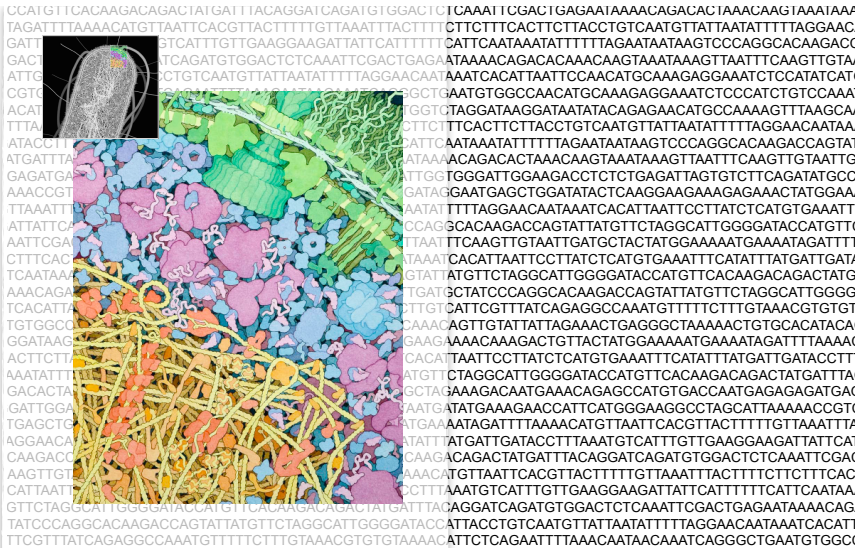


Genomics I

Biomedical Data Science: Mining and Modeling
CB&B 752 · MB&B 452
Matt Simon
Feb 2, 2021



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: (Monday'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

Cell Article

Spliceosome-targeted therapies trigger an antiviral immune response in triple-negative breast cancer

Graphical Abstract

Authors
Elizabeth A. Bowling, Jany H. Wang, Fude Gong, —, Kiyoko L. Pacifico, Xiang H.-F. Zhang, Thomas F. Westbrook

Correspondence
thomas@bcm.edu

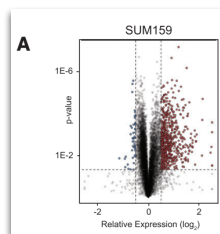
In Brief
Spliceosome-targeted therapies generate intron-retained, double-stranded RNAs that activate downstream antiviral signaling and extrinsic apoptosis in breast cancer.

Highlights

- Spliceosome-targeted therapies (STTs) induce widespread mis-spliced mRNA in cancer
- Mis-spliced, intron-retained mRNAs are an unexplored source of endogenous dsRNA
- STTs trigger antiviral signaling and extrinsic apoptosis in TNBCs via dsRNA sensors
- RNA mis-splicing in human breast cancers correlates with immune signatures

Bowling et al., 2021, Cell 184, 384–403
January 21, 2021 © 2020 Elsevier Inc.
<https://doi.org/10.1016/j.cell.2020.12.031>

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



Bowling et al., 2021, Cell 184, 384–403
January 21, 2021 © 2020 Elsevier Inc.
<https://doi.org/10.1016/j.cell.2020.12.031>

Figure 1. Spliceosome-targeted therapies stimulate antiviral signaling in MYC-driven triple-negative breast cancer
(A) Volcano plot of RNA-seq gene-expression changes due to spliceosome inhibition for two MYC-driven TNBC cell lines, SUM159 and LM2, treated with DSR or DMSO (n = 3 biological replicates).

Data and code availability
The datasets generated during this study are available in GEO [GEO: GSE163411, GSE163414, GSE163181, GSE163188, GSE163232].

Raw data can be found in genomics databases

Public on Dec 18, 2020

Titles SUM159 cells treated with suberoyl A6 (S6A)

Organism Homo sapiens

Experiment Date Experiments starting by high throughput sequencing

Summary The objective is to identify gene expression patterns that occur after spliceosome inhibition in triple-negative breast cancer.

General Design SUM159 cells were treated with vehicle (DMSO) or S6A for 72 hours, in biological triplicate.

Contributions Wang, J., Bowling, E., Westbrook, T.

Conflict of Interest We do not have any potential financial gain to declare or verify ODS.

