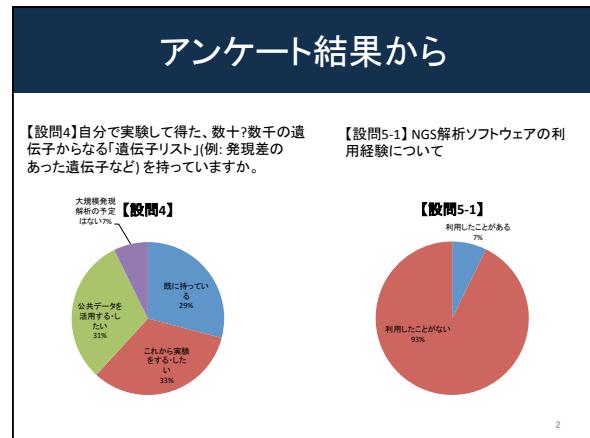


統合データベース講習会

In 群馬大学

Masaki Suimye Morioka



本日の講習内容

1. NGSについての簡単な説明 (illumina型中心)
 - （スライドで約1時間ほどの予定）
2. NGS解析のチュートリアル (ChIP-seq)

利用した画像について

Togo picture gallery by DBCLS is licensed under a Creative Commons Attribution 2.1 Japan license (c)

参考にさせていただいたweb site

DBCLS 大田さんの「GalaxyによるNGS解析」(<https://github.com/inutano/training/tree/master/ajaxs-advanced-01>)
株式会社ジナリス、竹田さんの「いまさら聞けないNGS超!入門」<http://www.slideshare.net/Genaris/ngs-46348501>

HiSeq X Ten System

1,000 genome sequencer, illumina HiSeq X ten

• Ten HiSeq X per one system.
• \$1,000/genome (including chemistry, personnel costs)

Jan. 2014. Press in illumina seminar

実物とは異なります

Togo picture gallery by DBCLS is licensed under a Creative Commons Attribution 2.1 Japan license (c)

ゲノムの塩基数とカバー率

Remember this....

1 billion = 100Mb

Byte Base!

A human genome: 30 billion base = 3Gb

Human whole genome sequencing:
commonly required over x20-30 coverage.

3Gb x 30 = 90 Gb / human

Performances of HiSeq X ten*

HiSeq 2500 (Rapid mode, v3)

- TP: 100Gb x3
- Reads: 最大10億
- Time: 27時間 x3

High coverage mode
11days...

320 peoples

HiSeq X ten system

- TP: 1.6Tb-1.8Tb
- Reads: maximum 60 billion
- Time: 3日間

1 week after...

* Based on Illumina web site information

HiSeq SBS Kit v4 and HiSeq Cluster Kit v4

HiSeq SBS Kit v4 (May 2014)

- 167Gb/Day
- 1Tb / 6 Day, 125 bp, max 40 billion reads (max 500Gb / flow-cell)
- Cluster density is increased 33% (High coverage mode)

From illumina web site information

7

NGSの最近の動向

より長く塩基配列を決定したい

Pacbio: SMRT-Seq
Single molecule real time sequencing

Each of the four nucleotides is labeled with a different colored fluorophore

Oxford NANOPORE

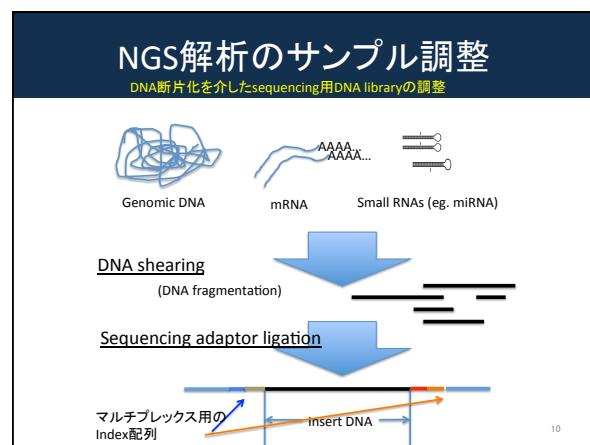
What could MinION be used for? Pathogen surveillance, metagenomic variant detection, selective sequencing and much more!

10X genomics

The compact, open Chromium™ Controller has been designed to enable the analysis of up to 100,000 or more samples in parallel. By using the Chromium™ Controller, users can analyze up to 100,000 samples in parallel, with no need for any of the traditional single-cell sequencing instruments such as flow cytometers or microfluidic devices.

