-Linux上でBowtie (ver.1)を実行する

- マッピング | [\(ESTレベルの長さの\)contig](#) (last modified 2014/06/19)
- マッピング | [基礎](#) (last modified 2013/06/19)
- マッピング | [single-end | ゲノム | basic aligner\(基礎\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- マッピング | [single-end | ゲノム | basic aligner\(応用\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/28)
- マッピング | [single-end | ゲノム | splice-aware aligner](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- マッピング | [paired-end | ゲノム | basic aligner\(基礎\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- マッピング | [paired-end | ゲノム | splice-aware aligner](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- マッピング | [paired-end | ゲノム | basic aligner\(応用\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- マッピング | [paired-end | ゲノム | splice-aware aligner](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- マップ後 | [出力](#)

マッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)

QuasRパッケージを用いて single-end RNA-seq データのリファレンスゲノム配列へのマッピングを行うやり方を示します。basic alignerの1つであるBowtie (Langmead et al., Genome Biol., 2009)を実装した Rbowtieパッケージを内部的に使っています。Bowtie自体は、複数個所にマップされるリードの取り扱い (uniquely mapped reads or multi-mapped reads) を "-m" オプションで指定したり、許容するミスマッチ数を指定する "-v" などの様々なオプションを利用可能ですが、「基礎」のところではやり方を示しませんでした。ここでは、マッピングのオプションをいくつか変更して挙動を確認したり、複数のRNA-seqファイルを一度にマッピングするやり方を示します。尚、出力ファイルは、"*_bam", "*_QC.pdf", "*_bed" の3つです。それ以外のファイルは基本無視で大丈夫です。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

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

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

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

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

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



①

①が、②オプション(-m 1 --best --strata -v 0)つきの
実行コマンド。③出力ファイル名はbowtie1_QuasR.sam

実行コマンド

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu 133 5月 17 21:54 agri.4.ebwt
-rwxrwxrwx 1 iu iu 4194810 5月 17 21:54 agri.rev.1.ebwt
-rwxrwxrwx 1 iu iu 72 5月 17 21:54 agri.rev.2.ebwt
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.2.bt2
-rwxrwxrwx 1 iu iu 53 5月 14 16:39 pigya.3.bt2
-rwxrwxrwx 1 iu iu 133 5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746 5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu 140 5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu 590 9月 29 2013 ref_genome.fa
-rwxrwxrwx 1 iu iu 396 10月 1 2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu 1616 5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v 0 agr
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
```



実行結果

①コマンド実行結果。②ぱっと見でQuasR上でのBowtie (ver.1)実行結果と同じだろうと安心する

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu      133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu      140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu      590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu      396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu     1616  5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v 0 agr
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] [ 3:29午後 ]
```

①

②

①この2リードが、`-m 1`オプション(1か所にのみマップされたリードを出力)という条件を満たさなかった…

実行結果

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu      133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu      140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu      590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu      396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu     1616  5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v 0 agr
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] █ [ 3:29午後 ]
```



①

複数個所にマップされるこの2つのリードなのだろう

QuasRのマッピング結果

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

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

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

実行結果

①この1リードが、`-v 0`オプション(許容する mismatches数は0)という条件を満たさなかった…

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu      133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu      140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu      590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu      396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu     1616  5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v 0 agr
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] █ [ 3:29午後 ]
```



①

(1か所にのみマップされるが)①の箇所でミスマッチがある、このリードなのだろう

QuasRのマッピング結果

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

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

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



オプション(-m 1 --best --strata -v 0)つきの、①
実行結果ファイル(bowtie1_QuasR.sam)の中身が...

実行結果

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso1]
-rwxrwxrwx 1 iu iu      133  5月 14 16:39 pigya.4.bt2
-rwxrwxrwx 1 iu iu 4194746  5月 14 16:39 pigya.rev.1.bt2
-rwxrwxrwx 1 iu iu      140  5月 14 16:39 pigya.rev.2.bt2
-rwxrwxrwx 1 iu iu      590  9月 29  2013 ref_genome.fa
-rwxrwxrwx 1 iu iu      396 10月  1  2013 sample_RNAseq1.fa
-rwxrwxrwx 1 iu iu     1616  5月 15 15:08 sample_RNAseq1.sam
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S
bowtie1_default.sam
# reads processed: 8
# reads with at least one reported alignment: 8 (100.00%)
# reads that failed to align: 0 (0.00%)
Reported 8 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v 0 agr
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
iu@bielinux[mapping_kiso1] █ [ 3:29午後 ]
```



bowtie1_QuasR.

① マップされなかったリードの違いを表しているのが赤枠部分。② -m 1 オプション (1か所にのみマップされたリードを出力) という条件を満たさなかったリードと、③ -v 0 オプション (許容するミスマッチ数は0) という条件を満たさなかったリードであることが既知の状態、④ XM: のところの数値の説明を読むとわかりやすい。逆に言えば、そういう実例をいくつか知ったうえでないと、いきなり説明を読んでもチンプンカンプンである場合が多い (個人の感想です)

Header	Field	Value
@HD	VN:1.0	SO:unsorted
@SQ	SN:chr1	LN:48
@SQ	SN:chr2	LN:160
@SQ	SN:chr3	LN:100
@SQ	SN:chr4	LN:123
@SQ	SN:chr5	LN:100
@PG	ID:Bowt:	VN:1.1 CL:"bowtie --wrapper basic-0 -m 1 --best --strata -v 0 agri -f sample_RNAseq1.
chr1_11_45	0	chr1 11 255 35M * 0 0 CC I I XA:i:0 MD:Z:35 NM:i:0
chr2_16_50	0	chr2 16 255 35M * 0 0 TT I I XA:i:0 MD:Z:35 NM:i:0
chr2_1_35	0	chr2 1 255 35M * 0 0 AC I I XA:i:0 MD:Z:35 NM:i:0
chr3_11_45	4	* 0 0 * * 0 0 TT I I XA:i:1 XM:i:1
chr3_15_49	4	* 0 0 * * 0 0 CC I I XA:i:1 XM:i:1
chr3_3_37	0	chr3 3 255 35M * 0 0 GC I I XA:i:0 MD:Z:35 NM:i:0
chr3_1_35	0	chr3 1 255 35M * 0 0 GC I I XA:i:0 MD:Z:35 NM:i:0
chr5_1_35	4	* 0 0 * * 0 0 GC I I XA:i:0 XM:i:0



Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`でgenenameの個数を確認

SAM形式の解説

①SAM形式のBowtie出力ファイルの説明。②の部分が12となるまでページ下部に移動

4. 平成30年05月29日

講義資料PDF

[bowtie1_default.sam](#)(Bowtieのデフォルトオプション実行結果)

[bowtie1_QuasR.sam](#)(-m 1 --best --strata -L 10での実行結果)

[Bowtieマニュアル中のSAM bowtie output](#)

SAM bowtie output

Following is a brief description of the SAM format as output by bowtie when the `-S/--sam` option is specified. For more details, see the [SAM format specification](#).

When `-S/--sam` is specified, bowtie prints a SAM header with @HD, @SQ and @PG lines. When one or more `--sam-RG` arguments are specified, bowtie will also print an @RG line that includes all user-specified `--sam-RG` tokens separated by tabs.

Each subsequent line corresponds to a read or an alignment. Each line is a collection of at least 12 fields separated by tabs; from left to right, the fields

1. Name of read that aligned
2. Sum of all applicable flags. Flags relevant to Bowtie are:
 - 1 The read is one of a pair
 - 2 The alignment is one end of a proper paired-end alignment
 - 4 The read has no reported alignments
 - 8 The read is one of a pair and has no reported alignments

SAM形式の解説

①12の、②のあたりにXM:についての説明があります。③赤下線部分が今回指定した-m 1 (1か所のみマップされたリードを出力)と関連しています

12. Optional fields. Fields are tab-separated. For descriptions of all possible optional fields, see the SAM format specification. bowtie outputs some of these optional fields for each alignment, depending on the type of the alignment:

NM:i:<N> Aligned read has an edit distance of <N>.

CM:i:<N> Aligned read has an edit distance of <N> in colorspace. This field is present in addition to the NM field in -C/--color mode, but is omitted otherwise.

MD:Z:<S> For aligned reads, <S> is a string representation of the mismatched reference bases in the alignment. See SAM format specification for details. For colorspace alignments, <S> describes the decoded nucleotide alignment, not the colorspace alignment.

XA:i:<N> Aligned read belongs to stratum <N>. See Strata for definition.

XM:i:<N> For a read with no reported alignments, <N> is 0 if the read had no alignments. If -m was specified and the read's alignments were suppressed because the -m ceiling was exceeded, <N> equals the -m ceiling 1, to indicate that there were at least that many valid alignments (but all were suppressed). In -M mode, if the alignment was randomly selected because the -M ceiling was exceeded, <N> equals the -M ceiling 1, to indicate that there were at least that many valid alignments (of which one was reported at random).

XM:i:<N>

For a read with no reported alignments, <N> is 0 if the read had no alignments. If -m was specified and the read's alignments were suppressed because the -m ceiling was exceeded, <N> equals the -m ceiling 1, to indicate that there were at least that many valid alignments (but all were suppressed). In -M mode, if the alignment was randomly selected because the -M ceiling was exceeded, <N> equals the -M ceiling 1, to indicate that there were at least that many valid alignments (of which one was reported at random).

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`で`genename`の個数を確認

マッピング結果の違い

Bio-Linux環境で行ったマッピングは、① bowtie2のデフォルト、②bowtieのデフォルト、③bowtieのオプション(-m 1 --best --strata -v 0)つきの計3通りであった

```
iu@bielinux[mapping_kiso1] bowtie2 -x pigya -f sample_RNAseq1.f  
a -S sample_RNAseq1.sam  
8 reads; of these:  
 8 (100.00%) were unpaired; of these:  
  0 (0.00%) aligned 0 times  
  3 (37.50%) aligned exactly 1 time  
  5 (62.50%) aligned >1 times  
100.00% overall alignment rate
```



```
iu@bielinux[mapping_kiso1] bowtie agri -f sample_RNAseq1.fa -S  
bowtie1_default.sam  
# reads processed: 8  
# reads with at least one reported alignment: 8 (100.00%)  
# reads that failed to align: 0 (0.00%)  
Reported 8 alignments to 1 output stream(s)
```



```
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v 0 agr  
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam  
# reads processed: 8  
# reads with at least one reported alignment: 5 (62.50%)  
# reads that failed to align: 1 (12.50%)  
# reads with alignments suppressed due to -m: 2 (25.00%)  
Reported 5 alignments to 1 output stream(s)
```



①bowtie2のデフォルトの結果は、②3リードが1個所に、③5リードが複数個所にマップされるというものであった

マッピング結果の違い

①
iu@bielinux[mapping_kiso1] bowtie2 -x pigya -f sample_RNAseq1.f
a -S sample_RNAseq1.sam
8 reads; of these:
8 (100.00%) were unpaired; of these:
0 (0.00%) aligned 0 times
3 (37.50%) aligned exactly 1 time ②
5 (62.50%) aligned >1 times ③
100.00% overall alignment rate

sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
②
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
③
>chr3_3_37
GGGGACTATTTCCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGGACTATTTCCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGCGGTCTATTTCCCCGCTTGCAGGAATCGTGTC

マッピング結果

①bowtie2のデフォルトの結果は、②3リードが1個所に、③5リードが複数個所にマップされるというものであった。④bowtieのオプション(-m 1 --best --strata -v 0)つきの結果は、⑤-v 0オプション(許容するミスマッチ数は0)を満たさなかった1リードと、

```
iu@bielinux[mapping_kiso1] bowtie2 -x pigya -f sample_RNAseq1.f
a -S sample_RNAseq1.sam
8 reads; of these:
 8 (100.00%) were unpaired; of these:
 0 (0.00%) aligned 0 times
 3 (37.50%) aligned exactly 1 time
 5 (62.50%) aligned >1 times
100.00% overall alignment rate
```

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

```
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
```

マッピング結果

① bowtie2のデフォルトの結果は、② 3リードが1個所に、③ 5リードが複数個所にマップされるというものであった。④ bowtieのオプション(-m 1 --best --strata -v 0) 付きの結果は、⑤ -v 0 オプション(許容するミスマッチ数は0)を満たさなかった1リードと、⑥ -m 1 オプション(1か所へのみマップ)を満たさなかったリード

```
iu@bielinux[mapping_kiso1] bowtie2
a -S sample_RNAseq1.sam
8 reads; of these:
 8 (100.00%) were unpaired; of these:
 0 (0.00%) aligned 0 times
 3 (37.50%) aligned exactly 1 time
 5 (62.50%) aligned >1 times
100.00% overall alignment rate
```

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

```
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
```

課題1の基礎情報

①6番目と7番目のリードは、②bowtie2では2回以上マップされ(aligned >1 times)、③bowtieでは少なくとも1回はマップ(at least one)されている

```
iu@bielinux[mapping_kiso1] bowtie2 -x pigya -f sample_RNAseq1.f
a -S sample_RNAseq1.sam
8 reads; of these:
 8 (100.00%) were unpaired; of these:
 0 (0.00%) aligned 0 times
 3 (37.50%) aligned exactly 1 time
 5 (62.50%) aligned >1 times
100.00% overall alignment rate
```

②

```
sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCCGCTTGCAGGAATCGTGTCAGTTGGTATA
xchr3_15_49
① CCGCTTGCAGGAATCGTGTCAGTTGGTATACAGG
>chr3_3_37
GGGGACTATTTCCCCGCTTGCAGGAATCGTGTCAG
>chr3_1_35
GGGGGGACTATTTCCCCGCTTGCAGGAATCGTGTC
>chr5_1_35
GCGCGGTCTATTTCCCCGCTTGCAGGAATCGTGTC
```

①

```
iu@bielinux[mapping_kiso1] bowtie -m 1 --best --strata -v
i -f sample_RNAseq1.fa -S bowtie1_QuasR.sam
# reads processed: 8
# reads with at least one reported alignment: 5 (62.50%)
# reads that failed to align: 1 (12.50%)
# reads with alignments suppressed due to -m: 2 (25.00%)
Reported 5 alignments to 1 output stream(s)
```

③

課題1

①6番目のリード(chr3_3_37)と②7番目のリード(chr3_1_35)について、③リファレンス配列上のどこにマップされたのか示せ。リファレンス配列名(例: chr1)とマップされた領域の左端の位置(例: 8番目の塩基)のみでよい



②



①



②

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

