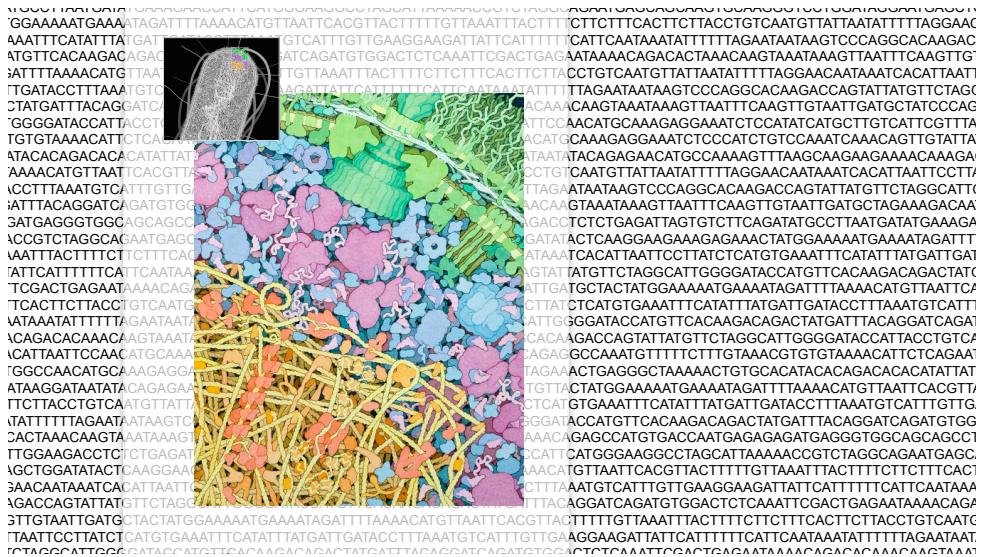


Genomics I

Biomedical Data Science: Mining and Modeling
CB&B 752 • MB&B 452

Matt Simon
Jan 31, 2022



What is genomics?

1. The **global** study of how biological **information** is encoded in genome sequence

Genes
Regulatory sequences
Genetic variation

2. How this information is **read out** to produce distinct **biological outcomes**

Gene expression and regulation
Cellular identity, differentiation and development
Phenotypic variation among individuals and species

In practice, many experiments that involve **deep sequencing** are considered genomics.

Overview

Genomics I: (today's lecture): Focus on sequencing technology and genomes.

Genomics II: (Wednesday's lecture): Focus on applications of sequencing technology.

Overview

- Sequencing data: from wet lab to fastq.
- Applications to studying genomes and much much more.
- * Sophisticated use of data from genomics requires an integrated understanding of the biological experiment, sample preparation and down stream computational analyses of the data.

Importance of genomics data: these data are central to most biomedical and biological

The screenshot shows the header of a Nature journal article titled 'Decoding gene regulation in the fly brain'. It includes author information: Jasper Janssens, Sara Albar, Ibrahim Ihsan Taskiran, Joy N. Irmaili, Alicia Estacio Gomez, Gabriel Aughey, Katalin I. Searle, Florian V. De Ros, Carmen Bravo Gonzalez-Blas, Marc Dione, Krista Grimes, Xiao-Jiun Qu, Dafni Pasavakzat, Gert Hulsemans, Samira Makhzani, Maxime De Weegeenier, Valérie Christiaens, Tony Southall & Stein Aerts. The article was published in Nature 601, 630-636 (2022). It also shows metrics: 17K Accesses | 231 Altmetric | Metrics. The abstract discusses the use of single-cell transcriptome analysis and electron microscopy mapping to reveal gene regulatory networks in the fly brain.

Methods: scRNA-seq, scATAC-seq, DamID, CUT&Tag, ChIP-seq

Data availability

The data generated for this study have been deposited in NCBI's Gene Expression Omnibus and are accessible through GEO Series accession numbers GSE163697 and GSE181494 (DGPR lines). We also provide a dedicated website to browse the results of the analyses and processed data (<https://flybrain.aertslab.org/>), which provides link-outs to the SCope session (<http://scope.aertslab.org/>) (Fly Brain), UCSC hub (<http://genome.ucsc.edu/cgi-bin/hgTracks?hgDb=Dm6HubU1!http://uscstracks.aertslab.org/papers/flybrain/hub.txt>), the eGRNs in NDX, the DeepExplainer plots of enhancers and other information. The following online databases were used: FlyBase (<https://flybase.org/>), FlyMine (<https://www.flymine.org/flymine>), icis-Target (<https://phioned.kuleuven.be/apps/icb-fcb-icisTarget/>), Flylight (<https://web.iapella.org/cgi-bin/flew.cgi>), CIS-BP (<http://cispb.cbr.uottawa.ca/>), ENCODE (<https://www.encodeproject.org/>). ENCF0704WCH. The following publicly accessible datasets were also used: GSE107451 (scRNA-seq adult brain), GSE157202 (scRNA-seq larval brain), GSE101581 (scATAC-seq embryo). The neural network is from Özel et al².

