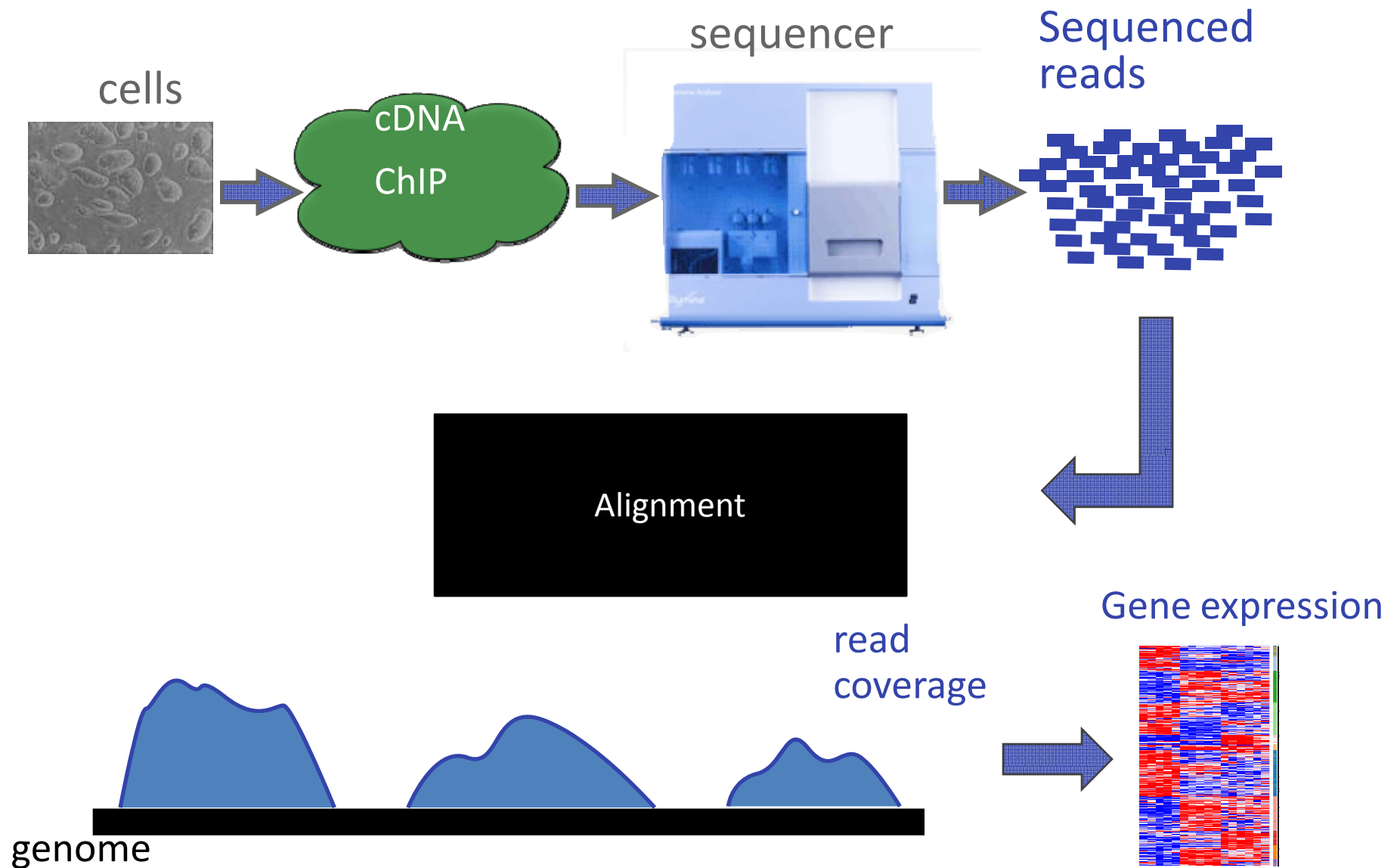# Data: databases, annotations, genomic resources