```
sample_RNAseq1.fa - メモ帳
ファイル(F) 編集(E) 書式(O) 表示(V) ヘルプ(H)
>chr1_11_45
CGCTTACGAGATCAGGCTAAGAGTGGATGCTGAGT
>chr2_16_50
TATCTATGGCCTAAAAACATAGACACCTTGAGGAG
>chr2_1_35
AGGGAGGGGGTCCAGTATCTATGGCCTAAAAACAT
>chr3_11_45
TTTCCCGCTTGCAGGAATCGTGTTCAGTTGGTATA
>chr3_15_49
CCCGCTTGCAGGAATCGTGTTCAGTTGGTATACAGG
>chr3_3_37
GGGGACTATTTCCCGCTTGCAGGAATCGTGTTCAG
>chr3_1_35
GGGGGACTATTTCCCGCTTGCAGGAATCGTGTTC
>chr5_1_35
GCGCGGTCTATTTCCCGCTTGCAGGAATCGTGTTC
```


課題2

①6番目のリード(chr3_3_37)と②7番目のリード(chr3_1_35)のマッピング結果「リファレンス配列名とマップされた領域の左端の位置」として、③「chr3上の3番目の塩基」と④「chr3上の1番目の塩基」が採用された理由について自由に考えを述べよ。

sample_RNAseq1.sam

Header	Field 1	Field 2	Field 3	Field 4	Field 5	Field 6	Field 7	Field 8	Field 9	Field 10	Field 11	Field 12	Field 13	Field 14	Field 15	Field 16	Field 17
@HD	VN:1.0	SO:unsorted															
@SQ	SN:chr1	LN:48															
@SQ	SN:chr2	LN:160															
@SQ	SN:chr3	LN:100															
@SQ	SN:chr4	LN:123															
@SQ	SN:chr5	LN:100															
@PG	ID:bowt:	PN:bow	VN:2	CL:"/usr/bin/../../lib/bowtie2/bin/bowtie2-align-s --wrapper basic-0 -x pigya -f sample_RNAS													
chr1_11_45	0	chr1	11	42	35M	*	0	0	CGIIAS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU	
chr2_16_50	0	chr2	16	42	35M	*	0	0	TIIAS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU	
chr2_1_35	0	chr2	1	42	35M	*	0	0	ACIIAS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU	
chr3_11_45	0	chr5	11	1	35M	*	0	0	TTIIAS:i:0	XS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU
chr3_15_45	0	chr3	15	35M	*	0	0	0	CCIIAS:i:0	XS:i:0	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU
chr3_3_37	0	chr3	3	31	35M	*	0	0	GCGIIAS:i:0	XS:i:-6	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU
chr3_1_35	0	chr3	1	35	35M	*	0	0	GCGIIAS:i:0	XS:i:-12	XN:i:0	XM:i:0	XO:i:0	XG:i:0	NM:i:0	MD:Z:35	YT:Z:UU
chr5_1_35	0	chr5	1	35M	*	0	0	0	GCGIIAS:i:-6	XS:i:-18	XN:i:0	XM:i:1	XO:i:0	XG:i:0	NM:i:1	MD:Z:3G31	YT:Z:UU

Contents

- マッピング (アラインメント) の続き
 - おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
 - マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
 - SAM形式の解説、マッピング結果の違い、課題
 - Linux環境以外でのBowtie2実行手段
- カウント情報取得
 - アノテーション情報がない場合: 単一サンプル、複数サンプル
 - アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - grepでgenenameの個数を確認

他のBowtie2実行手段

- DDBJ Pipeline (Nagasaki et al., DNA Res., 20: 383–90, 2013) ①
 - DDBJが提供するクラウド解析環境
- Galaxy (Goecks et al., Genome Biol., 11: R86, 2010) ②
 - Galaxy projectが提供するクラウド解析環境。Galaxy mainというサイトが有名。
- Illumina BaseSpace
 - Illumina社が提供するクラウド解析環境。
- ...

他のBowtie2実行手段

乳酸菌学会誌NGS連載の、①第6回がDDBJ Pipeline、②第11-12回がGalaxyについての解説。但し、マッピングについては書かれていない

(Rで)塩基配列解析

(last modified 2018/05/01, since 2010)

このウェブページのR関連部分は、[インストール](#)についての推奨手順(Windows2018.03.12版とMacintosh)済みであるという前提で記述しています。初心者の方は[基本的な利用法](#)(Windows2015.04.03版)的にまとめた[書籍](#)もあります。(2015/04/03)

What's new?

- Silhouetteスコアの新たな使い道提唱論文([Zhao et al.](#))
- Silhouetteスコアの新たな使い道提唱論文([Zhao et al.](#))
- 「平成29年度NGSハンズオン講習会」の[動画](#)が公開
- [門田からメール返信をもらえない場合は](#) (last modified 2015/03/31)
- [はじめに](#) (last modified 2015/03/31)
- 参考資料 | [書籍](#)、[学会誌](#) (last modified 2017/11/13)
- 参考資料 | [講習会](#)、[講義](#)、[講演資料](#) (last modified 2017/11/13)

- 書籍 | トランスクリプトーム解析 | [4.3.3 2群間比較](#) (last modified 2014/04/28)
- 書籍 | トランスクリプトーム解析 | [4.3.4 他の実験デザイン\(3群間\)](#) (last modified 2014/04/28)
- 書籍 | [日本乳酸菌学会誌](#) | [について](#) (last modified 2018/05/10) **NEW**
- 書籍 | [日本乳酸菌学会誌](#) | [第1回イントロダクション](#) (last modified 2016/12/22)
- 書籍 | [日本乳酸菌学会誌](#) | [第2回GUI環境からコマンドライン環境へ](#) (last modified 2015/11/26)
- 書籍 | [日本乳酸菌学会誌](#) | [第3回Linux環境構築からNGSデータ取得まで](#) (last modified 2017/07/02)
- 書籍 | [日本乳酸菌学会誌](#) | [第4回クオリティコントロールとプログラムのインストール](#) (last modified 2017/06/25)
- 書籍 | [日本乳酸菌学会誌](#) | [第5回アセンブル、マッピング、そしてQC](#) (last modified 2017/06/25)
- 書籍 | [日本乳酸菌学会誌](#) | [第6回ゲノムアセンブリ](#) (last modified 2017/06/21)
- 書籍 | [日本乳酸菌学会誌](#) | [第7回ロングリードアセンブリ](#) (last modified 2017/06/28)
- 書籍 | [日本乳酸菌学会誌](#) | [第8回アセンブリ後の解析](#) (last modified 2017/06/28)
- 書籍 | [日本乳酸菌学会誌](#) | [第9回ゲノムアノテーションとその可視化、DDBJへの登録](#) (last modified 2017/06/28)
- 書籍 | [日本乳酸菌学会誌](#) | [第10回DDBJへの塩基配列の登録\(後編\)](#) (last modified 2017/06/28)
- 書籍 | [日本乳酸菌学会誌](#) | [第11回統合データ解析環境Galaxy](#) (last modified 2017/11/13)
- 書籍 | [日本乳酸菌学会誌](#) | [第12回Galaxy:ヒストリーとワークフロー](#) (last modified 2018/03/23)
- イントロ | 一般 | [ランダムに行を抽出](#) (last modified 2014/07/17)
- イントロ | 一般 | [任意の文字列を行の最初に挿入](#) (last modified 2014/07/17)

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

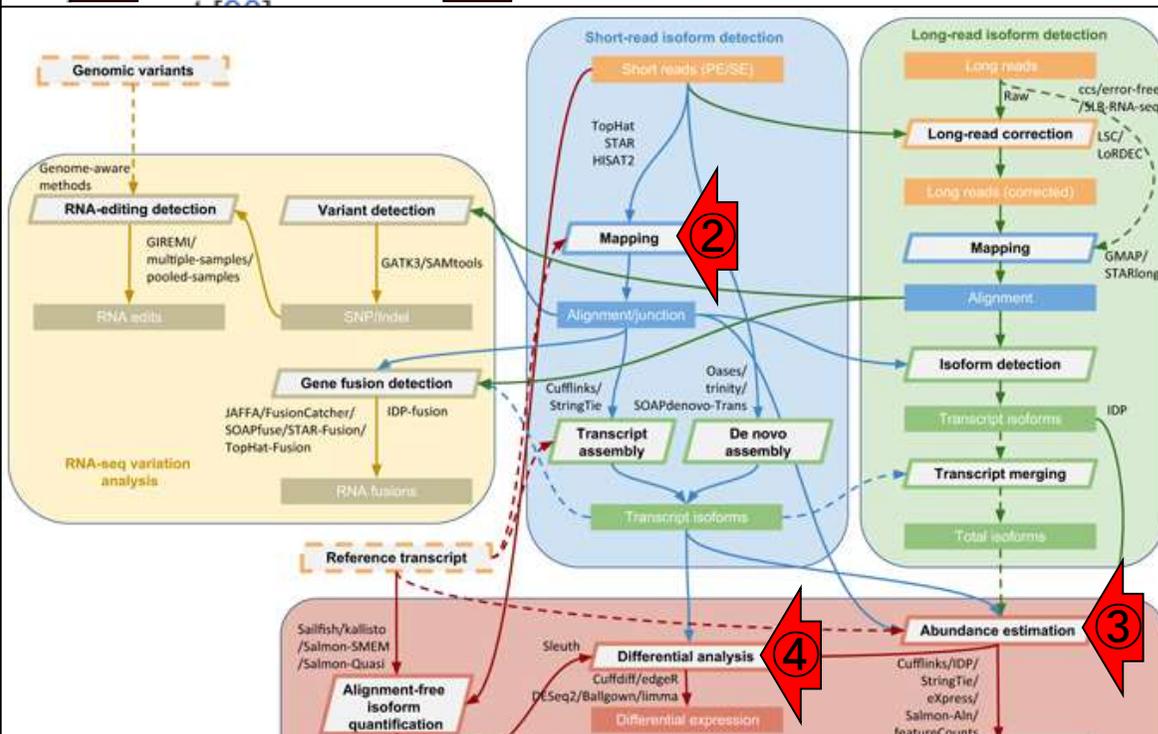
- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - grepでgenenameの個数を確認

全体像のおさらい

RNA-Seq data analysis

最近の総説 (Lowe et al., PLoS Comput. Biol., 13: e1005457, 2017)

processed to yield useful information. Data analysis usually requires a combination of **bioinformatics software** tools that vary according to the experimental design goals. The process can be broken down into the following four stages: quality control, alignment, quantification, and differential expression [89]. Most popular RNA-Seq programs are run from a command-line interface, either in a Unix environment or within the R/Bioconductor statistical



...ence needs to be
 ...ty scores for base
 ... representation of
 ... duplication rate [85].
 ... and FaQCs software
 ... tagged for special

RNACocktail (Sahraeian et al., Nat Commun., 8: 59, 2017)

単一サンプル

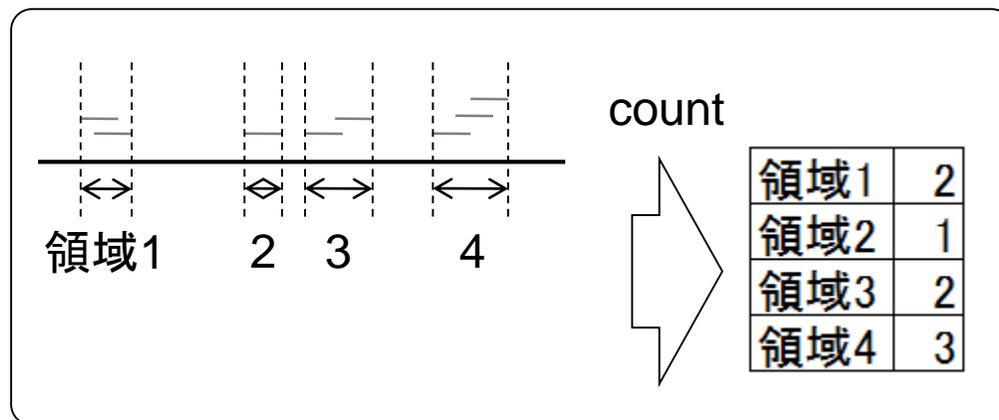
Rでアノテーション情報を利用する場合は、TxDbが基本。アノテーション情報がない場合は、マップされたリードの領域をたよりに転写領域を決める

■ アノテーション情報を利用する場合

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

■ アノテーション情報がない場合

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



複数サンプル

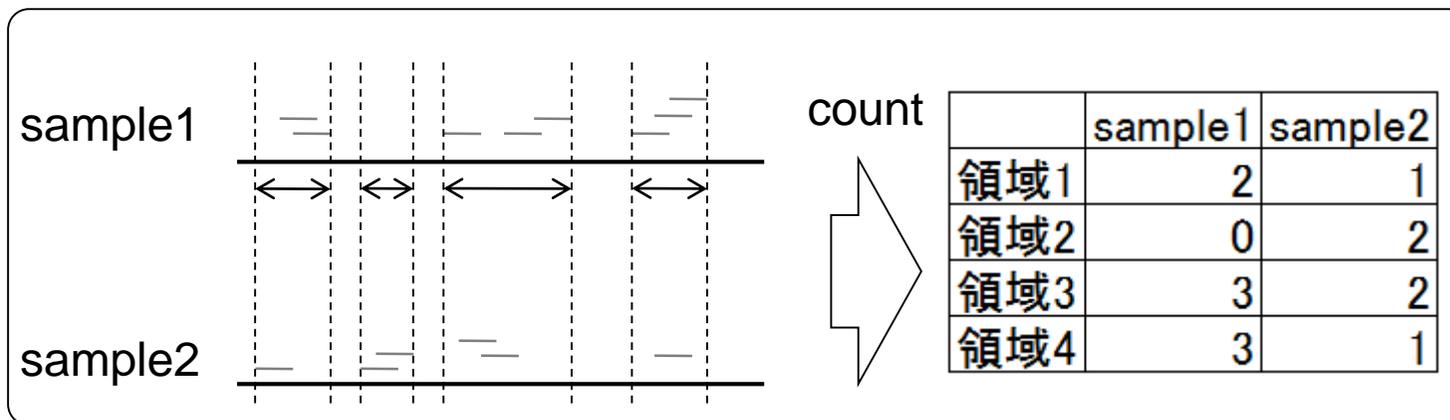
アノテーション情報がない場合の戦略は、複数サンプルの場合には領域が変わりうる。Cufflinks (最近ではStringTie)を知っているヒトはcuffmergeと同じイメージだと思えばよい

■ アノテーション情報を利用する場合

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

■ アノテーション情報がない場合

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



単一サンプル

利用可能なアノテーション情報がなく、単一サンプルで転写領域を定め、その領域にマップされるリード数をカウントするやり方を示します。①single-endのアノテーション無のところ

(Rで)塩基配列解析

(last modified 2018/05/01, since 2010)

このウェブページのR関連部分は、[インストール](#)についての推奨手順([Windows2018.03.12版](#)と[Macintosh](#))とツール済みであるという前提で記述しています。初心者の方は[基本的な利用法](#)([Windows2015.04.03版](#))と目的にまとめた[書籍](#)もあります。(2015/04/03)

What's new?

- [Silhouetteスコアの新たな使い道提唱論文](#)(2018/05/01)
- [Silhouetteスコアの新たな使い道提唱論文](#)(2018/05/01)
- 「平成29年度NGSハンズオン講習会」の[動画](#)

- [門田からメール返信をもらえない場合は](#) (last modified 2018/05/01)
- [はじめに](#) (last modified 2015/03/31)
- 参考資料 | [書籍、学会誌](#) (last modified 2015/07/04)
- 参考資料 | [講習会 講義 講演資料](#) (last modified 2015/07/04)

- [マッピング | 基礎](#) (last modified 2013/06/19)
- [マッピング | single-end | ゲノム | basic aligner\(基礎\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- [マッピング | single-end | ゲノム | basic aligner\(応用\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/28)
- [マッピング | single-end | ゲノム | splice-aware aligner](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- [マッピング | paired-end | ゲノム | basic aligner\(基礎\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- [マッピング | paired-end | ゲノム | basic aligner\(応用\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- [マッピング | paired-end | トランスクリプトーム | basic aligner\(基礎\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/10)
- [マッピング | paired-end | トランスクリプトーム | basic aligner\(応用\)](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/10)
- [マップ後 | について](#) (last modified 2013/06/19)
- [マップ後 | 出力ファイル形式について](#) (last modified 2013/11/05)
- [マップ後 | 出力ファイルの読み込み | BAM形式 | について](#) (last modified 2016/09/14)
- [マップ後 | 出力ファイルの読み込み | BAM形式 | rbamtools\(Kaisers 2015\)](#) (last modified 2016/09/14)
- [マップ後 | 出力ファイルの読み込み | BAM形式 | GenomicAlignments\(Lawrence 2013\)](#) (last modified 2016/09/14)
- [マップ後 | 出力ファイルの読み込み | Bowtie形式](#) (last modified 2013/06/18)
- [マップ後 | 出力ファイルの読み込み | SOAP形式](#) (last modified 2013/06/19)
- [マップ後 | 出力ファイルの読み込み | htSeqTools\(Planet 2012\)](#) (last modified 2013/06/19)
- [マップ後 | カウント情報取得 | について](#) (last modified 2017/01/11)
- [マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/02/26)
- [マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/22)
- [マップ後 | カウント情報取得 | paired-end | ゲノム | アノテーション有](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/13)
- [マップ後 | カウント情報取得 | paired-end | ゲノム | アノテーション無](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/07/02)
- [マップ後 | カウント情報取得 | paired-end | トランスクリプトーム](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/12)
- [マップ後 | カウント情報取得 | トランスクリプトーム | BEDファイルから](#) (last modified 2014/06/21)
- [マップ後 | 配列長とカウント数の関係](#) (last modified 2015/07/03)
- [正規化 | について](#) (last modified 2014/06/22)
- [正規化 | 基礎 | RPK or CPK \(配列長補正\)](#) (last modified 2015/07/04)

単一サンプル

①例題1をやります。②作業ディレクトリは「Desktop - hoge - mapping_kiso1」。③最低限この3つのファイルがあれば動きます

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無 | QuasR(Gaidatzis_2015)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

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

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

```
in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "ref_genome.fa" #入力ファイル
param_mapping <- "-m 1 --best --strata -v 0"#マッピング
```

#必要なパッケージをロード

```
library(QuasR) #パッケージの
library(GenomicAlignments) #パッケージの
```

#前処理(マッピング)

```
time_s <- proc.time() #計算時間を計
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #計算時間を計
time_e <- proc.time() #計算時間を計
time_e - time_s #計算時間を表
out #マッピングに
alignmentStats(out) #マッピング結
```

#本番(マップされたリードの和集合領域同定)

```
tmpfname <- out@alignments[,1] #ファイル名(
tmpsname <- out@alignments[,2] #サンプル名(
for(i in 1:length(tmpfname)){ #サンプル数(
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式ファ
```

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

①

②

③

単一サンプル

マップ後 | カウント情報取得 | single-end | ゲノム

①以前のマッピング結果が残っていても問題ありません。②このパラメータ(マッピングオプション)は前回指定したものと同じです。③マッピングを行う関数もコード内に記述されていますが、QuasRは③マッピングを実行する前に、④作業ディレクトリ内の⑤実行ログファイル(QuasR_log...)をまず見ます

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノムへのマッピングのフローを示します。アンテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

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

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

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

#必要なパッケージをロード
library(QuasR) #パッケージのインストール
library(GenomicAlignments) #パッケージのインストール

#前処理(マッピング)
time_s <- proc.time() #計算時間を計測
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッピングを実行
time_e <- proc.time() #計算時間を計測
time_e - time_s #計算時間を表示
out #マッピング結果
alignmentStats(out) #マッピング結果

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

```
R Console
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_146c2f3d4e64.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_146c6c6d54aa.bam"
[9] "sample_RNAseq1_146c6c6d54aa.bam.bai"
[10] "sample_RNAseq1_146c6c6d54aa.bam.txt"
[11] "sample_RNAseq1_146c6c6d54aa.bed"
[12] "sample_RNAseq1_146c6c6d54aa_QC.pdf"
> |
```

単一サンプル

マップ後 | カウント情報取得 | single-end | ゲノム |

そして、以前マッピングを行った結果が今回のものと同じ入力ファイルで②同じパラメータあれば、今回改めてマッピングを行うことはせずに以前の結果を流用してくれるというありがたい機能があります。マッピングに数時間以上かかるような場合は、便利さを実感します

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列の一連の流れを示します。アンテーション情報がない場合を想定しているため、[GenomicAlignments](#)パッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

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

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

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

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

```
#前処理(マッピング)
time_s <- proc.time()
```

```
R Console
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_146c2f3d4e64.txt"
[3] "ref_genome.fa"
```

```
C:\Users\kojik\Desktop\hoge\mapping_kiso1\QuasR_log_146c2f3d4e64.txt - 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)
QuasR_log_146c2f3d4e64.txt x
[1] "Executing bowtie on DESKTOP-3J8LKP8 using 1 cores. Parameters:"
[1] "%C:/Users/kojik/Desktop/hoge/mapping_kiso1/ref_genome.fa.Rbowtie/bowtieIndex%" %C:/Users/kojik/Desktop/hoge/mapping_kiso1/sample_RNAseq1
[1] "Converting sam file to sorted bam file on DESKTOP-3J8LKP8 : C:\Users\kojik\AppData\Local\Temp\Rtmpw5nBDI\sample_RNAseq1.fa2dd079d64
[1] "Genomic alignments for sample 1 (naeae) have been successfully created on DESKTOP-3J8LKP8"
598 バイト (598 バイト), 5 行。
k <- readGAlignments(tmpfname[i]) #BAM形式ノア
[12] "sample_RNAseq1_146c6c6d54aa_QC.pdf"
> |
```

コピー実行

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無 | QuasR(Gaidatzis_2015)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。
「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

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

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

```
in_f1 <- "mapping_single_genome1.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "ref_genome.fa" #入力ファイル
param_mapping <- "-m 1 --best --strata -v 0"#マッピング
```

#必要なパッケージをロード

```
library(QuasR) #パッケージの
library(GenomicAlignments) #パッケージの
```

#前処理(マッピング)