<https://www.nature.com/articles/s41586-021-04262-z.pdf>

Raw data can be found in genomics databases

The screenshot shows a search results page from the NCBI GEO database. The search query is 'COVID-19 is an emerging, rapidly evolving situation. Get the latest medical information from CDC. COVID-19 ClinicalTrials.gov' and 'Find NCBI SARS-CoV-2 Repository, sequence, and clinical context' with URL <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=ZS499>'. Below the search bar, it says 'More... GEO: View most recent entries for this dataset'. There are two main sections: 'Results (2)' and 'Summary'. The first result is for 'Public on Aug 10, 2021' and 'Decoding gene regulation in the fly brain'. It includes a 'View details' button and a 'GO to Experiment' link. The second result is for 'Status' and 'Title' (Decoding gene regulation in the fly brain) with a 'GO to Experiment' link. The 'Summary' section is expanded, showing detailed information about the experiment, including 'Detailed description', 'Contact', 'Experiment type', 'Sample', 'Series', 'Platform', 'Overall design', and 'Platforms'. It also lists 'GEO accessions' and 'Sample ID'. At the bottom, there is a 'Download' button.

- Most journals require authors to submit their data to a database (e.g., GEO) prior to publication.
- These databases entries contain raw data and processed data.
- These data can be used to examine the authors' claims, but also to test new hypotheses.

What is the output from an Illumina sequencing experiment?

One read (fastq format)

```
#HWI-ST1239:178:HOKPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCCTGTAGGCTGCCTAGCCCTCCCTGCAGGGTAAGTGGGAGGAAGAGCAGAGGACTTAGTGGGCCTCCCAGGG
+
#1=DDFFFFHHHHHIJIIJJIIJJJJJJJIIJ?FHIDGIJ=GIHIIHHIJIHEHHIHGFFFFEEDDDDDDDDDDDDDDD
```

1. Read identifier
2. **Sequence**
3. Quality score identifier “+”
4. Quality score

Central questions

Where do these data come from?

How does the way we collect it influence what we know?

Which best describes your experience with analysis of sequencing data?

(A) I have no relevant experience with DNA sequencing data.

(B) I've read/thought about DNA sequencing data but never worked with it.

(C) I've worked with some DNA sequencing data.

(D) I've worked with a lot of DNA sequencing data.

Which best describes your wet lab experience?

(A) I have never conducted research that requires molecular biology.

(B) I've done a lot of molecular biology (cloning, etc.) but only worked with Sanger sequencing.

(C) I've generated at least one deep sequencing data set.

(D) I've done a lot of deep sequencing.

Workflow

1. Isolation of sample.

e.g., Isolate DNA and shear.



2. Library preparation

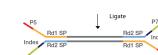
e.g., Add known sequences to the ends.

3. Sequencing

e.g., Illumina Novaseq

4. Analysis

e.g., Map to genome and interpret.



Metrics for evaluating sequencing technology

Throughput:

- Number of high quality bases per unit time
- Number of independent samples run in parallel
- Difficulty of sample preparation

Yield

- Number of useful reads per sample
- Read length

Cost

- Per run cost
- Per base cost
- Equipment
- Reagents
- Labor
- Analysis

Quality

- Accuracy per base

What is sequencing?

