## Slide 1

**King Mongkut's University of Technology Thonburi**
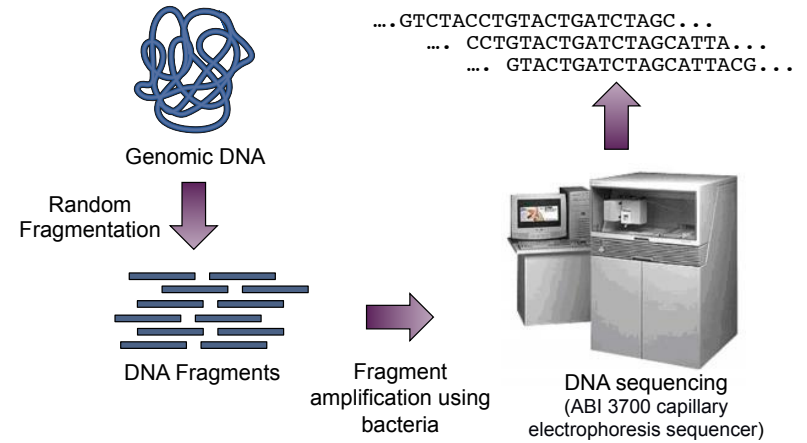
**Next-generation Sequencing Data Analysis**
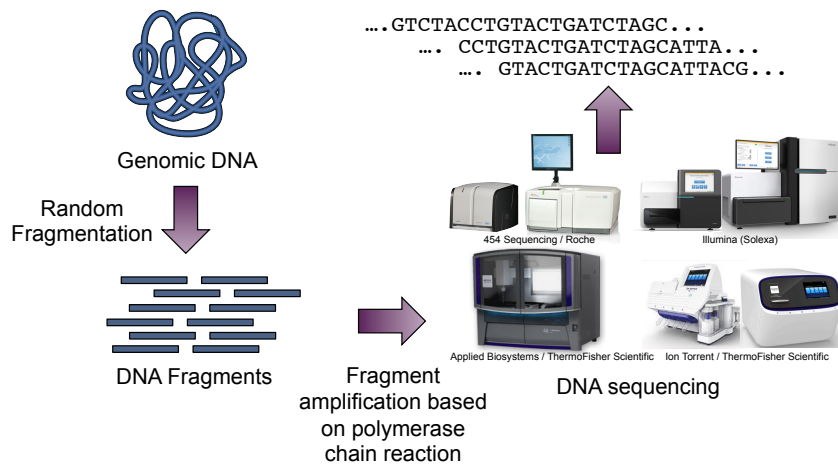**CSBio 2015 Pre-conference workshop: 22 November 2015**

Weerayuth Kittichotirat
weerayuth.kit@kmutt.ac.th

Systems biology and Bioinformatics Research Group
Pilot Plant Development and Training Institute
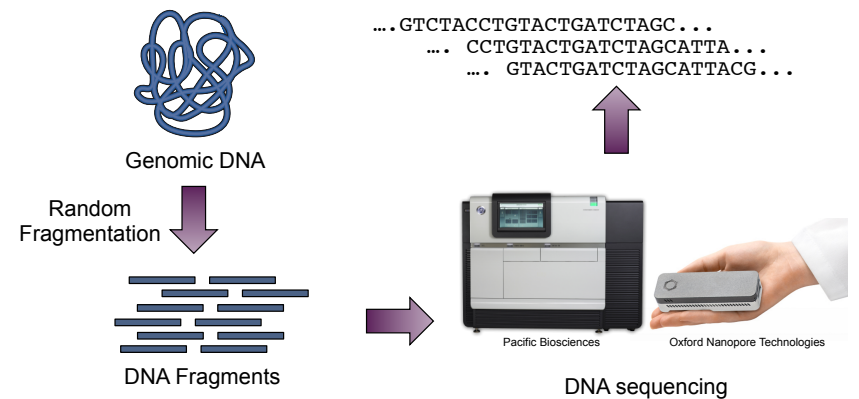King Mongkut's University of Technology Thonburi

CS Bio 2015 Bangkok, Thailand

## Slide 2

# Overview of the previous generation of DNA sequencing

….GTCTACCTGTACTGATCTAGC...
…. CCTGTACTGATCTAGCATTA...
…. GTACTGATCTAGCATTACG...

