

Intro to High-Throughput DNA Sequencing

Analysis of Next-Generation Sequencing Data

Friederike Dündar

Applied Bioinformatics Core

Slides at <https://bit.ly/2T3sjRg>¹

January 14, 2020



Weill Cornell Medicine

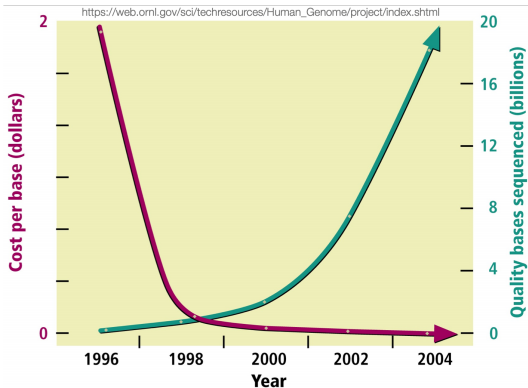
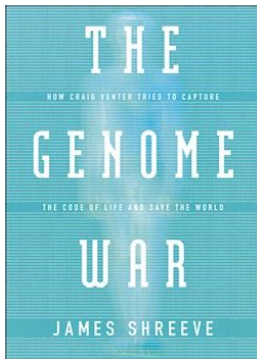
¹https://physiology.med.cornell.edu/faculty/skrabanek/lab/angsd/schedule_2020/

- 1 Why does sequencing data need bioinformatics?
- 2 What do we sequence?
- 3 How do we sequence?
- 4 Why do we sequence?
- 5 Experimental design
- 6 References

Why does sequencing data need bioinformatics?

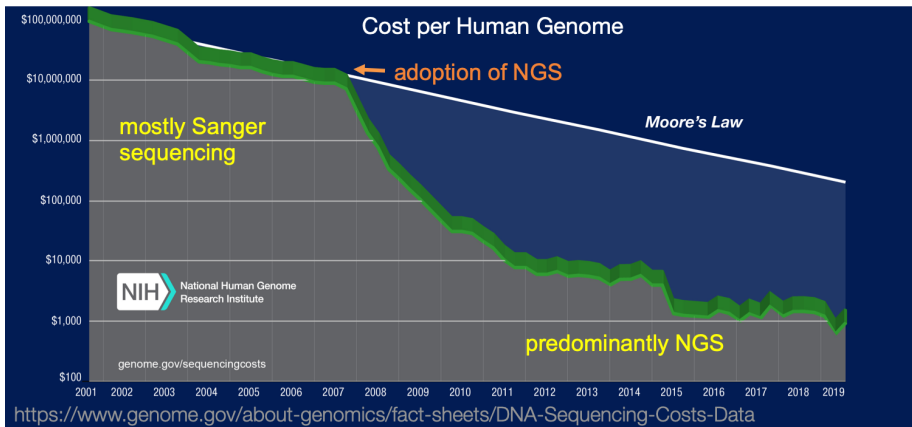
Evolution of sequencing data

The Human Genome Project (1990-2003) ushered in the era of "next"-generation sequencing.



Evolution of sequencing data: NGS

NGS = “next generation sequencing”



ANGSD relies on bioinformatics

Next-generation sequencing

High-throughput

millions of nucleotides can be sequenced at once

Relatively cheap*

experiments involving NGS have become abundant

relatively large data files
are being generated on a
regular basis



Bioinformatics

Processing

formatting, data wrangling

Alignment

Statistical analyses

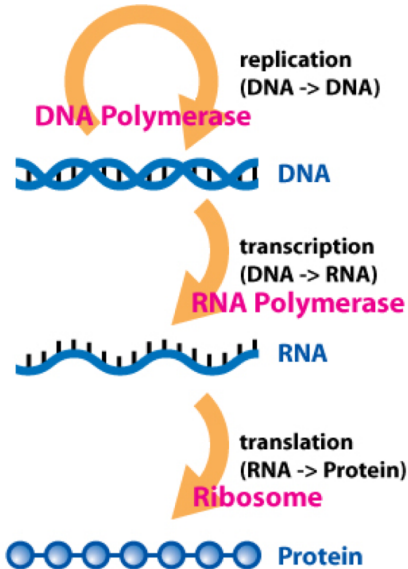
Interpretation

* The cost of analysis has
remained high and is difficult
to estimate!