Week 2

Presenter: Hennady Shulha

# Setting up

- Log in into ghpcc06.umassrc.org

```
$ ssh <username>@ghpcc06.umassrc.org
```

- Check that you are in your home directory:

```
$ pwd
$ cd
```

- Create "biocourse" directory

```
$ mkdir biocourse
$ cd biocourse
```

# Setting up

- Put into current folder file from web
  ```
  $ wget http://biocore.umassmed.edu/biocourse/week2.tar.gz
  ```

- Unpack compacted files
  ```
  $ tar -xvzf week2.tar.gz
  $ cd week2
  ```

- Unzip this given file
  ```
  $ gunzip reads.fastq.gz
  ```

alternative and more common commands for download: "sftp", "ftp"

sftp/ftp software: Filezilla. https://filezilla-project.org/

# File transfer software

## Used to connect your PC to any ftp/sftp location



- Quickly get output files
- Quickly load some stuff like software
- Back up priceless experimental results on external hardrive

# Protocols availability

**SFTP** – supported by Filezilla

SCP – not supported;

| | |
|---|---|
| Your/other computer in UMass | Cluster system in ghpcc06.umassrc.org |

~~FTP~~
**SFTP**
SCP
RSYNC
OTHER

**FTP**
**SFTP**
SCP
RSYNC
OTHER

~~FTP~~
**SFTP**
SCP
RSYNC
OTHER

World

VPN is required

6

# VPN

https://ssl.umassmed.edu

(contact UMass support if you do not have VPN account)

Used to connect your PC through UMass firewall if you are outside of UMass campus.



University of Massachusetts Medical School

Welcome to the UMass Medical School Intranet
**Secure Access SSL VPN**

Username [          ]     Please sign in to begin your secure session.
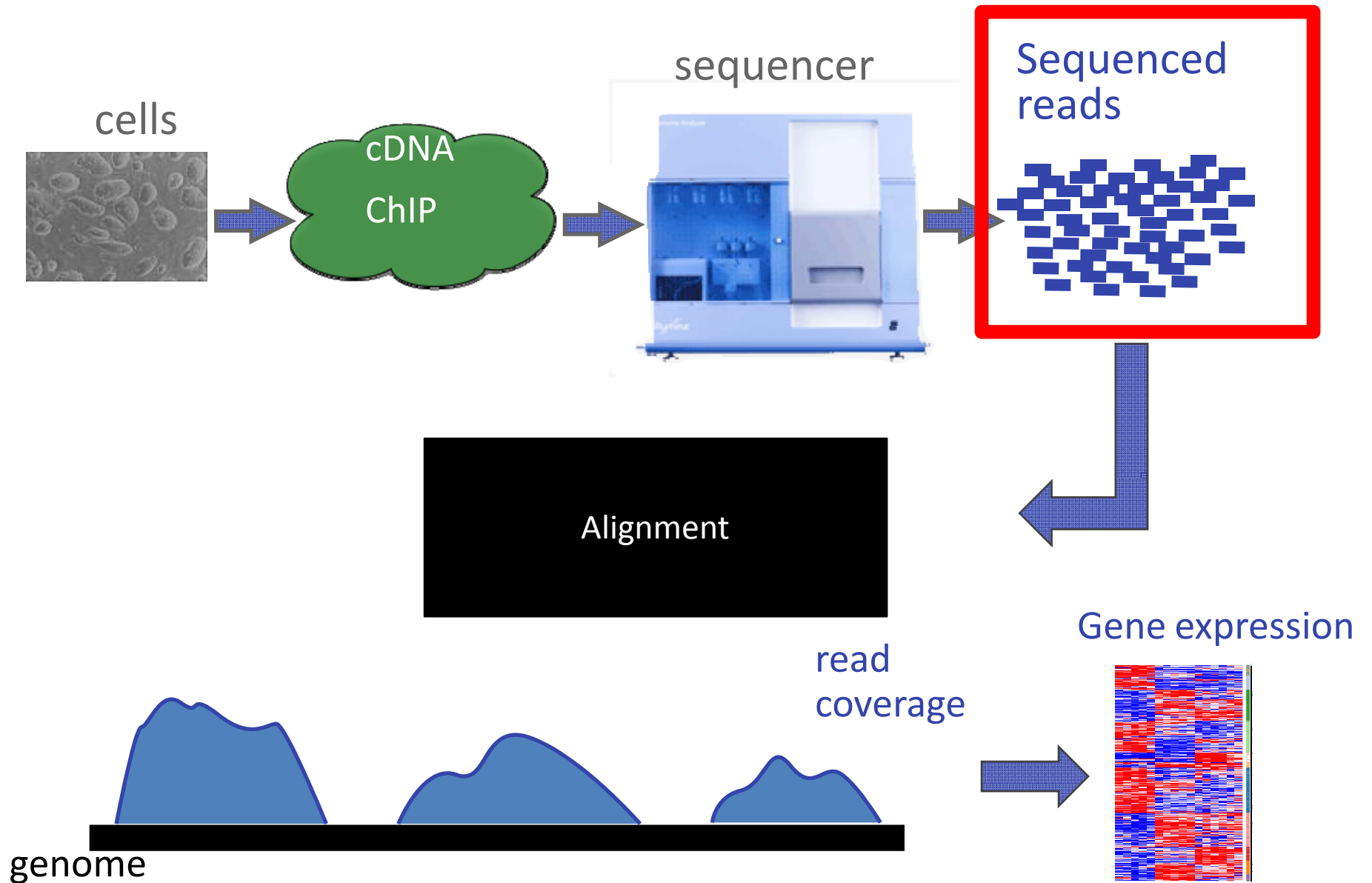Password [          ]     Forgot your password?