Genomic DNA

Random Fragmentation

DNA Fragments

Fragment amplification using bacteria

DNA sequencing (ABI 3700 capillary electrophoresis sequencer)

## Slide 3

# Overview of the next generation DNA sequencing

….GTCTACCTGTACTGATCTAGC...
…. CCTGTACTGATCTAGCATTA...
…. GTACTGATCTAGCATTACG...

Genomic DNA

Random Fragmentation

DNA Fragments

Fragment amplification based on polymerase chain reaction

DNA sequencing

454 Sequencing / Roche

Illumina (Solexa)

Applied Biosystems / ThermoFisher Scientific

Ion Torrent / ThermoFisher Scientific

## Slide 4

# Overview of the third generation DNA sequencing

….GTCTACCTGTACTGATCTAGC...
…. CCTGTACTGATCTAGCATTA...
…. GTACTGATCTAGCATTACG...

Genomic DNA

Random Fragmentation

DNA Fragments

Pacific Biosciences

Oxford Nanopore Technologies

DNA sequencing

# + Sequencing run types

- Single-end sequencing

Read 1

Target DNA

- Paired-end sequencing

Read 1

Target DNA

Read 2

Distance between reads is known

---

# + FASTQ file format

- FASTQ format is a <u>text-based format</u> for storing both
  - a biological sequence (usually nucleotide sequence)
  - and its corresponding quality scores[1].

Uses 4 lines for each sequence

Line 1 begins with a '@' character and is followed by a sequence identifier and an optional description

Line 2 is the raw sequence letters.

```
1  @MG00HS12:401:C3W95ACXX:4:1106:16550:82618/1
2  AGACTTACAATGATGATTCAAATGAAGGAAACTAAAAAGTAATGAAGCAAGGCAGAGGAAAA
3  +
4  @@@DDDDDDFHFDGFBFHGCG?B?:F@9:CCCH@?AEGG??FEGI@FHIGGGIIIEGIIIIB
5  @MG00HS12:401:C3W95ACXX:4:1207:20629:56402/1
6  ACCCGGCTAATGTTGTAGTTTTAGTAGAGACGGGGTTTCCCTATGTTGGTTAGGCTGGTCTC
7  +
8  @C@FFFFFFHHHIJHIJHIIIJIJHIIIJJGJJJJJ@FIJJJJGHIIIJIJIJIIIJJGGHHEH
```

Line 3 begins with a '+' character

Line 4 encodes the quality values for the sequence in Line 2

- Quality values are encoded using ASCII scheme

```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
```

↑ Lowest quality value

Highest quality value ↑

[1]https://en.wikipedia.org/wiki/FASTQ_format

---

# + Quality values are encoded differently for different platform

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...............................................
.....................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.......................
...........................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII..................
.....................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....................
LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.............................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                        |        |      |                                |          |
33                       59      64     73                               104        126
0.....................26...31.......40
                   -5...0.......9.......................40
                   0.......9.......................40
                        3.....9.......................40
0.2.....................26...31.......41
```

S - Sanger        Phred+33,  raw reads typically (0, 40)
X - Solexa        Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64,  raw reads typically (3, 40)
   with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
   (Note: See discussion above).
L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)

https://en.wikipedia.org/wiki/FASTQ_format

---

# + Quality value

- A quality value $Q$ is an integer mapping of $p$ (i.e., the probability that the corresponding base call is incorrect)[1]

$$Q_{sanger} = -10 \log_{10} p$$

| Quality score | Probability of incorrect bases | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 17 | 1 in 50 | 98% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1,000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |

[1]https://en.wikipedia.org/wiki/FASTQ_format

# + Quality control check using FastQC

- A quality control tool for high throughput sequence data.

- Developed by the Bioinformatics Group at the Babraham Institute, United Kingdom

- Available at http://www.bioinformatics.babraham.ac.uk/projects/fastqc/



**Babraham Bioinformatics**

About | People | Services | Projects | Training | Publications

**FastQC**

| Function | A quality control tool for high throughput sequence data. |
|---|---|
| Language | Java |
| Requirements | A suitable Java Runtime Environment |
| | The Picard BAM/SAM Libraries (included in download) |

---

# + Quality control check using FastQC

- Transfer raw sequencing data on to the Linux server

1) Open WinSCP

WinSCP.exe
WinSCP: SFTP, FTP
Martin Prikryl

2) Fill in Host Name, User name, Password, and click Login

`stud.sbi.kmutt.ac.th`



---

# + Quality control check using FastQC

- Transfer raw sequencing data on to the Linux server

3) Click and Drop raw sequence read files from the left to right panel



---

# + Quality control check using FastQC

- Log into the Linux server

1) Open Putty

putty.exe

2) Fill in Host Name, Click Open, and Enter your password when prompted

`userX@stud.sbi.kmutt.ac.th`

# + Quality control check using FastQC

- A few useful commands

  - `ls` : Show all files and directories (folders)

  - `mkdir` : Make a new directory (folder)

  - `cd` : Change directory (folder)

---

# + Quality control check using FastQC

- Download and install **FastQC**
  - `mkdir fastqc`
  - `cd fastqc`
  - `wget http://www.bioinformatics.babraham.ac.uk/projects/fastqc/fastqc_v0.11.4.zip`
  - `unzip fastqc_v0.11.4.zip`
  - `cd FastQC`
  - `chmod 755 fastqc`
  - `ls`