One-at-a-time methods		Short read deep sequencing				Long read deep sequencing			
		a. Maxam-Gilbert Sequencing b. Sanger Sequencing	a. Illumina Sequencing b. Ion Torrent	Read length (kb)	Read accuracy (%)	Throughput per flow cell (Gb)	Estimated cost per Gb (US\$)	Maximum throughput per year (Gb) ^a	
Sequencing technology	Platform	Data type	N50	Maximum	Mean	Maximum			
Pacific Biosciences (PacBio)	RS II ^b	CLR	5-15	>60	7.75-15	2	333-933 ^c	4,380	
	Sequel	CLR	25-50	>100	87-92	6-10	20	98-195 ^d	
	Sequel II	CLR	30-60	>200	50-100	160	12-26 ^e	17,520	
Oxford Nanopore Technologies (ONT)	HIFI	10-20	>20	>99	15-30	35	43-88 ^f	93,440	
	MinION/GridION	Long	10-60	>1,000	2-20	30	50-500 ^g	10,200	
	PromethION	Ultra-long	100-200	>1,500	87-98	0.5-2	2.5	21,950 (MinION) 108,500 (GridION)	
Illumina	NextSeq 550	Long	10-60	>1,000	50-100	180	21-42 ^h	3,153,600	
	NextSeq 550	Single-end	0.075-0.15	0.15	16-30	>30	50-63 ⁱ	>47,782	
	Novaseq 6000	Paired-end	0.075-0.15 (>2)	0.15 (>2)	32-120	>120	40-60 ^j	>70,080	
Illumina	Novaseq 6000	Single-end	0.05-0.25	0.25	>99.9	6E-3,000	>3,000	10-35 ^k	>1,194,545
	Novaseq 6000	Paired-end	0.05-0.25 (>2)	0.25 (>2)					

The technology will change, but your need to critically understand the input and output will not.

Logsdon (2020) *Nat Rev Genetics*

The steps of sequencing experiments

1. Sample preparation

- a. Isolation
- b. Library construction

2. Sequencing

- a. Flow cell loading
- b. Cluster generation
- c. Sequencing
- d. Processing image files
- e. De-multiplexing samples

3. Data analysis

- a. Read filtering
- b. Alignment to a genome
- c. Diverse analyses

[Yale Center for Genome Analysis \(YOGA\)](#)

Next-Gen Sequencing ▾ Bioinformatics Microarrays ▾ Services & Fees ▾ Mendelian Center ▾ About YCGA ▾

Illumina Next-Generation Sequencing

Oxford Nanopore Sequencing	FFPE RNA library prep	\$229	\$298
PeBio SMRT Sequencing	Analysis	\$402	\$505
Affymetrix Microarrays	Consultation per hour	\$265	\$341
Illumina Microarrays	HiSeq2500 paired-end 2x75 sequencing lane	\$1,970	\$2,386
DNA Extraction & RNA Isolation	HiSeq2500 single-end 1x75 sequencing lane	\$1,581	\$1,920
10X Genomics LD prep	HiSeq2500 paired-end 2x150 sequencing lane	\$2,620	\$3,166
MissionBio Nanostaining and Parallel Sequencing	MiSeq 500 cycle	\$1,692	\$2,052
Sample Submission	MiSeq 600 cycle	\$2,030	\$3,034
	NovaSeq S1 paired-end 2x100 sequencing lane	\$3,527	\$4,254
	NovaSeq S4 paired-end 2x100 sequencing lane	\$5,297	\$6,370
	NovaSeq SP paired-end 2x150 sequencing lane	\$2,609	\$3,154
	NovaSeq S4 paired-end 2x150 sequencing lane	\$5,297	\$6,379

Retrieved Jan 30, 2022:
<https://medicine.yale.edu/keck/yoga/services/illumina/prices/>

What is the most raw form of data recorded in an Illumina sequencing experiment?

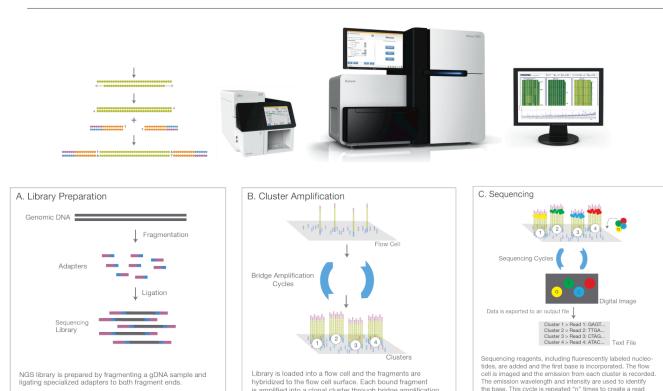
(A) A chromatogram.

(B) A string of letters.