Publication Date Dec 17, 2020

Last Update Date Dec 21, 2020

Contact Name Jany H. Wang

E-mail(s) jany@bcm.edu

Organization name Baylor College of Medicine

Street address 1 Baylor Plaza, 3048 Anderson Building

City Houston

State/province Texas

ZIP/postal code 77030

Country USA

Platforms (1) GPL13573 Illumina HiSeq500 (Illumina illumina)

Sample IDs GSE639113: SUM159 S06 Sample 1

File Name(s) GSE639113: SUM159 S06 Sample 2

Accession GSE639113: SUM159 S06 Sample 3

Annotations

Bioproject PRJNA58358

Study GSE639113

Downloaded Family

FASTQ format	FASTQ ID	FASTQ ID
FASTQ format family (FASTQ)	FASTQ ID	FASTQ ID
Series Matrix (Matrix)	FASTQ ID	FASTQ ID

Supplementary File **Size** **Download** **File type/extension**

GSE639113_M001 32.4 Mb (100%) (FASTQ) TAR (tar)

Raw data are available in SRA

Processed data provided as supplementary file

```

SRRT458868.1.HISEQ-250-CBM.MCX.3.1101.1497.1977.length51
NCATTTTCATGTTTCATGTACACTGATATACACTGTCACAA
SRRT458868.1.HISEQ-250-CBM.MCX.3.1101.1497.1980.length51
GTCTTFFRRHRH33333333333333333333333333333333
SRRT458868.2.HISEQ-250-CBM.MCX.3.1101.1497.1977.length51
HSKATAGGGGCGATCCAGCTCCAGCCCGGCTTCATCCCG
SRRT458868.2.HISEQ-250-CBM.MCX.3.1101.1497.1977.length51
AT:DDTDTGADGAGHGHC11113GGGAPR:GALG:34:258CC
SRRT458868.3.HISEQ-250-CBM.MCX.3.1101.1546.1574.length51
NCTTCATAGCAGCTTATTATATAGCCGAGGCTGGAAGAACAC
SRRT458868.3.HISEQ-250-CBM.MCX.3.1101.1546.1583.length51
AT:DDFFRRHRH33333333333333333333333333333333
SRRT458868.4.HISEQ-250-CBM.MCX.3.1101.1488.1958.length51
NSGAGAGAGACACACACACACACACACACACACACACACAC
SRRT458868.4.HISEQ-250-CBM.MCX.3.1101.1488.1958.length51
AT:DDFFRRHRH33333333333333333333333333333333
SRRT458868.5.HISEQ-250-CBM.MCX.3.1101.1546.1577.length51
NACCAACGAGGATGTTCTTTCAGAAAACAGAAAGATAGATAC
SRRT458868.5.HISEQ-250-CBM.MCX.3.1101.1546.1577.length51
AT:DDFFRRHRH33333333333333333333333333333333
SRRT458868.6.HISEQ-250-CBM.MCX.3.1101.1497.1989.length51
AT:DDFFRRHRH33333333333333333333333333333333
SRRT458868.7.HISEQ-250-CBM.MCX.3.1101.1497.1924.length51
HSKCTTTTCAGAGATATTAGTAGAGATTCCTTCTATTC
SRRT458868.7.HISEQ-250-CBM.MCX.3.1101.1497.1924.length51
AT:DDFFRRHRH33333333333333333333333333333333
SRRT458868.8.HISEQ-250-CBM.MCX.3.1101.1546.1574.length51
NCGAGTGAATATGGCAGAAACAGAAAGGCTGAAATTTAAAT
SRRT458868.8.HISEQ-250-CBM.MCX.3.1101.1546.1580.length51
AT:DDFFRRHRH33333333333333333333333333333333
  
```

- 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.

Central questions

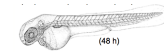
Where do these data come from?

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

Workflow

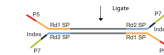
1. Isolation of sample.

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

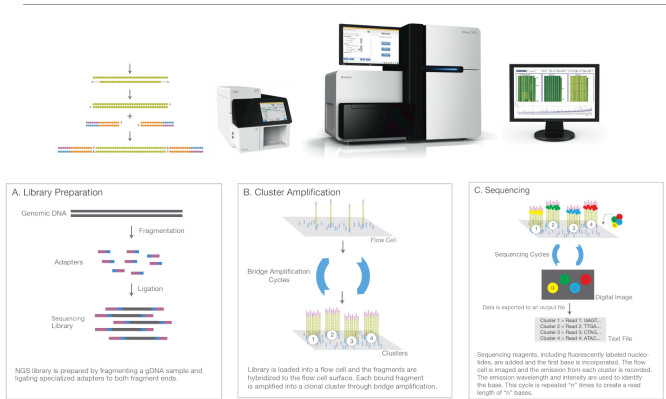
- Nanopore based
- Pacific Bioscience Sequencing

