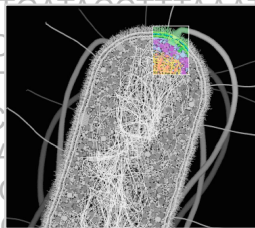


Genomics I

Biomedical Data Science: Mining and Modeling
CB&B 752 • MB&B 452
Matt Simon
Jan 27, 2025



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

Molecular Cell Supports

Meet the author

View all

Most read (last 30 days)

Article • Open Access

The ribotoxic stress response drives acute inflammation, cell death, and epidermal thickening in UV-irradiated skin *in vivo*

Vind et al.

Published online: November 25, 2024

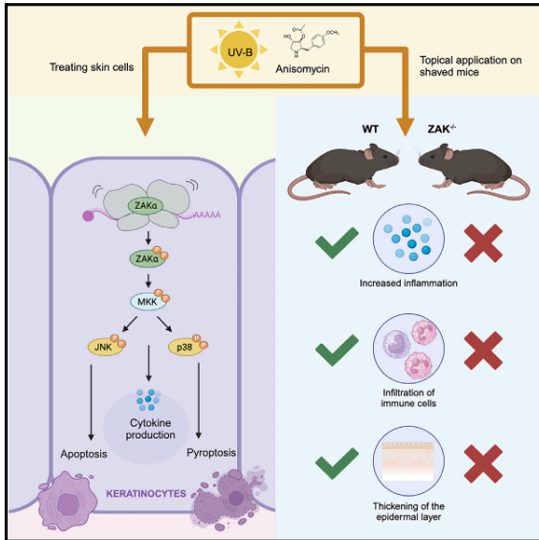


Article

Molecular Cell

The ribotoxic stress response drives acute inflammation, cell death, and epidermal thickening in UV-irradiated skin *in vivo*

Graphical abstract



Authors

Anna Constance Vind, Zhenzhen Wu, Muhammad Jasrie Firdaus, ..., Mads Gyrd-Hansen, Franklin L. Zhong, Simon Bekker-Jensen

Correspondence

vind@sund.ku.dk (A.C.V.), franklin.zhong@ntu.edu.sg (F.L.Z.), sbj@sund.ku.dk (S.B.-J.)

In brief

The acute skin reaction to sunburn, encompassing keratinocyte cell death, inflammation, and epidermal thickening, has traditionally been ascribed to DNA damage responses. Using a mouse model deficient for the *Zak* gene, Vind et al. demonstrate that these reactions depend on sensing of cytoplasmic RNA damage rather than nuclear DNA damage.

Data and code availability

- RNA sequencing raw data files have been deposited in NCBI's Gene Expression Omnibus (GEO) archive, with the accession code GEO: GSE251957.

Scope: Self | Format: HTML | Amount: Quick | GEO accession: GSE251957 | GEO

Series GSE251957 | Query DataSets for GSE251957

Status: Public on Jul 03, 2024

Title: Effects of ZAKalpa on UVB-dependent transcriptional changes in N/TERT-1 cells

Organism: Homo sapiens

Experiment type: Expression profiling by high throughput sequencing

Summary: This study investigates the effect of ZAKalpa kinase in UVB-triggered

Reads: 26,220,058 reads

Page: 10000 / 2623006 | quality scores | advanced options

Reads (separated)

SR82731488.99991.1 A00552.129.HSKH.LDS2:11172.24587.20957 Biological (Biological)

name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 24587, y: 32957

SR82731488.99992 name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 28903, y: 32957

SR82731488.99993 name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 28366, y: 32957

SR82731488.99994 name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 28930, y: 32957

SR82731488.99995 name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 32959, y: 32957

SR82731488.99996 name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 32344, y: 32957

SR82731488.99997 name: A00552.129.HSKH.LDS2:11172.24587.20957

member: GAGTCGTT-CTCCGCTAG

x: 2175, y: 32957

ZIP/Postal code: 308232

Country: Singapore

Platforms (1): GPL20301 Illumina HiSeq 4000 (Homo sapiens)

Samples (12): GSM7990569 N/TERT control Sham #1

More...: GSM7990570 N/TERT control Sham #2

GSM7990571 N/TERT control Sham #3

[https://www.cell.com/molecular-cell/fulltext/S1097-2765\(24\)00884-0](https://www.cell.com/molecular-cell/fulltext/S1097-2765(24)00884-0)

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

```
[ms2598@c16n06 ~]$ head -50 SRR13288692.fastq
@SRR13288692.1 1 length=86
GGATTATTTTACCATTTCCTTTTACGTGTGAAAACAGCAGACATCGCCAGTGTGGCCAACTCTNNNNNACTNNNNNNN
+SRR13288692.1 1 length=86
6AA/AEEEEEEEEAE/EAAB6EEEEEEEEEEEEEEEE/E/EE/EE<A/EEEE//EEEE#####6E6#####E#
@SRR13288692.2 2 length=86
GGGTGAACAGAGCCAAAGTTAATTTGACGCACGGCGAACACTGTCTCTTATACACATCTCCGAGCCANNNGAAGNNNNN
+SRR13288692.2 2 length=86
AAAAA6EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEA#####AEE#####E#
@SRR13288692.3 3 length=86
GACATTACATTACGAGTTTAGAAGCGGGCGCAGCAGCCACAACAGCAACAGCTGTCTTATACACANNNTTGGNNANNCN
+SRR13288692.3 3 length=86
A/AAEEAEAE6EE6EEE/EEEE/E//EEEEEEEEEEEEEE/EE/E6EA/E/EAEEEE/E/#####A#E#E#E#E#E#
@SRR13288692.4 4 length=86
ACGTGGTCTCTCTCCGATGGCAACTACGCCCTGTCGGATCTGGGATCAGACCTTCGCTGTGGANNNTTANNANNGAN
+SRR13288692.4 4 length=86
AAAA6EEEAEEEA/AEE/EEAE#####EEAAE</EAEE/EEAEA<<EAE/EEEA/EE//#####/E/#E#E#//
@SRR13288692.5 5 length=86
CCCGTACTCGCTCGCCGGGTGGCCCGATCAGCGGCAGCGGTACCCGTGACGATGGCGGCCNNNCGATNNCNGTN
+SRR13288692.5 5 length=86
A/AAAA/AEE/EAEE/</<EEEA//E//EE/A/EAE//A/EEA/E/A/<E/EEA/6/</66<#####/E/E#E#E#E#E#
@SRR13288692.6 6 length=86
AGTCCAACTGATGGTTCAAAGGCCGCAAGTTGAAATATCGATATCATGCAGATGATATCGANNNGATGNNCTNGAN
+SRR13288692.6 6 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####AEE#E#E#E#E#
@SRR13288692.7 7 length=86
GAAGGTAGTCCACAACACGCCACGGCTACGGTGTGACCACGGCTCTGGAACGGCTATGTACGANNNTTGNNTCNATN
+SRR13288692.7 7 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####AEE#E#E#E#E#
@SRR13288692.8 8 length=86
ATTGAGGTACAGATACAGATACACGGCAGGATACAAAAATACAGCTGTGGAATATTGTGTTTCAANNNTCATNNGGNTN
+SRR13288692.8 8 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEAEEAEEEEEEEEE#####AEE#E#E#E#E#
@SRR13288692.9 9 length=86
CTGTACACAAAGCACTTACTCTGTCTAGTGATGACCCGATAAGCTAGCTTGGTGTCTTAANNNGAAGNNAATTGN
+SRR13288692.9 9 length=86
AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#####AAAE#E#E#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  
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGGAGAGAGAGCAGAGGGACTTAGTGGGGCTCCCCAGGG  
+  
#1=DDFFFHHHHHIJIJJJIJJJJJIJJJ?FHIDGIJ=GIHGIIIHGIJIHEHIHHGFFFFEEEDDDDDDDDDDD
```

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 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.