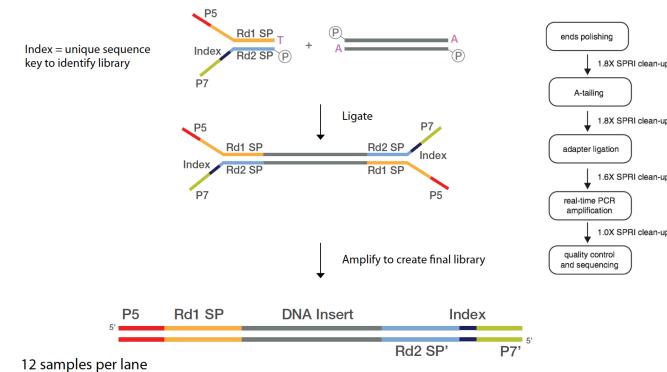
(C) A series of images.

(D) A readout of genomic locations.

Where do these reads come from?



Optional: Library preparation using ligation



Potential sources of bias:

1. Selective PCR amplification (issue of duplicates).
2. Size selection.
3. Enzyme specificities.

Optional: Library preparation using tagmentation

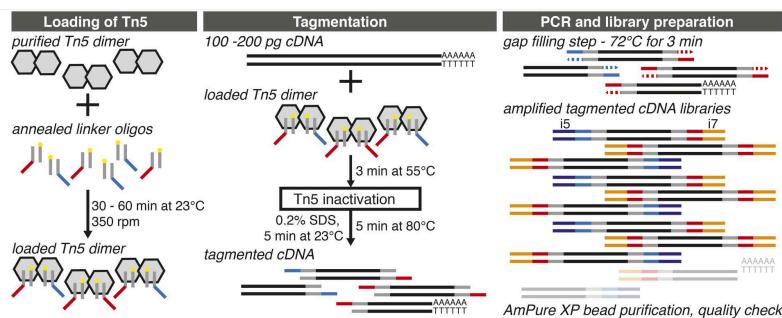
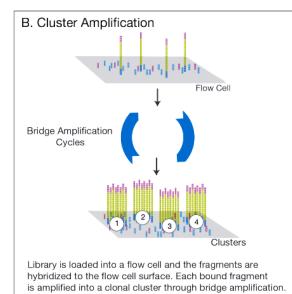


Figure from: Hennig BP, Velten L, Racke I, Tu CS, Thoms M, Rybin V, Besir H, Remans K, Steinmetz LM, G3 (Bethesda). 2018 Jan 4;8(1):79-89. PMID: 29118030.

Cluster amplification.

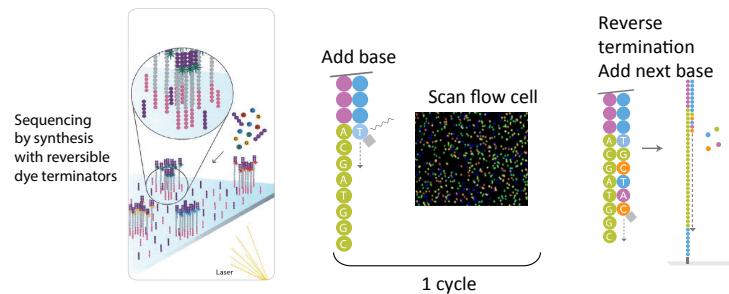


- Separate each individual molecule (randomly).
- Give each molecule an address (spatial location).
- Pack as many on as possible but avoid overlaps.



https://www.illumina.com/content/dam/illumina-marketing/documents/products/illumina_sequencing_introduction.pdf

Sequencing by synthesis



https://www.illumina.com/content/dam/illumina-marketing/documents/products/illumina_sequencing_introduction.pdf

How long are the reads?

TATTGCAATATGTTAACATCTAACAGGAAAAAATACCCACACAAAACAAACACAACCCCTTAGAACTGTGCTG
← →
75 nt

While there are other technologies that can give longer read lengths, Illumina reads are generally 50 nt - 250 nt

What limits the insert size and read length?

One read (fastq format)

```
@HWI-D00306:498:HBB89ADXX:1:1101:1180:1882 1:N:0:CGATGT
NCATCACTTCTGCAACCACCATGACGTCAATCTTCGCCAACCCAAACTCGAGATCGGAAGAGCACACGTCTG
+
#1EBBDDDFDFBFFFFIIIIIIIIIIIFEGIIIIFIGAGIIFIII=FEEEEFFFDDD=@9A@BBBBB=?BB<
```

- For each single end read: Incomplete incorporation of bases.
- For the size of the insert (especially for paired end analysis): Ability to get consistent clusters.

What is the output from an Illumina sequencing experiment?