```
time_s <- proc.time() #計算時間を計
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #計算時間を計
time_e <- proc.time() #計算時間を計
time_e - time_s #計算時間を表
out #マッピングに
alignmentStats(out) #マッピング結
```

#本番(マップされたリードの和集合領域同定)

```
tmpfname <- out@alignments[,1] #ファイル名(
tmpsname <- out@alignments[,2] #サンプル名(
for(i in 1:length(tmpfname)){ #サンプル数(
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式ファ
```

```
R Console
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_146c2f3d4e64.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_146c6c6d54aa.bam"
[9] "sample_RNAseq1_146c6c6d54aa.bam.bai"
[10] "sample_RNAseq1_146c6c6d54aa.bam.txt"
[11] "sample_RNAseq1_146c6c6d54aa.bed"
[12] "sample_RNAseq1_146c6c6d54aa_QC.pdf"
> |
```



コピペ実行後

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無 | QuasR(Gaidatzis_2015)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

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

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

```

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

#必要なパッケージをロード
library(QuasR) #パッケージのインストール
library(GenomicAlignments) #パッケージのインストール

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

#本番(マップされたリードの和集合領域同定)
tmpfname <- out@alignments[,1] #ファイル名(出力)
tmpsname <- out@alignments[,2] #サンプル名(出力)
for(i in 1:length(tmpfname)){ #サンプル数(出力)
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式ファイル読み込み
  }
}
m <- reduce(granges(k)) #GRangesオブジェクト
#本番(カウント情報取得)
tmp <- as.data.frame(m) #出力ファイル名
for(i in 1:length(tmpfname)){ #サンプル数
  tmpcount <- summarizeOverlaps(m, tmpfname[i]) #SummarizedExperimentオブジェクト
  count <- assays(tmpcount)$counts #行列オブジェクト
  colnames(count) <- tmpsname[i] #行列名
  tmp <- cbind(tmp, count) #保存したデータフレーム
}
#ファイルに保存
out_f <- sub(".bam", "_range.txt", tmpfname[i]) #出力ファイル名
write.table(tmp, out_f, sep="\t", append=F, quote=F)

```

*_range.txt

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無 | QuasR(Gaidatzis_2015)

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

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ18,190のRNA-seqデータ(sample_RNAseq1.fa)のref

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

```
in_f1 <- "mapping_single_genome1.txt" #入力ファイル
in_f2 <- "ref_genome.fa" #入力ファイル
param_mapping <- "-m 1 --best --strata -v 0"#マッピング

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

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

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

R Console

```
> out_f <- sub(".bam", "_range.txt", tmpfname[i])#$
> write.table(tmp, out_f, sep="\t", append=F, quot$
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_146c2f3d4e64.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_146c6c6d54aa.bam"
[9] "sample_RNAseq1_146c6c6d54aa.bam.bai"
[10] "sample_RNAseq1_146c6c6d54aa.bam.txt"
[11] "sample_RNAseq1_146c6c6d54aa.bed"
[12] "sample_RNAseq1_146c6c6d54aa_QC.pdf"
[13] "sample_RNAseq1_146c6c6d54aa_range.txt"
> |
```



* _range.txt

マップ後 | カウント情報取得 | single-end |

コピーしたコードの下部に移動。出力ファイルは何も指定していませんが、①*_range.txtという名前のファイルが作成されます。これは、②.bamという名前のファイルを内部的に入力として読み込み、その文字列中の.bamを_range.txtに置換したものを出力ファイル名として自動作成しているからそうなります

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」→「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ18,19のRNA-seqデータ(sample_RNAseq1.fa)のref

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

```
for(i in 1:length(tmpfname)){ #サンプル数(
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式ファ
  } else{
    k <- c(k, readGAlignments(tmpfname[i]))#BAM形式
  }
}
m <- reduce(granges(k)) #GRangesオブ

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

#ファイルに保存

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

```
R Console
> out_f <- sub(".bam", "_range.txt", tmpfname[i])#$
> write.table(tmp, out_f, sep="\t", append=F, quot$
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso1"
> list.files()
[1] "mapping_single_genome1.txt"
[2] "QuasR_log_146c2f3d4e64.txt"
[3] "ref_genome.fa"
[4] "ref_genome.fa.fai"
[5] "ref_genome.fa.md5"
[6] "ref_genome.fa.Rbowtie"
[7] "sample_RNAseq1.fa"
[8] "sample_RNAseq1_146c6c6d54aa.bam"
[9] "sample_RNAseq1_146c6c6d54aa.bam.bam"
[10] "sample_RNAseq1_146c6c6d54aa.bam.txt"
[11] "sample_RNAseq1_146c6c6d54aa.bed"
[12] "sample_RNAseq1_146c6c6d54aa_QC.pdf"
[13] "sample_RNAseq1_146c6c6d54aa_range.txt"
> |
```



*_range.txt

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

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無 | QuasR(Gaidatzis_2015)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ18,19のRNA-seqデータ(sample RNAseq1.fa)のref_genome.faへのマッピングの場合(mapping_single_genome1.txt) オプションを"-m 1 --best --strata -v 0"とした例です。

*.bed

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

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

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

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

*_range.txt

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

	A	B
1	FileName	SampleName
2	sample RNAseq1.fa	nameae

を計測するための表示(一番右側の列に用いたパラメータの結果(alignment

```
#本番(マッピング)の和集合領域取得
tmpfname <- out@alignments[,1] #ファイル名(in_f1の1列目に相当)をtmpfnameとして取り扱いたいだけです
tmpsname <- out@alignments[,2] #サンプル名(in_f1の2列目に相当)をtmpsnameとして取り扱いたいだけです
for(i in 1:length(tmpfname)){ #サンプル数(ファイル数)分だけループを回す
  if(i == 1){
    k <- readGAlignments(tmpfname[i]) #BAM形式ファイルを読み込んだ結果をkに格納(これはGAlignmentsオブジェクト)
```

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - grepでgenenameの個数を確認

複数サンプル

①例題5をやってみましょう。「デスクトップ - hoge - mapping_kiso2」フォルダを作成し、必要な入力ファイルを揃えてコピペ実行

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション無 | QuasR(Gaidatzis_2015)

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報がない場合を想定しているため、GenomicAlignmentsパッケージを利用して、マップされたリードの和集合領域(union range)を得たのち、領域ごとにマップされたリード数をカウントしています。「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

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

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

```
in_f1 <- "mapping_single_genome1.txt"
in_f2 <- "ref_genome.fa"
param_mapping <- "-m 1 --best --strata
```

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

#前処理(マッピング)

```
time_s <- proc.time()
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping)
time_e <- proc.time()
time_e - time_s
out
alignmentStats(out)
```

#本番(マップされたリードの和集合領域同定)

```
tmpfname <- out@alignments[,1]
tmpsname <- out@alignments[,2]
for(i in 1:length(tmpfname)){
  if(i == 1){
    k <- readGAlignments(tmpfname[i])
```

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

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

```
in_f1 <- "mapping_single_genome4.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqフ
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納(リファレン
out_f1 <- "hoge5_count.txt" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge5_genelength.txt" #出力ファイル名を指定してout_f2に格納
param_mapping <- "-m 1 --best --strata -v 1"#マッピング時のオプションを指定
```

#必要なパッケージをロード

```
library(QuasR) #パッケージの読み込み
library(GenomicAlignments) #パッケージの読み込み
```

#前処理(マッピング)

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

#本番(マップされたリードの和集合領域同定)

```
tmpfname <- out@alignments[,1] #ファイル名(in_f1の1列目に相当)をtmpfnameとして
tmpsname <- out@alignments[,2] #サンプル名(in_f1の2列目に相当)をtmpsnameとして
for(i in 1:length(tmpfname)){ #サンプル数(ファイル数)分だけループを回す
```

複数サンプル

①例題5をやってみましょう。「デスクトップ - hoge - mapping_kiso2」フォルダを作成し、必要な入力ファイルを揃えてコピー実行

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

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

```
in_f1 <- "mapping_single_genome4.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqフ
in_f2 <- "ref_genome.fa" #入力ファイル名を指定してin_f2に格納(リファレン
out_f1 <- "hoge5_count.txt" #出力ファイル名を指定してout_f1に格納
out_f2 <- "hoge5_genelength.txt" #出力ファイル名を指定してout_f2に格納
param_mapping <- "-m 1 --best --strata -v 1" #マッピング
```

#必要なパッケージをロード

```
library(QuasR) #パッケージの読
library(GenomicAlignments) #パッケージの読
```

#前処理(マッピング)

```
time_s <- proc.time() #計算時間を計測
out <- qAlign(in_f1, in_f2, alignmentParameter=param_ #計算時間を計測
time_e <- proc.time() #計算時間を計測
time_e - time_s #計算時間を表示
out #マッピングに用
alignmentStats(out) #マッピング結果
```

#本番(マップされたリードの和集合領域同定)

```
tmpfname <- out@alignments[,1] #ファイル名(in_
tmpsname <- out@alignments[,2] #サンプル名(in_
for(i in 1:length(tmpfname)){ #サンプル数(フ
```

```
R Console
詳しくは 'contributors()' と入力してください。
また、R や R のパッケージを出版物で引用する際の$
'citation()' と入力してください。

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

> getwd()
[1] "C:/Users/kadota/Desktop/hoge/mapping_kiso2"
> list.files()
[1] "mapping_single_genome4.txt"
[2] "ref_genome.fa"
[3] "sample_RNAseq1.fa"
[4] "sample_RNAseq2.fa"
> |
```

複数サンプル

①2つの出力ファイルのうち、主に取り扱うのはカウントデータを含むファイル(hoge5_count.txt)のほうです

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

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

```
in_f1 <- "mapping_single_genome4.txt" #入力ファイル名を
in_f2 <- "ref_genome.fa" #入力ファイル名を
out_f1 <- "hoge5_count.txt" #出力ファイル名を
out_f2 <- "hoge5_genelength.txt" #出力ファイル名を
param_mapping <- "-m 1 --best --stats --data -v 1" #マッピング

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

#前処理(マッピング)
time_s <- proc.time() #計算時間を計測する
out <- qAlign(in_f1, in_f2, alignmentParameter=param_mapping) #マッピングを実行
time_e <- proc.time() #計算時間を計測する
time_e - time_s #計算時間を表示する
out #マッピングに用いた出力
alignmentStats(out) #マッピング結果

#本番(マップされたリードの和集合領域同定)
tmpfname <- out@alignments[,1] #ファイル名(in_f1)
tmpsname <- out@alignments[,2] #サンプル名(in_f2)
for(i in 1:length(tmpfname)){ #サンプル数(ファイル名)をループして
```

```
R Console
> tmp <- cbind(tmp, h$width) #和集$
> write.table(tmp, out_f2, sep="\t", append=F, $
> getwd())
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso2"
> list.files()
[1] "hoge5_count.txt"
[2] "hoge5_genelength.txt"
[3] "mapping_single_genome4.txt"
[4] "QuasR_log_17ecaca51c5.txt"
[5] "ref_genome.fa"
[6] "ref_genome.fa.fai"
[7] "ref_genome.fa.md5"
[8] "ref_genome.fa.Rbowtie"
[9] "sample_RNAseq1.fa"
[10] "sample_RNAseq1_17ec48601f6d.bam"
[11] "sample_RNAseq1_17ec48601f6d.bam.bai"
[12] "sample_RNAseq1_17ec48601f6d.bam.txt"
[13] "sample_RNAseq2.fa"
[14] "sample_RNAseq2_17ec1390671e.bam"
[15] "sample_RNAseq2_17ec1390671e.bam.bai"
[16] "sample_RNAseq2_17ec1390671e.bam.txt"
> |
```

複数サンプル

①リストファイル中で指定した②サンプル名が、③のカウントデータ行列の④列名となります

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

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

```

in_f1 <- "mapping_single_genome4.txt"
in_f2 <- "ref_genome.fa"
out_f1 <- "hoge5_count.txt"
out_f2 <- "hoge5_genelength.txt"
param_mapping <- "-m 1 --best --strata -v 1"

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

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

```

① カファイル名を指定してin_f1に格納(RNA-seqフ
 ② カファイル名を指定してin_f2に格納(リファレン
 ③ #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
 #マッピング時のオプションを指定

#計算時間を計測するため
 #マッピングを行うqAlign
 #計算時間を計測するため
 #計算時間を表示(一番右側の数字
 #マッピングに用いたパラメー
 #マッピング結果をAlignment

FileName	SampleName
sample_RNAseq1.fa	sample1
sample_RNAseq2.fa	sample2

②
 ①列目に
 ②列目に
 (行数)分た

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

④

Contents

- マッピング (アラインメント) の続き
 - おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
 - マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
 - SAM形式の解説、マッピング結果の違い、課題
 - Linux環境以外でのBowtie2実行手段
- カウント情報取得
 - アノテーション情報がない場合: 単一サンプル、複数サンプル
 - アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`で`genename`の個数を確認

アノテーション情報あり

乳酸菌のアノテーション情報ファイル(.gff3)で定義されたgene領域にマップされるリード数をカウントするやり方を示します。①
single-endのアノテーション有のところ

(Rで)塩基配列解析

(last modified 2018/05/01, since 2010)

このウェブページのR関連部分は、[インストール](#)についての推奨手順([Windows2018.03.12版](#)と[Macintosh](#))済みであるという前提で記述しています。初心者の方は[基本的な利用法](#)([Windows2015.04.03版](#))的にまとめた[書籍](#)もあります。(2015/04/03)

What's new?

- [Silhouetteスコアの新たな使い道提唱論文](#)(2018/05/01)
- [Silhouetteスコアの新たな使い道提唱論文](#)(2018/05/01)
- 「平成29年度NGSハンズオン講習会」の[動画](#)

- [門田からメール返信をもらえない場合は](#) (last modified 2018/05/01)
- [はじめに](#) (last modified 2015/03/31)
- 参考資料 | [書籍](#)、[学会誌](#) (last modified 2018/05/01)
- 参考資料 | [講習会](#)、[講義](#)、[講演資料](#) (last modified 2018/05/01)

- [マッピング](#) | [基礎](#) (last modified 2013/06/19)
- [マッピング](#) | [single-end](#) | [ゲノム](#) | [basic aligner](#)(基礎) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- [マッピング](#) | [single-end](#) | [ゲノム](#) | [basic aligner](#)(応用) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/06/28)
- [マッピング](#) | [single-end](#) | [ゲノム](#) | [splice-aware aligner](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/21)
- [マッピング](#) | [paired-end](#) | [ゲノム](#) | [basic aligner](#)(基礎) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- [マッピング](#) | [paired-end](#) | [ゲノム](#) | [basic aligner](#)(応用) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/11)
- [マッピング](#) | [paired-end](#) | [トランスクリプトーム](#) | [basic aligner](#)(基礎) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/10)
- [マッピング](#) | [paired-end](#) | [トランスクリプトーム](#) | [basic aligner](#)(応用) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/10)
- [マップ後](#) | について (last modified 2013/06/19)
- [マップ後](#) | [出力ファイル形式](#)について (last modified 2013/11/05)
- [マップ後](#) | [出力ファイルの読み込み](#) | [BAM形式](#) | について (last modified 2016/09/14)
- [マップ後](#) | [出力ファイルの読み込み](#) | [BAM形式](#) | [rbamtools\(Kaisers 2015\)](#) (last modified 2016/09/14)
- [マップ後](#) | [出力ファイルの読み込み](#) | [BAM形式](#) | [GenomicAlignments\(Lawrence 2013\)](#) (last modified 2016/09/14)
- [マップ後](#) | [出力ファイルの読み込み](#) | [Bowtie形式](#) (last modified 2013/06/18)
- [マップ後](#) | [出力ファイルの読み込み](#) | [SOAP形式](#) (last modified 2013/06/19)
- [マップ後](#) | [出力ファイルの読み込み](#) | [htSeqTools\(Planet 2012\)](#) (last modified 2013/06/19)
- [マップ後](#) | [カウント情報取得](#) | について (last modified 2017/01/11)
- [マップ後](#) | [カウント情報取得](#) | [single-end](#) | [ゲノム](#) | [アノテーション有](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/02/26)
- [マップ後](#) | [カウント情報取得](#) | [single-end](#) | [ゲノム](#) | [アノテーション無](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2014/06/22)
- [マップ後](#) | [カウント情報取得](#) | [paired-end](#) | [ゲノム](#) | [アノテーション有](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/13)
- [マップ後](#) | [カウント情報取得](#) | [paired-end](#) | [ゲノム](#) | [アノテーション無](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2015/07/02)
- [マップ後](#) | [カウント情報取得](#) | [paired-end](#) | [トランスクリプトーム](#) | [QuasR\(Gaidatzis 2015\)](#) (last modified 2016/02/12)
- [マップ後](#) | [カウント情報取得](#) | [トランスクリプトーム](#) | [BEDファイルから](#) (last modified 2014/06/21)
- [マップ後](#) | [配列長とカウント数の関係](#) (last modified 2015/07/03)
- [正規化](#) | について (last modified 2014/06/22)
- [正規化](#) | [基礎](#) | [RPK or CPK](#) ([配列長補正](#)) (last modified 2015/07/04)

アノテーション情報あり

①例題10。マップされる側のリファレンス配列は、②乳酸菌ゲノム配列。③アノテーションファイル。④機能ゲノム学の第1回(5/8)で、転写物配列取得時に使ったものと同じです。

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis 2015)などを参考にしてください。

「ファイル」-「ディレクトリ」の変更で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ 10. mapping_single_genome7.txt 中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

mapping_single_genome7.txtのサンプル名(例: h1)です。hg19にマップした結果で、Entrez Gene IDに

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome_hg19"
out_f <- "hoge10.txt"
param_mapping <- "h"
param_txdb1 <- "hg19"
param_txdb2 <- "k"
param_reportlevel <- "gene"
```

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

#前処理(マッピング)
out <- qAlign(in_f1, in_f2, alignmentStats=out)

#前処理(TxDbオブジェクトの作成)
txdb <- makeTxDbFromGFF(in_f3, format="auto")
txdb

#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel)
dim(count)
head(count)

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample RNAseq4.fa)です。マップされる側のファイルは、Ensembl(Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名を指定してin_f2に格納
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #入力ファイル名を指定してin_f3に格納
out_f <- "hoge10.txt" #出力ファイル名を指定してout_fに格納
param_reportlevel <- "gene" #カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"
```

#必要なパッケージをロード

```
library(QuasR)
library(GenomicFeatures)
```

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="auto") #txdbオブジェクトの作成
txdb #確認してるだけです
```