[ Sign In ]

# How we represent genomic data?

- Genomic sequence
  - Genomes
  - PCR products
- Genomic annotations
  - Genes
  - miRNAs
- Experimental results
  - Sequencing experiment
  - Array hybridization
- Process data for visualization
  - How many reads per base?
  - Probes are on

# File formats

- ## Binary
  - Compressed to save space; not directly readable by human; advantage is that only part of file is needed to do visual display

- ## Text readable
  - Can be opened in any text editor; occupy large space (2-3x comparing with binary)

Good command to test what is inside is "head"

http://genome.ucsc.edu/FAQ/FAQformat.html

# File formats: Fasta

Used to keep sequences of genomes, proteins etc.

Example:

```
>Sequence_Name Basic description of the sequence
ATCGATCGATGCATGCGAGTCGTAGTCGTAGTCGT
TACGTACGTAGTCGTGTACGTGTAGTCGAGTCGTA
ATCGTACGAGCGAGTCGTTGATGCTGAGTCGTGTC
TACGATGCGAGGCTGTAGTCGTAGTCGTAGTGTCC
TACGACGTGTATGCGTACGGATCGCGATTCGTAGC
```

# File formats: Fastq

Used mainly to store reads and information about quality

## Example:

```
@EAS54_6_R1_2_1_413_324
CCCTTCTTGTCTTCAGCGTTTCTCC
+
;;3;;;;;;;;;;;7;;;;;;;88

@EAS54_6_R1_2_1_443_348
GTTGCTTCTGGCGTGGGTGGGGGGG
+EAS54_6_R1_2_1_443_348
;;;;;;;;;;;9;7;;.7;393333
```

# File formats: Fastq quality

- ## Quality values

```
@EAS54_6_R1_2_1_413_324
CCCTTCTTGTCTTCAGCGTTTCTCC
+
@;3;;;;;;;;;;;7;;;;;;;88

@EAS54_6_R1_2_1_443_348
GTTGCTTCTGGCGTGGGTGGGGGGG
+EAS54_6_R1_2_1_443_348
;;;;;;;;;;;9;7;;.7;393333
```