```
[weerayuth@stud FastQC]$ ls
cisd-jhdf5.jar   Help            LICENSE.txt    RELEASE_NOTES.txt  uk
Configuration    INSTALL.txt     net            run_fastqc.bat
fastqc           jbzip2-0.9.jar  org            sam-1.103.jar
fastqc_icon.ico  LICENSE_JHDF5.txt  README.txt  Templates
[weerayuth@stud FastQC]$
```

---

# + Quality control check using FastQC

- Run **FastQC** to check quality of sequencing data
  - `cd`
  - `mkdir output-fastqc`
  - `fastqc/FastQC/fastqc S2_L001_R1_001.fastq S2_L001_R2_001.fastq -o output-fastqc`
  - `cd output-fastqc`
  - `ls`

```
[weerayuth@stud output-fastqc]$ ls
S2_L001_R1_001_fastqc.html   S2_L001_R2_001_fastqc.html
S2_L001_R1_001_fastqc.zip    S2_L001_R2_001_fastqc.zip
[weerayuth@stud output-fastqc]$
```

---

# + Quality control check using FastQC

- Copy FastQC result back

  1) Open WinSCP

  WinSCP.exe
  WinSCP: SFTP, FTP
  Martin Prikryl

  2) Fill in Host Name, User name, Password, and click Login

# + Quality control check using FastQC

- Copy FastQC result back
  3) Click and Drop folder output-fastqc from the right to left panel



# + FastQC output files

- FastQC generates a HTML report and a zip file containing individual graphs for each input file

## Slide (FastQC Report)

**FastQC Report**

**Summary**

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

✓ **Per sequence quality scores**

Quality score distribution over all sequences

Average Quality per read

Mean Sequence Quality (Phred Score)

Produced by **FastQC** (version 0.11.4)

- Explanation of each quality check can be found at
  - http://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/

## Slide 22

# Coffee Break

## Slide 23

# *De novo* genome assembly

Genomic DNA

Randomly fragmented and sequenced

Sequence reads are stitched back based on overlapping sequences into longer <u>contig sequences</u>

ATGTTCCGATTAGGAAACCTATCTGTAACTGTTTCATTCAGTAAAAGGAGGAAA

http://people.mpi-inf.mpg.de/~sven/images/assembly.png

## Slide 24

# *De novo* assembly using Velvet

- Velvet is a *de novo* genomic assembler specially designed for short read sequencing technologies

- Developed by Daniel Zerbino and Ewan Birney at the European Bioinformatics Institute (EMBL-EBI), United Kingdom

- Available at https://www.ebi.ac.uk/~zerbino/velvet/

EMBL-EBI

**Velvet**

Sequence assembler for very short reads

- **Current version: 1.2.10**
- **Manual** and **extension for Columbus** in pdf format

# *De novo* assembly using Velvet

- Download and install **Velvet**
  - `cd`
  - `mkdir velvet`
  - `cd velvet`
  - `wget https://www.ebi.ac.uk/~zerbino/velvet/velvet_1.2.10.tgz`
  - `tar -xzf velvet_1.2.10.tgz`
  - `cd velvet_1.2.10`
  - `make 'BIGASSEMBLY=1' 'LONGSEQUENCES=1' 'MAXKMERLENGTH=151'`

---

# *De novo* assembly using Velvet