Paired read (fastq format)

```
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGCGCTGCCCTGCCAGGGTAAGTGGGAGGAGAGAGCAGAGGGACTTAGTGGGCCTCCCCAGGG
+
#1=DDFFFHHHHHIJIIJJJJJJJJJJ?FHIDGIJ=GIHIIHGIJIHEHIIHHGFFFFEEEDDDDDDDDDDD
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCTGCAGTCCTCGGCCCTGTGTTAGACCCAGAACTAGGGTGCCAGGTACCACTTAATCCTT
+
##4<@0000000000?@0@??????@0?@????????????????>?????????@>????@0@?@?????
```

- Read identifier
 - Instrument
 - Flow cell
 - Read ID
 - Coordinates
 - Which read from a paired end sample
 - Which index for multiplexed read
- Sequence
- Quality score identifier “+”
- Quality score

What is the output from an Illumina sequencing experiment?

Many reads...

```
@HWI-D00306:498:HB889ADXX:1:1101:1180:1882 1:N:0:CGATGT
NCATCACTTCTGCACCAGCCATGACGTCATACTCGTCCGAACCCCCAAACTCGAGATCGGAAGAGCACACGCTCTG
+
#11BBDDDFDFBFFFIIIIIIIIIIIIIFEGIIIIIFIGAGIIFIII=FEEEEFFFD@9A@BBBBB=?BB<
@HWI-D00306:498:HB889ADXX:1:1101:1167:1902 1:N:0:CGATGT
TATTGCAATATGTTAACAACTAAACAAGGAAAAAAATACCCACACAAAACACAACACAAACCTTAGACTGTCGTG
+
@HWI-D00306:498:HB889ADXX:1:1101:1190:1928 1:N:0:CGATGT
ACCAAGCCACAAATAAGTTAGTGTTCCTAGTACATGCTGAGTTATTTGATCCCCTATCTATACACTGCTACTGTC
+
@<@DDDDDBCDDDG>?2<AFFBCCEEHEIEGHIEGEIID@CDGFEEFIDGCFCDABFG>FBFGFGLIEIFFFD
@HWI-D00306:498:HB889ADXX:1:1101:1157:1931 1:N:0:CGATGT
CTGAGATTCTTGGCCATAGTCCTTAACCACCTACGCCAACCTGCCGTGGTTGCCCTCTCGATCG
+
CCFFFFFHIIHHHHHIIJJIIJJIIIGHHIIJGGJIGIJJJJJJIIIIIJJIIJJIIJJJJIIJJGJJHCHFBDDFFDDECB
```

Generally ~ 2,000,000,000 reads/sequencing lane

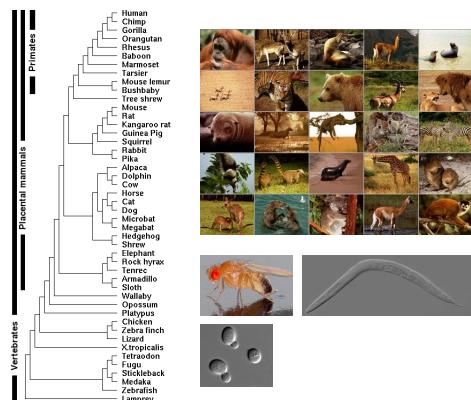
Note: This is for an Illumina NovaSeq with current chemistry, but this number changes

What do I do with my sequencing reads?



Source: Slate via Noonan

Many reference genomes are available



A 75 nt sequencing read matches to a reference genome perfectly, except for one mismatch. What might have caused this?

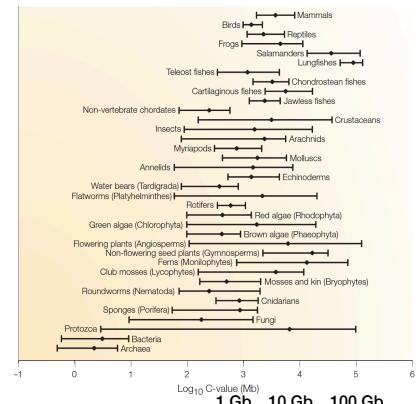
- (A)** An error introduced during library preparation.
- (B)** An error in a base call during sequencing.
- (C)** A single nucleotide difference between the genome of the biological sample and the reference genome.
- (D)** Any of the above.

There is a wide range of genome sizes.

kb = 1000 bp
 Mb = 1×10^6 bp
 Gb = 1×10^9 bp
 Tb = 1×10^{12} bp

Human haploid genome ~ 3 Gb

75 nt $\times 3 \times 10^8$ reads/lane is about the right scale, but the amount of **coverage** necessary depends on application.