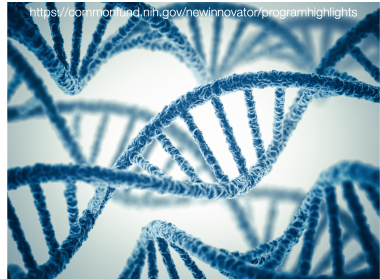
What do we sequence?

An essential macromolecule of life: DNA

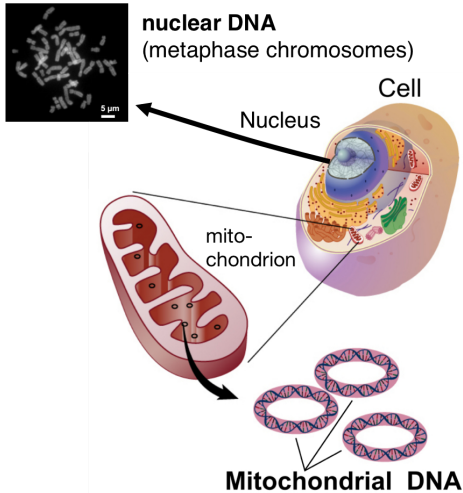
https://en.wikipedia.org/wiki/Central_dogma_of_molecular_biology



- the "hard drive" of all living organisms
- determines the traits of an organism
- contains the blueprint information for proteins

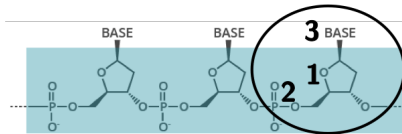


Two general types of eukaryotic DNA



- **GENOMIC (NUCLEAR) DNA**
 - ▶ contained and replicated within the **nucleus**
 - ▶ "linear"
 - ▶ multiple chromosomes, which are inherited from both parents
- **MITOCHONDRIAL DNA**
 - ▶ contained and replicated within **mitochondria**
 - ▶ circular
 - ▶ represents 1 chromosome
 - ▶ inherited (only!) from the mother

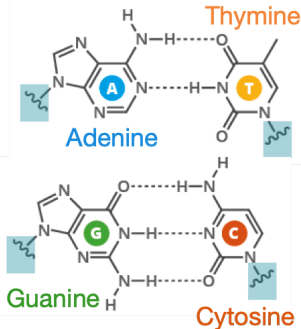
DNA: Deoxyribonucleic acid



sugar-phosphate backbone

+ nitrogenous bases

based off of <https://bit.ly/35OXXM>



• each **nucleotide**:

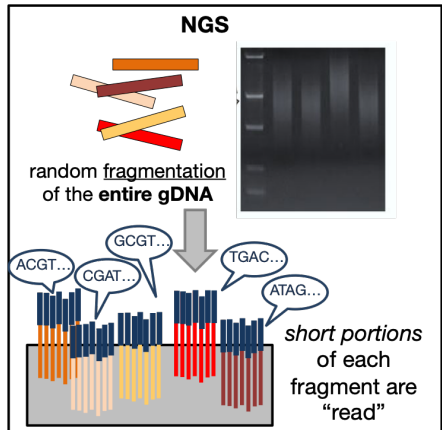
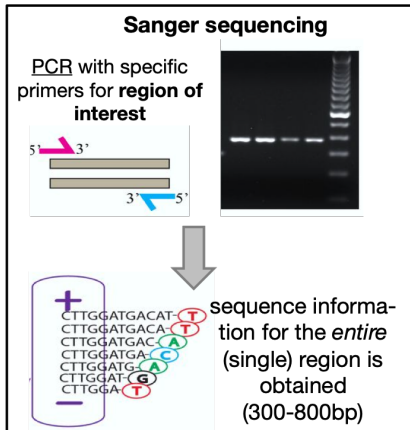
- ① **sugar**: 2'-deoxyribose (5 carbon atoms = pentose)
- ② **phosphate**: 1-3 linked phosphate units attached to the 5'-carbon of the sugar
- ③ **nitrogenous base**: either a single-ring pyrimidine (cytosine, thymine) or a double-ring purine (adenine, guanine)

