

Genomics I

Biomedical Data Science: Mining and Modeling
CB&B 752 · MB&B 452
Matt Simon
Jan 23, 2023

The image is a composite graphic. On the left side, there is a grid of small colored squares, resembling a DNA microarray or a sequencing chip. In the center, there is a 3D molecular model of a protein complex, with various subunits colored in shades of green, blue, pink, and yellow. On the right side, there is a DNA double helix structure, colored in shades of green and blue. The background of the entire image is a light blue pattern of overlapping circles.

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

Article | Published: 21 December 2022

BRD8 maintains glioblastoma by epigenetic reprogramming of the p53 network

[Xueqin Sun](#), [Olaf Klingbeil](#), [Bin Lu](#), [Caizhi Wu](#), [Carlos Ballon](#), [Meng Ouyang](#), [Xiaoli S. Wu](#), [Ying Jin](#), [Yon Hwangbo](#), [Yu-Han Huang](#), [Tim D. D. Somerville](#), [Kenneth Chang](#), [Jung Park](#), [Taemoon Chung](#), [Scott K. Lyons](#), [Junwei Shi](#), [Hannes Vogel](#), [Michael Schulner](#), [Christopher R. Vakoc](#) & [Alea A. Mills](#) 

Nature **613**, 195–202 (2023) | [Cite this article](#)

9178 Accesses | 216 Altmetric | [Metrics](#)

Abstract

Inhibition of the tumour suppressive function of p53 (encoded by *TP53*) is paramount for cancer development in humans. However, p53 remains unmutated in the majority of cases of glioblastoma (GBM)—the most common and deadly adult brain malignancy^{1,2}. Thus, how p53-mediated tumour suppression is countered in *TP53* wild-type (*TP53*^{WT}) GBM is unknown. Here we describe a GBM-specific epigenetic mechanism in which the chromatin regulator bromodomain-containing protein 8 (BRD8) maintains H2AZ occupancy at p53 target loci through the EP400 histone acetyltransferase complex. This mechanism causes a repressive chromatin state that prevents transactivation by p53 and sustains proliferation. Notably, targeting the bromodomain of BRD8 displaces H2AZ, enhances chromatin accessibility and

Methods: ChIP-seq, ATAC-seq, RNA-seq

Data availability

The ChIP-seq and RNA-seq data generated in this study is available at the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>) with accession number GSE158551. The RNA-seq data of GBM cells isolated from patient specimens are from datasets GSE84465 and GSE121720 in the GEO database.

<https://www.nature.com/articles/s41586-022-05551-x>

Raw data can be found in genomics databases

Experiment attributes:

GEO Accession: GSM4802676

Links:

NCBI link: [NCBI Entrez \(gds\)](#)