```
make 'BIGASSEMBLY=1' 'LONGSEQUENCES=1' 'MAXKMERLENGTH=151'
```

- Allow Velvet to handle more than 2.2 billion reads.
- This will cost more memory overhead.

- Allow Velvet to handle input sequences that are longer than 32kbp.
- This will cost more memory overhead.

- Allow Velvet to handle longer word length (default is 31bp)
- Longer word require more memory

---

# *De novo* assembly using Velvet

- Download and install **Velvet**
  - `cd`
  - `mkdir velvet`
  - `cd velvet`
  - `wget https://www.ebi.ac.uk/~zerbino/velvet/velvet_1.2.10.tgz`
  - `tar -xzf velvet_1.2.10.tgz`
  - `cd velvet_1.2.10`
  - `make 'BIGASSEMBLY=1' 'LONGSEQUENCES=1' 'MAXKMERLENGTH=151'`
  - `ls`

```
[weerayuth@stud velvet_1.2.10]$ ls
ChangeLog              debian                         Manual.pdf   third-party
Columbus_manual.pdf    doc                            obj          update_velvet.sh
contrib                For_MAC_or_SPARC_users.txt     README.txt   velvetg
CREDITS.txt            LICENSE.txt                    src          velveth
data                   Makefile                       tests
[weerayuth@stud velvet_1.2.10]$ 
```

---

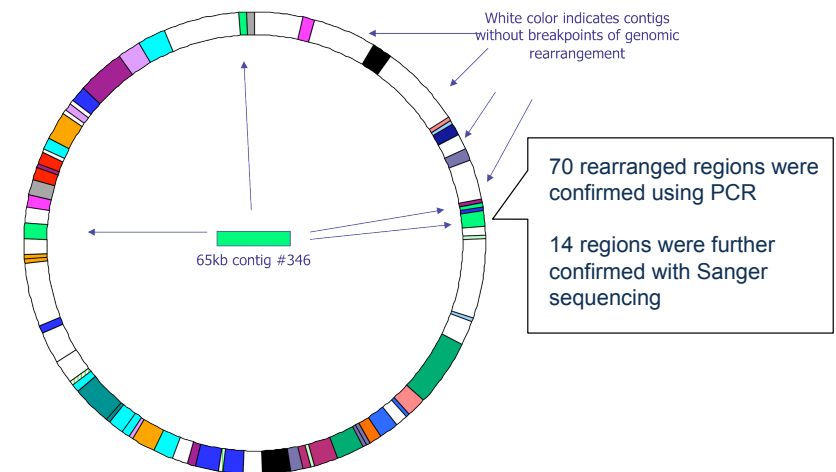# *De novo* assembly using Velvet

- Running *de novo* assembly
  - `cd`
  - `mkdir output-velvet`
  - `velvet/velvet_1.2.10/velveth output-velvet 83 -fastq -shortPaired -separate S2_L001_R1_001.fastq S2_L001_R2_001.fastq`

# *De novo* assembly using Velvet

`velvet/velvet_1.2.10/velveth`

- `velveth` produces hashtable and output files that are required for `velvetg`

---

# *De novo* assembly using Velvet

`velvet/velvet_1.2.10/velveth`

`output-velvet`

- `velveth` produces hashtable and output files that are required for `velvetg`
- Output directory

---

# *De novo* assembly using Velvet

`velvet/velvet_1.2.10/velveth`

`output-velvet`

`83`

- `velveth` produces hashtable and output files that are required for `velvetg`
- Output directory
- The word length in bp that are being hashed
- This can affect the assembly result

---

# Choice of hash length

- It must be an odd number. If an even number is entered, Velvet will just decrement it and proceed

- It must be below or equal to MAXKMERHASH length

- It must be shorter than the read length, otherwise you simply will not observe any overlaps between reads

- Longer hash length allows more specific overlap but fewer reads will be used in the assembly resulting in a decrease coverage

- Shorter hash length allows more reads to be used in the assembly and result in an increase in sensitivity and coverage but will also introduce more errors and higher computation overhead

- Choice of hash length can affect the assembly output and therefore tests with different lengths are usually carried out to find the length that works best

## Slide 33

# + Choice of hash length

*Chart: Number of contigs (y-axis, 100–155) vs Hash lengths (bp) (x-axis, 30 32 34 36 38 40 42 44 46 48 50 52 54 56 58 60 62 64 66 68 70 72 74 76 78 80 82 84 86 88 90 92 94 96 98 100)*

— No of contigs

---

## Slide 34

# + *De novo* assembly using Velvet

`velvet/velvet_1.2.10/velveth`

- `velveth` produces hashtable and output files that are required for `velvetg`

`output-velvet`
- Output directory

`83`
- The word length in bp that are being hashed
- This can affect the assembly result

`-fastq`
`-shortPaired`
- Specify type of input sequences

---

## Slide 35

# + Types of input sequences

| Supported file formats are: | Read categories are: |
|---|---|
| **fasta** (default) | **short** (default) |
| **fastq** | **shortPaired** |
| **fasta.gz** | **short2** (same as short, but for a separate insert-size library) |
| **fastq.gz** | **shortPaired2** (see above) |
| **sam** | **long** (for Sanger, 454 or even reference sequences) |
| **bam** | |
| **eland** | |
| **gerald** | |

https://www.ebi.ac.uk/~zerbino/velvet/Manual.pdf

---

## Slide 36

# + *De novo* assembly using Velvet

`velvet/velvet_1.2.10/velveth`

- `velveth` produces hashtable and output files that are required for `velvetg`

`output-velvet`
- Directory that contains `velveth` output files

`83`
- The word length in bp that are being hashed
- This can affect the assembly result

`-fastq`
`-shortPaired`
- Specify type of input sequences

`-separate`
- Specify that pair end input sequences are in two separate files

`S2_L001_R1_001.fastq`
`S2_L001_R2_001.fastq`
- Input files

## Slide 37

# + *De novo* assembly using Velvet

- Running *de novo* assembly
  - cd
  - mkdir output-velvet
  - velvet/velvet_1.2.10/velveth output-velvet 83 -fastq -shortPaired -separate S2_L001_R1_001.fastq S2_L001_R2_001.fastq
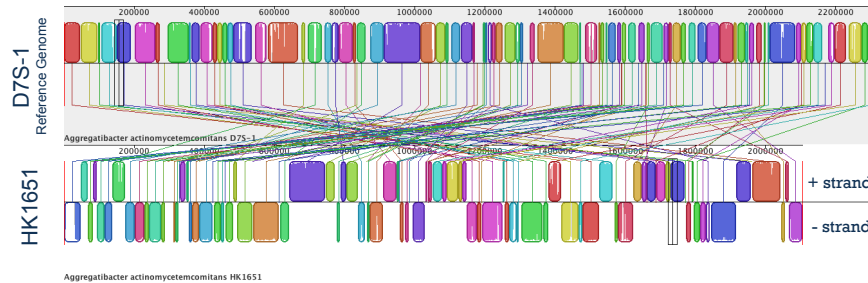  - velvet/velvet_1.2.10/velvetg output-velvet -exp_cov auto -min_contig_lgth 300

## Slide 38

# + *De novo* assembly using Velvet

velvet/velvet_1.2.10/velvetg → • velvetg makes use of velveth output files to create a sequence assembly

output-velvet → • Output directory

-exp_cov auto → • Let Velvet automatically determined the expected coverage
• This option is intended mainly for standard genomic sequencing

-min_contig_lgth 300 → • Specify the mininum contig length in the output contigs.fa file

## Slide 39

# + *De novo* assembly using Velvet

- Running *de novo* assembly
  - cd
  - mkdir output-velvet
  - velvet/velvet_1.2.10/velveth output-velvet 83 -fastq -shortPaired -separate S2_L001_R1_001.fastq S2_L001_R2_001.fastq
  - velvet/velvet_1.2.10/velvetg output-velvet -exp_cov auto -min_contig_lgth 300
  - cd output-velvet
  - ls