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.

Workflow

1. Isolation of sample.

e.g., Isolate DNA and shear.

2. Library preparation

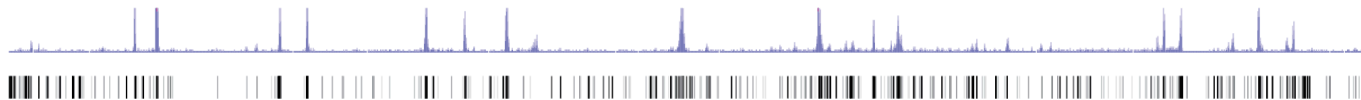
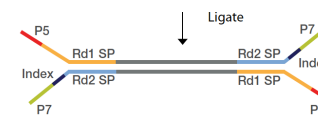
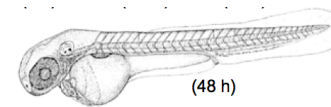
e.g., Add known sequences to the ends.

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

- a. Maxam-Gilbert Sequencing
- b. Sangar Sequencing

Short read deep sequencing

- a. **Illumina Sequencing**
- b. Ion Torrent

Long read deep sequencing

- a. **Nanopore based**
- b. **Pacific Bioscience Sequencing**

Spatial -omics

- a. **DNA probe based**
- b. **Antibody probe based**

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	98–195 ^d	17,520
	Sequel II	CLR	30–60	>200		50–100	160	13–26 ^e	93,440
		HiFi	10–20	>20	>99	15–30	35	43–86 ^e	10,220
Oxford Nanopore Technologies (ONT)	MinION/GridION	Long	10–60	>1,000	87–98	2–20	30	50–500 ^f	21,900 (MinION) 109,500 (GridION)
		Ultra-long	100–200	>1,500		0.5–2	2.5	500–2,000 ^f	913 (MinION) 4,563 (GridION)
	PromethION	Long	10–60	>1,000		50–100	180	21–42 ^f	3,153,600
Illumina	NextSeq 550	Single-end	0.075–0.15	0.15	>99.9	16–30	>30	50–63 ^g	>47,782
		Paired-end	0.075–0.15 (×2)	0.15 (×2)		32–120	>120	40–60 ^g	>70,080
	NovaSeq 6000	Single-end	0.05–0.25	0.25		65–3,000	>3,000	10–35 ^h	>1,194,545
		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*

Bressan (2023) *Science*

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

Sequencing

Service	Yale Rate	Non-Yale Rate
MiSeq 500 Cycle	\$1,905	\$2,485
NextSeq Usage	\$1,039	\$1,358
NovaSeq X Plus 25B 2x150	\$2,936	\$3,825
NovaSeq SP 2x150	\$2,617	\$3,410

...per 1.25 billion reads!

Retrieved Jan 27, 2025:

<https://medicine.yale.edu/genetics/research/ycga/service-fees/>

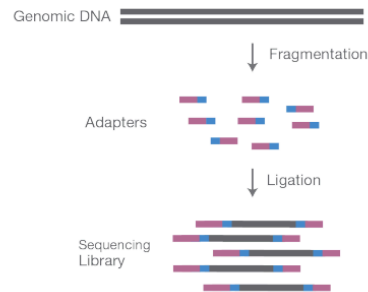
What is the most raw form of data recorded in an Illumina sequencing experiment?

- (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?

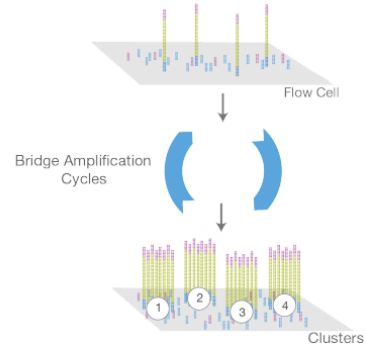


A. Library Preparation



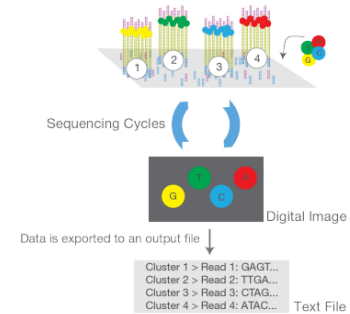
NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

B. Cluster Amplification



Library is loaded into a flow cell and the fragments are hybridized to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

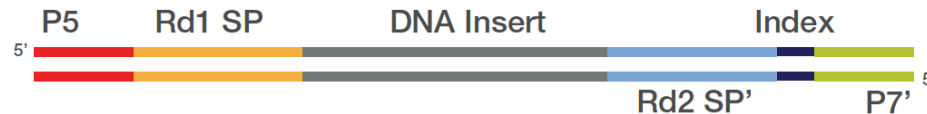
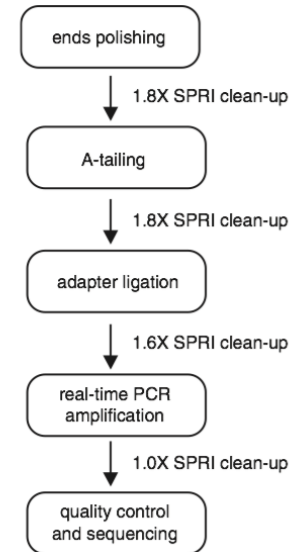
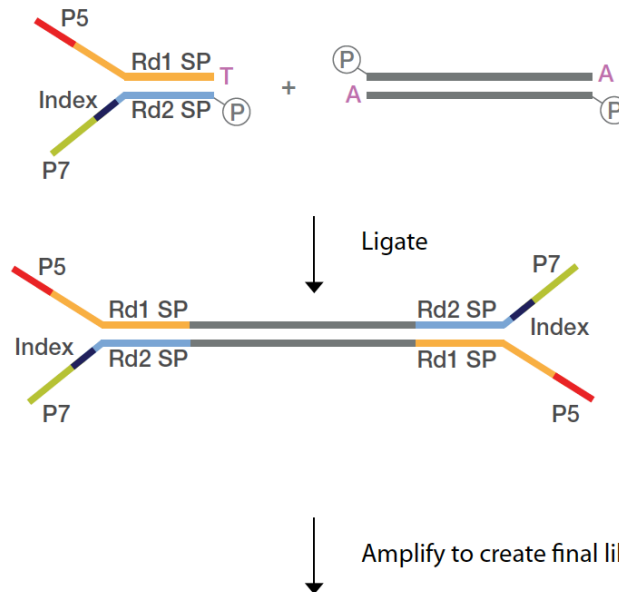
C. Sequencing



Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.