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2) #マッピングを行うqAlign関数を実行した結果
alignmentStats(out) #マッピング結果(alignment statistics)
```

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行別列別
dim(count) #行数と列数を表示
head(count) #確認してるだけです
```

講義日程 (平成30年度)

1. 平成30年05月08日

講義資料PDF

.gff3ファイル(約1.3MB)

.faファイル(約2.2MB)

(Rで)塩基配列解析

(Rで)マイクロアレイデータ解析

plasmid1.gff3(課題用)

plasmid2.gff3(課題用)

おさらい

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . + . ID=gc0829395.1.30.dna.chromosome.Chromosome.fa
Chromosome ena transcript 360 1676 . + . ID=tr0829395.1.30.dna.chromosome.Chromosome.gff3
Chromosome ena exon 360 1676 . + . Parent=gc0829395.1.30.dna.chromosome.Chromosome.fa
Chromosome ena CDS 360 1676 . + 0 ID=CDS0829395.1.30.dna.chromosome.Chromosome.gff3
###
Chromosome ena gene 1852 2991 . + . ID=gc0829395.1.30.dna.chromosome.Chromosome.fa
Chromosome ena transcript 1852 2991 . + . ID=tr0829395.1.30.dna.chromosome.Chromosome.gff3
Chromosome ena exon 1852 2991 . + . Parent=gc0829395.1.30.dna.chromosome.Chromosome.fa
Chromosome ena CDS 1852 2991 . + 0 ID=CDS0829395.1.30.dna.chromosome.Chromosome.gff3
###
Chromosome ena gene 3233 3457 . + . ID=gc0829395.1.30.dna.chromosome.Chromosome.fa
Chromosome ena transcript 3233 3457 . + . ID=tr0829395.1.30.dna.chromosome.Chromosome.gff3
Chromosome ena exon 3233 3457 . + . Parent=gc0829395.1.30.dna.chromosome.Chromosome.fa
Chromosome ena CDS 3233 3457 . + 0 ID=CDS0829395.1.30.dna.chromosome.Chromosome.gff3
###
Chromosome ena gene 3467 4588 . + . ID=gc0829395.1.30.dna.chromosome.Chromosome.fa
```

有 | QuasR(Gaidatzis_2015) NEW

るマッピングから、カウントデータ取得までの一連の
オブジェクトをネットワーク経由で取得するのを基本
の詳細については [マッピング | single-end | ゲノム](#)
をコピペ。

ファイルを乳酸菌ゲノムにマッピングする場合:
FASTA形式ファイル([sample RNAseq4.fa](#))です。マップされる側のファイルは、[Ensembl \(Zerbino Lactobacillus casei 12A\)](#)の multi-FASTA形式ゲノム配列ファイル
[jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3](#))を読み込んでカウン

#入力ファイル名を指定してin_flに格納(RNA-seqファイル)
"jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名
"jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #入力ファイル名
#出力ファイル名を指定してout_flに格納
#カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"

#パッケージの読み込み
#パッケージの読み込み

"auto") #txdbオブジェクトの作成
#確認してるだけです

#マッピングを行うqAlign関数を実行した結果をoutに格納
#マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納
dim(count) #行数と列数を表示
head(count) #確認してるだけです
```

geneレベルのカウントデータ

①例題10は、②GFF3ファイルから、③gene領域内にマップされたリード数をカウントするやり方

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」→「ディレクトリ」の変更で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ70 ① 10. mapping single genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

mapping_single_genome7.txtのサンプル名(例: hg19)にマップした結果で、Entrez Gene IDに対して

```
in_f1 <- "mapping  
in_f2 <- "BSgenom  
out_f <- "hoge1.t  
param_mapping <-  
param_txdb1 <- "h  
param_txdb2 <- "k  
param_reportlevel
```

```
#必要なパッケージを  
library(QuasR)  
library(GenomicFe
```

```
#前処理(マッピング)  
out <- qAlign(in_  
alignmentStats(ou
```

```
#前処理(TxDbオブジ  
txdb <- makeTxDbF  
txdb
```

```
#本番(カウントデー  
count <- qCount(o  
<
```

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample RNAseq4.fa)です。マップされる側のファイルは、Ensembl (Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus hokkaidonensis jcm 18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus hokkaidonensis jcm 18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)  
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名を指定してin_f2に格納  
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #入力ファイル名を指定してin_f3に格納  
out_f <- "hoge10.txt" #出力ファイル名を指定してout_fに格納  
param_reportlevel <- "gene" #カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"
```

#必要なパッケージをロード

```
library(QuasR)  
library(GenomicFeatures)
```

```
#パッケージの読み込み  
#パッケージの読み込み
```

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="auto") #txdbオブジェクトの作成  
txdb #確認してるだけです
```

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2) #マッピングを行うqAlign関数を実行した結果をoutに格納  
alignmentStats(out) #マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列
```

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納  
dim(count) #行数と列数を表示  
head(count) #確認してるだけです
```

geneレベルのカウント

geneレベルとはいっても、バクテリアなので例えば exon を指定しても同じ結果になります。gene を指定したときは、①のようなgeneと書かれた行の赤枠領域にマップされたリード数を調べることになります。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
```

Chromosome	ena	gene	360	1676	+	ID=gc
Chromosome	ena	transcript	360	1676	+	ID=tr
Chromosome	ena	exon	360	1676	+	Parent
Chromosome	ena	CDS	360	1676	+	ID=C
###						
Chromosome	ena	gene	1852	2991	+	ID=gc
Chromosome	ena	transcript	1852	2991	+	ID=tr
Chromosome	ena	exon	1852	2991	+	Parent
Chromosome	ena	CDS	1852	2991	+	ID=C
###						
Chromosome	ena	gene	3233	3457	+	ID=gc
Chromosome	ena	transcript	3233	3457	+	ID=tr
Chromosome	ena	exon	3233	3457	+	Parent
Chromosome	ena	CDS	3233	3457	+	ID=C
###						
Chromosome	ena	gene	3467	4588	+	ID=gc

あるマッピングから、カウントデータ取得までの一連のオブジェクトをネットワーク経由で取得するのを基本的に詳しくは [マッピング | single-end | ゲノム](#) をコピペ。

ファイルを乳酸菌ゲノムにマッピングする場合:

FASTA形式ファイル([sample RNAseq4.fa](#))です。マップされる側のファイルは、[Ensembl \(Zerbino Lactobacillus casei 12A\)](#)の multi-FASTA形式ゲノム配列ファイル ([0829395.1.30.dna.chromosome.Chromosome.fa](#))です。マッピング結果に対して、GFF3形式の ([jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3](#))を読み込んでカウン

#入力ファイル名を指定してin_f1に格納(RNA-seqファイル)

```
"jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名
"jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #入力ファイル名
```

#出力ファイル名を指定してout_flに格納

#カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"

#パッケージの読み込み

#パッケージの読み込み

"auto")#txdbオブジェクトの作成

#確認してるだけです

#マッピングを行うqAlign関数を実行した結果をoutに格納

#マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#カウントデータ行列を取得してcountに格納
```

```
dim(count)
```

#行数と列数を表示

```
head(count)
```

#確認してるだけです

Contents

- マッピング (アラインメント) の続き
 - おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
 - マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
 - SAM形式の解説、マッピング結果の違い、課題
 - Linux環境以外でのBowtie2実行手段
- カウント情報取得
 - アノテーション情報がない場合: 単一サンプル、複数サンプル
 - アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - grepでgenenameの個数を確認

マップする側のファイル

マップする側のファイルは、①のリストファイル内に記載されている、②sample_RNAseq4.fa。③が中身

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータのFASTA

mapping_single_genome7.txtのサンプル名(例: human)を、hg19にマップした結果で、Entrez Gene IDに

```
in_f1 <- "mapping
in_f2 <- "BSgenom
out_f <- "hoge1.t
param_mapping <-
param_txdb1 <- "h
param_txdb2 <- "k
param_reportlevel
```

#必要なパッケージを
library(QuasR)
library(GenomicFe

#前処理(マッピング)
out <- qAlign
alignmentStats

#前処理(TxDb)
txdb <- make
txdb

#本番(カウントデータ)
count <- qCount(o

<

10. mapping_single_genome7.txt中のFASTA形式ファイルを乳酸菌ゲノム

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample RNA et al., Nucleic Acids Res., 2018)から提供されている Lactobacillus casei 12Aの (Lactobacillus hokkaidonensis jcm 18461.GCA_000829395.1.30.dna.chromosc アノテーションファイル(Lactobacillus hokkaidonensis jcm 18461.GCA_000829 ト情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829
out_f <- "hoge10.txt"
param_reportlevel <- "gene"
```

#必要なパッケージをロード

library(QuasR)

library(GenomicFe

FileName	SampleName
sample_RNAseq4.fa	Lacto

#前処理(マッピング)
out <- qAlign(in_f1, in_f2)
alignmentStats(out)

#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel)#
dim(count)
head(count)

① カファイル名を指定
② #出力ファイル名を指定
③ #カウントデータ取得時

#パッケージの読み込み

#パッケージの読み込み

o")#txdbオブジェク
認してるだけです

#マッピングを行うqAli
#マッピング結果(alig

```
>Chromosome_361_400
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC
>Chromosome_1637_1676
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG
>Chromosome_1851_1890
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA
>Chromosome_1843_1882
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA
>Chromosome_1833_1872
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT
>Chromosome_1823_1862
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC
>Chromosome_1813_1852
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA
>Chromosome_3418_3457
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG
>Chromosome_3420_3459
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC
>Chromosome_3422_3461
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA
>Chromosome_3443_3482
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAGAAC
```

マップする側のファイル

①sample_RNAseq4.faは、②サンプルデータの例題47のコピペで作成しています。わざわざ見に行かなくてもよい

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

1. サンプルデータのFASTQ

mapping_single_genome7.txt
のサンプル名(例: hum)を
す。hg19にマップした結果
で、Entrez Gene IDに対

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome_homo_sapiens.GRCh38.dna.chromosome.1.fa"
out_f <- "hoge1.txt"
param_mapping <- list()
param_txdb1 <- "homo_sapiens.GRCh38.dna.chromosome.1.txdb"
param_txdb2 <- "k12"
param_reportlevel <- "gene"
```

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

#前処理(マッピング)
out <- qAlign(in_f1, in_f2, out_f, param_mapping, param_txdb1, param_txdb2, param_reportlevel)

#前処理(TxDbオブジェクトの作成)
txdb <- makeTxDbFromGFF(in_f3, format="auto")
txdb

#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel)
dim(count)
head(count)

#行数と列数を表示
#確認してるだけです

#確認してるだけです

10. mapping_single_genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample_RNAseq4.fa)です。マップされる側のファイルは、Ensembl (Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名を指定してin_f2に格納
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3" #入力ファイル名を指定してin_f3に格納
out_f <- "hoge10.txt" #出力ファイル名を指定してout_fに格納
param_reportlevel <- "gene" #カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"
```

#必要なパッケージをロード

```
library(QuasR)
library(GenomicFeatures)
```

#パッケージの読み込み
#パッケージの読み込み

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="auto") #txdbオブジェクトの作成
txdb
```

#確認してるだけです

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2, out_f, param_mapping, param_txdb1, param_txdb2, param_reportlevel)
```

#マッピングを行うqAlign関数を実行した結果をoutに格納
#マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納
dim(count) #行数と列数を表示
head(count) #確認してるだけです
```

マップする側のファイル

①マップする側(sample_RNAseq4.fa)のリードは、②アノテーションファイル(.gff3)中のgene領域を参考にしながら作成しています。

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ7のFASTA

mapping_single_genome7.txtのサンプル名(例: hum)です。hg19にマップした結果で、Entrez Gene IDに

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome_hg19"
out_f <- "hoge1.txt"
param_mapping <- "h"
param_txdb1 <- "hg19"
param_txdb2 <- "k"
param_reportlevel <- "gene"
```

```
library(QuasR)
library(GenomicFeatures)
```

```
out <- qAlign(in_f1, in_f2, txdb)
alignmentStats(out)
```

```
txdb <- makeTxDbFromGFF(in_f3, format="auto")
txdb
```

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)
dim(count)
head(count)
```

10. mapping_single_genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample_RNAseq4.fa)です。マップされる側のファイルは、Ensembl (Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名を指定してin_f2に格納
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #アノテーションファイル名を指定してin_f3に格納
out_f <- "hoge10.txt" #出力ファイル名を指定してout_fに格納
param_reportlevel <- "gene" #カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"
```

#必要なパッケージをロード

```
library(QuasR)
library(GenomicFeatures)
```

```
#パッケージの読み込み
#パッケージの読み込み
```

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="auto") #txdbオブジェクトの作成
txdb #確認してるだけです
```

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2)
alignmentStats(out)
```

```
#マッピングを行うqAlign関数を実行した結果をoutに格納
#マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列
```

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納
dim(count) #行数と列数を表示
head(count) #確認してるだけです
```

マップする側のファイル

①マップする側(sample_RNAseq4.fa)のリードは、②アノテーションファイル(.gff3)中の③gene領域を参考にしながら作成しています。全て完全一致でマップされるように設計。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . +
Chromosome ena transcript 360 1676 . +
Chromosome ena exon 360 1676 . +
Chromosome ena CDS 360 1676 . + 0
###
Chromosome ena gene 1852 2991 . +
Chromosome ena transcript 1852 2991 . +
Chromosome ena exon 1852 2991 . +
Chromosome ena CDS 1852 2991 . + 0
###
Chromosome ena gene 3233 3457 . +
Chromosome ena transcript 3233 3457 . +
Chromosome ena exon 3233 3457 . +
Chromosome ena CDS 3233 3457 . + 0
###
Chromosome ena gene 3467 4588 . +
```

有 | QuasR(Ga
るマッピングから、カウントデータ取得までの一連の
オブジェクトをネットワーク経由で取得するのを基本
の詳細については [マッピング | single-end | ゲノム](#)
をコピペ。

ファイルを乳酸菌ゲノム
STA形式ファイル([sample R](#)
[Lactobacillus casei 12Aの](#)
[0829395.1.30.dna.chromosc](#)
[s_jcm_18461.GCA_000829](#)

#入力ファイル名を指定
jcm_18461.GCA_000829
jcm_18461.GCA_000829
#出力ファイル名を指定
#カウントデータ取得時

#パッケージの読み込み
#パッケージの読み込み

"auto")#txdbオブジェク
#確認してるだけです

#マッピングを行うqAli
#マッピング結果(aligr

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#  
dim(count)  
head(count)  
<
```

#行数と列数を表示
#確認してるだけです

```
>Chromosome_361_400
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC
>Chromosome_1637_1676
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG
>Chromosome_1851_1890
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA
>Chromosome_1843_1882
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA
>Chromosome_1833_1872
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT
>Chromosome_1823_1862
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC
>Chromosome_1813_1852
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA
>Chromosome_3418_3457
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG
>Chromosome_3420_3459
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC
>Chromosome_3422_3461
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA
>Chromosome_3443_3482
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAAGAAC
```

マップする側のファイル

①最初の2リードは、②の領域内にマップされるように設計。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . + . ID=g...
Chromosome ena transcript 360 1676 . + . ID=tr...
Chromosome ena exon 360 1676 . + . Pare...
Chromosome ena CDS 360 1676 . + 0 ID=C...
###
Chromosome ena gene 1852 2991 . + . ID=g...
Chromosome ena transcript 1852 2991 . + . ID=tr...
Chromosome ena exon 1852 2991 . + . Pare...
Chromosome ena CDS 1852 2991 . + 0 ID=C...
###
Chromosome ena gene 3233 3457 . + . ID=g...
Chromosome ena transcript 3233 3457 . + . ID=tr...
Chromosome ena exon 3233 3457 . + . Pare...
Chromosome ena CDS 3233 3457 . + 0 ID=C...
###
Chromosome ena gene 3467 4588 . + . ID=g...
```

有 | QuasR(Gaidatzis_2015) NEW
 するマッピングから、カウントデータ取得までの一連の
 bオブジェクトをネットワーク経由で取得するのを基本
 の詳細については [マッピング | single-end | ゲノム](#)
 をコピペ。

ファイルを乳酸菌ゲノムにマ...
 STA形式ファイル([sam](#)...
[Lactobacillus casei 12](#)
[s_jcm_18461.GCA_000829](#)

#入力ファイル名を指定
[jcm_18461.GCA_000829](#)
[jcm_18461.GCA_000829](#)
 #出力ファイル名を指定
 #カウントデータ取得時

#パッケージの読み込み
 #パッケージの読み込み

"auto")#txdbオブジェク
 #確認してるだけです

#マッピングを行うqAli
 #マッピング結果(align

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#  
dim(count)  
head(count)  
<
```

#行数と列数を表示
 #確認してるだけです

```
>Chromosome_361_400  
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC  
>Chromosome_1637_1676  
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG  
>Chromosome_1851_1890  
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA  
>Chromosome_1843_1882  
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA  
>Chromosome_1833_1872  
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT  
>Chromosome_1823_1862  
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC  
>Chromosome_1813_1852  
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA  
>Chromosome_3418_3457  
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG  
>Chromosome_3420_3459  
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC  
>Chromosome_3422_3461  
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA  
>Chromosome_3443_3482  
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAAGAAC
```

マップする側のファイル

①3-7番目の5リードは、②の領域に一部がかかるように設計。領域内ではない。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena transcript 360 1676 + ID=transcript0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena exon 360 1676 + Parent=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena CDS 360 1676 + 0 ID=CDS0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
###
Chromosome ena gene 1852 2991 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena transcript 1852 2991 + ID=transcript0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena exon 1852 2991 + Parent=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena CDS 1852 2991 + 0 ID=CDS0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
###
Chromosome ena gene 3233 3457 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena transcript 3233 3457 + ID=transcript0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena exon 3233 3457 + Parent=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena CDS 3233 3457 + 0 ID=CDS0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
###
Chromosome ena gene 3467 4588 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
```

有 | QuasR(Gaidatzis_2015) NEW
 するマッピングから、カウントデータ取得までの一連の
 bオブジェクトをネットワーク経由で取得するのを基本
 の詳細については [マッピング | single-end | ゲノム](#)
 をコピペ。

ファイルを乳酸菌ゲノムにマ
 STA形式ファイル([sample R](#)
[Lactobacillus casei 12A](#)の
[0829395.1.30.dna.chromoso](#)
[s_jcm_18461.GCA_000829](#)

#入力ファイル名を指定
[jcm_18461.GCA_000829](#)
[jcm_18461.GCA_000829](#)
 #出力ファイル名を指定
 #カウントデータ取得時

#パッケージの読み込み
 #パッケージの読み込み
 "auto")#txdbオブジェク
 #確認してるだけです

#マッピングを行うqAli
 #マッピング結果(align