```
[weerayuth@stud output-velvet]$ ls
contigs.fa  Graph2  LastGraph  Log  PreGraph  Roadmaps  Sequences  stats.txt
[weerayuth@stud output-velvet]$
```

## Slide 40

**Initial results of genome sequencing of**
***Aggregatibacter actinomycetemcomitans* strain D7S-1**



White color indicates contigs without breakpoints of genomic rearrangement

70 rearranged regions were confirmed using PCR

14 regions were further confirmed with Sanger sequencing

65kb contig #346

Alignment of D7S contigs based on the genomic map of HK1651.

## Genomic rearrangement between strain D7S-1 and HK1651 of *Aggregatibacter actinomycetemcomitans*



- Whole genome sequence alignment created using the Mauve progressive alignment software

---

## Special challenges with next generation sequencing data

- Typically, only an incomplete genome is generated

- The cost of closing all gaps to produce a complete genome is still high
  - More incomplete genome sequences will be in public databases in the future

- Each technology is prone to making certain type of errors
  - Roach/454 and Ion Torrent tends to produce insertion/deletion in homopolymer regions

- Mapping to a reference genome may not be possible or even misleading

- Incomplete genome and sequence error produce "even greater" challenges in downstream analysis such as gene prediction and annotation

---

## + Whole-exome sequencing

- The targeted sequencing of the subset of the human genome that code for RNA or amino-acid.

- About 1% (30Mb) of the human genome.

- It is estimated that 85% of the disease-causing mutations are located in coding and functional regions of the genome

- Rabbani, B., Tekin, M., & Mahdieh, N. (2013). The promise of whole-exome sequencing in medical genetics. Journal of human genetics, 59(1), 5-15.
- van Dijk, E. L., Auger, H., Jaszczyszyn, Y., & Thermes, C. (2014). Ten years of next-generation sequencing technology. Trends in genetics, 30(9), 418-426.

---

## Whole-exome sequencing