Optional: Library preparation using ligation

Index = unique sequence
key to identify library



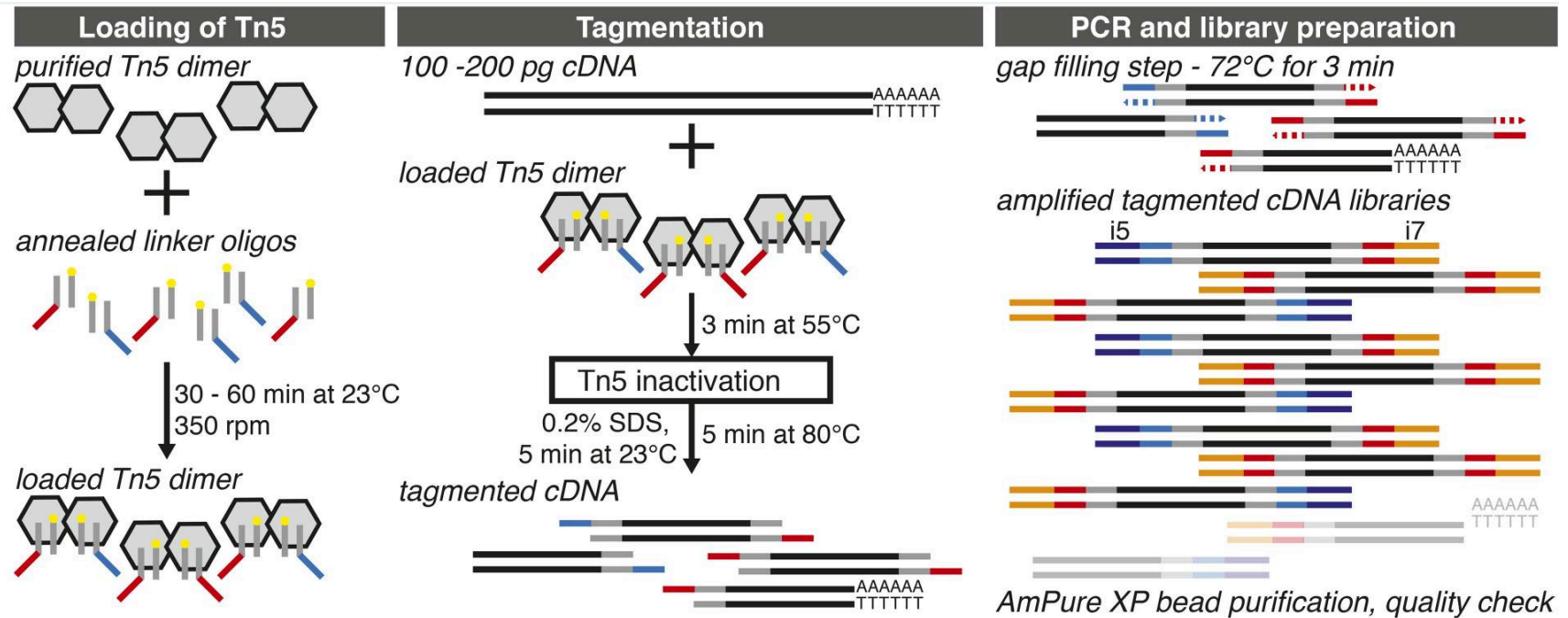
12 samples per lane

Potential sources of bias:

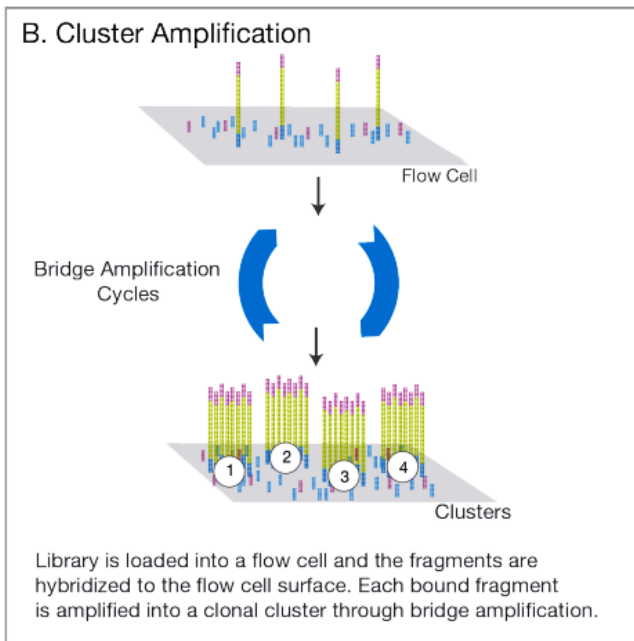
1. Selective PCR amplification (issue of duplicates).
2. Size selection.
3. Enzyme specificities.

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

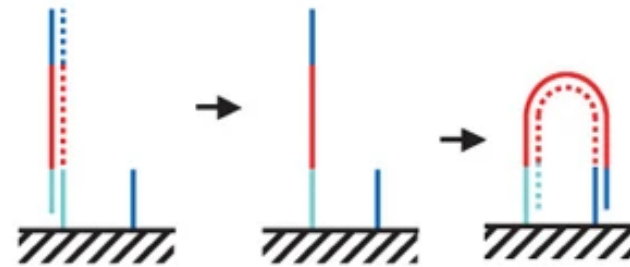
Optional: Library preparation using tagmentation



Cluster amplification.



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



Vol 456 | 6 November 2008 | doi:10.1038/nature07517 nature

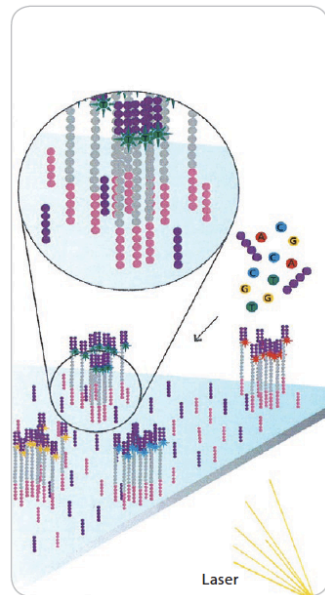
ARTICLES

Accurate whole human genome sequencing using reversible terminator chemistry

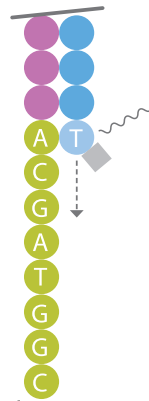
A list of authors and their affiliations appears at the end of the paper

Sequencing by synthesis

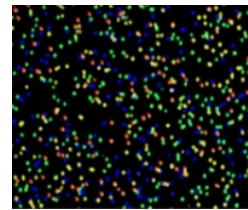
Sequencing by synthesis with reversible dye terminators



Add base

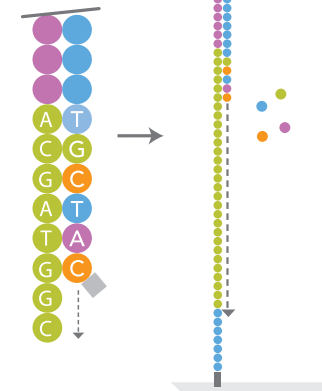


Scan flow cell



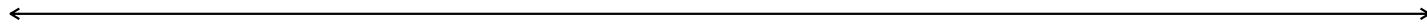
1 cycle

Reverse termination
Add next base



How long are the reads?