```
>Chromosome_361_400
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC
>Chromosome_1637_1676
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG
>Chromosome_1851_1890
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA
>Chromosome_1843_1882
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA
>Chromosome_1833_1872
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT
>Chromosome_1823_1862
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC
>Chromosome_1813_1852
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA
>Chromosome_3418_3457
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG
>Chromosome_3420_3459
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC
>Chromosome_3422_3461
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA
>Chromosome_3443_3482
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAAGAAC
```

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#
dim(count)
head(count)
<
```

マップする側のファイル

①8番目のリードは、②の領域内にマップされるように設計。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 + ID=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena transcript 360 1676 + ID=transcript_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena exon 360 1676 + Parent=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena CDS 360 1676 + 0 ID=CDS_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
###
Chromosome ena gene 1852 2991 + ID=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena transcript 1852 2991 + ID=transcript_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena exon 1852 2991 + Parent=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena CDS 1852 2991 + 0 ID=CDS_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
###
Chromosome ena gene 3233 3457 + ID=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena transcript 3233 3457 + ID=transcript_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena exon 3233 3457 + Parent=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
Chromosome ena CDS 3233 3457 + 0 ID=CDS_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
###
Chromosome ena gene 3467 4588 + ID=gene_0829395.1.30.dna.chromosome_jcm_18461.GCA_000829395.1
```

有 | QuasR(Gaidatzis_2015) NEW
 するマッピングから、カウントデータ取得までの一連の
 オブジェクトをネットワーク経由で取得するのを基本
 の詳細については [マッピング | single-end | ゲノム](#)
 をコピペ。

ファイルを乳酸菌ゲノムにマ
 STA形式ファイル([sample R](#)
[Lactobacillus casei 12A](#)の
[0829395.1.30.dna.chromoso](#)
[s_jcm_18461.GCA_000829](#)

#入力ファイル名を指定
[jcm_18461.GCA_000829](#)
[jcm_18461.GCA_000829](#)
 #出力ファイル名を指定
 #カウントデータ取得時

#パッケージの読み込み
 #パッケージの読み込み

"auto")#txdbオブジェク
 #確認してるだけで

#マッピングを行うqAli
 #マッピング結果(align

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#  
dim(count)  
head(count)  
<
```

#行数と列数を表示
 #確認してるだけです

```
>Chromosome_361_400  
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC  
>Chromosome_1637_1676  
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG  
>Chromosome_1851_1890  
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA  
>Chromosome_1843_1882  
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA  
>Chromosome_1833_1872  
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT  
>Chromosome_1823_1862  
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC  
>Chromosome_1813_1852  
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA  
>Chromosome_3418_3457  
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG  
>Chromosome_3420_3459  
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC  
>Chromosome_3422_3461  
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA  
>Chromosome_3443_3482  
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAAGAAC
```

マップする側のファイル

①9-10番目のリードは、②の領域に一部がかかるように設計。領域内ではない。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 + ID=gene
Chromosome ena transcript 360 1676 + ID=transcript
Chromosome ena exon 360 1676 + Parent=transcript
Chromosome ena CDS 360 1676 + 0 ID=CDS
###
Chromosome ena gene 1852 2991 + ID=gene
Chromosome ena transcript 1852 2991 + ID=transcript
Chromosome ena exon 1852 2991 + Parent=transcript
Chromosome ena CDS 1852 2991 + 0 ID=CDS
###
Chromosome ena gene 3233 3457 + ID=gene
Chromosome ena transcript 3233 3457 + ID=transcript
Chromosome ena exon 3233 3457 + Parent=transcript
Chromosome ena CDS 3233 3457 + 0 ID=CDS
###
Chromosome ena gene 3467 4588 + ID=gene
```

有 | QuasR(Gaidatzis_2015) NEW
 するマッピングから、カウントデータ取得までの一連の
 オブジェクトをネットワーク経由で取得するのを基本
 の詳細については [マッピング | single-end | ゲノム](#)
 をコピペ。

ファイルを乳酸菌ゲノムにマ
 STA形式ファイル([sample R](#)
[Lactobacillus casei 12A](#)の
[0829395.1.30.dna.chromoso](#)
[s_jcm_18461.GCA_000829](#)

#入力ファイル名を指定
[jcm_18461.GCA_000829](#)
[jcm_18461.GCA_000829](#)
 #出力ファイル名を指定
 #カウントデータ取得時

#パッケージの読み込み
 #パッケージの読み込み
 "auto")#txdbオブジェク
 #確認してるだけです

#マッピングを行うqLi
 #マッピング結果(af
 #行数と列数を表示
 #確認してるだけです

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#  
dim(count)  
head(count)  
<
```

```
>Chromosome_361_400  
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC  
>Chromosome_1637_1676  
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG  
>Chromosome_1851_1890  
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA  
>Chromosome_1843_1882  
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA  
>Chromosome_1833_1872  
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT  
>Chromosome_1823_1862  
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC  
>Chromosome_1813_1852  
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA  
>Chromosome_3418_3457  
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG  
>Chromosome_3420_3459  
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC  
>Chromosome_3422_3461  
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA  
>Chromosome_3443_3482  
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAGAAC
```



マップする側のファイル

①11番目のリードは、②と③の領域にまたがってマップされるように設計。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena transcript 360 1676 + ID=transcript0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena exon 360 1676 + Parent=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena CDS 360 1676 + 0 ID=CDS0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
###
Chromosome ena gene 1852 2991 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena transcript 1852 2991 + ID=transcript0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena exon 1852 2991 + Parent=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena CDS 1852 2991 + 0 ID=CDS0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
###
Chromosome ena gene 3233 3457 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena transcript 3233 3457 + ID=transcript0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena exon 3233 3457 + Parent=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
Chromosome ena CDS 3233 3457 + 0 ID=CDS0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
###
Chromosome ena gene 3467 4588 + ID=gene0829395.1.30.dna.chromosome.jcm.18461.GCA.000829395.1
```

有 | QuasR(Gaidatzis_2015) NEW
 するマッピングから、カウントデータ取得までの一連の
 bオブジェクトをネットワーク経由で取得するのを基本
 の詳細については [マッピング | single-end | ゲノム](#)
 をコピペ。

ファイルを乳酸菌ゲノムにマ
 STA形式ファイル([sample R](#)
[Lactobacillus casei 12A](#)の
[0829395.1.30.dna.chromoso](#)
[s_jcm_18461.GCA_000829](#)

#入力ファイル名を指定
[jcm_18461.GCA_000829](#)
[jcm_18461.GCA_000829](#)
 #出力ファイル名を指定
 #カウントデータ取得時

#パッケージの読み込み
 #パッケージの読み込み

"auto")#txdbオブジェク
 #確認してるだけです

#マッピングを行うqAli
 #マッピング結果(align

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)#  
dim(count)  
head(count)  
<
```

```
>Chromosome_361_400  
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC  
>Chromosome_1637_1676  
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG  
>Chromosome_1851_1890  
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA  
>Chromosome_1843_1882  
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA  
>Chromosome_1833_1872  
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT  
>Chromosome_1823_1862  
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC  
>Chromosome_1813_1852  
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA  
>Chromosome_3418_3457  
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG  
>Chromosome_3420_3459  
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC  
>Chromosome_3422_3461  
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA  
>Chromosome_3443_3482  
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAAGAAC
```

最低限3リードは...

全11リード中①これらの3リードは、②の領域内にマップされる。従って、得られるカウントの総和は、最低でも3はあるはず。という予想を立てて実際にマッピングを行う。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM829395.1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
```

Chromosome	ena	gene	360	1676	+	ID=gc0829395.1.30.dna.chromosome.jcm_18461.GCA_000829395.1
Chromosome	ena	transcript	360	1676	+	ID=tr0829395.1.30.dna.transcript.jcm_18461.GCA_000829395.1
Chromosome	ena	exon	360	1676	+	Parent=tr0829395.1.30.dna.transcript.jcm_18461.GCA_000829395.1
Chromosome	ena	CDS	360	1676	+	ID=cd0829395.1.30.dna.coding_sequence.jcm_18461.GCA_000829395.1
###						
Chromosome	ena	gene	1852	2991	+	ID=gc0829395.1.30.dna.chromosome.jcm_18461.GCA_000829395.1
Chromosome	ena	transcript	1852	2991	+	ID=tr0829395.1.30.dna.transcript.jcm_18461.GCA_000829395.1
Chromosome	ena	exon	1852	2991	+	Parent=tr0829395.1.30.dna.transcript.jcm_18461.GCA_000829395.1
Chromosome	ena	CDS	1852	2991	+	ID=cd0829395.1.30.dna.coding_sequence.jcm_18461.GCA_000829395.1
###						
Chromosome	ena	gene	3233	3457	+	ID=gc0829395.1.30.dna.chromosome.jcm_18461.GCA_000829395.1
Chromosome	ena	transcript	3233	3457	+	ID=tr0829395.1.30.dna.transcript.jcm_18461.GCA_000829395.1
Chromosome	ena	exon	3233	3457	+	Parent=tr0829395.1.30.dna.transcript.jcm_18461.GCA_000829395.1
Chromosome	ena	CDS	3233	3457	+	ID=cd0829395.1.30.dna.coding_sequence.jcm_18461.GCA_000829395.1
###						
Chromosome	ena	gene	3467	4588	+	ID=gc0829395.1.30.dna.chromosome.jcm_18461.GCA_000829395.1

有 | QuasR(Ga
るマッピングから、カウントデータ取得までの一連のオブジェクトをネットワーク経由で取得するのを基本的にコピー。
ファイルを乳酸菌ゲノムにマッピングする場合、FASTA形式ファイル(sample.R)を指定し、Lactobacillus casei 12462 (GenBank: JCM 18461.GCA_000829395.1)を指定する。
#入力ファイル名を指定
jcm_18461.GCA_000829395.1
#出力ファイル名を指定
jcm_18461.GCA_000829395.1
#カウントデータ取得時のオプション
#パッケージの読み込み
#パッケージの読み込み
"auto") #txdbオブジェクトを指定する
#確認してるだけでいい
#マッピングを行うqAligner
#マッピング結果(alignment)

```
>Chromosome_361_400
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC
>Chromosome_1637_1676
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG
>Chromosome_1851_1890
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA
>Chromosome_1843_1882
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA
>Chromosome_1833_1872
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT
>Chromosome_1823_1862
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC
>Chromosome_1813_1852
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA
>Chromosome_3418_3457
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG
>Chromosome_3420_3459
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC
>Chromosome_3422_3461
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA
>Chromosome_3443_3482
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAAGAAC
```

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ取得
dim(count) #行数と列数を表示
head(count) #確認してるだけです
<
```

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`でgenenameの個数を確認

マッピング実行

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」-「ディレクトリ」の変更で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ70 ① 10. mapping single genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

mapping_single_genome7.txtのサンプル名(例: hg19)をhg19にマップした結果で、Entrez Gene IDに

```
in_f1 <- "mapping
in_f2 <- "BSgenom
out_f <- "hoge1.t
param_mapping <-
param_txdb1 <- "h
param_txdb2 <- "k
param_reportlevel
```

#必要なパッケージを
library(QuasR)
library(GenomicFe

#前処理(マッピング)
out <- qAlign(in_
alignmentStats(ou

#前処理(TxDbオブジ
txdb <- makeTxDbF
txdb

#本番(カウントデー
count <- qCount(o

<

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample RNAseq4.fa)です。マップされる側のファイルは、Ensembl (Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus hokkaidonensis jcm 18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus hokkaidonensis jcm 18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名を指定してin_f2に格納
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #入力ファイル名を指定してin_f3に格納
out_f <- "hoge10.txt" #出力ファイル名を指定してout_fに格納
param_reportlevel <- "gene" #カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"
```

#必要なパッケージをロード

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

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="auto") #txdbオブジェクトの作成
txdb #確認してるだけです
```

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2) #マッピングを行うqAlign関数を実行した結果をoutに格納
alignmentStats(out) #マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列
```

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納
dim(count) #行数と列数を表示
head(count) #確認してるだけです
```

マッピング実行

①例題10をやってみましょう。②「デスクトップ - hoge - mapping_kiso3」フォルダを作成し、③必要な入力ファイルを揃えてコピペ実行

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」→「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピペ。

1. サンプルデータ70 ① 10. mapping_single_genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

mapping_single_genome7.txtのサンプル名(例: hmg)です。hg19にマップした結果で、Entrez Gene IDに對

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome_hg19"
out_f <- "hoge10.txt"
param_mapping <- list()
param_txdb1 <- "h"
param_txdb2 <- "k"
param_reportlevel <- "gene"
```

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

#前処理(マッピング)
out <- qAlign(in_f1, in_f2, alignmentStats=out)

#前処理(TxDbオブジェクト作成)
txdb <- makeTxDbFromGFF(in_f3, txdb)

#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel)

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample_RNAseq4.fa)です。マップされる側のファイルは、Ensembl(Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aのmulti-FASTA形式ゲノム配列ファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3" #入力ファイル名
out_f <- "hoge10.txt"
param_reportlevel <- "gene"
```

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

#前処理(アノテーション情報を取得)
txdb <- makeTxDbFromGFF(in_f3, txdb)

#前処理(マッピング)
out <- qAlign(in_f1, in_f2)
alignmentStats(out)

#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納
dim(count) #行数と列数を表示
head(count) #確認してるだけです

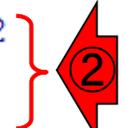
```
R Console
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso3"
> list.files()
[1] "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa"
[2] "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3"
[3] "mapping_single_genome7.txt"
[4] "sample_RNAseq4.fa"
```

途中経過

①アノテーションファイル(.gff3)読み込みのあたり。得られるカウントデータは、②2,000行以上からなるのだろう、などと妄想

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

R Console
> #前処理 (アノテーション情報を取得)
> txdb <- makeTxDbFromGFF(in_f3, format="auto")#$
Import genomic features from the file as a GRang$
Prepare the 'metadata' data frame ... OK
Make the TxDb object ... OK
> txdb                                     #確認し$
TxDb object:
# Db type: TxDb
# Supporting package: GenomicFeatures
# Data source: Lactobacillus_hokkaidonensis_jcm_$
# Organism: NA
# Taxonomy ID: NA
# miRBase build ID: NA
# Genome: NA
# transcript_nrow: 2262 }
# exon_nrow: 2262      }
# cds_nrow: 2194      }
# Db created by: GenomicFeatures package from Bi$
```



途中経過

①全11リードがマップされたのは妥当。そのように設計しているから

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

R Console
C:/Users/kojik/Desktop/hoge/mapping_kiso3\QuasR_$
Genomic alignments have been created successfully

> alignmentStats(out)
              seqlength mapped unmapped
Lacto:genome 2277985      11         0
>
> #本番 (カウントデータ取得)
> count <- qCount(out, txdb, reportLevel=param_r$
extracting gene regions from TxDb...done
counting alignments...done
collapsing counts by query name...done
> dim(count)
[1] 365  2
> head(count)
      width Lacto
accA   750    0
accB   369    0
```

#マップ\$
①

#行数と\$

#確認し\$

ん?!

①列数はともかく、なんで365行しかないの
だろう?!2000行以上ないとオカシイはずなの
に…と疑問を持ちつつも、最後まで眺める。

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

R Console
C:/Users/kojik/Desktop/hoge/mapping_kiso3\QuasR_ $
Genomic alignments have been created successfully

> alignmentStats(out) #マッピ$
      seqlength mapped unmapped
Lacto:genome 2277985      11      0
>
> #本番 (カウントデータ取得)
> count <- qCount(out, txdb, reportLevel=param_r$
extracting gene regions from TxDb...done
counting alignments...done
collapsing counts by query name...done
> dim(count) #行数と$
[1] 365  2
> head(count) #確認し$
      width Lacto
accA   750    0
accB   369    0
```



最後まで無事完了したら
、こんな感じになります

最後まで完了

```
extracting gene regions from TxDb...done
counting alignments...done
collapsing counts by query name...done
> dim(count) #行数と$
[1] 365 2
> head(count) #確認し$
      width Lacto
accA    750     0
accB    369     0
accC   1347     0
accD    789     0
ackA   1191     0
acpS    363     0
>
> #ファイルに保存
> tmp <- cbind(rownames(count), count) #保存し$
> write.table(tmp, out_f, sep="\t", append=F, qu$
> |
```

カウントデータの概要

①countオブジェクトが、得たいカウントデータ情報。②(ヘッダー部分を除く)行数は365で、列数は2。③Lactoという列名が…

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

R Console
extracting gene regions from TxDb...done
counting alignments...done
collapsing counts by query name...done
> dim(count) #行数と$
[1] 365 2
> head(count) #確認し$
  width Lacto
accA  750    0
accB  369    0
accC 1347    0
accD  789    0
ackA 1191    0
acpS  363    0
>
> #ファイルに保存
> tmp <- cbind(rownames(count), count) #保存し$
> write.table(tmp, out_f, sep="\t", append=F, qu$
> |
```

カウントデータの列名

①のリストファイル(mapping_single_genome7.txt)中に、②列名として記述したLactoです。

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ7のFASTA

mapping_single_genome7.txtのサンプル名(例: hum)です。hg19にマップした結果で、Entrez Gene IDに

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome_homo_sapiens.GRCh38.dna.chromosome.X.fa"
out_f <- "hoge10.txt"
param_mapping <- list()
param_txdb1 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa"
param_txdb2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3"
param_reportlevel <- "gene"
```

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

```
#前処理(マッピング)
out <- qAlign(in_f1, in_f2, out_f, param_mapping, param_txdb1, param_txdb2, param_reportlevel)
```

```
#前処理(TxDb)
txdb <- makeTxDbFromGFF(param_txdb2)
```

```
#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel)
```

```
dim(count)
head(count)
```

10. mapping_single_genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample_RNAseq4.fa)です。マップされる側のファイルは、Ensembl(Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa"
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3"
out_f <- "hoge10.txt"
param_reportlevel <- "gene"
```

```
#必要なパッケージをロード
```

```
library(QuasR)
```

```
library(GenomicFeatures)
```