```
AGGTCGTTACGTACGCTAC
GACCTACATCAGTACATAG
GCATGACAAAGCTAGGTGT
```
Mapping, alignment, variant calling

1. Bamshad, Michael J., et al. "Exome sequencing as a tool for Mendelian disease gene discovery." *Nature Reviews Genetics* 12.11 (2011): 745-755.
2. Human All Exon. (n.d.). http://www.genomics.agilent.com/en/SureSelect-DNA-Target-Enrichment-Baits-/Human-All-Exon/?cid=AG-PT-124&tabId=AG-PR-1308

# Slide 1

**+ Whole-exome sequencing data analysis**

- Preparation of a reference human genome sequence
  - `cd`
  - `wget http://hgdownload.cse.ucsc.edu/goldenpath/hg19/chromosomes/chr7.fa.gz`
    - More information can be found at http://hgdownload.cse.ucsc.edu/goldenpath/hg19/chromosomes/
  - `gunzip chr7.fa.gz`
  - `ls`



```
[weerayuth@stud ~]$ ls
bcftools    fastqc       S2_L001_R1_001.fastq   tabix                 vcftools
bwa         output-fastqc S2_L001_R2_001.fastq  tumor_chr7_1.fastq    velvet
chr7.fa     output-velvet samtools              tumor_chr7_2.fastq
[weerayuth@stud ~]$
```

# Slide 2

**+ Whole-exome sequencing data analysis**

- Preparation of the human genome annotation information
  - Go to https://genome.ucsc.edu/
  - Click Table Browser



# Slide 3

- Retrieve annotation data for Chromosome 7 as shown below



# Slide 4

- Select Exons and click 'get BED' as shown below



- Transfer file 'chr7.bed' to the Linux server using WinSCP



```
[weerayuth@stud ~]$ ls
bcftools    fastqc       S2_L001_R2_001.fastq   tumor_chr7_2.fastq
bwa         output-fastqc samtools              vcftools
chr7.bed    output-velvet tabix                 velvet
chr7.fa     S2_L001_R1_001.fastq tumor_chr7_1.fastq
[weerayuth@stud ~]$
```

**+ Whole-exome sequencing data analysis**

- Download and install Burrows-Wheeler Aligner http://bio-bwa.sourceforge.net/
  - cd
  - mkdir bwa
  - cd bwa
  - wget 'http://downloads.sourceforge.net/project/bio-bwa/ bwa-0.7.12.tar.bz2?r=http%3A%2F%2Fsourceforge.net %2Fprojects%2Fbio-bwa%2Ffiles %2F&ts=1447657569&use_mirror=jaist' -O bwa-0.7.12.tar.bz2
  - bzip2 -d bwa-0.7.12.tar.bz2
  - tar -xf bwa-0.7.12.tar
  - cd bwa-0.7.12/
  - make
  - ls

```
bntseq.o        bwaseqio.o   bwtsw2_aux.o    kopen.o    pemerge.c
bwa             bwashm.c     bwtsw2_chain.c  kseq.h     pemerge.o
bwa.1           bwashm.o     bwtsw2_chain.o  ksort.h    QSufSort.c
```

**+ Whole-exome sequencing data analysis**

- Index reference genome
  - cd
  - bwa/bwa-0.7.12/bwa index chr7.fa
    - More information http://bio-bwa.sourceforge.net/bwa.shtml
  - ls