Sequencing of the human genome

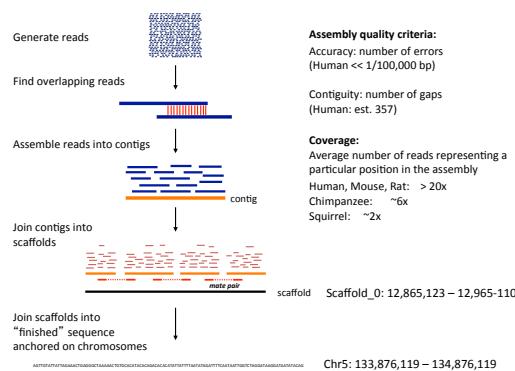
Victory declared 2003



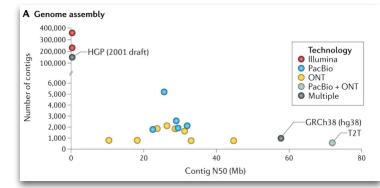
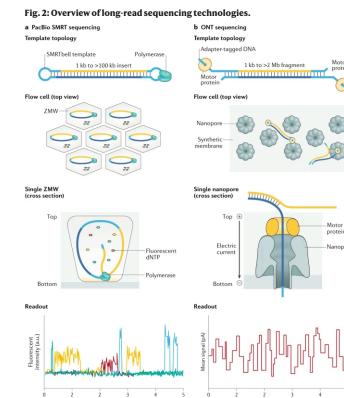
- Industrialization of Sanger sequencing, library construction, sample preparation, analysis, etc.
- \$3 billion total cost
- 1 Gb/month at largest centers (2005)

Novaseq 2 billion reads 2x150 bp. \$5000 \rightarrow <\$100/genome.

How to assemble a genome



The importance of long read sequencing



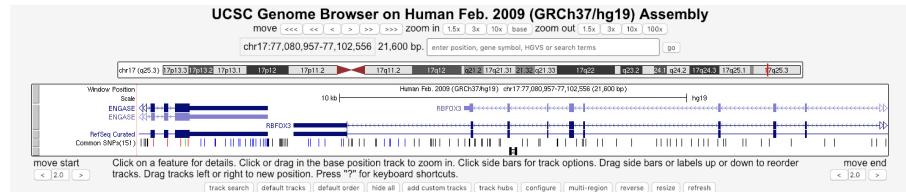
Logsdon (2020) *Nat Rev Genetics*

Example of a genome browser track (UCSC)

Chr5: 133,876,119 – 134,876,119

Our specific example:

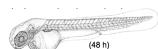
```
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGGGTAGCCCTCCCTGCAGGGTAAGTGGGAGGAGAGCAGAGGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHLLIJJJJIIJJJJJJJ#FHIDGIJ=GIGIIIIHGIIJIHEHHIHGFFFFEEEDDDDDDDDDDD
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCGTCAGTCCTCGGTCCCTGTTAGACCAAGACTAGGTGCCAGGCCAGGTACCACTAATCCTT
+
##4<????????????????????????????????????????????@>????@?@???????
```



Workflow

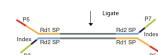
1. Isolation of sample.

e.g., Isolate DNA and shear.



2. Library preparation

e.g., Add known sequences to the ends.



3. Sequencing

e.g., Illumina Novaseq



4. Analysis

e.g., Map to genome and interpret.



Using sequencing to annotate the genome

1. Where are the cis-acting regulatory elements in DNA?

- A. DNase I hyper-sensitivity mapping (DNase-Seq).
- B. FAIRE to map regulatory elements.

2. Where do transcription factors bind?

- C. ChIP-seq of transcription factors (or in high res, ChIP-exo)
- D. Nucleosome mapping (MNase-Seq).

3. Where are different histone modifications found?

- E. ChIP-Seq of histone modifications.
- F. ChIP-Seq of chromatin writers, readers and erasers.

4. Where is RNA polymerase transcribing?

- G. ChIP-Seq of polymerase.
- H. GRO-Seq, NET-Seq and TT-Seq to measure RNA in the polymerase active site..

5. How is the genome organized in 3D?

- I. 4C/5C/Hi-C to measure chromatin conformation.

Applications of sequencing technology next week.

Conclusions

- Sequencing technology is central to our understanding of biology.
 - The decrease in cost and increase in throughput make sequencing data increasingly ubiquitous.
- * Sophisticated use of data from genomics requires an integrated understanding of the biological experiment, sample preparation and down stream computational analyses of the data.