SEQUENCING = IDENTIFYING THE ORDER OF THE BASES

How do we sequence?

Next-generation sequencing

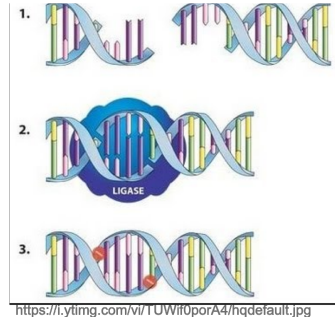
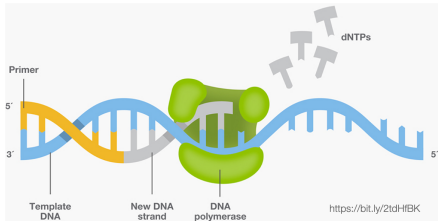
- refers to **highly parallelized sequencing** of millions of DNA fragments at the same time (in contrast to the traditional one-region-at-a-time approach)



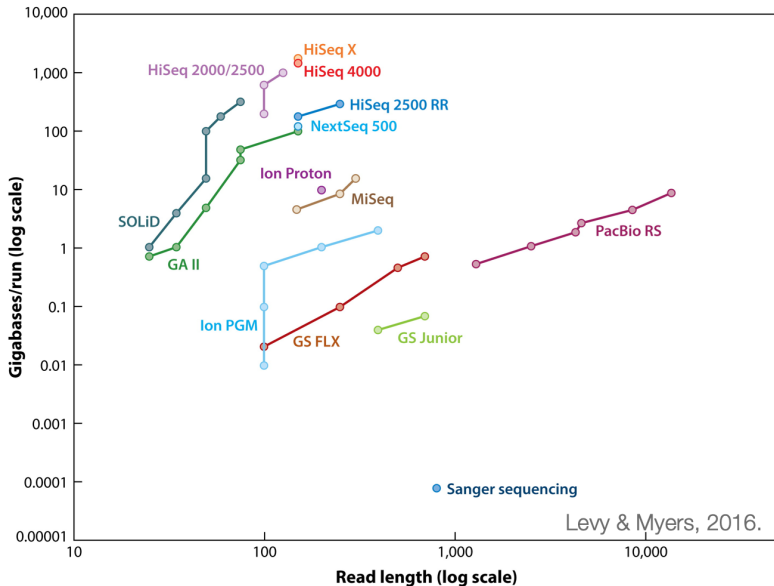
Decoding the DNA

DNA sequencing mostly relies on **enzymes** that are DNA “readers”:

- **DNA Polymerase**: synthesizes a new strand of DNA
 - ▶ **sequencing by synthesis** platforms: Solexa/Illumina, Ion Torrent
- **DNA Ligase**: joins the “sticky” ends of two strands of DNA together
 - ▶ **sequencing by ligation** NGS platforms: SOLiD, Complete Genomics



Next-generation sequencing platforms



Next-generation sequencing platforms

Unifying characteristics of the different NGS platforms:

- **short fragments (250-1000 bp)** are assessed via **short reads (50-250 bp)**
- require **clonal amplification** of every single DNA fragment
- markedly **higher error rates** than Sanger sequencing of the 1980s-1990s (0.1–15%)

