バイオインフォマティクスへの招待 ～高速シーケンサーと RNA-Seq～（講義・実習）
Introduction of Bioinformatics - Next Generation Sequencer and RNA-seq
（Lecture/Practice in Japanese）

門田 陽介・芦原 貴司（情報総合センター・医療情報部）
Yosuke Kadota, Takashi Ashihara（Information Technology and Management Center）

この講義・実習では、高速シーケンサーによる塩基配列の読み取りの原理を概説し、得られた膨大な塩基配列データから遺伝子の発現解析（RNA-Seq）を行う手法について、実際にマルチメディアセンターの PC を各自が使ってハンズオンで学ぶ。

We will show you outline the principles of sequence reading by high-speed sequencers, and let you know how to analysis gene expression (RNA-Seq) from the huge amount of sequence data obtained from high-speed sequencers, by using the PCs in the Multimedia Center by yourselves.

# バイオインフォマティクスへの招待

～高速シーケンサーとRNA-Seq～

国立大学法人
滋賀医科大学
SHIGA UNIVERSITY OF MEDICAL SCIENCE

滋賀医科大学は
開学50周年

50th
ANNIVERSARY

---

## What is RNAseq?

- Sequencing of RNA (mainly mRNA) by High-speed sequencer

- Quantitative analysis of gene expression

- Compare with qPCR
  - gene expression can be analyzed comprehensively.

  - Splicing analysis is also available

---

## Examples of analysis

- Comparison of differences in gene expression in normal and cancerous tissues
  - Identification of up- or down-regulated genes and elucidation of disease mechanisms

- Comparison of splicing differences between normal and cancerous tissues
  - Identification of used variants and elucidation of disease mechanisms

---

## Flow of RNAseq analysis

**Sample Preparation** → **High-speed sequencing** → **Expression Analysis**

- **Sample Preparation**
  - mRNA extraction
  - Reverse transcription to cDNA

- **High-speed sequencing**
  - Fastq file
    - Sequencing data

- **Expression Analysis**
  - quality check
  - Mapping
  - Expression Analysis
  - visualization

---

## Preparation of the environment required for RNAseq

- Installation of software required for analysis
  - Free UNIX-based software
    - Linux or Mac

---

## WSL(Windows Subsystem for Linux)

- Windows Subsystem for Linux (WSL) is a feature of Windows that allows you to run a Linux environment on your Windows machine, without the need for a separate virtual machine or dual booting. WSL is designed to provide a seamless and productive experience for developers who want to use both Windows and Linux at the same time.
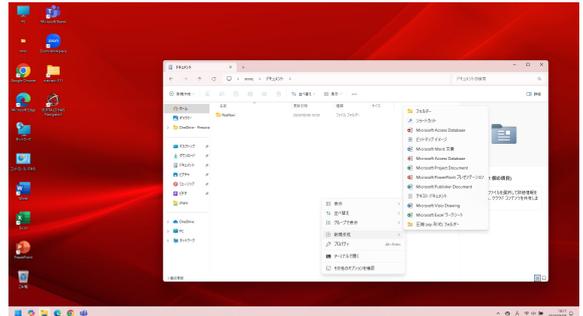
ubuntu

**Microsoft Learn**
https://learn.microsoft.com/en-us/windows/wsl/about
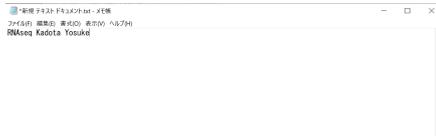
## CUI (Character User Interface)

- A method of operating a computer through character input using a keyboard, which allows the user to interact with the computer in a character-based manner.

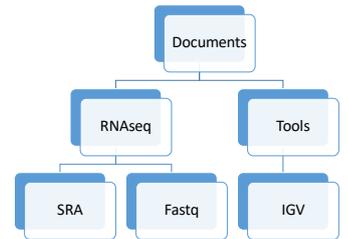## GUI（Graphical User Interface）

## GUIとは



Newtext.txt

## Let's operate by CUI

- Directory (folder) operations
  - cd : Moving Directories
  - pwd : Show current directory
  - ls : Listing of folders and files
  - mkdir : create folder

## Pwd : Show current directory

$ pwd

```
kadota@Kadota-Lenovo-WS2:~$ pwd
/home/kadota
(base) kadota@Kadota-Lenovo-WS2:~$
```