TATTGCAATATGTTAACAATCTAACAAGGAAAAAATACCCACACAAAACAAAACACAACCCTTAGAACTGTGCTG



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
NCATCACTTTCTGCACCAGCCATGACGTCAATCTTCGTCCGAACCCCAAACCTCGAGATCGGAAGAGCACACGTCTG
+
#11BBDDDFDFBFFFIIIIIIIIIIIIIIIFEGIIIIIFIGAGIIFIII=FFFFFFFDDDD=@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
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGGAGAGAGAGCAGAGGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHIJIJJJIJJJJJIJJJ?FHIDGIJ=GIHGIIIHGIJIHEHIHHGFFFFEEDDDDDDDDDDDDD

@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCTGCAGTCCTCGGTCTGTGTTAGACCAGAACTAGGTGCCCAGGCCAGGTACCACCTAATCCTT
+
##4<@@@@@@@@@?@@@@?@@?????@@??@????????????????????>???????????@>????@@@@?@@??????
```

1. Read identifier
 - a. Instrument
 - b. Flow cell
 - c. Read ID
 - d. Coordinates
 - e. Which read from a paired end sample
 - f. Which index for multiplexed read
2. Sequence
3. Quality score identifier “+”
4. Quality score

What is the output from an Illumina sequencing experiment?

Many reads...

```
@HWI-D00306:498:HBB89ADXX:1:1101:1180:1882 1:N:0:CGATGT
NCATCACTTTCTGCACCAGCCATGACGTCAATCTTCGTCCGAACCCCAAACCTCGAGATCGGAAGAGCACACGTCTG
+
#11BBDDDFDFBFFFIIIIIIIIIIIIFEGIIIIIFIGAGIIFIII=FEFFFFFFDDDD=@9A@BBBBB=?BB<