FileName	SampleName
sample_RNAseq4.fa	Lacto

```
#前処理(マッピング)
```

```
out <- qAlign(in_f1, in_f2, out_f, param_mapping, param_txdb1, param_txdb2, param_reportlevel)
```

```
alignmentStats(out)
```

```
#本番(カウントデータ取得)
```

```
count <- qCount(out, txdb, reportLevel=param_reportlevel)
```

```
dim(count)
```

```
head(count)
```

① カファイル名を指定してin_f1に格納(RNA-seqファイル)

#入力ファイル名を指定してin_f2に格納

#入力ファイル名を指定してin_f3に格納

#出力ファイル名を指定してout_fに格納

#カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"

#パッケージの読み込み

#パッケージの読み込み

o")#txdbオブジェクトの作成
認してるだけです

#マッピングを行うqAlign関数を実行した結果をoutに格納

#マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列

#行数と列数を表示

#確認してるだけです

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`で`genename`の個数を確認

カウントデータ

マッピング実行後は、こんな感じになります。
① hoge10.txt が目的のカウントデータです

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

R Console
> tmp <- cbind(rownames(count), count) #保存し$
> write.table(tmp, out_f, sep="\t", append=F, qu$
> getwd()
[1] "C:/Users/kojik/Desktop/hoge/mapping_kiso3"
> list.files()
[1] "hoge10.txt"
[2] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[3] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[4] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[5] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[6] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[7] "mapping_single_genome7.txt"
[8] "QuasR_log_3b6c12e745f9.txt"
[9] "sample_RNAseq4.fa"
[10] "sample_RNAseq4_3b6c652a602a.bam"
[11] "sample_RNAseq4_3b6c652a602a.bam.bai"
[12] "sample_RNAseq4_3b6c652a602a.bam.txt"
> |
```



	width	Lacto
accA	750	0
accB	369	0
accC	1347	0
accD	789	0
ackA	1191	0
acpS	363	0
addA	3711	0
adh	1056	0
adhE	2607	0
adk	660	0
alaS	2655	0
aldA	1464	0
aldB	714	0
amt	1323	0

カウントデータ

①hoge10.txtの中身は、② countオブジェクトと同じです

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

R Console
[6] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[7] "mapping_single_genome7.txt" $
[8] "QuasR_log_3b6c12e745f9.txt" $
[9] "sample_RNAseq4.fa" $
[10] "sample_RNAseq4_3b6c652a602a.bam" $
[11] "sample_RNAseq4_3b6c652a602a.bam.bai" $
[12] "sample_RNAseq4_3b6c652a602a.bam.txt" $
> head(count)
  width Lacto
accA   750    0
accB   369    0
accC  1347    0
accD   789    0
ackA  1191    0
acpS   363    0
> dim(count)
[1] 365  2
> |
```



	width	Lacto
accA	750	0
accB	369	0
accC	1347	0
accD	789	0
ackA	1191	0
acpS	363	0
addA	3711	0
adh	1056	0
adhE	2607	0
adk	660	0
alaS	2655	0
aldA	1464	0
aldB	714	0
amt	1323	0

1行目のカウントが0

①最初におやっ?!と思うのは、1行目のカウントが0であるという点です。この理由は…

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

R Console
[6] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[7] "mapping_single_genome7.txt" $
[8] "QuasR_log_3b6c12e745f9.txt" $
[9] "sample_RNAseq4.fa" $
[10] "sample_RNAseq4_3b6c652a602a.bam" $
[11] "sample_RNAseq4_3b6c652a602a.bam.bai" $
[12] "sample_RNAseq4_3b6c652a602a.bam.txt" $
> head(count)
  width Lacto
accA   750    0
accB   369    0
accC  1347    0
accD   789    0
ackA  1191    0
acpS   363    0
> dim(count)
[1] 365  2
> |
```



	width	Lacto
accA	750	0
accB	369	0
accC	1347	0
accD	789	0
ackA	1191	0
acpS	363	0
addA	3711	0
adh	1056	0
adhE	2607	0
adk	660	0
alaS	2655	0
aldA	1464	0
aldB	714	0
amt	1323	0

1行目のカウントが0

①GFF3ファイル内に出現する最初のgeneのカウントは、②2だと思っていたからです

##gff-version 3						
##sequence-region Chromosome 360 2277853						
#!genome-build European Nucleotide Archive ASM82939						
#!genome-version GCA_000829395.1						
#!genome-date 2014-11						
#!genome-build-accession GCA_000829395.1						
#!genebuild-last-updated 2014-11						
Chromosome	ena	gene	360	1676	1	ID=ge
Chromosome	ena	transcript	360	1676	+	ID=tr
Chromosome	ena	exon	360	1676	+	Pare
Chromosome	ena	CDS	360	1676	+	0 ID=C
###						
Chromosome	ena	gene	1852	2991	+	ID=ge
Chromosome	ena	transcript	1852	2991	+	ID=tr
Chromosome	ena	exon	1852	2991	+	Pare
Chromosome	ena	CDS	1852	2991	+	0 ID=C
###						
Chromosome	ena	gene	3233	3457	+	ID=ge
Chromosome	ena	transcript	3233	3457	+	ID=tr
Chromosome	ena	exon	3233	3457	+	Pare
Chromosome	ena	CDS	3233	3457	+	0 ID=C
###						
Chromosome	ena	gene	3467	4588	+	ID=ge

②

```
>Chromosome_361_400
TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC
>Chromosome_1637_1676
AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG
>Chromosome_1851_1890
CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA
>Chromosome_1843_1882
TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA
>Chromosome_1833_1872
CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT
>Chromosome_1823_1862
CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC
>Chromosome_1813_1852
AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA
>Chromosome_3418_3457
GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG
>Chromosome_3420_3459
TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC
>Chromosome_3422_3461
GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA
>Chromosome_3443_3482
ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAGAAC
```

geneレベルのカウントデータ

geneレベルのカウントデータを
得ているので、①が遺伝子名
(genename)のアルファベット順
にソートされているのだろう。

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

R Console
[6] "Lactobacillus_hokkaidonensis_jcm_18461.GCA$
[7] "mapping_single_genome7.txt" $
[8] "QuasR_log_3b6c12e745f9.txt" $
[9] "sample_RNAseq4.fa" $
[10] "sample_RNAseq4_3b6c652a602a.bam" $
[11] "sample_RNAseq4_3b6c652a602a.bam.bai" $
[12] "sample_RNAseq4_3b6c652a602a.bam.txt" $
> head(count)
width Lacto
accA 750 0
accB 369 0
accC 1347 0
accD 789 0
ackA 1191 0
acpS 363 0
> dim(count)
[1] 365 2
> |
```

 ①	width	Lacto
accA	750	0
accB	369	0
accC	1347	0
accD	789	0
ackA	1191	0
acpS	363	0
addA	3711	0
adh	1056	0
adhE	2607	0
adk	660	0
alaS	2655	0
aldA	1464	0
aldB	714	0
amt	1323	0

GFFファイル内にgenenameの情報があるはずだという確信のもと、①のあたりをよく見る

GFFファイルをよく見る

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . - ① ID=gene:LOOC260_100010;Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . - ① ID=gene:LOOC260_100020;Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome ena CDS 1852 2991 . + 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###
Chromosome ena gene 3233 3457 . - ① ID=gene:LOOC260_100030;biotype=protein_coding;description=S
Chromosome ena transcript 3233 3457 . + ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome ena exon 3233 3457 . + . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome ena CDS 3233 3457 . + 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###
Chromosome ena gene 3467 4588 . - ① ID=gene:LOOC260_100040;Name=recF;biotype=protein_coding;de
```

GFFファイルをよく

①赤枠内のName=の右側の文字がgenenameのようですね。②この遺伝子領域にはName=genenameがないこともわかる。これらが原因で2000行超にはならず365行となってしまったのかも…と考える。そして、②の領域[3233, 3457]内には、確か少なくとも1リード完璧にマップされたような…といううろ覚えの記憶をたどる。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM8
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
```

Chromosome	ena	gene	360	1676	.	①	ID=gene:LOOC260_100010	Name=dnaA	biotype=protein_coding;description=
Chromosome	ena	transcript	360	1676	.	+	ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=dnaA		
Chromosome	ena	exon	360	1676	.	+	Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;exon_id=1		
Chromosome	ena	CDS	360	1676	.	+	ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP84581		
###									
Chromosome	ena	gene	1852	2991	.	①	ID=gene:LOOC260_100020	Name=dnaN	biotype=protein_coding;description=
Chromosome	ena	transcript	1852	2991	.	+	ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=dnaN		
Chromosome	ena	exon	1852	2991	.	+	Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;exon_id=1		
Chromosome	ena	CDS	1852	2991	.	+	ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP84582		
###									
Chromosome	ena	gene	3233	3457	.	②	ID=gene:LOOC260_100030		biotype=protein_coding;description=S
Chromosome	ena	transcript	3233	3457	.	+	ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=		
Chromosome	ena	exon	3233	3457	.	+	Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;exon_id=1		
Chromosome	ena	CDS	3233	3457	.	+	ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP84583		
###									
Chromosome	ena	gene	3467	4588	.	①	ID=gene:LOOC260_100040	Name=recF	biotype=protein_coding;de

dnaAのカウント数

少なくとも2リードは完全に領域内にマップされているはずの、①dnaAのカウント数を調べる。

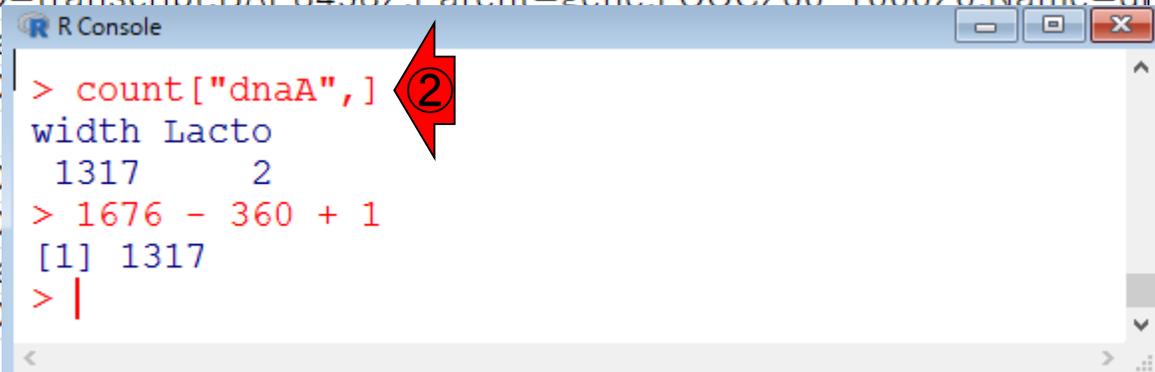
##gff-version 3						
##sequence-region Chromosome 360 2277853						
#!genome-build European Nucleotide Archive ASM82939v1						
#!genome-version GCA_000829395.1						
#!genome-date 2014-11						
#!genome-build-accession GCA_000829395.1						
#!genebuild-last-updated 2014-11						
Chromosome	ena	gene	360	1676	.	+ . ID=gene:LOOC260_100010;Name= <u>dnaA</u> ;biotype=protein_coding;d
Chromosome	ena	transcript	360	1676	.	+ . ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome	ena	exon	360	1676	.	+ . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome	ena	CDS	360	1676	.	+ 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###						
Chromosome	ena	gene	1852	2991	.	+ . ID=gene:LOOC260_100020;Name=dnaN;biotype=protein_coding;d
Chromosome	ena	transcript	1852	2991	.	+ . ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome	ena	exon	1852	2991	.	+ . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome	ena	CDS	1852	2991	.	+ 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###						
Chromosome	ena	gene	3233	3457	.	+ . ID=gene:LOOC260_100030;biotype=protein_coding;description=S4
Chromosome	ena	transcript	3233	3457	.	+ . ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome	ena	exon	3233	3457	.	+ . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome	ena	CDS	3233	3457	.	+ 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###						
Chromosome	ena	gene	3467	4588	.	+ . ID=gene:LOOC260_100040;Name=recF;biotype=protein_coding;de



dnaAのカウント数

少なくとも2リードは完全に領域内にマップされているはずの、①dnaAのカウント数を調べる。
②カウント数は2ですね。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . + . ID=gene:LOOC260_100010;Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + . ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . + . ID=gene:LOOC260_100020;Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + . ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Pa
Chromosome ena CDS 1852 2991 . + 0 ID
###
Chromosome ena gene 3233 3457 . + . ID
Chromosome ena transcript 3233 3457 . + . ID
Chromosome ena exon 3233 3457 . + . Pa
Chromosome ena CDS 3233 3457 . + 0 ID
###
Chromosome ena gene 3467 4588 . + . ID=gene:LOOC260_100040;Name=recF;biotype=protein_coding;de
```



```
R Console
> count["dnaA",]
width Lacto
1317 2
> 1676 - 360 + 1
[1] 1317
> |
```

dnaAのカウント数

少なくとも2リードは完全に領域内にマップされているはずの、①dnaAのカウント数を調べる。
②カウント数は2ですね。ちなみに、③1317は、dnaAの領域の長さです。

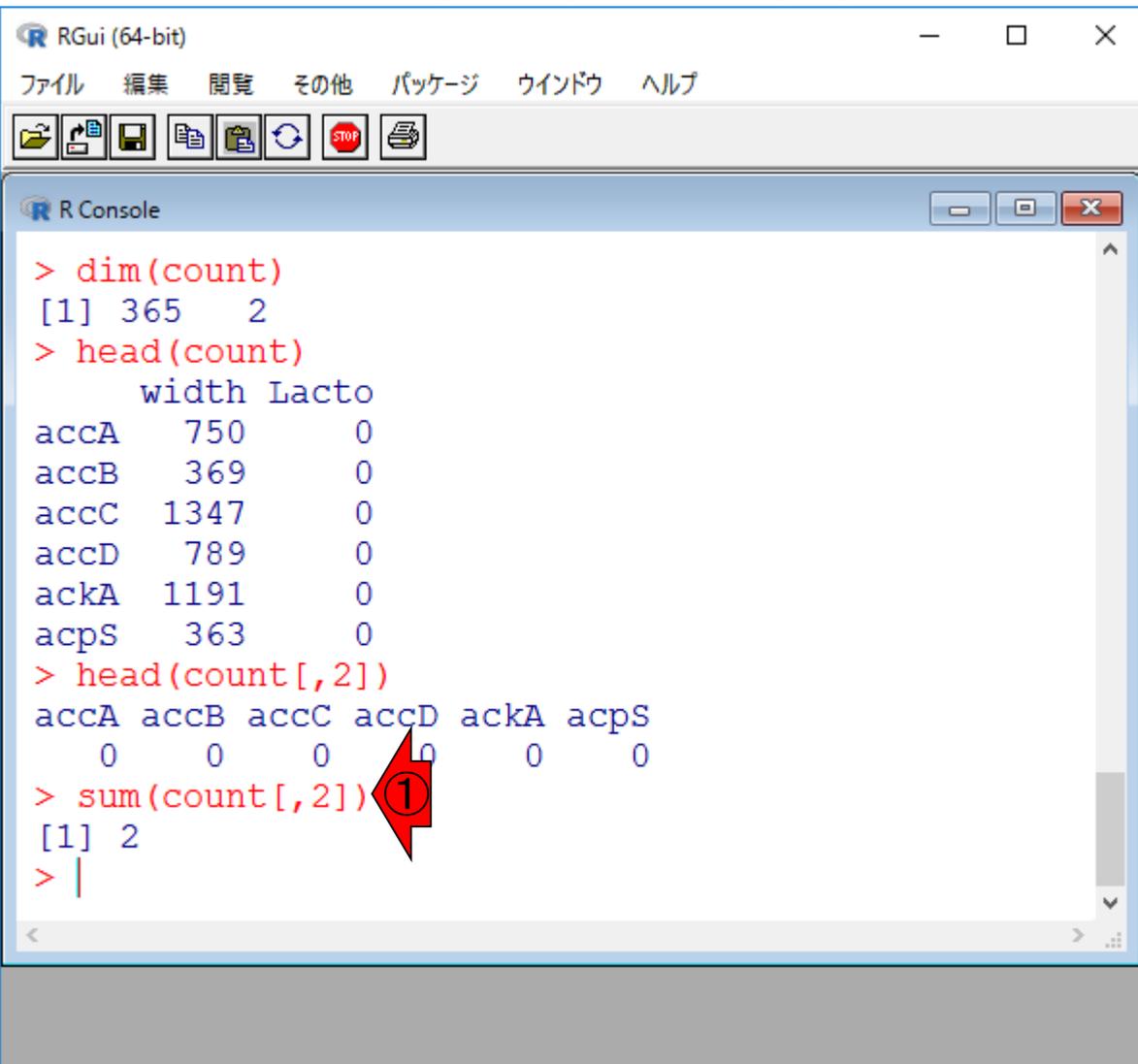
```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
```

Chromosome	ena	gene	360	1676	.	+	.	ID=gene:LOOC260_100010;Name=dnaA;biotype=protein_coding;d
Chromosome	ena	transcript	360	1676	.	+	.	ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome	ena	exon	360	1676	.	+	.	Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome	ena	CDS	360	1676	.	+	0	ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###								
Chromosome	ena	gene	1852	2991	.	+	.	ID=gene:LOOC260_100020;Name=dnaN;biotype=protein_coding;d
Chromosome	ena	transcript	1852	2991	.	+	.	ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome	ena	exon	1852	2991	.	+	.	Pa
Chromosome	ena	CDS	1852	2991	.	+	0	ID
###								
Chromosome	ena	gene	3233	3457	.	+	.	1317 2
Chromosome	ena	transcript	3233	3457	.	+	.	> 1676 - 360 + 1
Chromosome	ena	exon	3233	3457	.	+	.	[1] 1317
Chromosome	ena	CDS	3233	3457	.	+	0	>
###								
Chromosome	ena	gene	3467	4588	.	+	.	ID=gene:LOOC260_100040;Name=recF;biotype=protein_coding;de

```
R Console
> count["dnaA",]
width Lacto
1317 2
> 1676 - 360 + 1
[1] 1317
> |
```

全365遺伝子の領域にマップされたリードの総数は、①2でした

カウント総数の確認