## cd

$ cd ..

```
(base) kadota@Kadota-Lenovo-WS2:~$ cd ..
(base) kadota@Kadota-Lenovo-WS2:/home$ pwd
/home
```
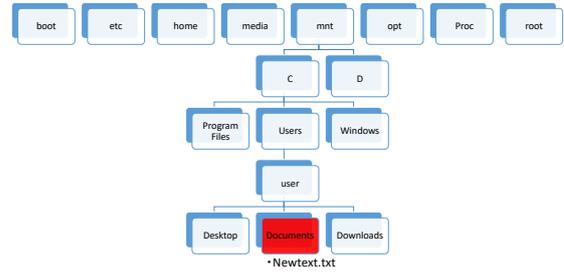
## Ls : Listing of folders and files

$ ls

```
(base) kadota@Kadota-Lenovo-WS2:/home$ cd ..
(base) kadota@Kadota-Lenovo-WS2:/$ pwd
/
(base) kadota@Kadota-Lenovo-WS2:/$ ls
bin    etc    lib    libx32      mnt    root    snap    tmp
boot   home   lib32  lost+found  opt    run     srv     usr
dev    init   lib64  media       proc   sbin    sys     var
(base) kadota@Kadota-Lenovo-WS2:/$ _
```

$ ls -l

$ ll

## Let's operate by CUI

## ls

```
(base) kadota@Kadota-Lenovo-WS2:~/Documents$ ls
Newtext.txt
```

## Cat : View all of the text files

$ cat Newtext.txt

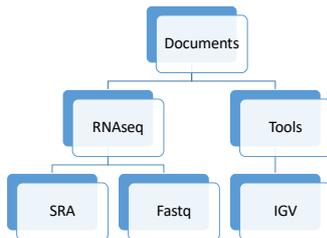## mkdir : create folder
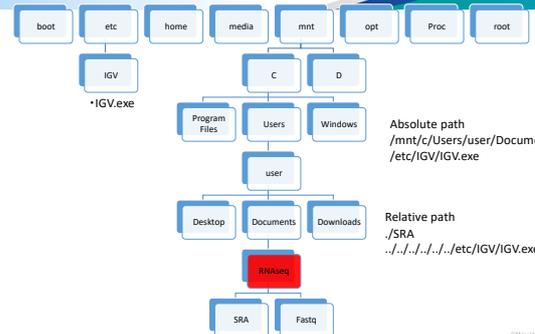
$ mkdir RNAseq

```
~/Documents$ mkdir RNAseq
~/Documents$ mkdir Tools
~/Documents$ cd RNAseq/
~/Documents/RNAseq$ mkdir SRA Fastq
~/Documents/RNAseq$ cd ..
~/Documents$ cd Tools/
~/Documents/Tools$ mkdir IGV
```

## absolute path and relative path



Absolute path
/mnt/c/Users/user/Documents/RNAseq/SRA
/etc/IGV/IGV.exe

Relative path
./SRA
../../../../../../etc/IGV/IGV.exe

# CUIで操作してみよう

- ファイルの操作
  - mv : Moving files and folders
  - cp : Copying files and folders
  - cat : View all of the text files
  - less : View some of the text files

---

# absolute path and relative path



boot | etc | home | media | mnt | opt | Proc | root

IGV

・IGV.exe

C | D

Program Files | Users | Windows

user

Desktop | Documents | Downloads

RNAseq

SRA | Fastq

Absolute path
/mnt/c/Users/user/Documents/RNAseq/SRA
/etc/IGV/IGV.exe

Relative path
./SRA
../../../../../../etc/IGV/IGV.exe

---

# less : View some of the text files

$ less ../Newtext.txt

---

# RNAseq flowchart

| | Analysis Software | File extension to be used |
|---|---|---|
| Download FASTQ file | SRA toolkit | sra, dra, fastq |
| Mapping | HISAT2、STAR | sam, bam |
| Expression Comparison | FeatureCount | txt, xlsx |
| Visualization | IGV | bam.bai |

---

# mapping

GCTATGAAAGGTT    CCGGAATTGGAC

ACGGTAACCGTAGCTATGAAAGGTT CCGTAAGTACGTTTAACCGGAATTGGACCAGTCAGTC

**Genome sequence**

---

# mapping



Exon1 | Exon2 | Exon3 | Exon4

**Genome sequence from ensembl**

## RNAseq flowchart

| | Analysis Software | File extension to be used |
|---|---|---|
| Download FASTQ file | SRA toolkit | sra, dra, fastq |
| Mapping | HISAT2、STAR | sam, bam |
| Expression Comparison | FeatureCount | txt, xlsx |
| Visualization | IGV | bam.bai |