Runs: 1 run, 43.5M spots, 3.3G bases, 1.1Gb

Run	# of Spots	# of Bases	Size	Published
SRR15101033	43,549,074	3.3G	1.1Gb	2022-11-10

```
[me2598@c16n06 ~]$ head -50 SRR13288692.fastq
@SRR13288692.1.1 length=86
GGATTATTTTACCAATTTTCTTTTACGTGCTGAAACAGCAGACAGTCCCGAGTGTGGCCAATCTNNNNACTNNNNNNN
+SRR13288692.1.1 length=86
GAA/AEEEEEEEEEE/EAEEEEEEEEEEEEEEEE/ABEE/-EEE/EE-/E/EEEE//EEEE#####G#####E#
@SRR13288692.2.2 length=86
GGGTACACAGACCAAGTTAATTTCAGCGACGCGGAAACACTGTCTCTTATACACATCTCCGAGCCANNNNAGANNNNN
+SRR13288692.2.2 length=86
AAAAEEEEEEEEEEEAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####E#####E#
@SRR13288692.3.3 length=86
GACATTACACATTACAGGTTTGAAGCGGCGCAGCAGCACACACACACAGCTGTCTTATACANNNTTGGNNNNCGN
+SRR13288692.3.3 length=86
A/AAAEAEAEAEAEAEAEAE/EEEE/E//EEEEEEEEEEEEEE/EE/E6EA/E/EAEEEE//#####EA#####E#
@SRR13288692.4.4 length=86
ACGTGGGCTCTCTCCGATGGCACTACGCCCTGTCGGATCTCGGATCAGACCTTCCGCTGGGANNNTCTANNANNAN
+SRR13288692.4.4 length=86
AAAAEEEEEEEA/AAEE/EAEEEE<EEAAE</EABAE/EE/EAEEA<EAE/EEEA/EE//#####E/#####E#
@SRR13288692.5.5 length=86
CCGTACTCTCTCTCCCGGGTGGCGCCGATCAGCGGCGAGCGGTACCCGTGACATGGCGGCCNNNCGATNCCNCGTH
+SRR13288692.5.5 length=86
A/AAAE/AAE/EAEE/<E/EEEA//E//E/N//EAE//N/EEA/E/N/<E/EEA/6/<66<#####E/#####E#
@SRR13288692.6.6 length=86
ACTCCACTGATGTTTCAAAGGCCCAAGTTGAAATATCGATATCATCGATGATATATCGANNNGATGNNCTNGAN
+SRR13288692.6.6 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####E#####E#
@SRR13288692.7.7 length=86
GAGGTACTCAGACTGGCCCGCGCTACGGTGTGACACCGCTCCGGAACCGCTATGTACCANNNTTGNNTCNATN
+SRR13288692.7.7 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####E#####E#
@SRR13288692.8.8 length=86
NFTCAGTACAGTACGATACGGCGACGATACAAAATAACAGCTGTGGAATATGTTTCAAANNNTCATNNGNFTN
+SRR13288692.8.8 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####E#####E#
@SRR13288692.9.9 length=86
CTCCACACAAAGCTTACACTCTGTCTAGTGTGATGACCCGATAGCTACTGTGTTCTTAAANNAGANNNAATTN
+SRR13288692.9.9 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####E#####E#
```