The INSTRUMENT Compact, sleek, efficient.

各社のwebサイトより転載させていただきました。



Illumina社のonline学習

Online Courses

- Sequencing
- Array
- Expert Video Type

Online Sequencing Courses

- Training Videos

ONLINE COURSE TITLE DESCRIPTION LENGTH

- Sequencing: Illumina Technology This course provides a general overview of the Illumina sequencing platform, including the function of nucleic acids at the completion of a sequencing run. Languages: Spanish, Chinese (Simplified), Japanese 30 min
- Sequencing: Illumina's Dual Indexing Strategy The goal of dual index sequencing is to increase the multiplex level of a sequencing run so that more samples can be sequenced on the same flow cell. Languages: Spanish, Chinese (Simplified), Japanese 10 min
- HiSeq 1500/2500: Best Practices This course covers best practices for your HiSeq 1500/2500. It includes tips for sample loading, instrument set-up and washes, best practices for performing a HiSeq sequencing run, BaseSpace Sequence Hub and data compression options for your sequencing runs, how to store your data cars, and how to prepare the instrument for cleaning or shutdown. 25 min
- HiSeq: Rapid Run Mode This course describes the Rapid Run Mode available for HiSeq 2500 or HiSeq 1500. It identifies when it is needed to perform a Rapid Run, describes the options for clustering the 2-lane flow cell, and lists ways to reduce data storage requirements for your HiSeq runs. 15 min

<http://support.illumina.com/training/online-courses/sequencing.html>

11

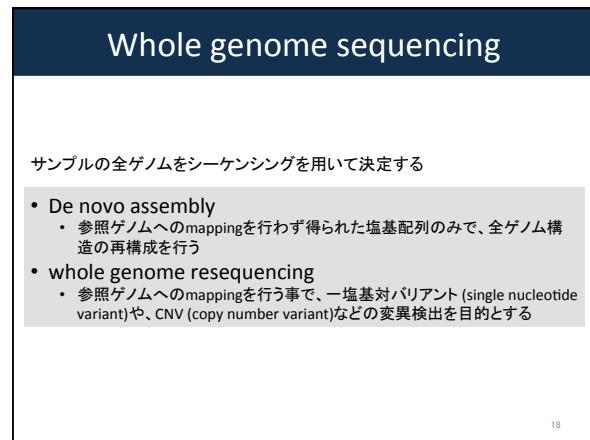
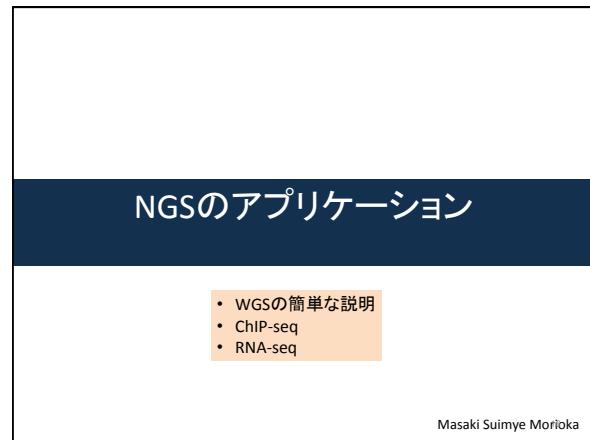
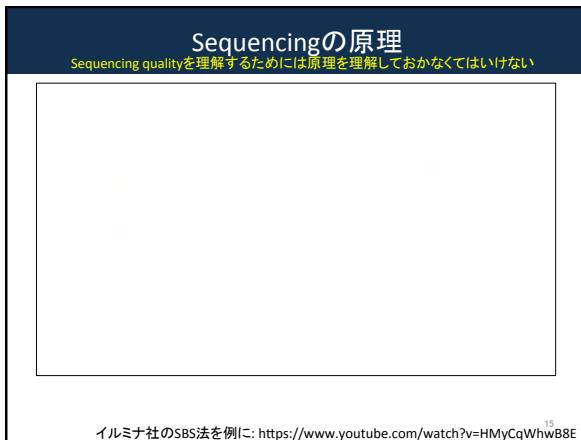
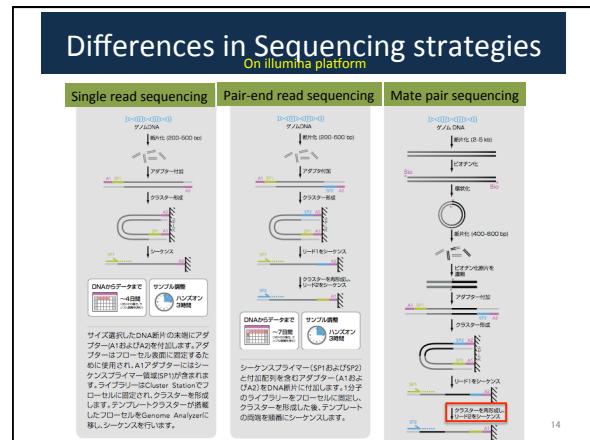
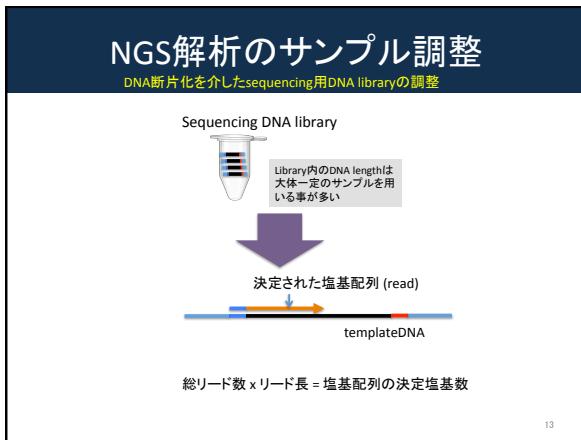
Illumina社のonline学習