```
> dim(count)
[1] 365  2
> head(count)
      width Lacto
accA    750     0
accB    369     0
accC   1347     0
accD    789     0
ackA   1191     0
acpS    363     0
> head(count[,2])
accA accB accC accD ackA acpS
  0    0    0    0    0    0
> sum(count[,2])
[1] 2
> |
```

Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`で`genename`の個数を確認

?qCount

カウントデータ取得部分では、qCount関数を利用しているので、他にどのようなオプションがあるのかを調べてみる。

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「ファイル」-「ディレクトリの変更」で解析したいファイルを置いてあるディレクトリに移動し以下をコピー。

1. サンプルデータ7のFASTA

mapping_single_genome7.txtのサンプル名(例: hum)です。hg19にマップした結果で、Entrez Gene IDに

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome_hg19.fa"
out_f <- "hoge1.txt"
param_mapping <- list()
param_txdb1 <- "h"
param_txdb2 <- "k"
param_reportlevel <- "gene"
```

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

```
#前処理(マッピング)
out <- qAlign(in_f1, in_f2, out_f, param_mapping)
```

```
#前処理(TxDbオブジェクトの作成)
txdb <- makeTxDbFromGFF(in_f3, format="auto")
txdb
```

```
#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel=param_reportlevel)
```

```
dim(count)
head(count)
```

10. mapping_single_genome7.txt中のFASTA形式ファイルを乳酸菌ゲノムにマッピングする場合:

マップする側のファイルは、サンプルデータ47のFASTA形式ファイル(sample_RNAseq4.fa)です。マップされる側のファイルは、Ensembl(Zerbino et al., Nucleic Acids Res., 2018)から提供されているLactobacillus casei 12Aの multi-FASTA形式ゲノム配列ファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa)です。マッピング結果に対して、GFF3形式のアノテーションファイル(Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.gff3)を読み込んでカウント情報を取得しています。

```
in_f1 <- "mapping_single_genome7.txt" #入力ファイル名を指定してin_f1に格納(RNA-seqファイル)
in_f2 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.dna.chromosome.Chromosome.fa" #入力ファイル名を指定してin_f2に格納
in_f3 <- "Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3" #入力ファイル名を指定してin_f3に格納
out_f <- "hoge10.txt" #出力ファイル名を指定してout_fに格納
param_reportlevel <- "gene" #カウントデータ取得時のレベルを指定: "gene", "exon", "promoter", "junction"
```

#必要なパッケージをロード

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

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="auto") #txdbオブジェクトの作成
txdb #確認してるだけです
```

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2, out_f, param_mapping) #マッピングを行うqAlign関数を実行した結果をoutに格納
alignmentStats(out) #マッピング結果(alignment statistics)の表示。seqlength: リファレンス配列
```

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel=param_reportlevel) #カウントデータ行列を取得してcountに格納
dim(count) #行数と列数を表示
head(count) #確認してるだけです
```

1

?qCount

```

RGui (64-bit)
ファイル 編集 閲覧 その他 パッケージ ウィンドウ
[Icons]
R Console
> dim(count)
[1] 365  2
> head(count)
      width Lacto
accA    750    0
accB    369    0
accC   1347    0
accD    789    0
ackA   1191    0
acpS    363    0
> head(count[,2])
accA accB accC accD ackA acpS
  0    0    0    0    0    0
> sum(count[,2])
[1] 2
> ?qCount|

```

①

qCount {QuasR}

R Documentation

Quantify alignments

Description

Quantify alignments from sequencing data.

Usage

```

qCount(proj,
        query,
        reportLevel=c(NULL, "gene", "exon", "promoter", "junction"),
        selectReadPosition=c("start", "end"),
        shift=0L,
        orientation=c("any", "same", "opposite"),
        useRead=c("any", "first", "last"),
        auxiliaryName=NULL,
        mask=NULL,
        collapseBySample=TRUE,
        includeSpliced=TRUE,
        includeSecondary=TRUE,
        mapqMin=0L,
        mapqMax=255L,
        absIsizeMin=NULL,
        absIsizeMax=NULL,
        maxInsertSize=500L,
        clobj=NULL)

```

orientation

①orientation部分の説明。②デフォルトはanyのようだ。これは、昔のRNA-seqが二本鎖の状態ではsequenceされていたためである。つまりリードの方向性(orientation)はわからなかったということ。リファレンス配列上の領域はわかるが、+鎖側(same strand)か-鎖側(opposite strand)かは不明だった、ということ。

qCount {QuasR}

Quantify alignments

Description

Quantify alignments from sequencing data

Usage

```
qCount (proj,  
        query,  
        reportLevel=c(NULL, "gene"),  
        selectReadPosition=c("start", "end"),  
        shift=0L,  
        orientation=c("any", "same", "opposite"),  
        useRead=c("any", "first", "last"),  
        auxiliaryName=NULL,  
        mask=NULL,  
        collapseBySample=TRUE,  
        includeSpliced=TRUE,  
        includeSecondary=TRUE,  
        mapqMin=0L,  
        mapqMax=255L,  
        absIsizeMin=NULL,  
        absIsizeMax=NULL,  
        maxInsertSize=500L,  
        clobj=NULL)
```

orientation

sets the required orientation of the alignments relative to the query region in order to be counted, one of:

- any (default): count alignment on the same and opposite strand
- same : count only alignment on the same strand
- opposite : count only alignment on the opposite strand

orientation

①今回のマップする側のリードは、+鎖側(same strand)で設計したので、おそらくanyでもsameでも同じ結果になる。しかし、おそらくoppositeを指定したら(orientation = "opposite"), カウントの総和は0になるだろう。未確認

qCount {QuasR}

Quantify alignments

Description

Quantify alignments from sequencing data

Usage

```
qCount (proj,  
        query,  
        reportLevel=c(NULL, "gene"),  
        selectReadPosition=c("start", "end", "middle"),  
        shift=0L,  
        orientation=c("any", "same", "opposite"),  
        useRead=c("any", "first", "last"),  
        auxiliaryName=NULL,  
        mask=NULL,  
        collapseBySample=TRUE,  
        includeSpliced=TRUE,  
        includeSecondary=TRUE,  
        mapqMin=0L,  
        mapqMax=255L,  
        absIsizeMin=NULL,  
        absIsizeMax=NULL,  
        maxInsertSize=500L,  
        clobj=NULL)
```

orientation

sets the required orientation of the alignments relative to the query region in order to be counted, one of:

- ① • any (default): count alignment on the same and opposite strand
- ② • same : count only alignment on the same strand
- ③ • opposite : count only alignment on the opposite strand

少しずつれたリード

①1塩基くらいずれていても、領域内の大部分にマップされたリードということでカウント情報として加えるにはどうすればよいのか？という視点でオプション名を眺める。①shiftとかのオプションをshift = 1などとすればいいのかな…などと妄想しながら説明文を読む。

qCount {QuasR}

Quantify alignments

Description

Quantify alignments from sequencing data.

Usage

```
qCount (proj,  
        query,  
        reportLevel=c(NULL, "gene", "exon", "promoter", "junction"),  
        selectReadPosition=c("start", "end"),  
        shift=0L,  
        orientation=c("any", "same", "opposite"),  
        useRead=c("any", "first", "last"),  
        auxiliaryName=NULL,  
        mask=NULL,  
        collapseBySample=TRUE,  
        includeSpliced=TRUE,  
        includeSecondary=TRUE,  
        mapqMin=0L,  
        mapqMax=255L,  
        absIsizeMin=NULL,  
        absIsizeMax=NULL,  
        maxInsertSize=500L,  
        clobj=NULL)
```

①

少しずつれたリード

①shiftの説明部分。なんで②3'-endに限定するの
か不明だが、もしかしたら③のところのデフォルトが
④startということと関連するのかな?!などと想像する

qCount {QuasR}

R Documentation

Quantify alignments

Description

Quantify alignments from sequencing data

Usage

```
qCount(proj,  
  query,  
  reportLevel=c(NULL, "gene",  
  selectReadPosition=c("start", "end", "middle", "all"),  
  shift=0L,  
  orientation=c("any", "same", "opposite"),  
  useRead=c("any", "first", "last"),  
  auxiliaryName=NULL,  
  mask=NULL,  
  collapseBySample=TRUE,  
  includeSpliced=TRUE,  
  includeSecondary=TRUE,  
  mapqMin=0L,  
  mapqMax=255L,  
  absIsizeMin=NULL,  
  absIsizeMax=NULL,  
  maxInsertSize=500L,  
  clobj=NULL)
```

③

selectReadPosition

defines the part of the alignment that has to be contained within a query region to produce an overlap (see Details). Possible values are:

- ④ • start (default): start of the alignment
- end: end of the alignment

①

shift

controls the shifting alignments towards their 3'-end before quantification. shift can be one of:

- an “integer” vector of the same length as the number of alignment files
- a single “integer” value
- the character string "halfInsert" (only available for paired-end experiments)

The default of 0 will not shift any alignments.

②

少しずつれたリード

ただ、こういう場合は、プログラムのバグの可能性など様々な要因が考えられるので、まずは③ selectReadPositionには手をつけず、① shift = 1を付けてどのような結果になるかを見てみる。

qCount {QuasR}

Quantify alignments

Description

Quantify alignments from sequencing data

Usage

```
qCount (proj,  
        query,  
        reportLevel=c(NULL, "gene",  
                        selectReadPosition=c("start", "end", "middle"),  
                        shift=0L,  
        orientation=c("any", "same", "opposite"),  
        useRead=c("any", "first", "last"),  
        auxiliaryName=NULL,  
        mask=NULL,  
        collapseBySample=TRUE,  
        includeSpliced=TRUE,  
        includeSecondary=TRUE,  
        mapqMin=0L,  
        mapqMax=255L,  
        absIsizeMin=NULL,  
        absIsizeMax=NULL,  
        maxInsertSize=500L,  
        clobj=NULL)
```

③ selectReadPosition

defines the part of the alignment that has to be contained within a query region to produce an overlap (see Details). Possible values are:

- ④ • start (default): start of the alignment
- end: end of the alignment

① shift

② controls the shifting alignments towards their 3'-end before quantification. shift can be one of:

- an “integer” vector of the same length as the number of alignment files
- a single “integer” value
- the character string "halfInsert" (only available for paired-end experiments)

The default of 0 will not shift any alignments.

1塩基のずれのみでカウントされるかもしれないのは、①のリード。②の遺伝子領域が対応します。

少しずつたリード

##gff-version 3						
##sequence-region Chromosome 360 2277853						
#!genome-build European Nucleotide Archive ASM82939v1						
#!genome-version GCA_000829395.1						
#!genome-date 2014-11						
#!genome-build-accession GCA_000829395.1						
#!genebuild-last-updated 2014-11						
Chromosome	ena	gene	360	1676	.	+ . ID=gene:LOOC260_10
Chromosome	ena	transcript	360	1676	.	+ . ID=transcript:BAP845
Chromosome	ena	exon	360	1676	.	+ . Parent=transcript:BA
Chromosome	ena	CDS	360	1676	.	+ 0 ID=CDS:BAP84581;Pa
###						
Chromosome	ena	gene	1852	2991	.	+ . ID=gene:LOOC260_10
Chromosome	ena	transcript	1852	2991	.	+ . ID=transcript:BAP845
Chromosome	ena	exon	1852	2991	.	+ . Parent=transcript:BA
Chromosome	ena	CDS	1852	2991	.	+ 0 ID=CDS:BAP84582;Pa
###						
Chromosome	ena	gene	3233	3457	.	+ . ID=gene:LOOC260_10
Chromosome	ena	transcript	3233	3457	.	+ . ID=transcript:BAP845
Chromosome	ena	exon	3233	3457	.	+ . Parent=transcript:BA
Chromosome	ena	CDS	3233	3457	.	+ 0 ID=CDS:BAP84583;Pa
###						
Chromosome	ena	gene	3467	4588	.	+ . ID=gene:LOOC260_10

>Chromosome_361_400	TGACTGATTTAGAAACACTTTGGGACACAATTAAGAATC
>Chromosome_1637_1676	AGAAGATGTCCAAAACCTTAAAATGGAGCTAAAGCCATAG
>Chromosome_1851_1890	CATGAAATTTACAATTAGTCGTGCAACTTTTACAGCCAAA
>Chromosome_1843_1882	TAACCAATCATGAAATTTACAATTAGTCGTGCAACTTTTA
>Chromosome_1833_1872	CTTCAAGGAGTAACCAATCATGAAATTTACAATTAGTCGT
>Chromosome_1823_1862	CAAATTCAACCTTCAAGGAGTAACCAATCATGAAATTTAC
>Chromosome_1813_1852	AAATTAAGACAAATTCAACCTTCAAGGAGTAACCAATCA
>Chromosome_3418_3457	GATTGCAGATAATGGGACATTTGTCATTCAAATGAGTAG
>Chromosome_3420_3459	TTGCAGATAATGGGACATTTGTCATTCAAATGAGTAGGC
>Chromosome_3422_3461	GCAGATAATGGGACATTTGTCATTCAAATGAGTAGGCAA
>Chromosome_3443_3482	ATTCAAATGAGTAGGCAACTTAAATGATTTTAAAGAAC

少しずつたリード

1塩基のずれのみでカウントされるかもしれないのは、①のリード。②の遺伝子領域が対応します。③その遺伝子名はdnaN。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . + . ID=gene:LOOC260_100010;Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + . ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . + . ID=gene:LOOC260_100020;Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + . ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome ena CDS 1852 2991 . + 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###
Chromosome ena gene 3233 3457 . + . ID=gene:LOOC260_100030;biotype=protein_coding;description=S4
Chromosome ena transcript 3233 3457 . + . ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome ena exon 3233 3457 . + . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome ena CDS 3233 3457 . + 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###
Chromosome ena gene 3467 4588 . + . ID=gene:LOOC260_100040;Name=recF;biotype=protein_coding;de
```



コピペ実行後

①例題11のコピペ実行後は、②のような感じになります。

マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「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

1. サンプルデータ7011. mapping single genome7.txt中のFASTA形式

のサンプル名(例: hg19)にマッピングした結果で、Entrez Gene IDに

例題10と基本的に同じですが、カウント情報を

```
in_f1 <- "mapping_single_genome7.t
in_f2 <- "Lactobacillus_hokkaido
in_f3 <- "Lactobacillus_hokkaido
out_f <- "hoge11.txt"
param_mapping <-
param_txdb1 <- "h
param_txdb2 <- "k
param_reportlevel
```

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

#前処理(マッピング)
out <- qAlign(in_
alignmentStats(ou

#前処理(TxDbオブジ
txdb <- makeTxDbF
txdb

#本番(カウントデー
count <- qCount(o

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

#前処理(アノテーション情報を取得)
txdb <- makeTxDbFromGFF(in_f3, for
txdb

#前処理(マッピング)
out <- qAlign(in_f1, in_f2)
alignmentStats(out)

#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLe
shift=param_shift)
```

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

R Console
call dbDisconnect() when finished working with a $
> dim(count)
[1] 365 2
> head(count)
      width Lacto
accA    750    0
accB    369    0
accC   1347    0
accD    789    0
ackA   1191    0
acps    363    0
>
> #ファイルに保存
> tmp <- cbind(rownames(count), count) #保存し$
> write.table(tmp, out_f, sep="\t", append=F, quo$
> |
```

全365遺伝子の領域にマップされたリードの総数は、①3になりました!

カウント総数の確認



マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「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

1. サンプルデータのFASTQ形式ファイル(SRR037430.fastq)のBSgenome.Hsapiens

mapping_single_genome7.txt中のFASTA形式のサンプル名(例: hum)です。hg19にマップした結果で、Entrez Gene IDに対

11. mapping_single_genome7.txt中のFASTA形式

例題10と基本的に同じですが、カウント情報を

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "Lactobacillus_hokkaidonei"
in_f3 <- "Lactobacillus_hokkaidonei"
out_f <- "hoge11.txt"
param_mapping <-
param_txdb1 <- "h"
param_txdb2 <- "k"
param_reportlevel
```

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

```
#前処理(アノテーション情報を取得)
txdb <- makeTxDbFromGFF(in_f3, format="gff", txdb=txdb)
```

```
#前処理(マッピング)
out <- qAlign(in_f1, in_f2, txdb=txdb, alignmentStats=out)
```

```
#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel="gene", shift=param_shift)
```

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "BSgenome.Hsapiens.hg19"
out_f <- "hoge1.txt"
param_mapping <-
param_txdb1 <- "hg19"
param_txdb2 <- "k"
param_reportlevel
```

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

```
#前処理(マッピング)
out <- qAlign(in_f1, in_f2, txdb=txdb, alignmentStats=out)
```

```
#前処理(TxDbオブジェクト取得)
txdb <- makeTxDbFromGFF(in_f3, format="gff", txdb=txdb)
```

```
#本番(カウントデータ取得)
count <- qCount(out, txdb, reportLevel="gene", shift=param_shift)
```

```

[1] 365  2
> head(count)
      width Lacto
accA    750    0
accB    369    0
accC   1347    0
accD    789    0
ackA   1191    0
acpS    363    0
>
> #ファイルに保存
> tmp <- cbind(rownames(count), count) #保存し$
> write.table(tmp, out_f, sep="\t", append=F, quo$
> sum(count[,2])
[1] 3
> |

```

dnaNのカウント数



マップ後 | カウント情報取得 | single-end | ゲノム | アノテーション有 | QuasR(Gaidatzis_2015) NEW

QuasRパッケージを用いたsingle-end RNA-seqデータのリファレンスゲノム配列へのBowtieによるマッピングから、カウントデータ取得までの一連の流れを示します。アノテーション情報は、GenomicFeaturesパッケージ中の関数を利用してTxDbオブジェクトをネットワーク経由で取得するのを基本としつつ、TxDbパッケージを読み込むやり方も示しています。マッピングのやり方やオプションの詳細についてはマッピング | single-end | ゲノム | basic aligner(応用) | QuasR(Gaidatzis_2015)などを参考にしてください。

「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

1. サンプルデータのFASTQ形式ファイル(SRR037430.fastq)のBSgenome.Hsapiens

mapping_single_genome7.txt中のFASTA形式のサンプル名(例: hum)です。hg19にマップした結果で、Entrez Gene IDに対