- 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 use 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:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGCGTAGCCTCCCTGCAGGTAAGTGGGAGGAGAGAGACAGAGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFHHHHHIJJIJJJIJJJIJJJ?FHIDGIJ=GIGHI I IHGIJIHEHIHHGFFFEEDDDDDDDDDDD
```

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

Central questions

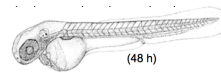
Where do these data come from?

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

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

Cost

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

Yield

- Number of useful reads per sample
- Read length

Quality

- Accuracy per base

What is sequencing?

One-at-a-time methods

- Maxam-Gilbert Sequencing
- Sanger Sequencing

Short read deep sequencing

- Illumina Sequencing
- Ion Torrent

Long read deep sequencing

- Nanopore based
- Pacific Bioscience Sequencing

Sequencing technology	Platform	Data type	Read length (kb)		Read accuracy (%)	Throughput per flow cell (Gb)		Estimated cost per Gb (US\$)	Maximum throughput per year (Gb) ^a
			N50	Maximum		Mean	Maximum		
Pacific Biosciences (PacBio)	RS II ^b	CLR	5–15	>60	87–92	0.75–1.5	2	333–933 ^c	4,380
	Sequel	CLR	25–50	>100		5–10	20	98–196 ^d	17,520
	Sequel II	CLR	30–60	>200	50–100	160	13–26 ^e	93,440	
Oxford Nanopore Technologies (ONT)	MiniION/GridION	Long	10–60	>1,000	87–98	2–20	30	50–500 ^f	21,900 (MiniION) 109,500 (GridION)
		Ultra-long	100–200	>1,500		0.5–2	2.5	500–2,000 ^f	913 (MiniION) 4,563 (GridION)
	PromethION	Long	10–60	>1,000		50–100	180	21–42 ^g	3,153,600
Illumina	NextSeq 550	Single-end	0.075–0.15	0.15	>99.9	16–30	>30	50–63 ^h	>47,782
		Paired-end	0.075–0.15 (x2)	0.15 (x2)		32–120	>120	40–60 ^h	>70,080
	NovaSeq 6000	Single-end	0.05–0.25	0.25		65–3,000	>3,000	10–35 ^h	>1,194,545

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

- Isolation
- Library construction

2. Sequencing

- Flow cell loading
- Cluster generation
- Sequencing
- Processing image files
- De-multiplexing samples

3. Data analysis

- Read filtering
- Alignment to a genome
- Diverse analyses

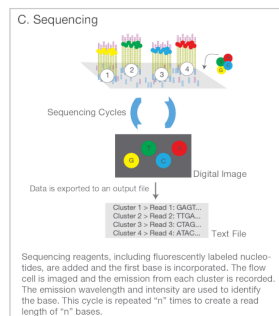
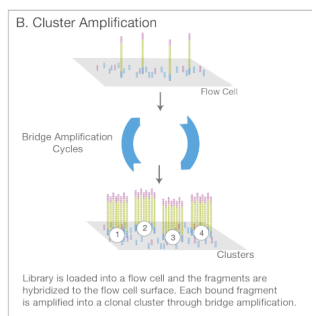
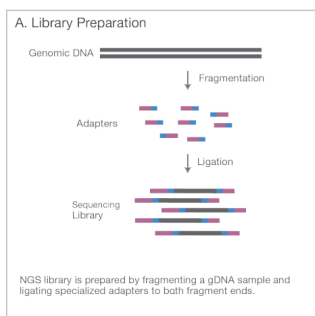
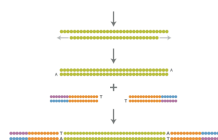
Yale Center for Genome Analysis (YCGA) [MENU](#)

RNA Library prep	2019	2017
cDNA extraction	\$56	\$78