@HWI-D00306:498:HBB89ADXX:1:1101:1167:1902 1:N:0:CGATGT
TATTGCAATATGTTAACAATCTAACAAGGAAAAAATACCCACACAAAACAAAACACAACCCTTAGAACTGTGCTG
+
B@@FFDFHFHHHJJJIJIGIIJJJJJIJJHFIJJJJJIJJJEHHJJIJJJJIJJJJJJGHHHHFBDFFFE>CEEC
@HWI-D00306:498:HBB89ADXX:1:1101:1190:1928 1:N:0:CGATGT
ACCAAGCCACAATAAGTTAGTGTTCATAGTACATGCTGAGTTATTTGATCCCGTATCTATACTGCTACTGTC
+
@<@DDDDD8CDDDGE?2<AFFBCCEEHEIEGHIIEGEIDD@CDGFFFEFIDGCFCDABFG>FBFGFGIEIFFDDDD
@HWI-D00306:498:HBB89ADXX:1:1101:1157:1931 1:N:0:CGATGT
CTGAGATTCTTTGCCATAGTCCTTAACCACTACGCAACTGCAACCAACCACCTTCCGTGGTTTGCCTCTCGATCG
+
CCFFFFFFHHHHHHIJJJIJJJIIGHHIJGGJIGIJJJJJJJIJJIIJJJIJJJIJJJJJHCHFBDFFFDDECB
```

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

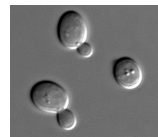
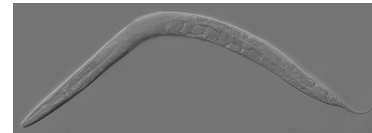
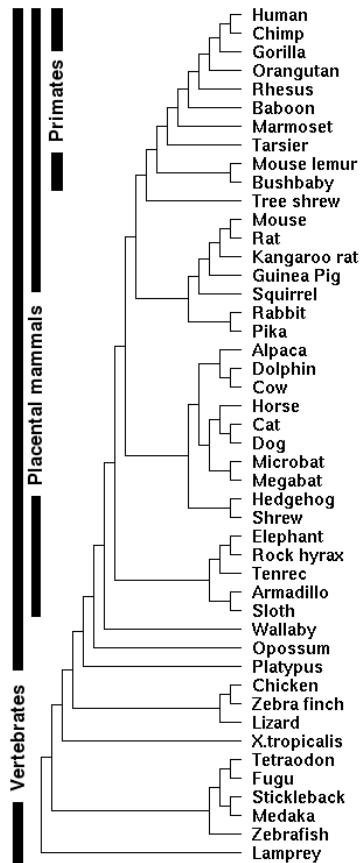
Note: This is for an Illumina NovaSeq with current chemistry, but this number changes

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?

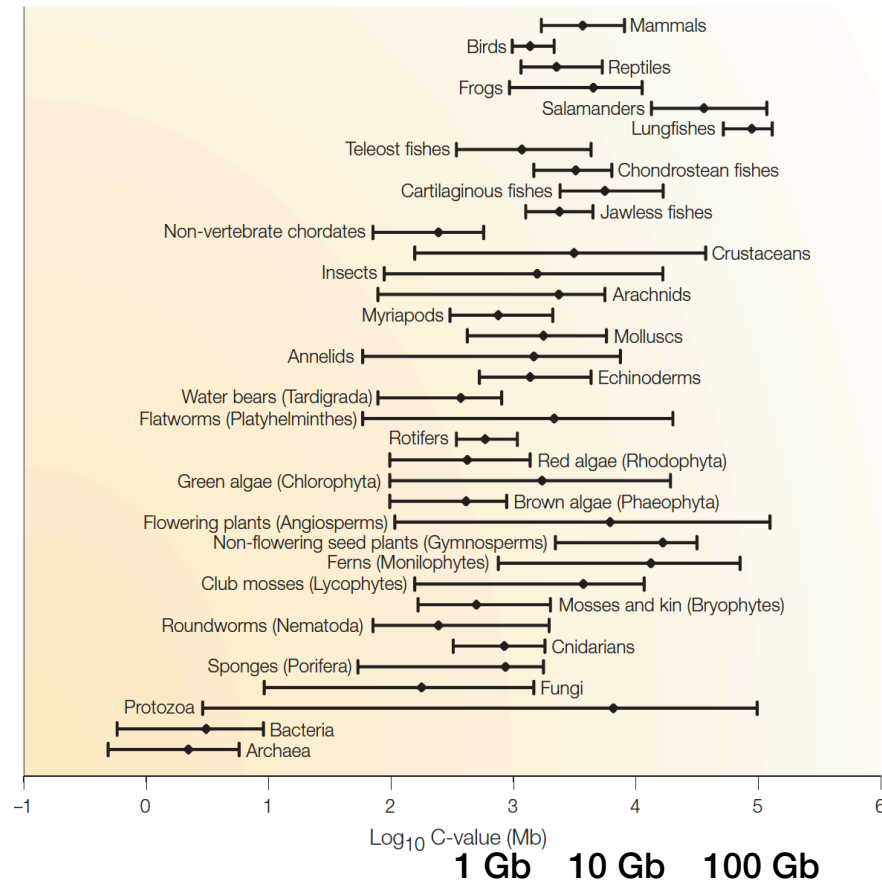
- (A)** An error introduced during library preparation.
- (B)** An error in a base call during sequencing.
- (C)** A single nucleotide difference between the genome of the biological sample and the reference genome.
- (D)** Any of the above.

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 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**



National Human
Genome Research
Institute

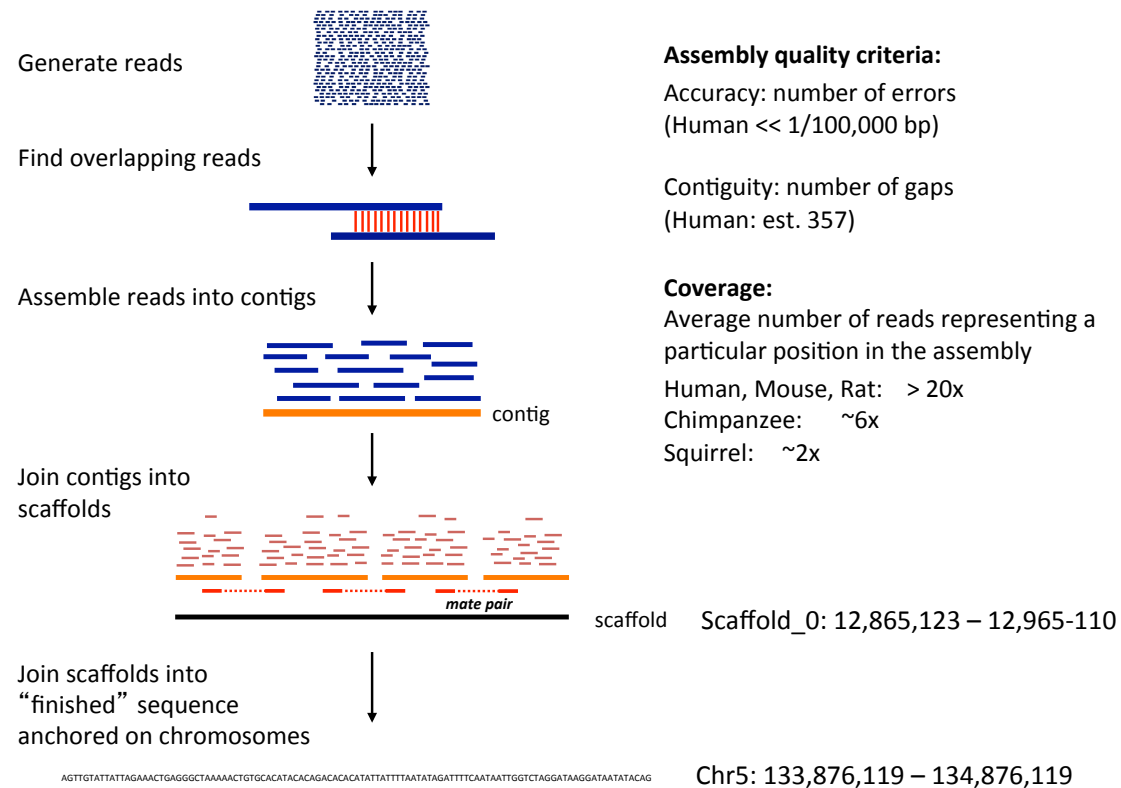


- Industrialization of Sanger sequencing, library construction, sample preparation, analysis, etc.
- \$3 billion total cost
- 1 Gb/month at largest centers (2005)



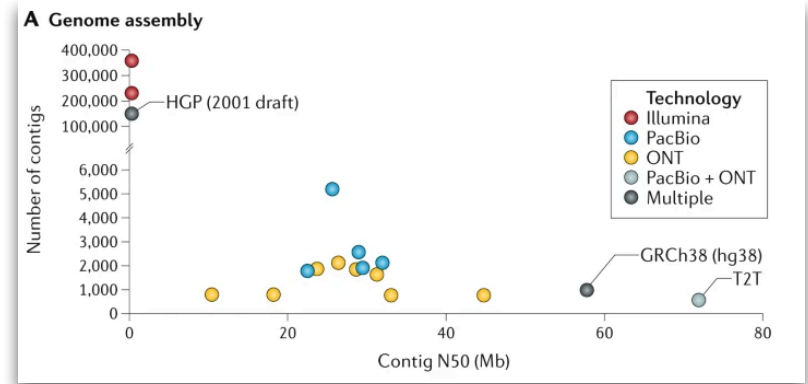
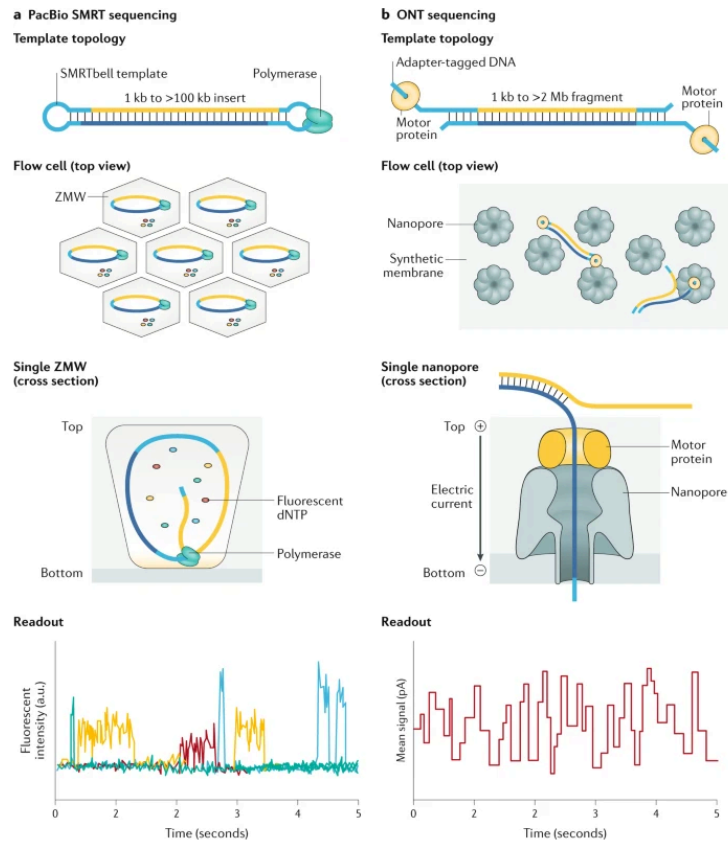
Novaseq 1 billion reads 2x150 bp. \$3000 -> <\$100/genome.

How to assemble a genome



The importance of long read sequencing

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



Logsdon (2020) *Nat Rev Genetics*

The importance of long read sequencing

SPECIAL SECTION COMPLETING THE HUMAN GENOME

RESEARCH ARTICLE

HUMAN GENOMICS

The complete sequence of a human genome

Sergey Nurk¹, Sergey Korotki¹, Arang Rhie¹, Mikko Rautiainen¹, Andrey V. Bakard¹, Alla Mikhonenko¹, Mitchell R. Vollger¹, Nicolas Altshuler¹, Lee Ukhabay¹, Ariel Gershenson¹, Sergey Agapov¹, Swarnal J. Hoyt¹, Mark Diekhans¹, Glenn A. Lopez¹, Michael Almgren¹, Sylvain E. Antonarakis¹, Matthew Borchers¹, Gerard G. Bouffard¹, Sholise Y. Brooks¹, Gina V. Caldas¹, Nae-Chyun Chen¹, Haoyi Cheng¹, Chen-Shan Chien¹, William Chow¹, Leonardo G. de Lima¹, Philip C. DiStefano¹, Richard Durbin¹, Tatsuro Dookhan¹, Ian T. Fiddler¹, Orlan Fomena¹, Robert S. Fulton¹, Arkarachi Fungtammasan¹, Erik Garntorp¹, Patrick G. S. Grady¹, Tessa A. Graves-Lindsay¹, Ira M. Hall¹, Nancy F. Hansen¹, Gabrielle A. Hartley¹, Marine Haudouze¹, Kerstin Howe¹, Michael W. Hunkapiller¹, Cheng Jian¹, Milton Jasi¹, Erich D. Jarvis¹, Peter Kerpedjiev¹, Melanie Kirchoff¹, Mikhail Kolmogorov¹, Jesse Kordas¹, Milos Kravtsov¹, Heng Li¹, Valerie V. Madar¹, Tobias Marschall¹, Ann M. McCartney¹, Jennifer McDaniel¹, Danny E. Miller¹, James C. Mullikin¹, Eugene W. Myers¹, Nathan D. Olson¹, Benedict Paten¹, Paul Peluso¹, Pavel A. Pevzner¹, David Porubsky¹, Tamara Potapova¹, Evgeny I. Rogov¹, Jeffrey A. Rosenblatt¹, Steven I. Sabido¹, Valerie A. Schneider¹, Fritz J. Sedlaczek¹, Kishwar Shafiq¹, Colin J. Sheehy¹, Alana Shumate¹, Ying Sim¹, Adrian F. A. Smith¹, Daniela C. Soto¹, Ivan Sovic¹, Jessica M. Storey¹, Aaron Streets¹, Beth A. Sullivan¹, Françoise Thibaud-Nissen¹, James Torrance¹, Justin Wagner¹, Brian P. Walenz¹, Aaron Wenger¹, Jonathan M. D. Wood¹, Chunlin Xiao¹, Stephanie M. Yan¹, Alice C. Young¹, Samantha Zaretski¹, Urovasi Sankar¹, Rajiv G. McCoo¹, Megha Y. Desai¹, Ivan A. Alexandrov¹, Jennifer L. Gerken¹, Rachel J. O'Neill¹, Winston Tang¹, Juejin M. Zou¹, Michael C. Schatz¹, Evan E. Eichler¹, Karen H. Miga¹, Adam M. Phillippy¹

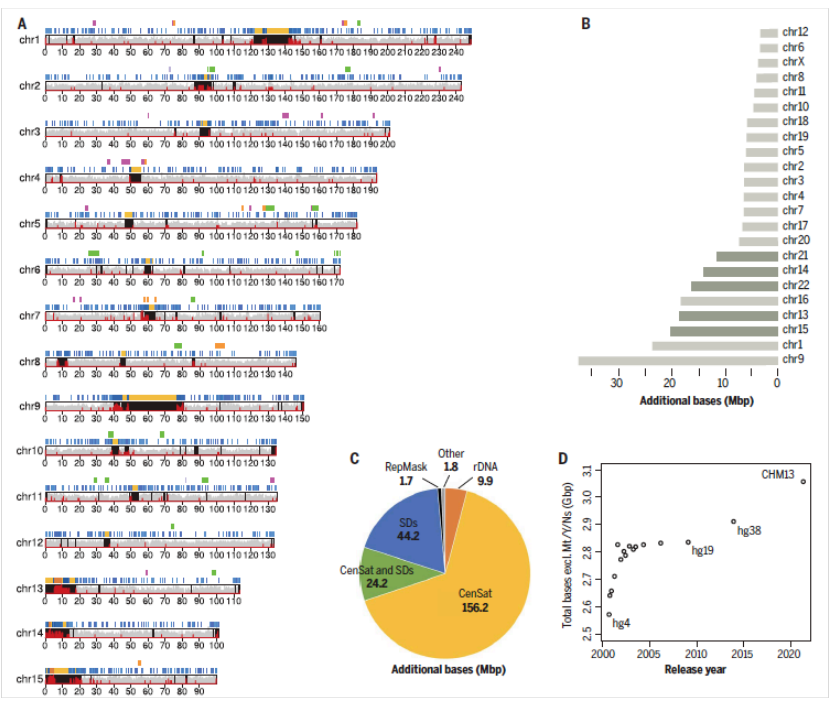
Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.056 billion-base pair sequence of a human genome, T2T-CHM3, that includes gapless assemblies for all chromosomes except Y, corrects errors in the prior references, and introduces nearly 200 million base pairs of sequence containing 1956 gene predictions, 59 of which are predicted to be protein coding. The completed regions include all centromeric satellite arrays, recent segmental duplications, and the short arms of all five acrocentric chromosomes, unlocking these complex regions of the genome to variational and functional studies.

The current human reference genome was released by the Genome Reference Consortium (GRC) in 2003 and most recently published in 2019 (GRCh38.p13) (1). This reference traces its origin to the publicly funded Human Genome Project (2) and has been continually improved over the past two decades. Unlike the competing Celera effort (3) and most modern sequencing projects based on “shotgun” sequence assembly (4),

the GRC assembly was constructed from sequenced bacterial artificial chromosomes (BACs) that were ordered and oriented along the human genome by means of radiation hybrid, genetic linkages, and fingerprint maps. However, limitations of BAC cloning led to an underrepresentation of repetitive sequences, and the opportunistic assembly of BACs derived from multiple individuals resulted in a mosaic of haplotypes. As a result, several GRC assembly gaps are unusable because of incompatible structural polymorphisms on their flanks, and many other repetitive and polymorphic regions were left unfinished or incorrectly assembled (5).

The GRCh38 reference assembly contains 151 mega-base pairs (Mbp) of unknown sequence distributed throughout the genome, including pericentromeric and subtelomeric regions, recent segmental duplications, amplicon gene arrays, and ribosomal DNA (rDNA) arrays, all of which are necessary for fundamental cellular processes (Fig. 1A). Some of the largest reference gaps include human satellite (HSAT) repeat arrays and the short arms of all five acrocentric chromosomes, which are represented in GRCh38 as mating-enzyme stretches of unknown bases (Fig. 1, B and C). In addition to these apparent gaps, other regions of GRCh38 are artificial or are otherwise incorrect. For example, the centromeric alpha satellite arrays are represented as computationally generated models of alpha satellite monomers to serve as decoys for resequencing analyses (6), and sequence assigned to the short arm of chromosome 21 appears falsely duplicated and poorly assembled (7). When compared with other human genomes, GRCh38 also shows a genome-wide deletion bias that is indicative of incomplete assembly (8). Despite finishing efforts from both the Human Genome Project (9) and GRC (1) that improved the quality of the reference, there was limited

¹Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ²Graduate Program in Bioinformatics and Systems Biology, University of California, San Diego, La Jolla, CA, USA. ³Center for Algorithmic Biotechnology, Institute of Translational Biomedicine, Saint Petersburg State University, Saint Petersburg, Russia. ⁴Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA. ⁵Department of Bioengineering, University of California, Berkeley, Berkeley, CA, USA. ⁶Yonsei University of Science and Technology, South Korea. ⁷Yandex Institute of General Genomics, Moscow, Russia. ⁸Department of Molecular Biology and Genetics, Johns Hopkins University, Baltimore, MD, USA. ⁹Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA. ¹⁰Institute for Systems Genomics and Department of Molecular and Cell Biology, University of Connecticut, Storrs, CT, USA. ¹¹USC Santa Cruz Genomes Institute, University of California, Santa Cruz, CA, USA. ¹²University of Geneva Medical School, Geneva, Switzerland. ¹³Medical Research Service, Kansas City, MO, USA. ¹⁴NIH Intramural Sequencing Center, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ¹⁵Department of Molecular and Cell Biology, University of California, Berkeley, Berkeley, CA, USA. ¹⁶Department