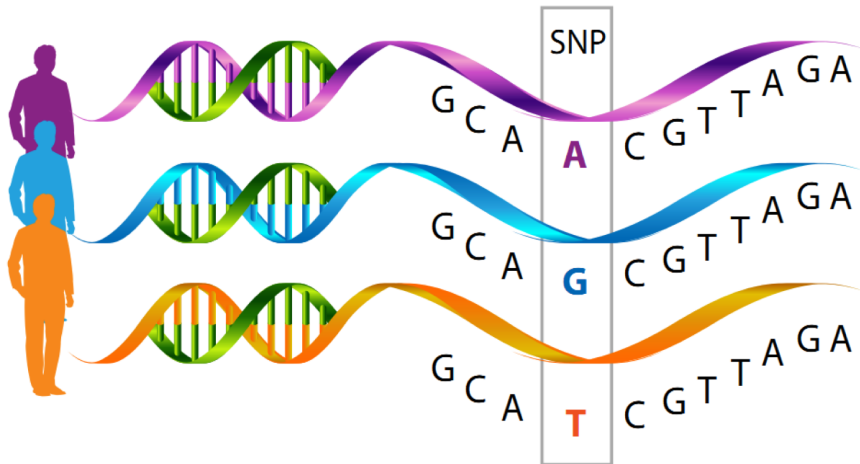


See Goodwin et al. [2016] for detailed descriptions of NGS platforms.

Why do we sequence?

Why do we sequence?

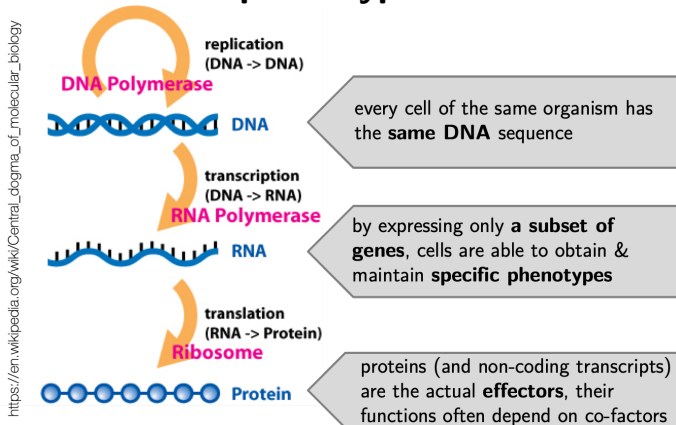
...to identify individuals

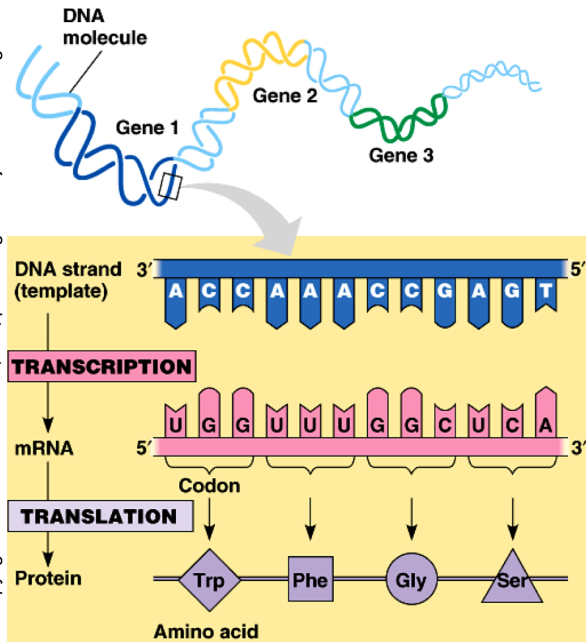


<https://neuroendocrine.wordpress.com/2014/03/27/dna-rna-snp-alphabet-soup-or-an-introduction-to-genetics/>

Why do we sequence?