Whole genome methylation(NEB enzymatic)	\$70	\$96
RNA Library prep (poly A selection)	\$77	\$105
RNA library prep (smRNA)	\$164	\$218
RNA Library prep (ribosomal depletion)	\$164	\$219
RNA Library prep (low input)	\$178	\$236
RNA Library prep (FFPE)	\$250	\$331
Analysis	\$508	\$666
Consultation per hour	\$527	\$690
MiSeq 500 Cycle	\$1708	\$2,227
NovaSeq S1 2x100	\$3,412	\$4,442
NovaSeq S1 2x150	\$4,025	\$5,238
NovaSeq SP 2x150	\$2,526	\$3,290
NovaSeq S4 2x150	\$4,568	\$5,944

Retrieved Jan 23, 2023:

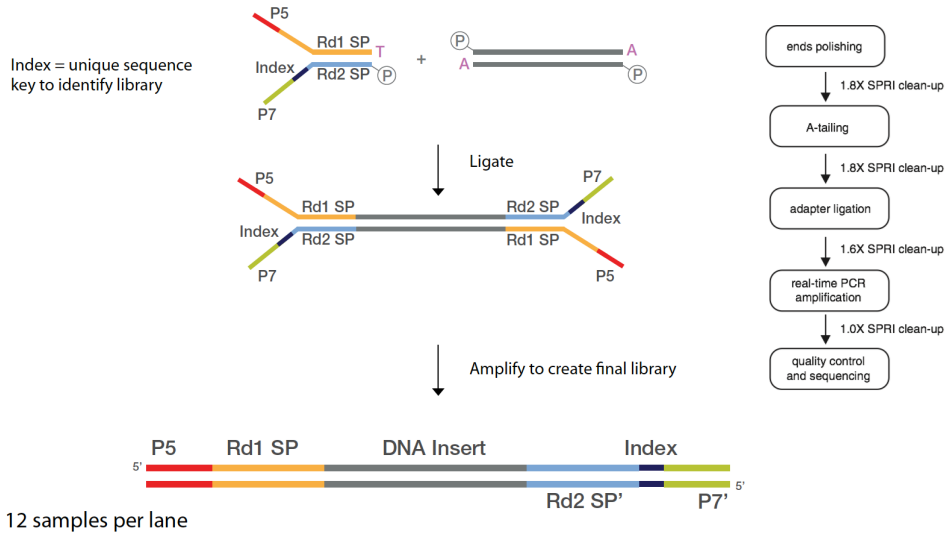
<https://medicine.yale.edu/keck/ycga/services/illuminaprices/>

Where do these reads come from?



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

Optional: Library preparation using ligation

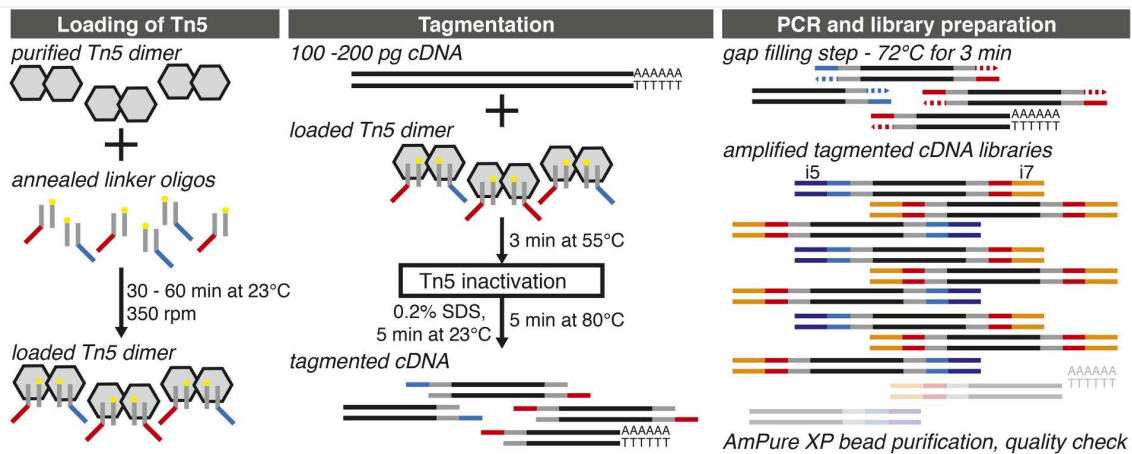


Potential sources of bias:

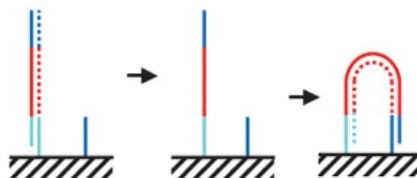
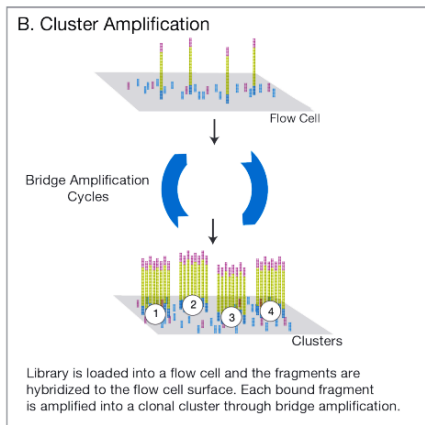
1. Selective PCR amplification (issue of duplicates).
2. Size selection.
3. Enzyme specificities.

Challenging but possible to analyze pg quantities of DNA. (In humans, ~6 pg DNA/cell).

Optional: Library preparation using tagmentation



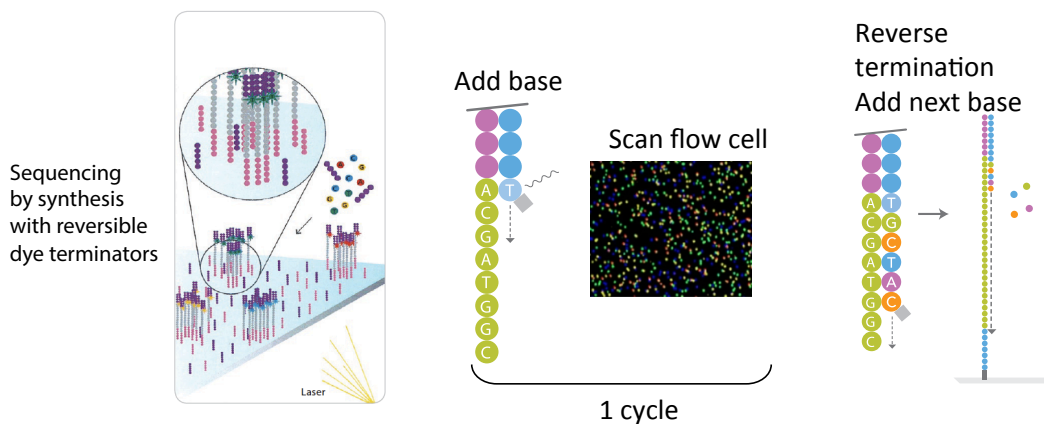
Cluster amplification.



