

RNA-seq演習

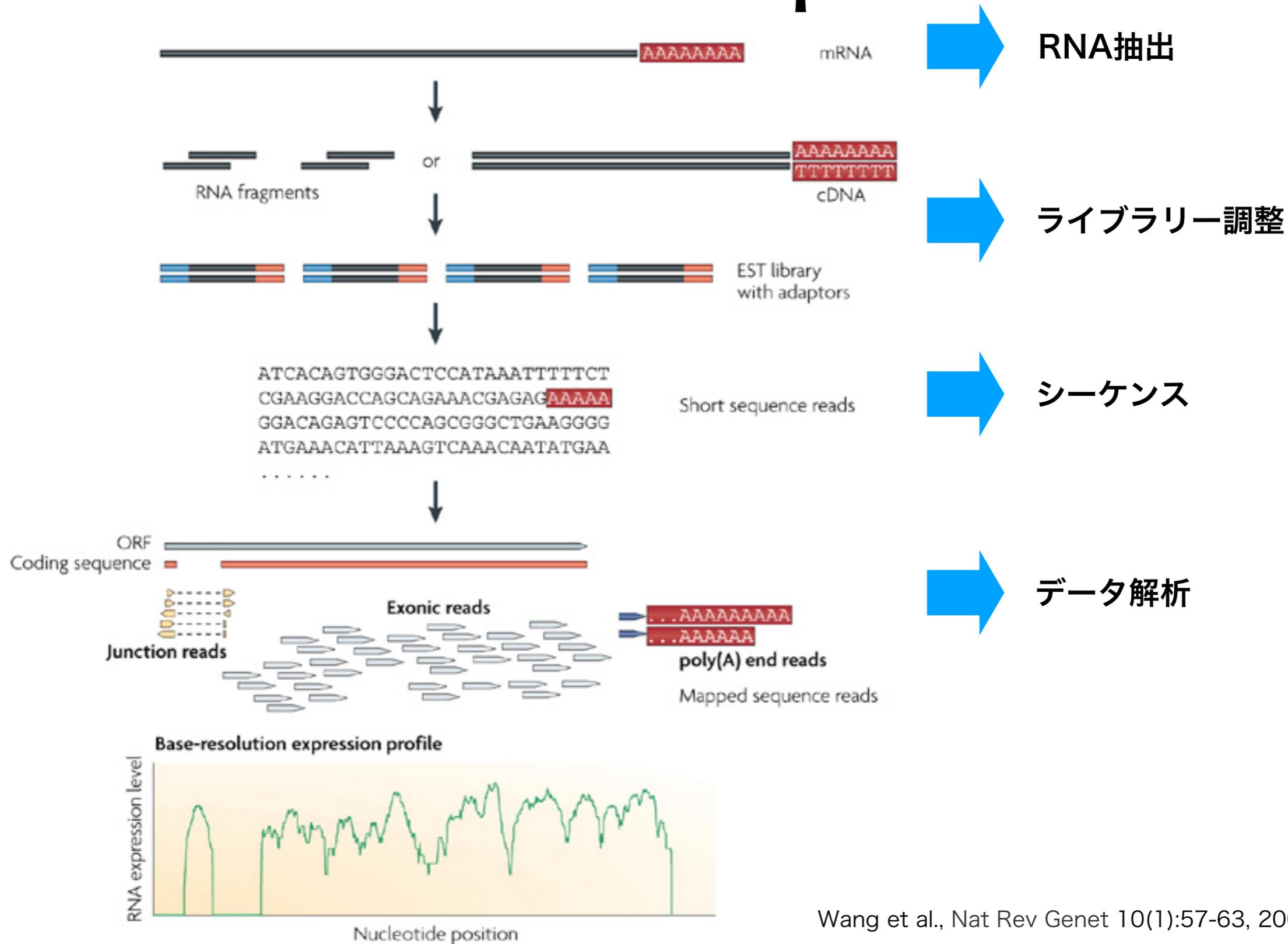
千葉大学真菌医学研究センター

高橋 弘喜

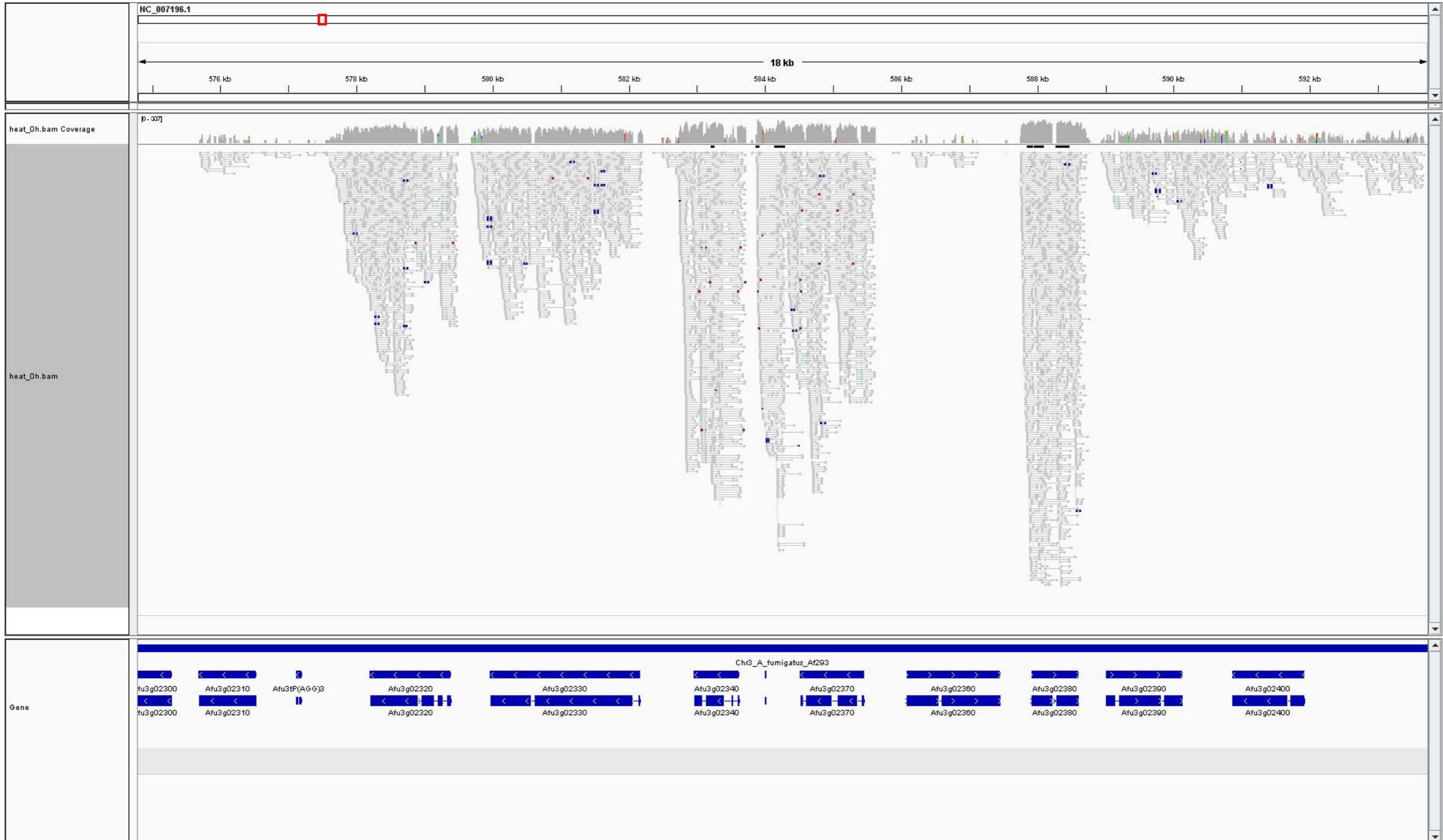
演習の内容

- テストデータを用いて、リードのマッピングから発現量の算出までの解析を遺伝研スパコン上でやってみる。

RNA-seq



結果例



本日の内容

リードのマッピング
(TopHat2)

発現量の算出
(cufflinks)

有意な発現変化を示す遺
伝子の探索
(cummeRbund)

スパコン上で計算

ローカルコンピュー
ター上で計算

スパコン使用方法(イメージ)