of Biomedical Informatics, Boston, MA, USA. ¹⁷Department of Biomedical Informatics, Boston, MA, USA. ¹⁸Department of Biomedical Informatics, Boston, MA, USA. ¹⁹Wellcome Sanger Institute, Cambridge, UK. ²⁰Department of Genetics, University of Cambridge, Cambridge, UK. ²¹Phosphor, Boulder, CO, USA. ²²Laboratory of Neurogenetics of Language and the Vertebrate Genome Lab, The Rockefeller University, New York, NY, USA. ²³Howard Hughes Medical Institute, The Rockefeller University, New York, NY, USA. ²⁴Department of Genetics, Washington University School of Medicine, St. Louis, MO, USA. ²⁵University of Tennessee Health Science Center, Memphis, TN, USA. ²⁶Genome Reference Consortium, Washington University in St. Louis, St. Louis, MO, USA. ²⁷Department of Genetics, Yale University School of Medicine, New Haven, CT, USA. ²⁸Comparative Genomics Analysis Unit, Cancer Genetics and Comparative Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ²⁹Phosphor Biocomputing, Merck, Kenilworth, NJ, USA. ³⁰Department of Computational and Systems Biology, Indian Institute of Science, Bangalore, India. ³¹Research Genomics LLC, Oakland, CA, USA. ³²Department of Computer Science and Engineering, University of California, San Diego, La Jolla, CA, USA. ³³Highland Genome Program, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ³⁴Phosphor Biocomputing, Merck, Kenilworth, NJ, USA. ³⁵Department of Medical Biometry and Biostatistics, Düsseldorf, Germany. ³⁶Biogenetics and Biomaterials Division, National Institute of Standards and Technology, Gaithersburg, MD, USA. ³⁷Department of Pediatrics, Division of Genetic Medicine, University of Washington and Seattle Children's Hospital, Seattle, WA, USA. ³⁸Max Planck Institute of Molecular Cell Biology and Genetics, Dresden, Germany. ³⁹Department of Psychiatry, University of Massachusetts Medical School, Worcester, MA, USA. ⁴⁰Faculty of Biology, University of Wrocław, Wrocław, Poland. ⁴¹Center for Genome Sciences, State University of New York, Stony Brook, NY, USA. ⁴²Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA. ⁴³National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA. ⁴⁴Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA. ⁴⁵Phosphor Genomics, Merck, Kenilworth, NJ, USA. ⁴⁶Department of Biochemistry and Molecular Medicine, University of California, Davis, CA, USA. ⁴⁷Institute for Systems Biology, Seattle, WA, USA. ⁴⁸Digital Biologics, Inc., Irvine, CA, USA. ⁴⁹Chen Zuckerman Biotech, San Francisco, CA, USA. ⁵⁰Department of Molecular Genetics and Microbiology, Duke University School of Medicine, Durham, NC, USA. ⁵¹Department of Biology, Johns Hopkins University, Baltimore, MD, USA. ⁵²Department of Pathology, University of Pittsburgh, Pittsburgh, PA, USA. ⁵³Research Center of Biotechnology of the Russian Academy of Sciences, Moscow, Russia. ⁵⁴Department of Biochemistry and Molecular Biology, University of Kansas Medical School, Kansas City, MO, USA. ⁵⁵Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA. ⁵⁶Department of Biomedical Engineering, University of California, Santa Cruz, Santa Cruz, CA, USA. ⁵⁷Corresponding author. Email: nurk@nhgri.nih.gov (S.N.), mihaylova@nhgri.nih.gov (A.M.P.). ⁵⁸These authors contributed equally to this work. Present address: Oxford Nanopore Technologies Inc., Lexington, MA, USA.



<https://www.science.org/doi/10.1126/science.abj6987>

TCCAAATCAAACAGTTGTATTATTAGAAACTGAGGGCTAAAACTGTGCACATACACAGACACACATATTATTTAATATAGATTTTCAATAATTGGTCTAGGATAAGG/
TAAGCAAGAAGAAAACAAAGACTGTTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACTTCTTA/
CAATAAATCACATTAATTCCTTATCTCATGTGAAATTCATATTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTT
CCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCAACAAGACAGACTATGTTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTA/
GTAATTGATGCTAGAAAAGACAATGAAACAGAGCCATGTGACCAATGAGAGAGATGAGGGTGGCAGCAGCCTGTTTTAGATAAAGGTACCTGATTGGTGGGATTGGA
ATATGCCTTAATGATATGAAAGAACCATTATGTTGGAAGGCCTAGCATTAAAAACCGTCTAGGCAGAATGAGCAGCAAGTGAAGGGTCTGGATAGGAATGAGCTC
TATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACTTCTTACCTGTCAATGTTATTAATTTTTAGGAAC
TGAAATTTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTAGAAATAAAGTCCCAGGCACAAGAC/
CCATGTTCAACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGT
TAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACTTCTTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAAT
GATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTAGAAATAAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGC
GACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTATCCCAG
ATTGGGGATACCATTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAATTCACACATGCAAAGAGGAAATCTCCATATCATGCTTGTCAATTCGTTTA
CGTGTGTA AACATTCTCAGAATTTTAAACAATAACAAATCAGGGCTGAATGTGGCCAACATGCAAAGAGGAAATCTCCATCTGTCCAAATCAAACAGTTGATTA
ACATACACAGACACACATATTATTTAATATAGATTTTCAATAATTTGGTCTAGGATAAGGATAATATACAGAGAACATGCCAAAAGTTAAGCAAGAAGAAAACAAAGAC
TTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACTTCTTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAATTCCTTA
ATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTAGAAATAAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATT
ATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTAGAAAAGACA
GAGATGAGGGTGGCAGCAGCCTGTTTTAGATAAAGGTACCTGATTGGTGGGATTGGAAGACCTCTCTGAGATTAGTGTCTTCAGATATGCCTTAATGATATGAAAGA
AAACCGTCTAGGCAGAATGAGCAGCAAGTGAAGGGTCTGGATAGGAATGAGCTGGATATACTCAAGGAAGAAAAGAGAAAATATGGAAAAATGAAAATAGATTT
TTAAATTTACTTTTTCTTTCACTTCTTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTTCATATTTATGATTGAT
ATTATTCATTTTTTCATTCAATAAATATTTTTAGAAATAAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCAACAAGACAGACTATC
AATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTACTATGAAAAATGAAAATAGATTTTAAAACATGTTAATTC
CTTTCACTTCTTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTTCATATTTATGATTGATACCTTTAAATGTCATTT
TCAATAAATTTTTTAGAATAAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCAACAAGACAGACTATGTTTACAGGATCAGAT
AAACAGACACAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTATCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATTACCTGTCA
TCACATTAATTTCAACATGCAAAGAGGAAATCTCCATATCATGCTTGTCAATTCGTTTATCAGAGGCCAAATGTTTTCTTTGTAACCGTGTGTA AACATTCTCAGAAT
TGTGGCCAACATGCAAAGAGGAAATCTCCATCTGTCCAAATCAAACAGTTGATTTATTAGAAACTGAGGGCTAAAACTGTGCACATACACAGACACACATATTAT
GGATAAGGATAATATACAGAGAACATGCCAAAAGTTAAGCAAGAAGAAAACAAAGACTGTTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTT/
ACTTCTTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTG/
AAATATTTTTTAGAATAAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCAACAAGACAGACTATGTTTACAGGATCAGATGTGG/
GACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTAGAAAGACAATGAAACAGAGCCATGTGACCAATGAGAGAGATGAGGGTGGCAGCAGCCT
GATTGGAAGACCTCTCTGAGATTAGTGTCTTCAGATATGCCTTAATGATATGAAAGAACCATTATGTTGGAAGGCCTAGCATTAAAAACCGTCTAGGCAGAATGAGC/
TGAGCTGGATATACTCAAGGAAGAAAAGAGAAAATATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACT
AGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAA
CAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCAACAAGACAGACTATGTTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACA
AAGTTGTAATTGATGCTACTATGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACTTCTTACCTGTCAATG
CATTAATTCCTTATCTCATGTGAAATTTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTAGAAATAA/
GTTCTAGGCATTGGGGATACCATGTTCAACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACAACAAGTAAAT
TATCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATTACCTGTCAATGTTATTAATTTTTAGGAACAATAAATCACATTAATTCACACATGCAA/
TTCGTTTTATCAGAGGCCAAATGTTTTCTTTGTAACCGTGTGTA AACATTCTCAGAATTTTAAACAATAACAAATCAGGGCTGAATGTGGCCAACATGCAAAGAGG/
TTGATTTATTAGAACTGAGGGCTAAAACTGTGCACATACACAGACACACATATTATTTAATATAGATTTTCAATAATTTGGTCTAGGATAAGGATAATATACAGAGAA/
ACAAAGACTGTTACTATGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTAAATTTACTTTTTCTTTCACTTCTTACCTGTCAATGTTATT/
ATTCCTTATCTCATGTGAAATTTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTAGAAATAAAGTC

What types of annotation do we have/want?

~3 billion bp