---

## (filename) extension

- SRA
  - Archive of FASTQ files
- FASTQ
  - Sequence data (text format)
- SAM
  - Mapping data (text format)
- BAM
  - Mapping data (binary format)
- GZ
  - compressed file

---

## Updated package list

$ sudo apt update

---

## Upgrade Package

$ sudo apt upgrade

---

## RNAseq flowchart

| | Analysis Software | File extension to be used |
|---|---|---|
| Download FASTQ file | SRA toolkit | sra, dra, fastq |
| Mapping | HISAT2、STAR | sam, bam |
| Expression Comparison | FeatureCount | txt, xlsx |
| Visualization | IGV | bam.bai |

---

## International Nucleotide Sequence Database

- **INSDC(International Nucleotide Sequence Database Collaboration)**

- INSDC is a global cooperative of independent government agencies or non-profit organizations that manage sequence databases, collect sequence information and annotations, preserve them in the scientific record, and create a comprehensive collection of such data that can be widely shared.

## International Nucleotide Sequence Database Collaboration

- The National Library of Medicine(NLM)
- National Center of Biotechnology Information(NCBI)
- NCBI Sequence Read Archive(SRA)

- The European Molecular Biology Laboratory(EMBL)
- European Bioinformatics Institute(EBI)
- EBI Sequence Read Archive(ERA)

- The Research Organization of Information and System(大学共同利用機関法人 情報・システム研究機構)
- National Institute of Genetics(日本国立遺伝学研究所)
- DDBJ Sequence Read Archive(DRA)

---

## International Nucleotide Sequence Database

- All sequence data registered in DDBJ are issued with DDBJ accession numbers. The accession number is assigned to each nucleotide sequence data registered in the database and is unique to that sequence data. Data registered in DDBJ are sent to NCBI and EBI at the time of publication, and are common worldwide. More than 99% of all registrations from Japan are made through DDBJ.

---

## Find an accession number

- Pubmed
  - Srsf7 Establishes the Juvenile Transcriptome through Age-Dependent Alternative Splicing in Mice

  - DRA009510
  - DRA009537
  - DRA009538

---

---

## Download Fastq file

- SRA-Toolkit

Download archive file
    $ prefetch [accession number]

Convert archive files to Fastq format
    $ fastq-dump [archive file]

---

## gzファイル

To compress
    $ gzip [input file]

To decompress
    $ gzip –d [input file]

## FASTQの中身を確認

```
$ less L-1_P1.fastq.gz

@MG00HS09:723:C9A15ACXX:7:1101:1483:1919 1:N:0:CGATGT
NGACCCGCTGAATTTAAGCATATTAGTCAGCGGAGGAAAAGAAACTAACCA
+
#11B?D@8DAF?DGGECFHDF?4ACFB?GEGGB6)?69DE;FGE=4@F)=C
```

- 1 : @Sequence ID and additional information
- 2 : Nucleotide sequence
- 3 : + Sequence ID and additional information
- 4 : sequence quality

## ASC2

| | | | | |
|---|---|---|---|---|
| 33! | 45- | 57  9 | 69E | 81Q |
| 34" | 46. | 58: | 70F | 82R |
| 35# | 47/ | 59; | 71G | 83S |
| 36$ | 48  0 | 60< | 72H | 84T |
| 37% | 49  1 | 61= | 73I | 85U |
| 38& | 50  2 | 62> | 74J | 86V |
| 39' | 51  3 | 63? | 75K | 87W |
| 40( | 52  4 | 64@ | 76L | 88X |
| 41) | 53  5 | 65A | 77M | 89Y |
| 42* | 54  6 | 66B | 78N | 90Z |
| 43+ | 55  7 | 67C | 79O | |
| 44, | 56  8 | 68D | 80P | |