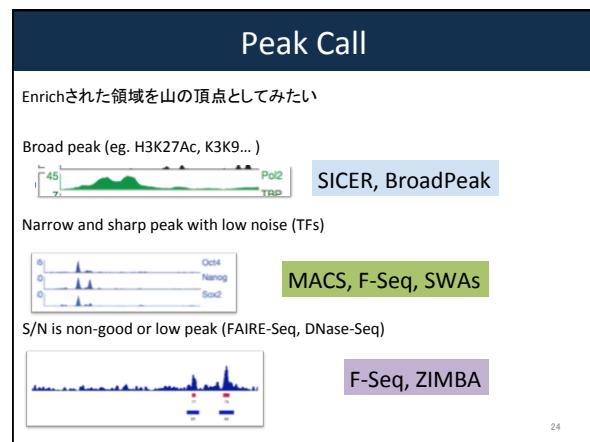
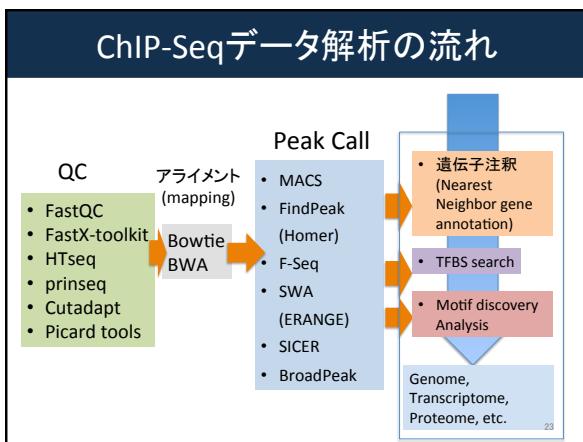
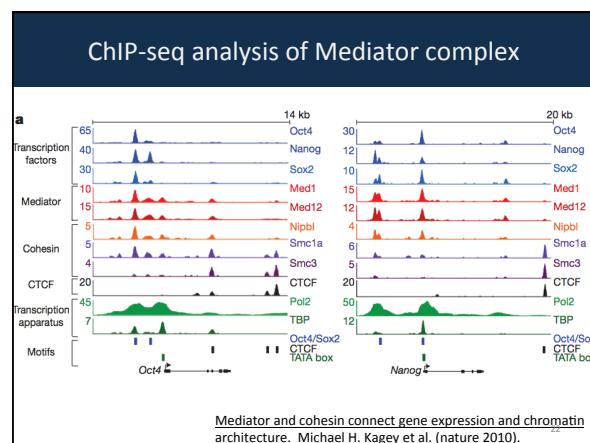
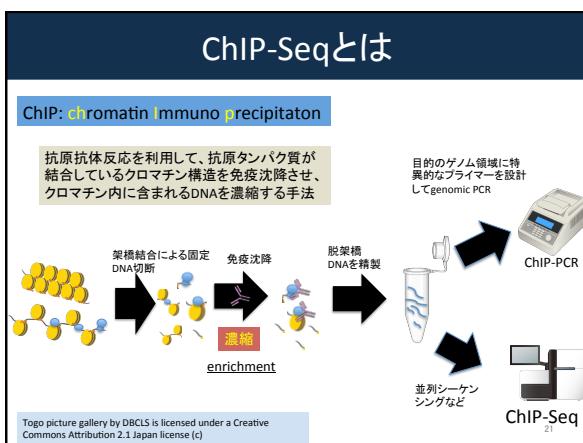
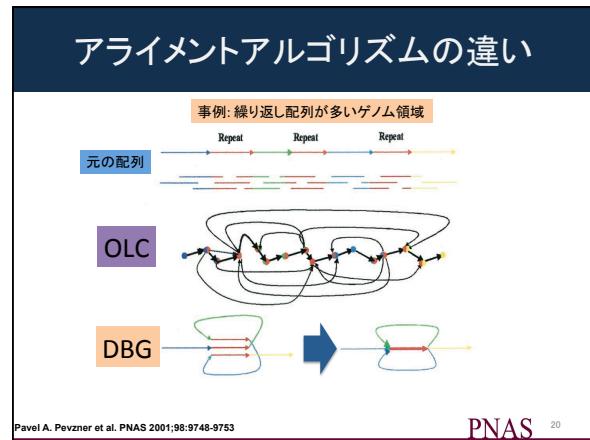
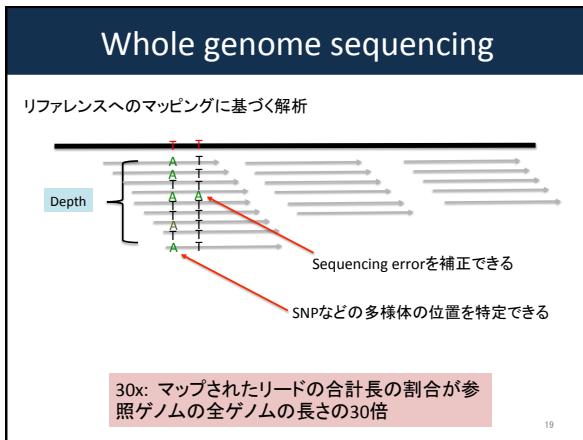
SBS (Sequencing by Synthesis) の概要