**FASTQ (Phred)**

Q = NUMERICS_ID_of_SYMBOL − 33

$Q = -10 \ \log_{10} P$   where P is probability of incorrect base calling

**Solexa**

Q = 10 * log(1 + 10 ** (NUMERICS_ID_of_SYMBOL - 64) / 10.0)) / log(10)

$Q_{\text{solexa-prior to v.1.3}} = -10 \ \log_{10} \dfrac{p}{1-p}$   where P is probability of incorrect base calling

(later they had different  Phred related versions (Phred+33; Phred+64))

# File formats: Fastq versions

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.................................

...............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII

 !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefgh.
 |                             |    |         |                            |
33                            59   64        73                          104
 0........................26...31.......40

                              0........9...............................40
```

S - Sanger          Phred+33,  raw reads typically (0, 40)

I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)

# File formats: Fastq versions

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.....................................
...........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX...............
.....................................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.........
..................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ........
..LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL..................................
!"#$%&'()*+,-./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)
```

I.e. if you are going to use quality scores (by aligners) – would be better to know a source that generated it!!!

# Fastq quality

- Alternative: check if general picture is fine and not use it.



You get:
median/mean;
25-75% quartile range;
10-90% points

Very useful software: "FastQC". It would do basic analysis for your FASTQ files.

```
$ module load fastqc/0.10.1          $ module avail
$ fastqc -help
$ fastqc  reads.fastq
```

# BAM/SAM: Read alignment format

- Source of a bam file
  - Fastq → aligner → BAM: Alignment results
- Aligners:
  - Bowtie: Aligns contiguous reads
  - Tophat: Aligns spliced reads

To load bowtie:

```
$  module load bowtie/1.0.0
```

# File formats: SAM/BAM

Information about genomic locations, plus other useful things like base qualities, number of mismatches etc.

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1 113 chr1 497 37 37M 15 100338662 0
CGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAG
0;==-==9;>>>>>=>>>>>>>>>>=>>>>>>>>>
XT:A:U NM:i:0 SM:i:37 AM:i:0 X0:i:1 X1:i:0 XM:i:0
```

BAM contains the same info but in compressed, binary format.

http://samtools.sourceforge.net/SAMv1.pdf

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
```

**Header:** basic description; sizes of chromosomes; other

http://samtools.sourceforge.net/SAMv1.pdf

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1
```

**Sequence name:** whatever the read was named in fastq file

http://samtools.sourceforge.net/SAMv1.pdf

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1 113 chr1 497
```

**Genomic position:** position how it was mapped by a mapper

http://samtools.sourceforge.net/SAMv1.pdf

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1 113 chr1 497 37 37M
```

**CIGAR string:** contains information about mapping like gaps.

http://samtools.sourceforge.net/SAMv1.pdf

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1 113 chr1 497 37 37M 15 100338662 0
CGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAG
```

**Sequence by itself:** mapped sequence

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1 113 chr1 497 37 37M 15 100338662 0
CGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAG
0;==-==9;>>>>>=>>>>>>>>>>=>>>>>>>>>>
```

**Quality:** quality from FASTQ file

http://samtools.sourceforge.net/SAMv1.pdf

# File formats: SAM/BAM

```
@HD VN:1.0
@SQ SN:1 LN:249250621
Sequence1 113 chr1 497 37 37M 15 100338662 0
CGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAG
0;==-==9;>>>>>=>>>>>>>>>>>=>>>>>>>>>
XT:A:U NM:i:0 SM:i:37 AM:i:0 X0:i:1 X1:i:0 XM:i:0
```

**Other info:** for example, NM:i:0 means 0 mismatches

http://samtools.sourceforge.net/SAMv1.pdf

# File conversion

- ## SAMTOOLS

```
$ module load samtools/0.0.19
```

```
Program: samtools (Tools for alignments in the SAM format)
Version: 0.1.18 (r982:295)