| ASCII文字 | ASCIIコード | クオリティスコア | 塩基が間違っている確率 (P) | ASCII文字 | ASCIIコード | クオリティスコア | 塩基が間違っている確率 (P) |
|---|---|---|---|---|---|---|---|
| ! | 33 | 0 | 1.0 (100%) | 6 | 54 | 21 | 0.00794 (0.794%) |
| " | 34 | 1 | 0.794 (79.4%) | 7 | 55 | 22 | 0.00631 (0.631%) |
| # | 35 | 2 | 0.631 (63.1%) | 8 | 56 | 23 | 0.00501 (0.501%) |
| $ | 36 | 3 | 0.501 (50.1%) | 9 | 57 | 24 | 0.00398 (0.398%) |
| % | 37 | 4 | 0.398 (39.8%) | : | 58 | 25 | 0.00316 (0.316%) |
| & | 38 | 5 | 0.316 (31.6%) | ; | 59 | 26 | 0.00251 (0.251%) |
| ' | 39 | 6 | 0.251 (25.1%) | < | 60 | 27 | 0.00199 (0.199%) |
| ( | 40 | 7 | 0.199 (19.9%) | = | 61 | 28 | 0.00158 (0.158%) |
| ) | 41 | 8 | 0.158 (15.8%) | > | 62 | 29 | 0.00126 (0.126%) |
| * | 42 | 9 | 0.126 (12.6%) | ? | 63 | 30 | 0.00100 (0.1%) |
| + | 43 | 10 | 0.100 (10%) | @ | 64 | 31 | 0.000794 (0.0794%) |
| , | 44 | 11 | 0.0794 (7.94%) | A | 65 | 32 | 0.000631 (0.0631%) |
| - | 45 | 12 | 0.0631 (6.31%) | B | 66 | 33 | 0.000501 (0.0501%) |
| . | 46 | 13 | 0.0501 (5.01%) | C | 67 | 34 | 0.000398 (0.0398%) |
| / | 47 | 14 | 0.0398 (3.98%) | D | 68 | 35 | 0.000316 (0.0316%) |
| 0 | 48 | 15 | 0.0316 (3.16%) | E | 69 | 36 | 0.000251 (0.0251%) |
| 1 | 49 | 16 | 0.0251 (2.51%) | F | 70 | 37 | 0.000199 (0.0199%) |
| 2 | 50 | 17 | 0.0199 (1.99%) | G | 71 | 38 | 0.000158 (0.0158%) |
| 3 | 51 | 18 | 0.0158 (1.58%) | H | 72 | 39 | 0.000126 (0.0126%) |
| 4 | 52 | 19 | 0.0126 (1.26%) | I | 73 | 40 | 0.000100 (0.01%) |
| 5 | 53 | 20 | 0.0100 (1.0%) | | | | |

## RNAseq flowchart

| Analysis Software | File extension to be used |
|---|---|
| Download FASTQ file → SRA toolkit | sra, dra, fastq |
| Mapping → HISAT2、STAR | sam, bam |
| Expression Comparison → FeatureCount | txt, xlsx |
| Visualization → IGV | bam.bai |

## Mapping

- HISAT2
- hisat2 -p [number of CPU] -x [genome index] -U [input file] -S [output file]

- Download Index file of the genome
  - http://daehwankimlab.github.io/hisat2/download/#m-musculus

- Mapping by using hisat2. Make sam file.

```
$ hisat2 -p 8 -x HISATindex/genome -U fastq.gz/L-1_P1.fastq.gz -S sam/L1.sam
```

## SAMファイル

- QNAME: Query name. Read identifier (usually the read ID).
- FLAG: Bit flag indicating alignment information. It has binary information such as read mapping status and pair information.
- RNAME: Reference name. The name of the reference sequence to which the read is aligned (e.g., chromosome name).
- POS: Leftmost position of 1-base. The starting position of the read mapping.   MAPQ: Mapping quality. A score indicating the reliability of the alignment.
- CIGAR: CIGAR string. Indicates the alignment pattern of the reads (e.g., '76M' means that 76 bases were matched).
- RNEXT: The reference name to which the next read is aligned. In the case of paired-end sequencing, the reference sequence to which the other read is aligned.
- PNEXT: The leftmost position of one base of the next read. The starting position of the partner of the paired-end read.
- TLEN: Template Length. Insertion size between paired-end reads.
- SEQ: Read Sequence. Actual nucleotide sequence.
- QUAL: Quality score. Read sequence quality expressed in ASCII characters.

## Convert sam file to bam file

- samtools

Convert sam file to bam file
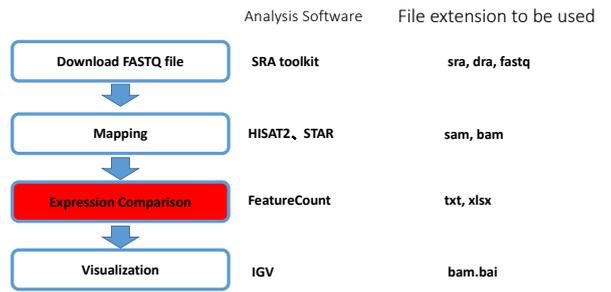    $ samtools view -@ 8 -bS [input sam file] > [output bam file]
Sorting bam files
    $ samtools sort -@ 8 -o [output sorted bam file] [input bam file]
Create an index of bam files
    $ samtools index [input sorted bam file]