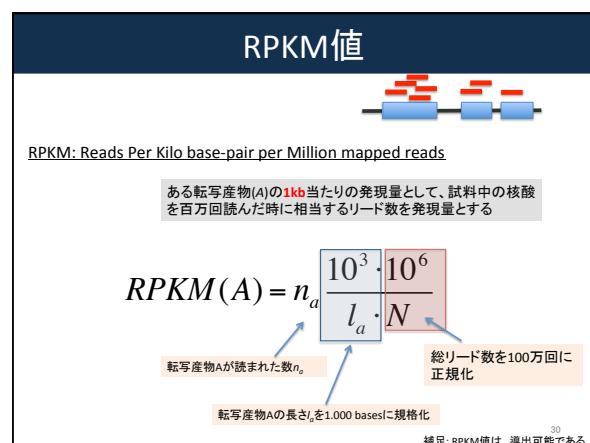
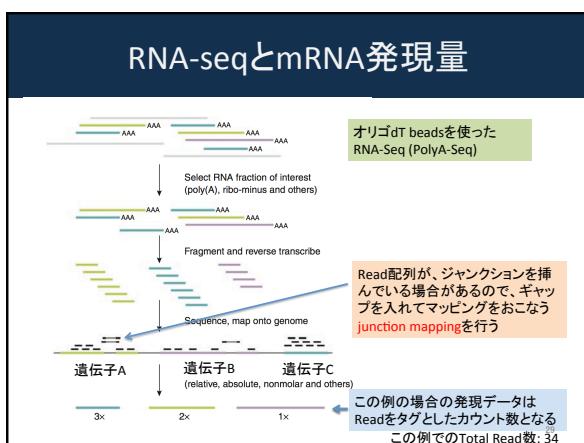
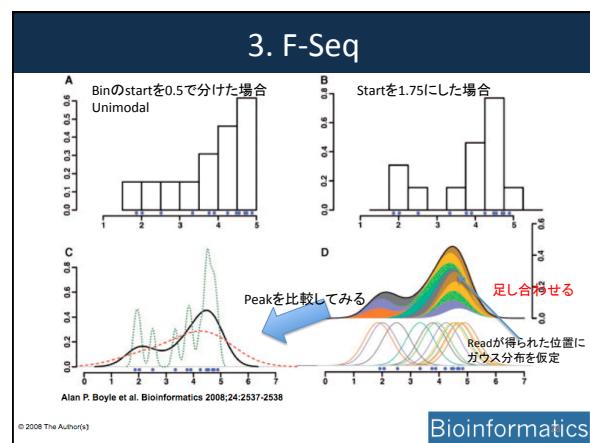
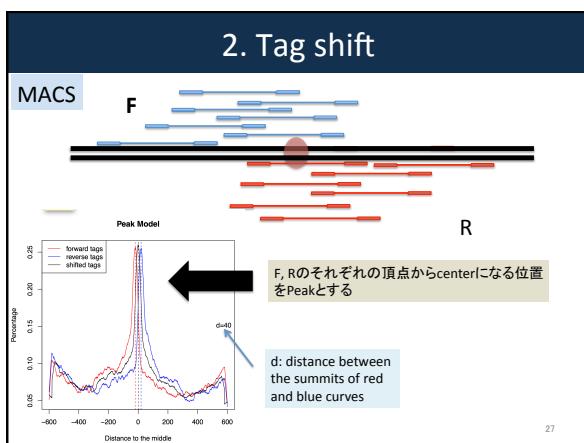
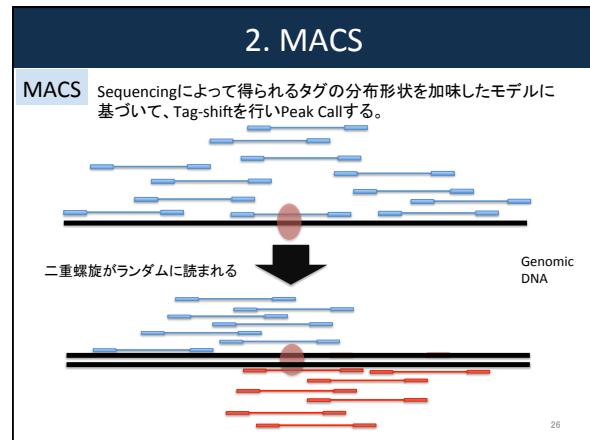
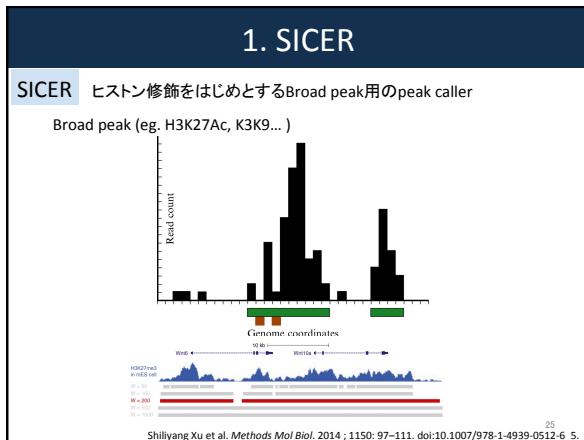
Illuminaのユニークなシーケンスアプライバーにより、どんな生物、組織、または分子構造でもDNAレベルあるいはRNAレベルで研究できるようになりました。Illuminaでは、このプロセスをSequencing by Synthesis、略してSBSと呼んでいます。この例では、DNAが使用されています。詳細については、各ステップをクリックして下さい。