```
ACAATAATCACATTAATTCCTTATCTCATGTGAAATTCATATTTATGATTG
ATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTGAAT
AAATATTTTTAGAATAAAGTCCCAGGCACAAGACCAGTATTATGTTCT
AGGCATTGGGGATACCATTGTTCAACAAGACAGACTATGATTACAGGATC
AGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAAG
TAAATAAAGTTAATTTCAAAGTTGTAATGATGCTAGAAAGACAATGAAACA
GAGCCATGTACCAATGAGAGAGATGAGGGTGGCAGCAGCCTGTTTTA
GATAAGGTACCTGATTGGTGGGATTGGAAGACCTCTGAGATTAGTGT
CTTCAGATATGCCTTAATGATATGAAAGAACCATTGATGGAAGGCCTAG
CATTAAAACCGTCTAGGCAGAATGAGCAGCAAGTGCAAGGGTCTCTGG
ATAGGAATGAGCTGATATACTCAAGGAAGCAAGGAAACTATGSAAAA
ATGAAAATAGATTTAAAACATGTAATTCACGTTACTTTTTGTAATTTA
CTTTCTCTTTCACCTCTTACCTGTCAATGTTTAAATTTTTAGGAACA
ATAAATCACATTAATTCCTTATCTCATGTGAAATTTCAATTTATGATTGATA
CCITTAATGTCATTTGTGAAGGAAGATTATTCATTTTTTCATTCAATAAA
TATTTTTAGAAATAAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGG
CATTGGGGATACCATTGTTCAACAAGACAGACTATGATTACAGGATCAAGT
GTGGACTCTCAAATTCGACTGAGAATAAACAAGACACTAAACAAGTAAT
AAAGTTAATTTCAAAGTTGTAATTGATGCTACTATGGAATAAATAAATAA
TTTTAAACATGTAATTCACGTACTTTTTGTTAAATTTACTTTTTCTCTTT
CACCTCTTACCTGTCAATGTTAATTTTTTGAACACATAAATCACATT
AATTCCTTATCTCATGTGAAATTTCAATTTATGATTGATACCTTTAATGT
CATTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATTTTTAGA
ATAAATAGTCCCAGGCACAAGACCAGTATTGTTCTAGGCATTGGGGAT
ACCATGTTCAACAAGACAGACTATGATTACAGGATCAGATGTGGACTCTC
AAATTCGACTGAGAATAAACAAGACACAACAAGTAAATAAAGTTAATTT
CAAGTTGTAATTTGATGCTATCCAGGCACAAGACA....
```

Genes:

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

Genetic variation:

- SNPs and CNVs

Sequence conservation

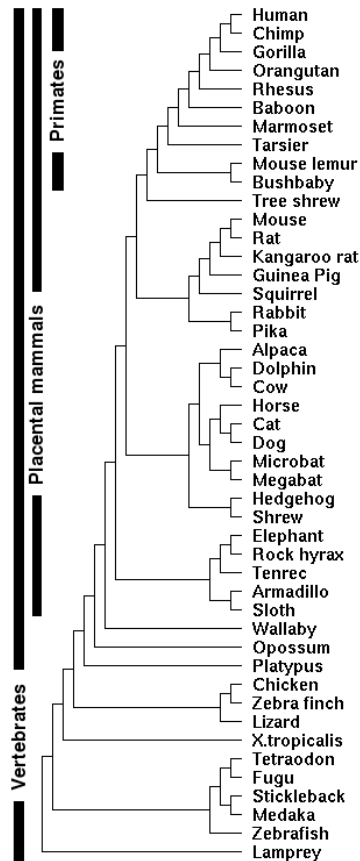
Regulatory sequences:

- 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

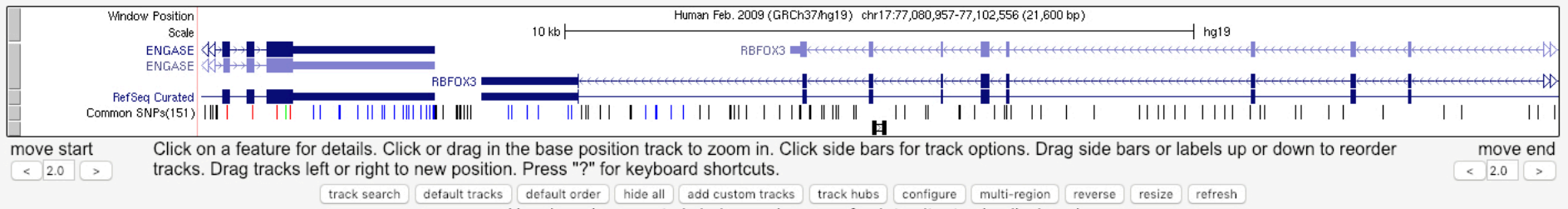
Our specific example:

```
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGGAGAGAGAGACAGAGGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHIJIJJJJJJJJJJJJJJ?FHIDGIJ=GIHGIIIHGIJIHEHIHGGFFFFEEEDDDDDDDDDDDDDDD

@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCTGCAGTCCCTCGGTCCGTGTGTAGACCAGAAGTAGGTGCCCAGGCCAGGTACCACCTAATCCTT
+
##4<@&&&&&&&&&&@&&&&?@&&&&?@&?????@&??@????????????????????>????????????@>????@&&?@&???????
```

UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x
chr17:77,080,957-77,102,556 21,600 bp. go



Workflow

1. Isolation of sample.

e.g., Isolate DNA and shear.

2. Library preparation

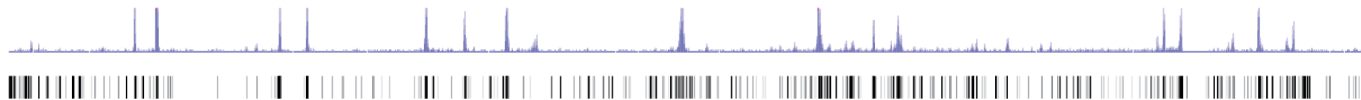
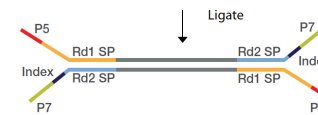
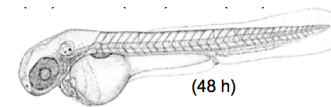
e.g., Add known sequences to the ends.

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 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.