...to understand the molecular basis of different cellular **phenotypes**

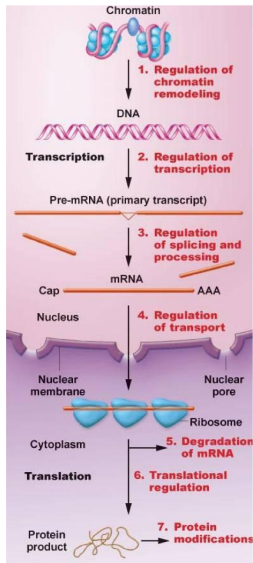
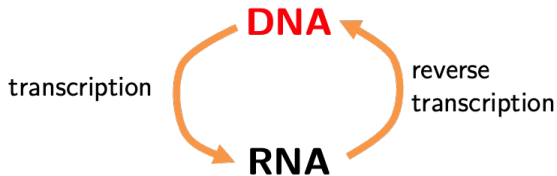




Central dogma:
the genetic code
serves as a manual for
building different
proteins.

Understanding the genetic code and its interpretation

Both RNA and DNA molecules can be assessed through sequencing in a high-throughput manner.

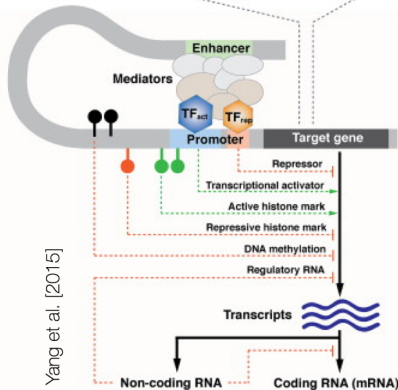


Understanding the genetic code and its interpretation

NGS can be used for (A) **qualitative** as well as (B) **quantitative** approaches.

(A) “Reading” the actual sequence

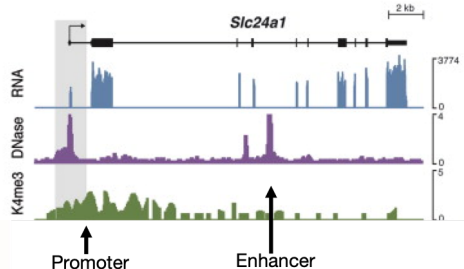
WT Rho: ...CCCTTCTCCAATGCGACGGGTGGTACGCAGCCCTTCGAGTACCCACAG...
 ADRP: ...CCCTTCTCCAATGCGA**T**GGGTGGTACGCAGCCCTTCGAGTACCCACAG...
 ...CCCTTCTCCAATGCGACGGGTGGTACGCAG**C**CTTCGAGTACCCACAG...



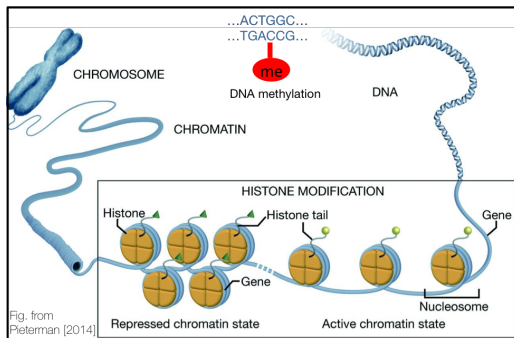
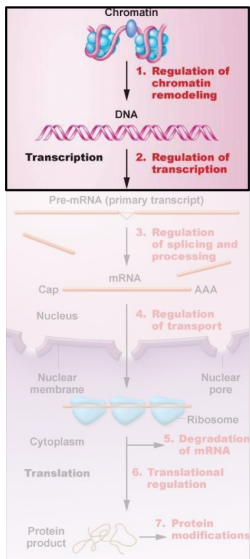
Yang et al. [2015]