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	98-195 ^d	17,520
	Sequel II	CLR	30-60	>200		50-100	160	13-26 ^e	93,440
Oxford Nanopore Technologies (ONT)	MiniON/GridION	Long	10-60	>1,000	>99	15-30	35	43-86 ^f	10,220
		Ultra-long	100-200	>1,500	87-98	0.5-2	2.5	500-2,000 ^g	21,900 (MiniON) 109,500 (GridION)
	PromethION	Long	10-60	>1,000		50-100	180	21-42 ^h	3,153,600
Illumina	NextSeq 550	Single-end	0.075-0.15	0.15	>99.9	16-30	>30	50-63 ⁱ	>47,782
		Paired-end	0.075-0.15 (x2)	0.15 (x2)		32-120	>120	40-60 ^j	>70,080
	NovaSeq 6000	Single-end	0.05-0.25	0.25		65-3,000	>3,000	10-35 ^k	>1,194,545

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

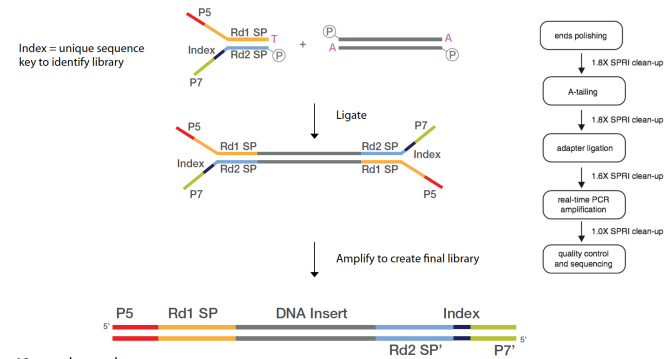
Logsdon (2020) Nat Rev Genetics

Where do these reads come from?



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

Optional: How do you make a sequencing library?



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).

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
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGAGAGACAGAGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHJJJJJJJJJJJJ?FHIDGJ=GIHGIIIHGIJIHEHIHGGFFFEEDDDDDDDDDDD
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCTGCAGTCTCGTCTGTGTAGACCAGAACTAGGTGCCAGCCAGGATACCACCTAATCCTT
+
##4<@@@0000000?000?000?????@000?000????????????????????>?????????@??000?000?????
```

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 limits the insert size and read length?

One read (fastq format)

```
@HWI-D0306:498:HBB89ADXX:1:1101:1180:1882 1:N:0:CGATGT
NCATCACTTTCGACCAGCCATGACGTAATCTTCGTCGAAACCCAAACTCGAGATCGGAAGACACAGCTCTG
+
#11BBDDDFDFBFFFIIIIIIIIIIIIIFEGIIIIIFIGAGIIFIIIFEEEEFFFD@9A@BBBBB=?BB<
```

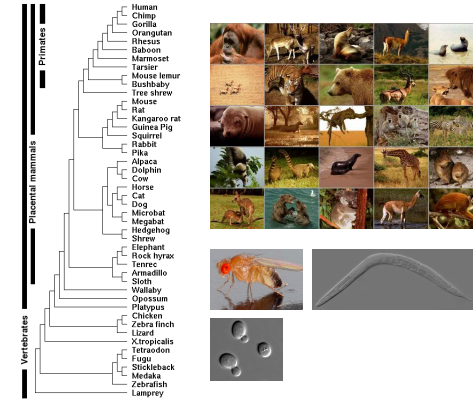
- For each single end read: Incomplete incorporation of bases.
- For the size of the insert (especially for paired end analysis): Ability to get consistent clusters.

What do I do with my sequencing reads?



Source: Slate via Noonan

Many reference genomes are available

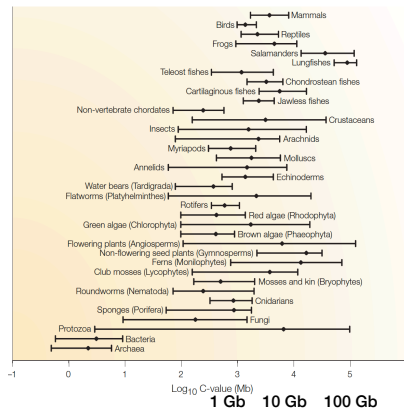


There is a wide range of genome sizes.

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

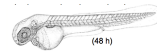
• 1 Gb/month at largest centers (2005)

Novaseq 20 billion reads 2x150 bp. \$1000 → \$100/genome.

Workflow

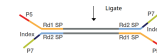
1. Isolation of sample.

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.