```
[weerayuth@stud ~]$ ls
bcftools   chr7.fa.amb   chr7.fa.sa     S2_L001_R1_001.fastq  tumor_chr7_1.fastq
bwa        chr7.fa.ann   fastqc         S2_L001_R2_001.fastq  tumor_chr7_2.fastq
chr7.bed   chr7.fa.bwt   output-fastqc  samtools              vcftools
chr7.fa    chr7.fa.pac   output-velvet  tabix                 velvet
```

**+ Whole-exome sequencing data analysis**

- Align sequence reads to reference genome
  - cd
  - bwa/bwa-0.7.12/bwa mem chr7.fa tumor_chr7_1.fastq tumor_chr7_2.fastq > alignment.sam
    - More information http://bio-bwa.sourceforge.net/bwa.shtml
    - SAM format specification https://samtools.github.io/hts-specs/SAMv1.pdf
  - ls

```
[weerayuth@stud ~]$ ls
alignment.sam  chr7.fa.amb  fastqc         samtools
bcftools       chr7.fa.ann  output-fastqc  tabix
bwa            chr7.fa.bwt  output-velvet  tumor_chr7_1.fastq
chr7.bed       chr7.fa.pac  S2_L001_R1_001.fastq  tumor_chr7_2.fastq
chr7.fa        chr7.fa.sa   S2_L001_R2_001.fastq  vcftools
```

**+ Whole-exome sequencing data analysis**

- Download and install Samtools (http://www.htslib.org/)

- Samtools is a suite of programs for interacting with high-throughput sequencing data
  - cd
  - mkdir samtools
  - cd samtools
  - wget https://github.com/samtools/samtools/releases/ download/1.2/samtools-1.2.tar.bz2
  - bzip2 -d samtools-1.2.tar.bz2
  - tar -xf samtools-1.2.tar
  - cd samtools-1.2/
  - make
  - ls

```
bam_color.c   bam_rmdupse.o   examples   samtools
bam_color.o   bamshuf.c       faidx.c    samtools.1
```

## Slide 1

**+ Whole-exome sequencing data analysis**

- "Clean" alignment result using Samtools
  - `cd`
  - `samtools/samtools-1.2/samtools fixmate -O bam alignment.sam alignment.fixmate.bam`
    - BWA can sometimes leave unusual FLAG information on SAM records, it is helpful when working with many tools to first clean up read pairing information and flags
  - `samtools/samtools-1.2/samtools sort -O bam -o alignment.fixmate.sorted.bam -T alignment.fixmate.temp alignment.fixmate.bam`
    - Sort records from name order into coordinate order
  - `ls`

```
[weerayuth@stud ~]$ ls
alignment.fixmate.bam         chr7.fa.ann      S2_L001_R2_001.fastq
alignment.fixmate.sorted.bam  chr7.fa.bwt      samtools
alignment.sam                 chr7.fa.pac      tabix
```

## Slide 2

**+ Whole-exome sequencing data analysis**

- Download and install BCFtools - utilities for variant calling and manipulating VCFs and BCFs (http://www.htslib.org/doc/bcftools.html)
  - `cd`
  - `mkdir bcftools`
  - `cd bcftools`
  - `wget https://github.com/samtools/bcftools/releases/download/1.2/bcftools-1.2.tar.bz2`
  - `bzip2 -d bcftools-1.2.tar.bz2`
  - `tar -xf bcftools-1.2.tar`
  - `cd bcftools-1.2/`
  - `make`
  - `ls`

```
AUTHORS      HMM.c      plugins      vcfcnv.c
bcftools     HMM.h      polysomy.c   vcfcnv.o
bcftools.h   HMM.o      prob1.c      vcfconcat.c
```

## Slide 3

**+ Whole-exome sequencing data analysis**

- Call sequence variants from alignment data
  - `cd`
  - `samtools/samtools-1.2/samtools mpileup -go variant.bcf -f chr7.fa -Q 30 -l chr7.bed alignment.fixmate.sorted.bam`
    - Use mpileup to produce a BCF file that contains all of the locations in the genome.
    - http://www.htslib.org/doc/samtools.html
  - `bcftools/bcftools-1.2/bcftools call -vmO v -o variant.vcf variant.bcf`
    - Call genotypes and reduce our list of sites to those found to be variant by passing this file into bcftools call
    - http://www.htslib.org/doc/bcftools.html
  - `ls`