- ①ゲートウェイノード(gw.ddbj.nig.ac.jp)にログインする
- ②qloginを実行しインタラクティブノードにログインする
- ③qloginしたホストからジョブをUGEに投入する
- ④UGEは負荷の低いノードでジョブを実行する
- ⑤ジョブ実行結果をlustreのホームディレクトリに出力する
- ⑥ジョブ実行結果を確認する



fastq

```
1 @SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
2 GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
3 +SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
4 IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```

- 1行目： '@'に続き配列ID
- 2行目： 配列の文字列
- 3行目： '+'に続いて， 配列IDか， 改行
- 4行目： asciiコードで表現したquality value

fastq

33:! 34:" 35:# 36:\$ 37:% 38:& 39:' 40:(41:) 42:* 43:+ 44:,
45:- 46:. 47:/ 48:0 49:1 50:2 51:3 52:4 53:5 54:6 55:7 56:8
57:9 58:: 59:; 60:< 61:= 62:> 63:? 64:@ 65:A 66:B **67:C** 68:D
69:E 70:F **71:G** 72:H **73:I** 74:J 75:K 76:L 77:M 78:N 79:O 80:P
81:Q 82:R 83:S 84:T 85:U 86:V 87:W 88:X 89:Y 90:Z 91:[92:\
93:] 94:^ 95:_ 96:` 97:a 98:b 99:c 100:d 101:e 102:f 103:g
104:h 105:i 106:j 107:k 108:l 109:m 110:n 111:o 112:p 113:q
114:r 115:s 116:t 117:u 118:v 119:w 120:x 121:y 122:z 123:
{ 124:| 125:} 126:~

fastq

```
2 GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
4 IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```

S Q P (塩基の信頼性)