- 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



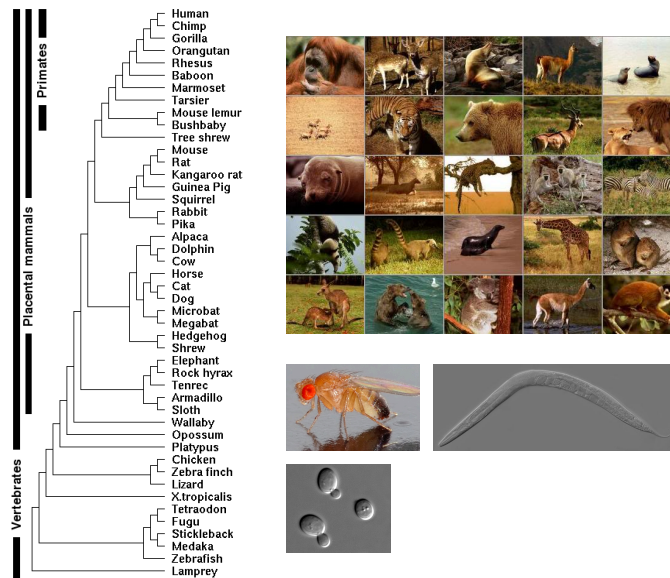
https://www.illumina.com/content/dam/illumina-marketing/documents/products/illumina_sequencing_introduction.pdf

What do I do with my sequencing reads?



Source: Slate via Noonan

Many reference genomes are available

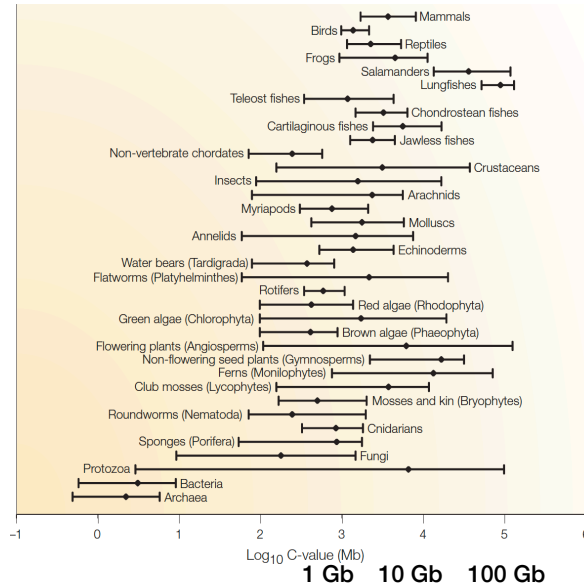


There is a wide range of genome sizes.

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 x 3×10^8 reads/lane is about the right scale, but the amount of **coverage** necessary depends on application.



Sequencing of the human genome

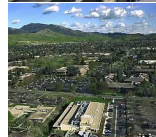
Victory declared **2003**



- Industrialization of Sanger sequencing, library construction, sample preparation, analysis, etc.

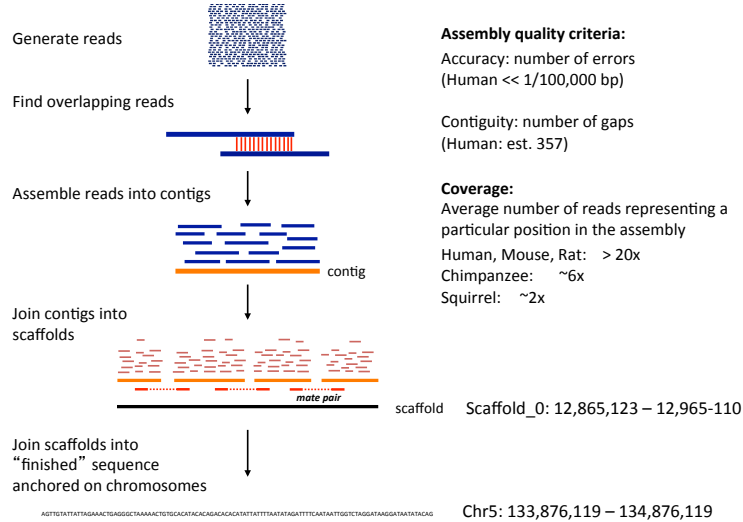
- \$3 billion total cost

- 1 Gb/month at largest centers (2005)



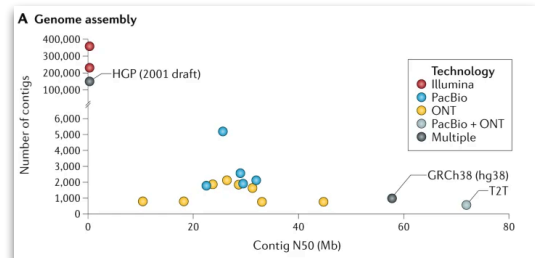
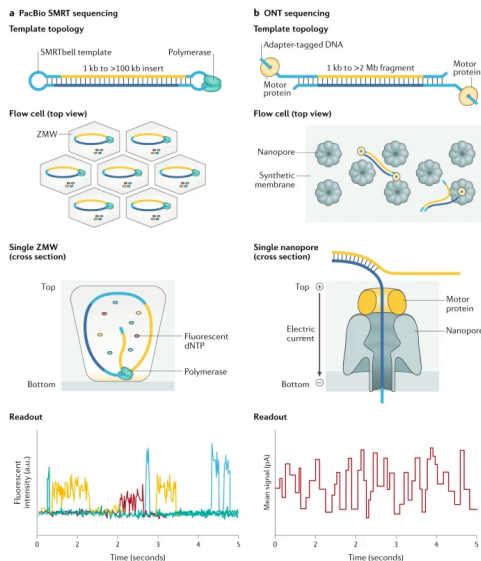
Novaseq 2 billion reads 2x150 bp. \$5000 -> <\$100/genome.

How to assemble a genome



The importance of long read sequencing

Fig. 2: Overview of long-read sequencing technologies.



What types of annotation do we have/want?

~3 billion bp