```
bcftools                chr7.fa.sa       tumor_chr7_2.fastq
bwa                     fastqc           variant.bcf
chr7.bed                output-fastqc    variant.vcf
```

## Slide 4

**+ Whole-exome sequencing data analysis**

- Download and install VCFtools (https://vcftools.github.io/index.html)
- VCFtools provides easily accessible methods for working with complex genetic variation data in the form of VCF files
  - `cd`
  - `mkdir vcftools`
  - `cd vcftools`
  - `wget 'http://downloads.sourceforge.net/project/vcftools/vcftools_0.1.13.tar.gz?r=http%3A%2F%2Fsourceforge.net%2Fprojects%2Fvcftools%2Ffiles%2F&ts=1448009724&use_mirror=jaist' -O vcftools_0.1.13.tar.gz`
  - `tar -xzf vcftools_0.1.13.tar.gz`
  - `cd vcftools_0.1.13`
  - `make`
  - `export PERL5LIB=/home/weerayuth/vcftools/vcftools_0.1.13/perl`

## + Whole-exome sequencing data analysis

- Filter variant result using vcftools
  - `cd`
  - `cat variant.vcf | vcftools/vcftools_0.1.13/bin/vcf-annotate --filter MinDP=20/RefN -H > variant.filtered.vcf`
    - Filter out variants that are supported by less than 20 reads
    - Filter out variants where reference sequence is N
    - https://vcftools.github.io/perl_module.html#vcf-annotate
  - `ls`

```
bwa              fastqc            variant.bcf
chr7.bed         output-fastqc     variant.filtered.vcf
chr7.fa          output-velvet     variant.vcf
chr7.fa.amb      S2_L001_R1_001.fastq   vcftools
```

---

## + Whole-exome sequencing data analysis

- More information on vcf format can be found at http://www.1000genomes.org/wiki/analysis/variant%20call%20format/vcf-variant-call-format-version-41

- The resulting vcf file can be further annotated to add more functional information using variant annotation tools

- This can be done by using command-line or web-based variant annotation tools

- An example of a web-based variant annotation tool is wANNOVAR by Wang Genomics Lab at University of Southern California http://wannovar.usc.edu/

---

## + Final comments

- This workshop only introduce open-source software for doing NGS data analysis

- Advantages of open-source software
  - Free
  - Clear methods
  - Most run on Linux platforms (stable, can easily make your own pipelines)

- Disadvantage of open-source software
  - Most run on Linux platforms (Requires knowledge in Linux system and command line)
  - Lots of small software that only do one specific job
  - Can become obsolete very quickly

---

## + Final comments

- A good collection of software packages for next generation sequence data analysis can be found at http://seqanswers.com/wiki/Software

| | Bioinformatics method | Biological technology | Operating system | Language | Maintained | Licence |
|---|---|---|---|---|---|---|
| .NET BIO | Programming Library | | Windows Linux | C# | Yes | |
| 4peaks | Sequence analysis | Sanger | Mac OS X | | Yes | Freeware |
| A5 | Sequence assembly | Illumina | Linux Mac OS X | | Yes | GPLv3 |
| AB Large Indel Tool | Mapping | ABI SOLiD | Linux 64 | Perl | No | GPL |
| AB Small Indel Tool | Read mapping Alignment | ABI SOLiD | Linux 64 | Perl C++ | Maybe | GPL |
| ABBA | Sequence assembly Scaffolding | | Linux | | Maybe | Artistic License |
| ABMapper | Read mapping Alignment | Illumina | Linux | C++ Perl | Yes | GPLv3 |
| ABySS | Sequence assembly De Bruijn graph | Illumina 454 ABI SOLiD Sanger | POSIX Linux Mac OS X | C++ | Yes | Commercial Freeware |
| Adapter Removal (software) | Adapter Removal (software) | Illumina 454 | Linux 64 Windows Mac OS X | Java | Yes | Custom Licence |
| ADTEx | Hidden Markov Model Expectation Maximization | Illumina | GNU/Linux | Python R | Yes | GPLv3 |
| AGE | Alignment Gap extension | Illumina | | | Maybe | Creative Commons license (Attribution-NonCommerical) |

- This lecture is merely a small introduction to the big world of next generation sequence data analysis

- The endless possibility and new discovery is waiting for you

# + Acknowledgements



My students from the Bioinformatics and Systems biology program at KMUTT

+

## *Thank you for your attention...*