I: 73 → 40 → 1.0×10^{-4}
9: 57 → 24 → 3.9×10^{-3}
G: 71 → 38 → 1.6×10^{-4}
C: 67 → 34 → 4.0×10^{-4}

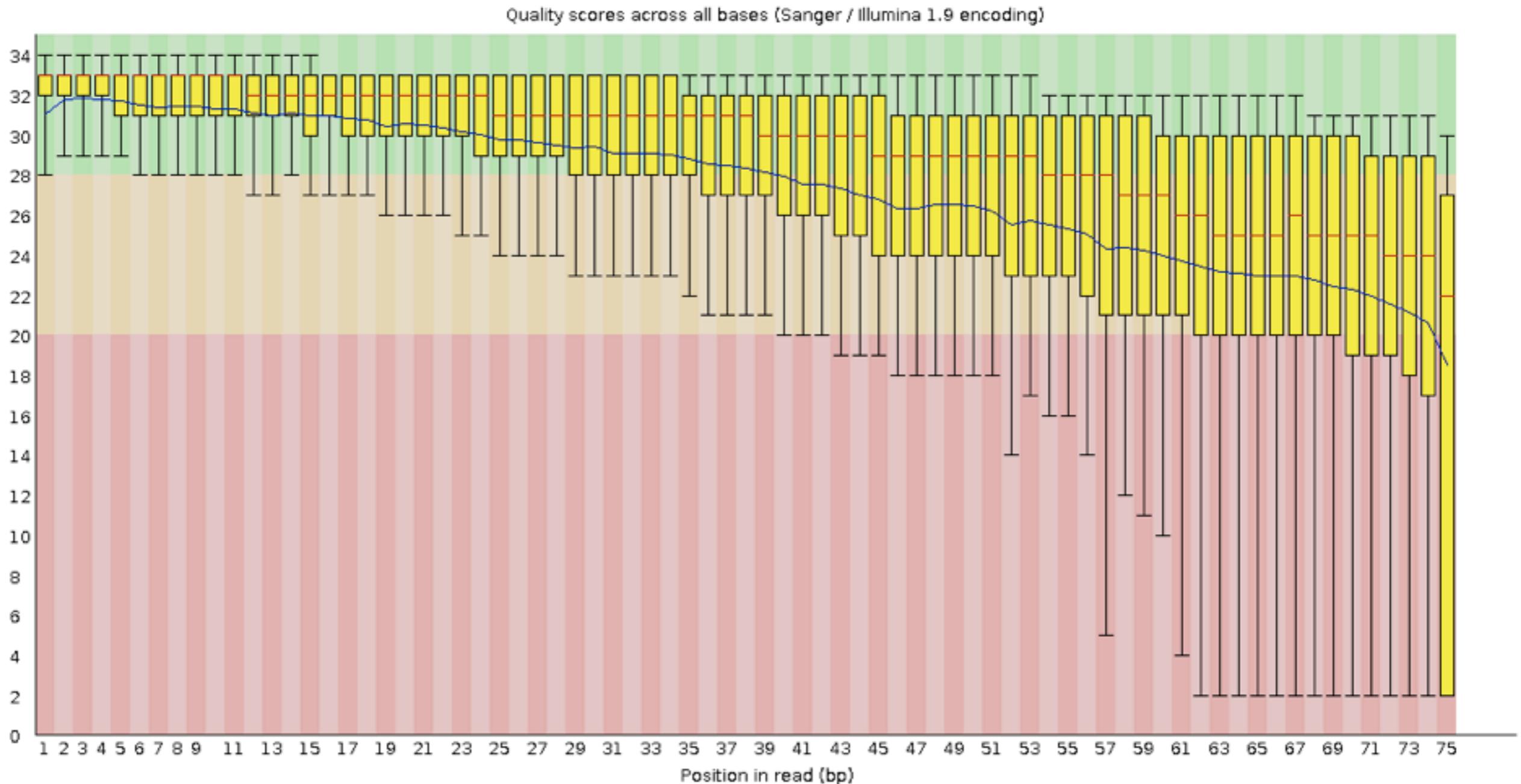
$$Q = -10 \log_{10} p$$

$$40 = -10 \log_{10} P$$

$$P = 10^{-4}$$

データクオリティ

シーケンスが進むにつれて信頼性が落ちていく



sam/bam

- NGSデータをリファレンスへマッピングすると、sam形式（テキストファイル）で出力される。（DNA-seq, RNA-seq, ChIP-seqなど）
- 11列のデータ（タブ区切り）で記載されている。

sam/bam

QNAME	Query template NAME
FLAG	bitwise FLAG
RNAME	References sequence NAME
POS	1- based leftmost mapping POSition
MAPQ	MAPping Quality
CIGAR	CIGAR String
RNEXT	Ref. name of the mate/next read
PNEXT	Position of the mate/next read
TLEN	observed Template LENgth
SEQ	segment SEQuence
QUAL	ASCII of Phred-scaled base QUALity+33