Usage:     samtools <command> [options]

Command:  view          SAM<->BAM conversion
          sort          sort alignment file
          mpileup       multi-way pileup
          depth         compute the depth
          faidx         index/extract FASTA
          tview         text alignment viewer
          index         index alignment
          idxstats      BAM index stats (r595 or later)
          fixmate       fix mate information
          flagstat      simple stats
          calmd         recalculate MD/NM tags and '=' bases
          merge         merge sorted alignments
          rmdup         remove PCR duplicates
          reheader      replace BAM header
          cat           concatenate BAMs
          targetcut     cut fosmid regions (for fosmid pool only)
          phase         phase heterozygotes
```

# File conversion

- SAM<->BAM conversion

```
$ samtools view -b -S reads.sam > reads.bam
```
[samopen] SAM header is present: 25 sequences.

```
$ samtools view -h reads.bam | less
```

```
$ samtools sort
```
Usage: samtools sort [-on] [-m <maxMem>] <in.bam> <out.prefix>

```
$ samtools sort reads.bam sorted
```

cells

cDNA
ChIP

sequencer

Sequenced reads

Alignment

genome

read coverage

Gene expression

# Aggregation



**Genomic position**

**Genomic position**

# File formats: WIG/bigWIG

Aggregated information about genomic locations.

**VARIABLE Step**

variableStep chrom=chr2
300701 12.5
300702 12.5
300703 12.5
300704 12.5
300705 12.5


is equivalent to:

variableStep chrom=chr2 span=5
300701 12.5

bigWIG contains the same info but in compressed, binary format.

# File formats: WIG/bigWIG

Aggregated information about genomic locations.

**FIXED Step**

fixedStep chrom=chr3 start=400601 step=100
11
22
33

displays the values 11, 22, and 33 as single-base regions on chromosome 3 at positions 400601, 400701, and 400801, respectively. Adding span=5 to the declaration line:

fixedStep chrom=chr3 start=400601 step=100 span=5
11
22
33

# Tools

- Bedtools (BED manipulation but BAM support is available)

https://bedtools.googlecode.com/files/BEDTools-User-Manual.v4.pdf

```
$ module load bedtools/2.17.0

$ head a.bed
chr1   100   200
chr1   1000   2000

$ head b.bed
chr1   150   250

$ intersectBed -a a.bed -b b.bed
chr1   150   200
```

# Tools

- Bedtools

| Utility | Description |
|---|---|
| intersectBed | Returns overlapping features between two BED/GFF/VCF files. *Also supports BAM format as input and output.* |
| windowBed | Returns overlapping features between two BED/GFF/VCF files within a "window". *Also supports BAM format as input and output.* |
| closestBed | Returns the closest feature to each entry in a BED/GFF/VCF file. |
| coverageBed | Summarizes the depth and breadth of coverage of features in one BED/GFF file (e.g., aligned reads) relative to another (e.g., user-defined windows). *Also supports BAM format as input and output.* |
| genomeCoverageBed | Histogram or a "per base" report of genome coverage. *Also supports BAM format as input and output.* |
| pairToBed | Returns overlaps between a BEDPE file and a regular BED/GFF/VCF file. *Also supports BAM format as input and output.* |
| pairToPair | Returns overlaps between two BEDPE files. |
| bamToBed | Converts BAM alignments to BED and BEDPE formats. *Also supports BAM format as input and output.* |
| bedToBam | Converts BED/GFF/VCF features (both blocked and unblocked) to BAM format. |
| bedToIgv | Creates a batch script to create IGV images at each interval defined in a BED/GFF/VCF file. |
| bed12ToBed6 | Splits BED12 features into discrete BED6 features. |
| subtractBed | Removes the portion of an interval that is overlapped by another feature. |
| mergeBed | Merges overlapping features into a single feature. |
| fastaFromBed | Creates FASTA sequences from BED/GFF intervals. |
| maskFastaFromBed | Masks a FASTA file based upon BED/GFF coordinates. |
| shuffleBed | Permutes the locations of features within a genome. |
| slopBed | Adjusts features by a requested number of base pairs. |
| sortBed | Sorts BED/GFF files in useful ways. |
| linksBed | Creates an HTML links from a BED/GFF file. |
| complementBed | Returns intervals not spanned by features in a BED/GFF file. |
| overlap | Computes the amount of overlap (positive values) or distance (negative values) between genome features and reports the result at the end of the same line. |
| groupBy | Summarizes a dataset column based upon common column groupings. Akin to the SQL "group by" command. |
| unionBedGraphs | Combines multiple BedGraph files into a single file, allowing coverage/other comparisons between them. |
| annotateBed | Annotates one BED/VCF/GFF file with overlaps from many others. |

# File formats: bedGraph

| CHR | Start | End | Value |
|-----|-------|-----|-------|
| chr1 | 1 | 100 | 100 |
| chr1 | 200 | 300 | 20 |