---

## RNAseq flowchart

| | Analysis Software | File extension to be used |
|---|---|---|
| Download FASTQ file | SRA toolkit | sra, dra, fastq |
| Mapping | HISAT2、STAR | sam, bam |
| Expression Comparison | FeatureCount | txt, xlsx |
| Visualization | IGV | bam.bai |

---

## Expression Comparison

- FeatureCounts

featureCounts -T [number of CPU] --extraAttributes gene_name -a [GTF file] -o [output file] [input file1] [input file2] [input file3]

- Download GTF (gene feature format)
  - https://www.ensembl.org/info/data/ftp/index.html

---

## Expression Comparison

- FeatureCounts

- Count the number of reads mapped to each gene

$ featureCounts -T 8 --extraAttributes gene_name -a GTF/Mus_musculus.GRCm38.93.sorted.gtf -o counts_result.txt bam_sort/L1.sort.bam

$ featureCounts -T 8 --extraAttributes gene_name -a GTF/Mus_musculus.GRCm38.93.sorted.gtf -o counts_result.txt bam_sort/L1.sort.bam bam_sort/L2.sort.bam bam_sort/L3.sort.bam bam_sort/L4.sort.bam bam_sort/L5.sort.bam bam_sort/L6.sort.bam bam_sort/L7.sort.bam bam_sort/L8.sort.bam bam_sort/L9.sort.bam

---

## Expression Comparison

- Gene expression ≠ Number of reads mapped on the gene

- Longer genes have more reads mapped (intergenic bias)

- The more samples(amount of cDNA), the more reads mapped (sample bias)

- These biases need to be corrected before comparing expression levels.

---

## Expression Comparison
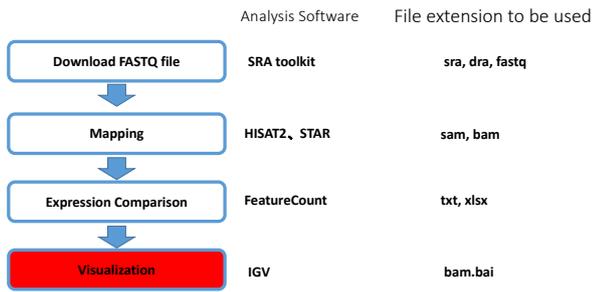
- Excel
  - TPM（transcripts per million）

  - qi is the number of reads mapped
  - li is the transcript length

$$A_i = \frac{q_i}{l_i} * 10^3$$

$$TPM_i = A_i * \frac{1}{\sum_j A_j} * 10^6$$

**RNAseq flowchart**

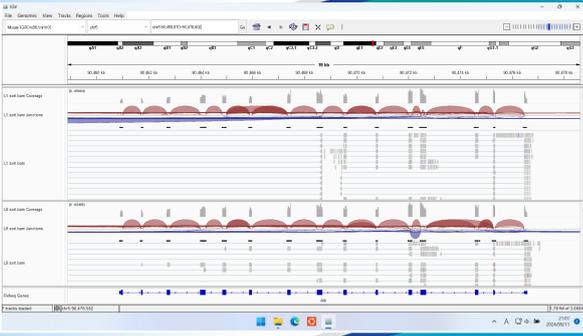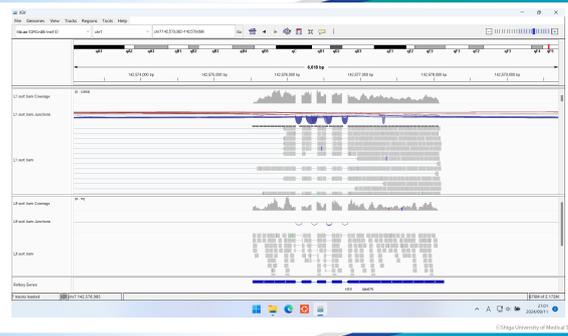| | Analysis Software | File extension to be used |
|---|---|---|
| Download FASTQ file | SRA toolkit | sra, dra, fastq |
| Mapping | HISAT2、STAR | sam, bam |
| Expression Comparison | FeatureCount | txt, xlsx |
| Visualization | IGV | bam.bai |

**Visualization**

- IGV（Integrative Genomics Viewer）

**Visualization**

**Visualization**

**Visualization**