例

```
Coord      12345678901234 5678901234567890123456789012345
ref        AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1    TTAGATAAAGGATA*CTG
+r002      aaaAGATAA*GGATA
+r003      gcctaAGCTAA
+r004      ATAGCT.....TCAGC
-r003      ttagctTAGGC
-r001/2    CAGCGGCAT
```

```
@HD VN:1.5 SO:coordinate
```

```
@SQ SN:ref LN:45
```

```
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

sam \leftrightarrow bam

- bamファイルは、samファイルをバイナリーに変換した
もの。それにより、データサイズの圧縮、データへの高速
なアクセスが可能になる。
- sam \leftrightarrow bamは、samtoolsを使用することで相互に変換可
能。

samtools

- sam→bam
\$ samtools view -bS test.sam > test.bam
- bam→sam
\$ samtools view -h test.bam > test.sam
- ソート (test.sort.bamが作成)
\$ samtools sort test.bam test.sort
- indexの作成 (test.sort.bam.baiが作成)
\$ samtools index test.sort.bam

テストデータ

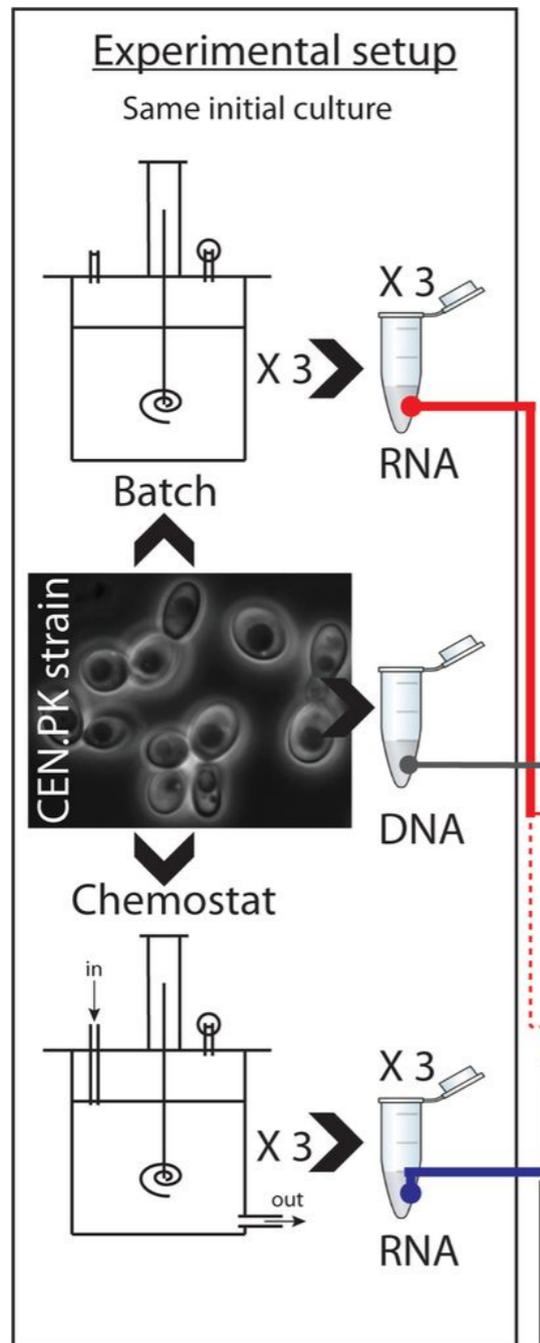
SRA051410  [FTP](#)

Submission Detail	
Alias	S. cerevisiae CENPK RNA-seq
Submission ID	
Submission Date	2012-04-04
Center Name	Chalmers University of Technology
Lab Name	

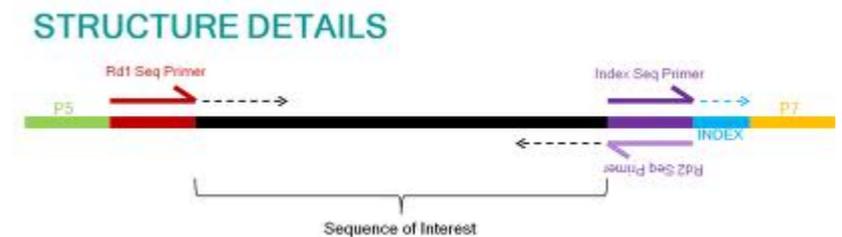
Navigation	
Study	SRP012047
Experiment	SRX135198  FASTQ  SRA
	SRX135710  FASTQ  SRA
	SRX135711  FASTQ  SRA
	SRX135712  FASTQ  SRA
	SRX135713  FASTQ  SRA
Sample	SRS307298
	SRS308058
Run	SRR453566  FASTQ  SRA