```

AGCAATTAATCAGTAAATTCCTTATCTCAGTGTGAATTTGADATTTTATGATTC
ATACCTTTAAAGTGCATTTGTTGAGGAGAGATATTTCATTTTTCATTCGAT
AAATATTTTAAAGTAAATAGTCCAGGACAAAGACAGATATATGTTCT
AAGCATTGGGATACCACTTCCAGAGAGAGATATGATTTAGAGAGAC
AGATGTGGACTCTCAAAATTCGAGTGGAGATAAAGACAGACTAAGCAG
TAAATTAAGGTTTAAATTCAGCTTGTATTTGATGTCGGAAGACATGAAACA
GAGCCATGTGACCAATGAGAGAGATGAGGGTGGACAGCCTGTTTAA
GATAGAGTCCCTGTTGGTGGATGAGAGACCTCTCTGAGATAGTGT
CTTCAGATGCCTTAAATGATGAAAGAACCATTCATGGGAAGCCCTAG
CATTAAAGCCCTTAGGCGATGAGAGAGAGGTCAGAGGCTGCTGG
ATAGGATGAGCTGGATATACCAAGGAAAGAAAGAACTATGGAAA
ATGAATAGATTTTAAAGACATTTTAACTAGCTTACTTTTGTAAATTA
CTTTCTCTTCACTTCTTACCTGTCAATGTTATTAATTTTAAAGACA
ATAAACGATTAATTCCTTACCTGATGTAAATTCATTTATATGATGATA
GCTTAAATGTCATTTGAGGAGATATTCATTTTTCATTAAGAAA
TATTTTAAAGATAAGTCCAGGACAAAGACAGATATATGTTCTAGG
CATTGGGATACCACTTCCAGAGAGAGATATTTAGAGAGACAGAT
GTGACTCTCAAAATTCGACTGAGATAAAGACAGACTAAGACAGATAAT
AAGCTTAAATTCAGGTTTAAATGATGATGAGTGGTAAAGATAAGATA
TTTAAAGATGTAATTCAGCTTACTTTTGTAAATGATTTTCTCTT
CACTTCTTACCTGCTGATTTATTAATTTTAAAGATAAATGAGAT
AATTCCTATCTCATGTGAATTTCAATTTATGATGATACCTTAAATGT
GTTTGTGAGGAGAGATATTCATTTTTCATTAATTTTAAAGATA
ATAAAGTCCAGGACAAAGACAGATATATGTTTAGGCAATGGGAT
AGCATGTCAGAGAGAGATATTCAGATTCAGATGATGATGATGATGAT
AAATTCGACTGAGATAAAGACAGACAAAGATAAATGATTAATTT
CAAGTTGATTTGATGCTATCCAGGACAGACCA...
    
```

Genes:

- Coding, noncoding, miRNA, etc.
- Isoforms
- Expression

Genetic variation:

- SNPs and CNVs

Sequence conservation

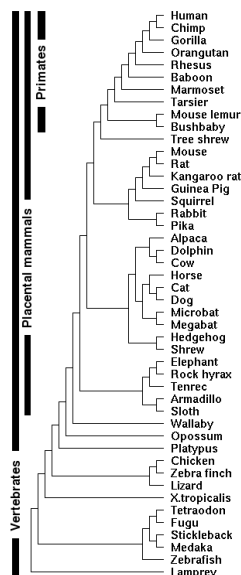
Regulatory sequences:

- Promoters
- Enhancers
- Insulators

Epigenetics:

- DNA methylation
- Chromatin

Degrees of genomic annotation vary widely



ENCODE and modENCODE

Human, Mouse (Fly, Worm, Yeast):

- Chromosome assemblies
- Dense gene and regulatory maps, variation, etc.

Other models (Dog, Chicken, Zebrafish):

- Chromosome assemblies
- Partial gene maps; variation; little regulatory data

Low coverage vertebrate genomes:

- Scaffold assemblies
- Few annotated genes
- Used for comparative purposes

Where do you look for existing annotations?

UCSC Genome Browser (genome.ucsc.edu):

Visualization, data recovery, simple analysis
(also <http://genome-preview.ucsc.edu/>)

ENSEMBL (ensembl.org):

Visualization, data recovery, simple analysis

Integrative Genomics Viewer

(broadinstitute.org/software/igv/):

Local genome viewer (visualize local and remote data)

Galaxy (main.g2.bx.psu.edu):

Complex data analysis and workflows

Example of a genome browser track (UCSC)

Chr5: 133,876,119 – 134,876,119

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 class.

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.