11. mapping_single_genome7.txt中のFASTA形式

例題10と基本的に同じですが、カウント情報を

```
in_f1 <- "mapping_single_genome7.txt"
in_f2 <- "Lactobacillus_hokkaidonei"
in_f3 <- "Lactobacillus_hokkaidonei"
out_f <- "hoge11.txt"
param_mapping <-
param_txdb1 <- "h"
param_txdb2 <- "k"
param_reportlevel
```

#必要なパッケージをロード

```
library(QuasR)
library(GenomicFeatures)
```

#前処理(アノテーション情報を取得)

```
txdb <- makeTxDbFromGFF(in_f3, format="gff",
txdb
```

#前処理(マッピング)

```
out <- qAlign(in_f1, in_f2)
alignmentStats(out)
```

#本番(カウントデータ取得)

```
count <- qCount(out, txdb, reportLevel="gene",
shift=param_shift)
```

RGui (64-bit)

ファイル 編集 閲覧 その他 パッケージ ウィンドウ ヘルプ Vignettes



R Console

```
accA 750 0
accB 369 0
accC 1347 0
accD 789 0
ackA 1191 0
acpS 363 0
>
> #ファイルに保存
> tmp <- cbind(rownames(count), count) #保存し$
> write.table(tmp, out_f, sep="\t", append=F, quo$
> sum(count[,2])
[1] 3
> count["dnaN",]
width Lacto
1140 1
> |
```



Contents

■ マッピング (アラインメント) の続き

- おさらい: 入力ファイル (マップする側、される側)、QuasRの結果、Bowtie2の結果
- マップされなかったリード: Bowtie (デフォルト)、Bowtie (QuasRと同じオプション)
- SAM形式の解説、マッピング結果の違い、課題
- Linux環境以外でのBowtie2実行手段

■ カウント情報取得

- アノテーション情報がない場合: 単一サンプル、複数サンプル
- アノテーション情報がある場合
 - 概要
 - マップする側のファイルの説明
 - マッピング実行
 - 結果の解釈
 - カウント情報取得時のオプション
 - `grep`でgenenameの個数を確認

残る問題

(他にも沢山あるが)なぜ得られたカウントデータの行数が2000行超にならずに365行となってしまったのか、について考える

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . - ① ID=gene:LOOC260_100010 Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . - ① ID=gene:LOOC260_100020 Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome ena CDS 1852 2991 . + 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###
Chromosome ena gene 3233 3457 . - ② ID=gene:LOOC260_100030;biotype=protein_coding;description=S4
Chromosome ena transcript 3233 3457 . + ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome ena exon 3233 3457 . + . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome ena CDS 3233 3457 . + 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###
Chromosome ena gene 3467 4588 . - ① ID=gene:LOOC260_100040 Name=recF;biotype=protein_coding;de
```

残る問題

Linux環境で、genenameの個数が本当に365個だったのかを検証する。

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . - ① ID=gene:LOOC260_100010 Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . - ① ID=gene:LOOC260_100020 Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome ena CDS 1852 2991 . + 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###
Chromosome ena gene 3233 3457 . + . ID=gene:LOOC260_100030;biotype=protein_coding;description=S4
Chromosome ena transcript 3233 3457 . + . ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome ena exon 3233 3457 . + . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome ena CDS 3233 3457 . + 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###
Chromosome ena gene 3467 4588 . - ① ID=gene:LOOC260_100040 Name=recF;biotype=protein_coding;de
```

お約束の①pwdと②ls。このような状況で③gff3ファイルのみを取り扱う

Bio-Linuxのターミナル画面

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3] 15:20
① iu@bielinux[mapping_kiso3] pwd [ 3:19午後 ]
/home/iu/Desktop/mac_share/mapping_kiso3
② iu@bielinux[mapping_kiso3] ls [ 3:19午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chrom
osome.Chromosome.gff3
iu@bielinux[mapping_kiso3] [ 3:19午後 ]
```



Tips: ワイルドカード

①*.gff3と書くことで、.gff3で終わる全てのファイルのみをリストアップすることができます。今回の場合は、ファイルが1つしかないなので、タブ補完で直打ちしてもよいといえはよい。

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
iu@bielinux[mapping_kiso3] pwd [ 3:19午後 ]
/home/iu/Desktop/mac_share/mapping_kiso3
iu@bielinux[mapping_kiso3] ls [ 3:19午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chrom
osome.Chromosome.gff3
① iu@bielinux[mapping_kiso3] ls *.gff3 [ 3:19午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chrom
osome.Chromosome.gff3
iu@bielinux[mapping_kiso3] [ 3:27午後 ]
```

残る問題

(他にも沢山あるが)なぜ得られたカウントデータの行数が2000行超にならずに365行となってしまったのか、について考える

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . - ① ID=gene:LOOC260_100010 Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . - ① ID=gene:LOOC260_100020 Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome ena CDS 1852 2991 . + 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###
Chromosome ena gene 3233 3457 . - ② ID=gene:LOOC260_100030;biotype=protein_coding;description=S4
Chromosome ena transcript 3233 3457 . + ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome ena exon 3233 3457 . + . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome ena CDS 3233 3457 . + 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###
Chromosome ena gene 3467 4588 . - ① ID=gene:LOOC260_100040 Name=recF;biotype=protein_coding;de
```

ID=geneを含む行数

Geneの領域数をカウントする。全体をざっと眺めて、①ID=geneを含む行数をカウントすればよいだろう

```
##gff-version 3
##sequence-region Chromosome 360 2277853
#!genome-build European Nucleotide Archive ASM82939v1
#!genome-version GCA_000829395.1
#!genome-date 2014-11
#!genome-build-accession GCA_000829395.1
#!genebuild-last-updated 2014-11
Chromosome ena gene 360 1676 . + . ID=gene:LOOC260_100010;Name=dnaA;biotype=protein_coding;d
Chromosome ena transcript 360 1676 . + . ID=transcript:BAP84581;Parent=gene:LOOC260_100010;Name=d
Chromosome ena exon 360 1676 . + . Parent=transcript:BAP84581;Name=BAP84581-1;constitutive=1;e
Chromosome ena CDS 360 1676 . + 0 ID=CDS:BAP84581;Parent=transcript:BAP84581;protein_id=BAP8
###
Chromosome ena gene 1852 2991 . + . ID=gene:LOOC260_100020;Name=dnaN;biotype=protein_coding;d
Chromosome ena transcript 1852 2991 . + . ID=transcript:BAP84582;Parent=gene:LOOC260_100020;Name=d
Chromosome ena exon 1852 2991 . + . Parent=transcript:BAP84582;Name=BAP84582-1;constitutive=1;e
Chromosome ena CDS 1852 2991 . + 0 ID=CDS:BAP84582;Parent=transcript:BAP84582;protein_id=BAP8
###
Chromosome ena gene 3233 3457 . + . ID=gene:LOOC260_100030;biotype=protein_coding;description=S4
Chromosome ena transcript 3233 3457 . + . ID=transcript:BAP84583;Parent=gene:LOOC260_100030;biotype=
Chromosome ena exon 3233 3457 . + . Parent=transcript:BAP84583;Name=BAP84583-1;constitutive=1;e
Chromosome ena CDS 3233 3457 . + 0 ID=CDS:BAP84583;Parent=transcript:BAP84583;protein_id=BAP8
###
Chromosome ena gene 3467 4588 . + . ID=gene:LOOC260_100040;Name=recF;biotype=protein_coding;de
```

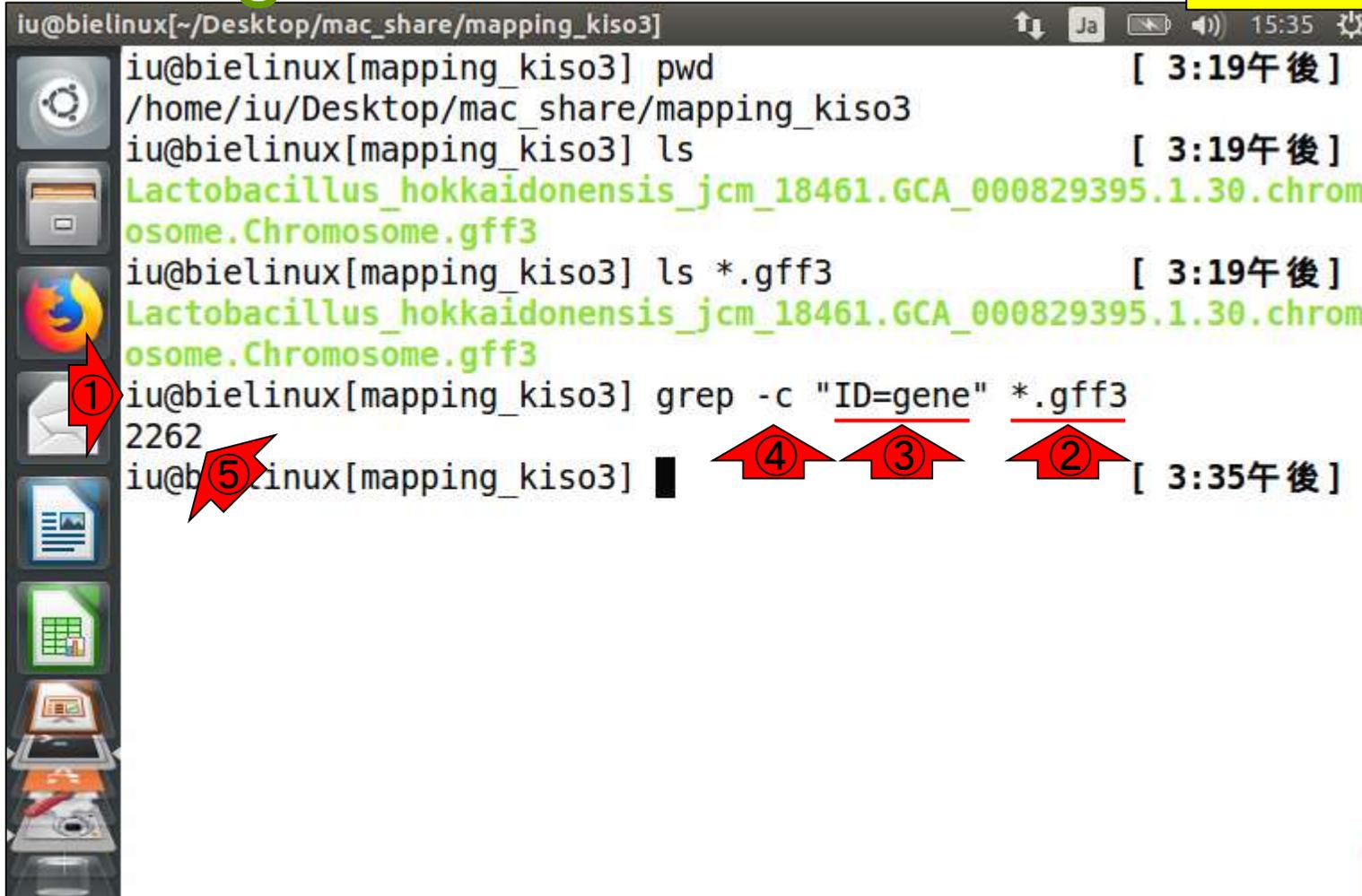


①

ID=geneを含む行数

①grepコマンドで、②*.gff3というファイルに対して、③ID=geneという文字列を含む、④行数は、⑤2262行

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
iu@bielinux[mapping_kiso3] pwd [ 3:19午後 ]
/home/iu/Desktop/mac_share/mapping_kiso3
iu@bielinux[mapping_kiso3] ls [ 3:19午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3
iu@bielinux[mapping_kiso3] ls *.gff3 [ 3:19午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chromosome.Chromosome.gff3
① iu@bielinux[mapping_kiso3] grep -c "ID=gene" *.gff3
2262
iu@bielinux[mapping_kiso3] [ 3:35午後 ]
```



ID=geneを含む行を表示

①を実行すると、ID=geneという文字列を含む行がそのまま表示される。2262行分あるので、画面がざっと流れる。

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
iu@bielinux[mapping_kiso3] pwd [ 3:42午後 ]
/home/iu/Desktop/mac_share/mapping_kiso3
iu@bielinux[mapping_kiso3] ls [ 3:42午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chrom
osome.Chromosome.gff3
iu@bielinux[mapping_kiso3] ls *.gff3 [ 3:42午後 ]
Lactobacillus_hokkaidonensis_jcm_18461.GCA_000829395.1.30.chrom
osome.Chromosome.gff3
iu@bielinux[mapping_kiso3] grep -c "ID=gene" *.gff3
2262
① iu@bielinux[mapping_kiso3] grep "ID=gene" *.gff3 [ 3:42午後 ]
```

ID=geneを含む行を表示

画面がざっと流れた結果。この結果をざっと眺めても、genenameがあるのは5個中①2個と少ないですね。

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
in;gene_id=L00C260_122640;logic_name=ena;version=1
Chromosome      ena      gene      2273924 2275312 .      -      .
ID=gene:L00C260_122650;Name=trmE;biotype=protein_coding;description=tRNA modification GTPase;gene_id=L00C260_122650;logic_name=ena;version=1
Chromosome      ena      gene      2275488 2276288 .      -      .
ID=gene:L00C260_122660;biotype=protein_coding;description=single-stranded DNA-binding protein;gene_id=L00C260_122660;logic_name=ena;version=1
Chromosome      ena      gene      2276455 2277288 .      -      .
ID=gene:L00C260_122670;biotype=protein_coding;description=membrane protein;gene_id=L00C260_122670;logic_name=ena;version=1
Chromosome      ena      gene      2277304 2277648 .      -      .
ID=gene:L00C260_122680;biotype=protein_coding;description=ribonuclease P;gene_id=L00C260_122680;logic_name=ena;version=1
Chromosome      ena      gene      2277719 2277853 .      -      .
ID=gene:L00C260_122690;Name=rpmH;biotype=protein_coding;description=50S ribosomal protein L34;gene_id=L00C260_122690;logic_name=ena;version=1
iu@bielinux[mapping_kiso3] [ 3:45午後 ]
```

ID=geneを含む行をファイル出力

①のようにするとターミナル画面上に表示するのではなく、`uge.txt`に保存できます

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
in;gene_id=L00C260_122640;logic_name=ena;version=1
Chromosome      ena      gene      2273924 2275312 .      -      .
ID=gene:L00C260_122650;Name=trmE;biotype=protein_coding;descrip
tion=tRNA modification GTPase;gene_id=L00C260_122650;logic_name
=ena;version=1
Chromosome      ena      gene      2275488 2276288 .      -      .
ID=gene:L00C260_122660;biotype=protein_coding;description=singl
e-stranded DNA-binding protein;gene_id=L00C260_122660;logic_nam
e=ena;version=1
Chromosome      ena      gene      2276455 2277288 .      -      .
ID=gene:L00C260_122670;biotype=protein_coding;description=membr
ane protein;gene_id=L00C260_122670;logic_name=ena;version=1
Chromosome      ena      gene      2277304 2277648 .      -      .
ID=gene:L00C260_122680;biotype=protein_coding;description=ribon
uclease P;gene_id=L00C260_122680;logic_name=ena;version=1
Chromosome      ena      gene      2277719 2277853 .      -      .
ID=gene:L00C260_122690;Name=rpmH;biotype=protein_coding;descrip
tion=50S ribosomal protein L34;gene_id=L00C260_122690;logic_nam
e=ena;version=1
iu@bielinux[mapping_kiso3] grep "ID=gene" *.gff3 > uge.txt
```

①wcでuge.txtの行数を確認。確かに2262行ですね。

行数を確認

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
tion=tRNA modification GTPase;gene_id=L00C260_122650;logic_name
=ena;version=1
Chromosome      ena      gene      2275488 2276288 .      -      .
ID=gene:L00C260_122660;biotype=protein_coding;description=singl
e-stranded DNA-binding protein;gene_id=L00C260_122660;logic_nam
e=ena;version=1
Chromosome      ena      gene      2276455 2277288 .      -      .
ID=gene:L00C260_122670;biotype=protein_coding;description=membr
ane protein;gene_id=L00C260_122670;logic_name=ena;version=1
Chromosome      ena      gene      2277304 2277648 .      -      .
ID=gene:L00C260_122680;biotype=protein_coding;description=ribon
uclease P;gene_id=L00C260_122680;logic_name=ena;version=1
Chromosome      ena      gene      2277719 2277853 .      -      .
ID=gene:L00C260_122690;Name=rpmH;biotype=protein_coding;descrip
tion=50S ribosomal protein L34;gene_id=L00C260_122690;logic_nam
e=ena;version=1
iu@bielinux[mapping_kiso3] grep "ID=gene" *.gff3 > uge.txt
iu@bielinux[mapping_kiso3] wc uge.txt                                [ 3:52午後 ]
 2262 24154 398272 uge.txt
iu@bielinux[mapping_kiso3] █                                        [ 3:52午後 ]
```

①

②

Name=を含む行数を表示

①grepコマンドで、②uge.txtというファイルに対して、③Name=という文字列を含む、④行数は、⑤457行。大分365個に近づいてきました。本当はここで365行になって一件落着のつもりでしたが…。Name=はあるがgenename部分が空っぽとか、同じ遺伝子名のものがあるとか…

```
iu@bielinux[~/Desktop/mac_share/mapping_kiso3]
Chromosome   ena      gene    2275488 2276288 .
ID=gene:L00C260_122660;biotype=protein_coding;description=e-stranded DNA-binding protein;gene_id=L00C260_122660;logic_name=ena;version=1
Chromosome   ena      gene    2276455 2277288 .
ID=gene:L00C260_122670;biotype=protein_coding;description=membrane protein;gene_id=L00C260_122670;logic_name=ena;version=1
Chromosome   ena      gene    2277304 2277648 .
ID=gene:L00C260_122680;biotype=protein_coding;description=ribonuclease P;gene_id=L00C260_122680;logic_name=ena;version=1
Chromosome   ena      gene    2277719 2277853 .
ID=gene:L00C260_122690;Name=rpmH;biotype=protein_coding;description=50S ribosomal protein L34;gene_id=L00C260_122690;logic_name=ena;version=1
iu@bielinux[mapping_kiso3] grep "ID=gene" *.gff3 > uge.txt
iu@bielinux[mapping_kiso3] wc uge.txt [ 3:52午後 ]
  2262  24154 398272 uge.txt
iu@bielinux[mapping_kiso3] grep -c "Name=" uge.txt [ 3:52午後 ]
457
iu@bielinux[mapping_kiso3] █ [ 3:56午後 ]
```

①

⑤

④

③

②