	SRR453567  FASTQ  SRA
	SRR453568  FASTQ  SRA
	SRR453569  FASTQ  SRA
	SRR453570  FASTQ  SRA
	SRR453571  FASTQ  SRA
	SRR453572  FASTQ  SRA
	SRR453573  FASTQ  SRA
	SRR453574  FASTQ  SRA
	SRR453575  FASTQ  SRA
	SRR453576  FASTQ  SRA
	SRR453577  FASTQ  SRA
	SRR453578  FASTQ  SRA

<http://trace.ddbj.nig.ac.jp/DRASearch/>

テストデータ



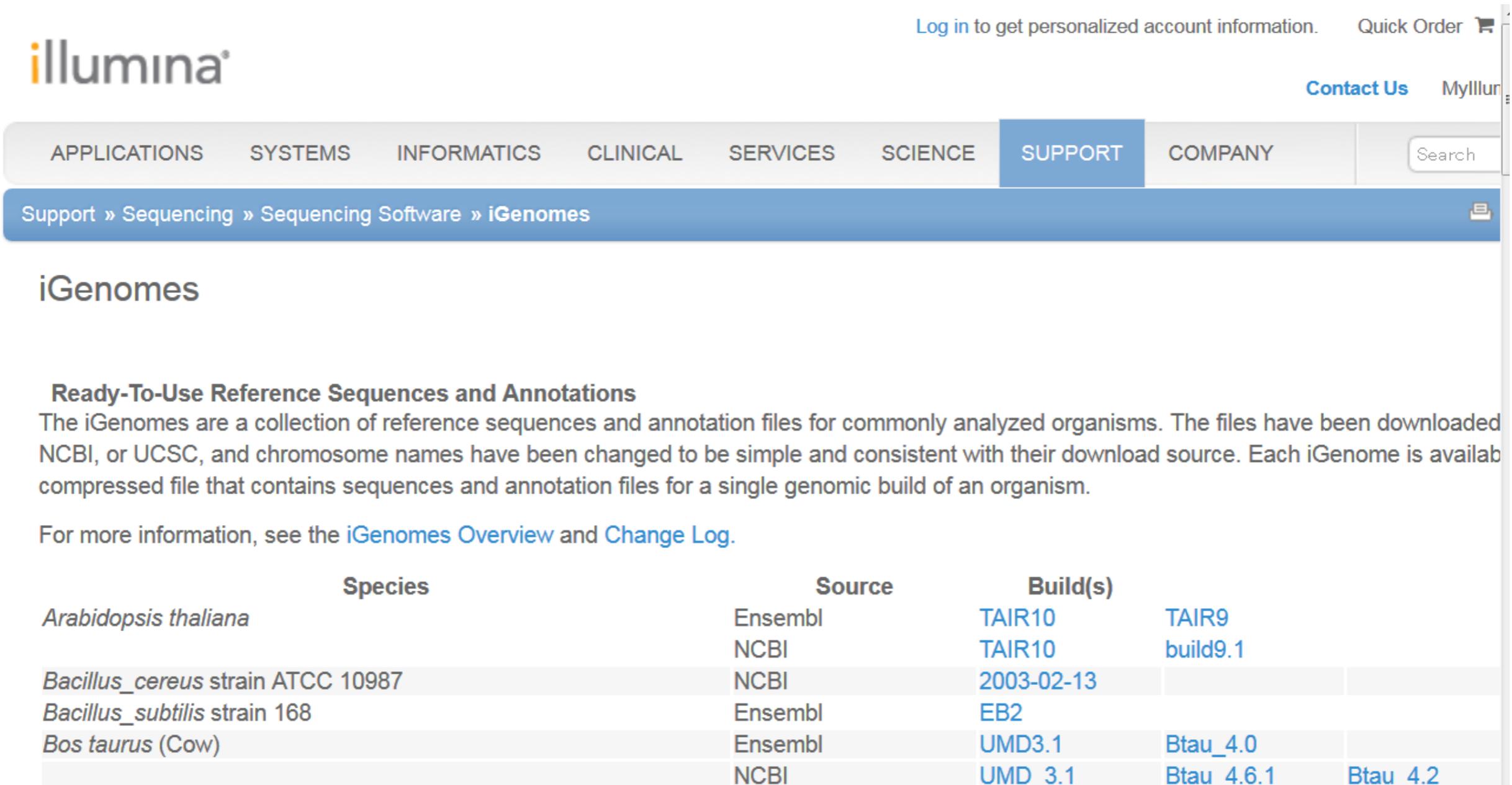
Illumina HiSeq2000
100bp PE



サンプル	Run	リード数
batch1	SRR453566	5,725,730
batch2	SRR453567	7,615,732
batch3	SRR453568	5,565,734
chemo1	SRR453569	4,032,514
chemo2	SRR453570	6,745,975
chemo3	SRR453571	6,163,396

アノテーションファイル

iGenomesから取得



illumina®

Log in to get personalized account information. Quick Order

Contact Us MyIllumina

APPLICATIONS SYSTEMS INFORMATICS CLINICAL SERVICES SCIENCE **SUPPORT** COMPANY

Support » Sequencing » Sequencing Software » iGenomes

iGenomes

Ready-To-Use Reference Sequences and Annotations

The iGenomes are a collection of reference sequences and annotation files for commonly analyzed organisms. The files have been downloaded from NCBI, or UCSC, and chromosome names have been changed to be simple and consistent with their download source. Each iGenome is available as a compressed file that contains sequences and annotation files for a single genomic build of an organism.