```
$ genomeCoverageBed -bg -ibam sorted.bam -g mm10.chrom.sizes >
out.bg
```

track type=bedGraph name=test_track

# File formats: BED

Information about genomic locations. TAB delimited.

- Simple case:

chr1     5          10



Area covered by this BED line

5          6          7          8          9          10

Position on chr1

- Advanced (12 columns, details in the link on slide9):

chr22 100 500 cloneA 96 + 100 500 0 2 4,5, 0,35
chr22 200 600 cloneB 90 - 200 600 0 2 4,5, 0,36

# File formats: GTF

Gene information (usually from some external databases)

TAB separated information

<seqname> <source> <feature> <start> <end> <score> <strand> <frame> [attributes] [comments]

AB000381 Twinscan  CDS         700  707  .   +  2  gene_id "001"; transcript_id "001.1";
AB000381 Twinscan  start_codon 380  382  .   +  0  gene_id "001"; transcript_id "001.1"

# Visualization

- UCSC browser

  http://genome.ucsc.edu

  http://biocore.umassmed.edu/ucsc.html

- ENSEMBL

  http://www.ensembl.org/index.html

- IGV browser

  http://www.broadinstitute.org/igv

# Visualization

- UCSC – handles bunch of different things

# Visualization

**http://biocore.umassmed.edu/ucsc.html**

       to use for  visualization, any other activities where you would expect heavy competition on ucsc.edu.

**http://genome.ucsc.edu/index.html**

       to use for tools that are not under heavy competition like BLAT, datatables download.

# Visualization

# Visualization

# Visualization

# Visualization

# Visualization

**http://genome.ucsc.edu/index.html**



AGCGAATTGGAATGACCTAACATTTCTGTGACATCT

# Visualization/Databases



Try to download Human hg19 RefSeq genes with "fields selection".

# Databases

- NCBI (http://www.ncbi.nlm.nih.gov)

# Databases

- ## NCBI download     The Gene Expression Omnibus (GEO):    GSE44690

# Databases

- ## NCBI download

# Databases

- ## NCBI download

| Supplementary file | Size | Download | File type/resource |
|---|---|---|---|
| GSM1087281_201.mm.wild.type.42.rep1.chipAMYB.peaks.bed.gz | 1.1 Mb | (ftp)(http) | BED |
| SRX/SRX244/SRX244353 | | (ftp) | SRA Experiment |

http://www.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software

```
$ module load sratoolkit/2.3.4-2
$ illumina-dump data.sra
```

BED format is already described

www.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software

NCBI    Site map    All databases    Search

Sequence Read Archive

Main  Browse  Search  Download  Submit  Documentation  Software

Software  XML Schema  Toolkit Documentation

**SRA Toolkit**

# Questions?

# Homework

Use information from today's seminar and get instructions from this location:

Server: ghpcc06.umassrc.org

Folder/file: /project/umw_biocore/seminar/Step2.docx