<http://support.illumina.com/training/online-courses/sequencing.html>

ページ 5 / 9







FPKM値

FPKM: Fragments Per Kilo base-pair per Million mapped reads
ペアエンドのデータに対応したRPKM値と考えてください

ペアエンドのデータを1つのフラグメント(DNAライブリ作成時のインサートフラグメント)として捉えてカウントする

31

NGSデータの特徴のまとめ

Introduction of NGS analysis

1. 塩基配列である
2. Read(読み取った配列)が体積した山になる
3. 試料のプールからシーケンサでランダムに読まれていく
4. シーケンシングされるread数は、決まったread数にはならない(試料ごとにTotalのread数が異なる)

1. アライメントする、アッセンブリ	1. Genome mapping, Genome assembler
2. 山をみつける	2. Peak Caller
3. ランダムサンプリングされていることを前提にする	3. Poisson Distribution
4. 全体のリード数を整えてから、NGSデータ間の比較	4. RPKM FPKM

32

NGSのファイル形式

- Sequencing quality
- Mapping quality

33

NGS解析から得られるデータ

Raw data

例	Illumina型: 画像 PacBio型: 動画
---	------------------------------

配列データ

例	Fastq形式 (.fastq) HDF5形式 (.h5)
---	----------------------------------

シーケンサから出力された情報をもとに塩基が決定され、その塩基配列情報の品質評価値とともに文字情報として出力される。

シーケンサから出力された塩基情報のもともらしさ
Phred Scale

34

Fastqファイル (sanger形式)

head sample.fastq
↓ ファイルの中身を先頭から数行抜き出すunixコマンド

```
@SRR445816.18948743 HWUSI-EAS366_145:3:39:5673:1812/1
ACTCCTGCTTCAGGTGATCCATCCGGCTCAGGCCCTCT
+
BDGGG@GDGGBG@DDE@ABFGGGB<CFEAE<:???:=
@SRR445816.28703832 HWUSI-EAS366_145:3:61:8637:6699/1
ACTGACTCAAATGTTAACCTCCTTGGCAACACTCT
+
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
@SRR445816.43211059 HWUSI-EAS366_145:3:94:3572:7097/1
TAAGCCCCCTCTTAGGATTATAACCTCATCACT
```

35

Phred Score

General NGS data format

$$Q_{phred} = -10 \log_{10}(P)$$

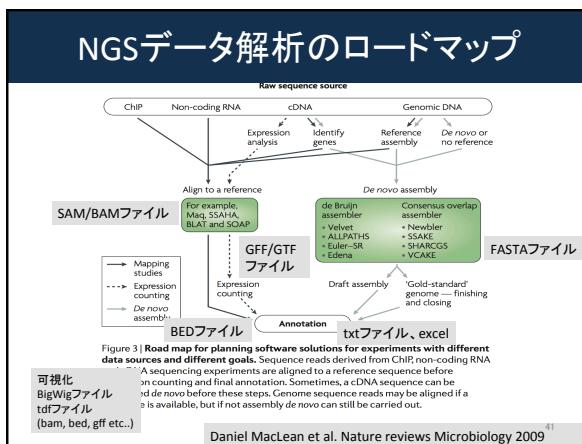
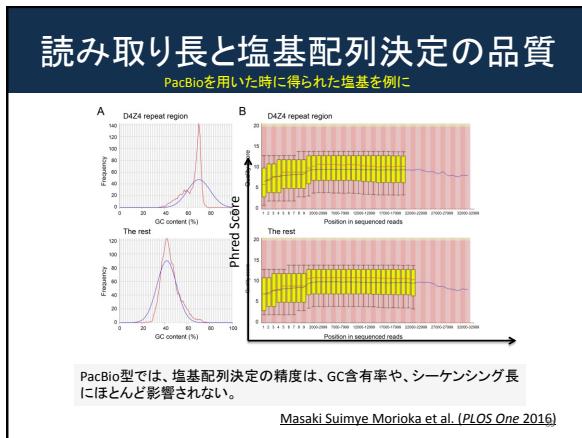
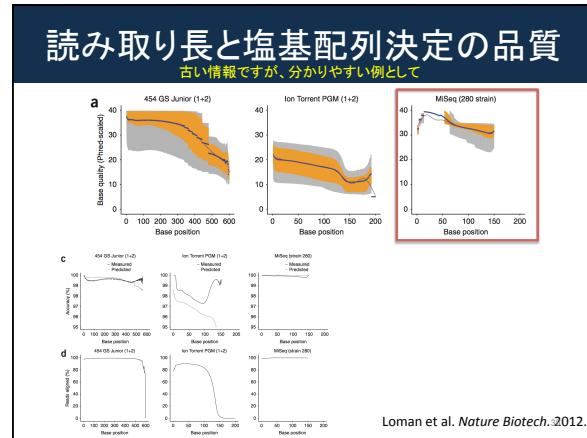
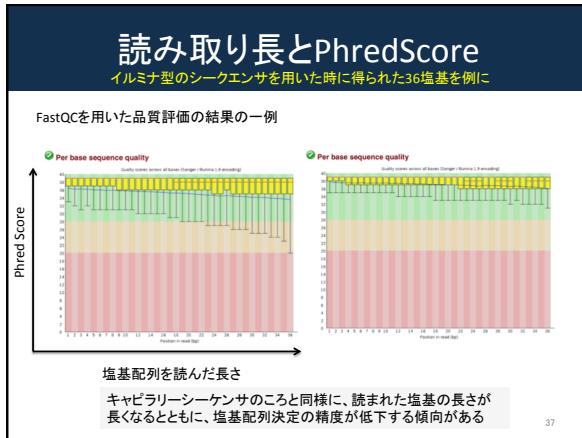
Phred Score	塩基が違う確率(P)	正確性
10	1/10	90%
20	1/100	99%
30	1/1000	99.9%
40	1/10000	99.99%