(B) Characterizing DNA regions with certain properties

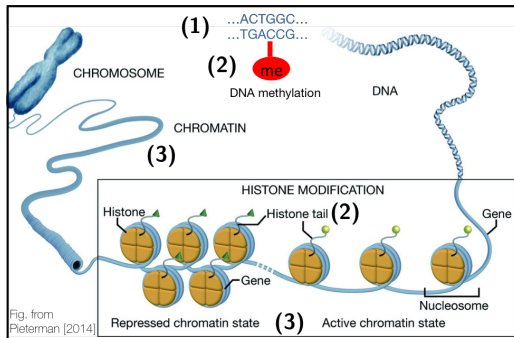
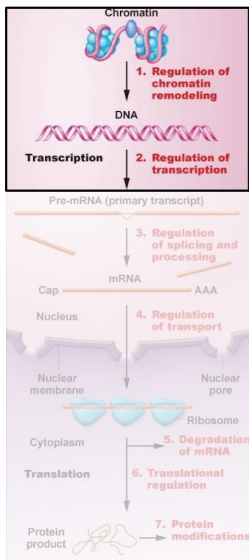
that are first biochemically enriched in a sample of interest – their DNA sequence is only needed to identify the locus of origin, the information of interest is based on the abundance (= enrichment) of seq. reads



Understanding DNA: it's not just about the letters



Understanding DNA: it's not just about the letters



(1) DNA sequence: genome assembly, variant detection

- whole genome (WGS), whole exome (WES), amplicons

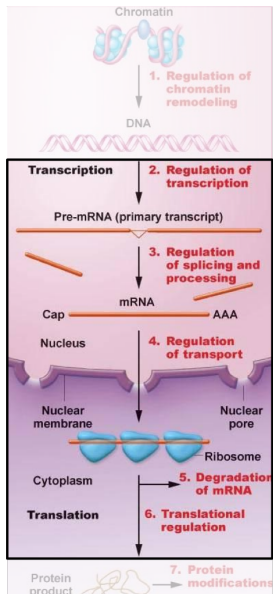
(2) Interrogating epigenetic marks

- histone marks, TF binding: ChIP-seq
- DNA methylation: bisulfite sequencing

(3) Understanding the chromatin structure

- active/repressed: ATAC-, DNase, MNase-seq, ...
- 3D interactions: Hi-C, ChIA-PET

Understanding RNA



(1) Gene expression: sequencing transcripts

- transcript identification & quantification (including non-coding transcripts): RNA-seq
- nascent transcripts: PRO-, GRO-seq

(2) Identifying RNA-binding proteins

- RIP-seq, CLIP-seq

(3) Determining RNA structures

- PARS, Structure-seq

(4) Assessing RNA modifications

- meRIP-seq, ICE (adenosine-inosine editing)

(5) Understanding translation

- Ribo-seq

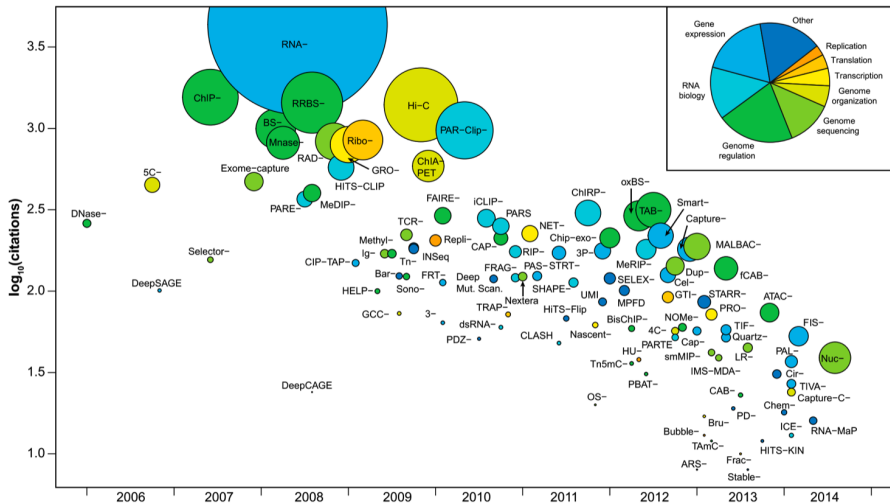


Fig. from Reuter, J. A., Spacek, D. V., & Snyder, M. P. (2015). *Molecular Cell*, 58(4), 586–597.