For more information, see the [iGenomes Overview](#) and [Change Log](#).

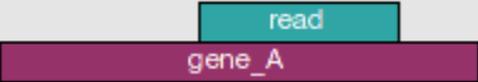
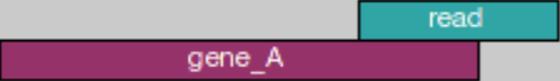
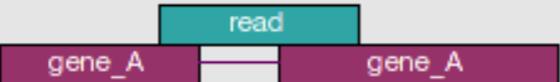
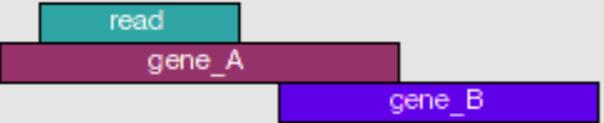
Species	Source	Build(s)
<i>Arabidopsis thaliana</i>	Ensembl	TAIR10 TAIR9
	NCBI	TAIR10 build9.1
<i>Bacillus_cereus</i> strain ATCC 10987	NCBI	2003-02-13
<i>Bacillus_subtilis</i> strain 168	Ensembl	EB2
<i>Bos taurus</i> (Cow)	Ensembl	UMD3.1 Btau_4.0
	NCBI	UMD_3.1 Btau_4.6.1 Btau_4.2

アノテーションファイル

- iGenomesでは、マッピングに必要なインデックスファイル、アノテーションファイルが用意されている。
 - BWAIndex
 - Bowtie2Index
 - BowtieIndex
 - genes.gtf

その他

どのリードを転写物由来とするか

	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

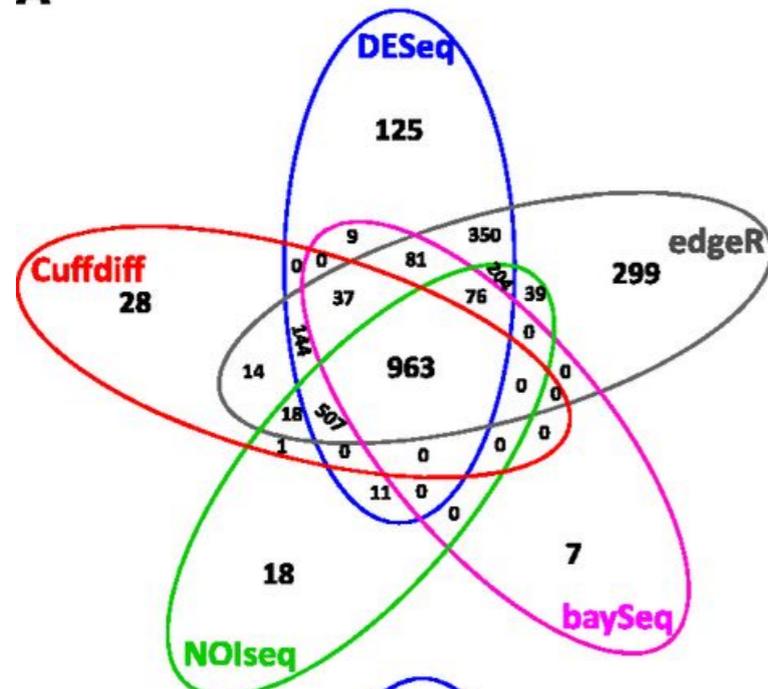
その他

- 解析手法（マッピング含む）による違い

Table S4 Number of DGE (Q-values < 10e-5) from different methods (for microarray the number DGE is 1603)

Method	Gsnap	N.Gsnap	Stampy	N.Stampy	TopHat	N.TopHat	De novo
Cuffdiff	2061	2172	1712	1741	1671	1726	1623
DESeq	2690	2731	2507	2503	2412	2432	2197
edgeR	3087	3161	2732	2742	2649	2673	2385
baySeq	1785	1807	1173	1198	1092	1133	1175
NOISeq*	2097	2070	1837	1784	1804	1754	1595

A



その他

- 別の解析手法
 - HTSeq
 - DESeq2
 - edgeR
- Viewer
 - IGV
- データQC
 - FastQC