Phred Scoreは、アスキーコードで書かれており、1文字に対応する数字が記載される。
例えば、文字'a'は97に相当する。
サンガーフォーマットのクオリティ値は、アスキーコードの数字から33引いた値(文字'a'は65に相当する)

```
HWUSI-EAS1869:80:FC665TUAXX:4:1:1329:1035 1:N:0:
GTCCTTAAGTACTGTCATAGGGCTCTGTATCCA
GG>>EEEEEEBEG@GGDDEEEED=>GGGBBBGEGGGG<
```

大文字"G"は、71なので、71 - 33 = 38
文字">"は、62なので、62 - 33 = 29

36



SAM/BAMファイル

SAM (Sequence Alignment/Map) ファイルは、シーケンサから得られたリード配列がゲノムDNAのどの領域に由来するのかを調べるために、既に読まれたゲノム配列をリファレンスとしたアライメントを行った時に得られるアライメント結果のファイル。

```
iu@bio[tutorial150806] head -n5 sample.sam [12:12午前]
@HD VN:1.0 SO:unsorted
@SQ SN:chr17 LN:81195210
@PG ID:bowtie2 PN:bowtie2 VN:2.2.4 CL:"/usr/bin/../lib/bowtie2/bin/bowtie2-align-s --wr
SRR44516.15448471 0 chr17 5102645 1 36M * 0 CTCCCAAGCTGAGATGCAATCTGGTC
SRR44516.28783832 0 chr17 18949619 1 36M * 0 ACTGACTCAAATGTTAACATCCTTTGGC
iu@bio[tutorial150806]
```

例では、36塩基基準で
がmatchという意味

Header

Op	HD	Informations
R	0	alignment switch (can be a sequence pack < v mismatch)
C	1	clipping start from the reference
D	2	clipping end from the reference
F	4	soft clipping (clipped sequence present in \$2)
H	5	padding character (padding character for the \$2, padding taken from packed reference)
I	6	sequence length
P	8	sequence rank

Read Name
SAM tag
chromosome (if read is has no alignment, there will be a "*" here)
position (1-based index, "left end of read")
MAPQ (マッピングクオリティ, 0=non-unique, >10 probably unique)
CIGAR (カイガル) (read length, intrepidリードの長さにあらかじめなどに利用)
junctionCountなどによる情報
Name of mate (mate pair information)
Position of mate (mate pair information)
Template length (always zero for me)

http://samtools.github.io/hts-specs/SAMv1.pdf

42

GFFファイル

リファレンスにアライメントされていないリードの情報も含む

General Feature Format (GFF)

ゲノムの特徴について一行で記載されているファイル。現行は、GFF3。GTF (General Transfer Format)は、GFF2と同一。

- Fields
- Track lines
- More information

大きく3つの構成要素

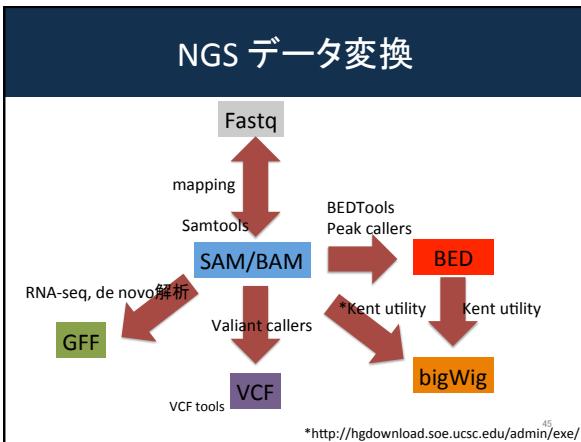
Sample GTF output from Ensembl data dump:

```
1 transcripted unprocessed_pseudogene gene    11869 14409 . + . gene_id "ENSG000000223972"; gene_name "DOX11L1"
1 processed _transcript transcript 11869 14409 . + . gene_id "ENSG000000223972"; transcript_id "ENSED
```

Sample GFF output from Ensembl export:

X	Ensembl Repeat	2419108	2419128	42	.	.	.	hid-trif; hid=trif-1; hid=trif-2;
X	Ensembl Repeat	2419108	2419128	2502	.	.	.	hid=AllRepeat; hid=repeat; hid=303
X	Ensembl Repeat	2419108	2419128	2502	.	.	.	hid=dust; hid=trif; hid=2419108; hid=2419128
X	Ensembl Pred.trans.	2416676	2418760	450..19	-	2	.	genomic::GENSCAN0000001935
X	Ensembl Variation	2431425	2431425	42	+	.	.	.
X	Ensembl Variation	2413803	2413805	42	+	.	.	.

<http://www.ensembl.org/info/website/upload/gff.html>



解析環境の選択

Illumina BaseSpace

<https://basespace.illumina.com/apps/>

BaseSpace SEQUENCE HUB DASHBOARD PREP RUNS PROJECTS APPS PUBLIC DATA Q ? 🌐 MASAKI MISHO... ⓘ | illumina

Applications



16S Metagenomics
Illumina, Inc.



Amplicon DS
Illumina, Inc.



BWA Aligner
BaseSpace Labs



BWA Enrichment
Illumina, Inc.



RNA Seq
Illumina, Inc.



Small RNA
Illumina, Inc.



ChIP-Seq (1)



Differential Expression (1)



Gene Fusion Detection (1)



HiFiA (1)



Metagenomics (12)



Methyl Seq (8)



Proteomics (10)



Quality (9)



Resequencing (23)



RNA Seq (14)



Small RNA (7)

Search Apps

Categories

ChIP-Seq (1)	De Novo Assembly (7)
Differential Expression (1)	Gene Fusion Detection (1)
HiFiA (1)	Metagenomics (12)
Methyl Seq (8)	Proteomics (10)
Quality (9)	Resequencing (23)
RNA Seq (14)	Small RNA (7)

GalaxyからChIP-seq専用に作成されたnebula (<http://nebula.curie.fr>)

The screenshot shows the Nebula web interface. On the left, there's a sidebar with links like 'UPLOAD YOUR DATA', 'Get Data', 'FILES MANIPULATION', 'Filter and Sort', 'Compare Enrichments', 'NCBI BLAST', 'NCBI BLAST+ Discovery', 'NCBI BLAST+ Tools', 'NCBI BLAST+ Training', 'NCBI Park Annotation', 'NCBI ASAT Tools', 'NCBI JBrowse Tools'. The main area has a purple header 'Read carefully before starting' with a note about server use policy. Below it is a file list table:

ヒストリー	オプション	Using OK
tutorial	455.3 Mb	
20 MACS on data 37.1mm	0 / 2	
REPORTS		
18 MACS on data 37.1mm	0 / 2	
16 Item to Rel on data 37.1mm	0 / 2	
17 seq7ch19.3mm	0 / 2	
18 seq7ch19.3mm	0 / 2	

At the bottom right of the main area, it says '49'.

ChIP-seqデータをとにかく可視化

http://www.devbio.med.kyushu-u.ac.jp/sra_tailor/

Srataller

[presenting] 現存のChIP-seqデータを全て可視化する
本日のセミナー、2014年12月22日にはライカサイエンス結合データベースセンター（DBCLS）にて行われた「第3回 助教（九州大学大9号）セミナー」を主催者よりご了承を得て、現存のChIP-seqデータを全て可視化する」をお送りします。ChIP-seqデータの可視化ツールとして、Asperaを利用しているので、proxyのsettingを忘れずに

YouTubeはこちで

http://togotv.dbcls.jp/20150106.html

統合解析環境UGENE

Unipro U GENE

Home Documentation Learn Community Downloads Support Contacts Donate

Unipro UGENE 1.22 last update

UGENE is free open-source cross-platform bioinformatics software. It works perfectly on Windows, Mac OS and Linux and requires only a few clicks to install.

[Download PDF](#) [Download UGNE](#)

Cite Us Support and Services

NGS解析には別途、BWA, Tuxedo Toolsなどのinstallが必要であるが、pipelineを作れる