Main steps of typical NGS experiments

TEMPLATE PREP

**Obtaining the
molecules of
interest:**

DNA, RNA,
nucleotide-protein
complexes



**Library
preparation:**

fragmentation and
ligation of
sequencing adapters



Amplification

SEQUENCING

Sequencing by
Synthesis

vs.

Sequencing by
Ligation

short reads vs. long
reads

BIOINFORMATICS

Base calling



Alignment

Identifying loci of
the sequenced
fragments



**Additional
processing**



Interpretation

Experimental design

Where to sequence at WCM?

Genomics and Epigenomics Sequencing Services

- highly experienced staff
- nevertheless: know the issues you need to discuss with them!

Jenny Xiang, M.D.

Director of Genomics Services

WCM CLC Genomics and Epigenomics Core Facility

(212) 746-4258

jzx2002@med.cornell.edu

Alicia Alonso, Ph.D.

Director of Epigenomics Services

WCM CLC Genomics and Epigenomics Core Facility

(212) 746-3260

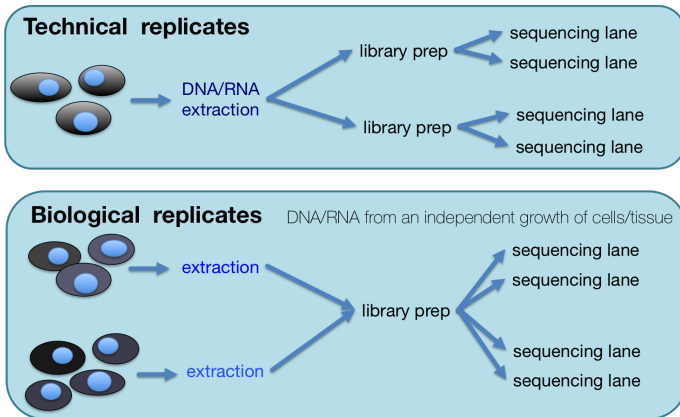
ala2035@med.cornell.edu

Experimental design considerations

- **How many replicates?**
- How to avoid batch effects?
- How many reads?

Why do we need replicates?

- replicates are needed to understand the **level of noise**

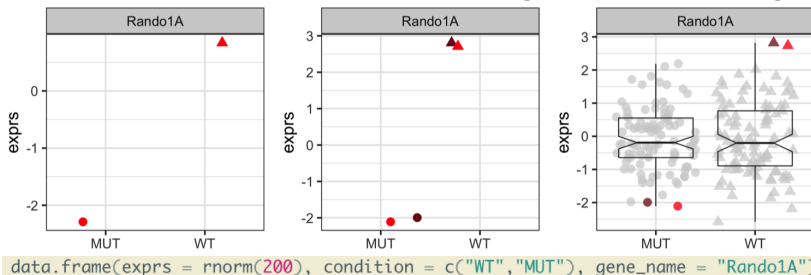


Cross-platform replicates sometimes may make sense, too.

Why do we need replicates?

"Samples are our windows to the population." [Krzywinski and Altman, 2013]

- definitely needed for quantitative assessments, e.g. RNA-seq for determining expression level differences [Schurch et al., 2016]



- qualitative approaches (e.g. variant calling) also benefit from technical replicates [Robasky et al., 2014, Derryberry et al., 2016]

Experimental design considerations

- How many replicates?
- **How to avoid batch effects?**
 - ▶ Understanding typical sources of noise and artifacts
- How many reads?

General problems for NGS

Problems = sources of technical noise

Sample preparation

- DNA/RNA extraction with varying degrees of degradation
- contaminations
- mislabelling, mishandling

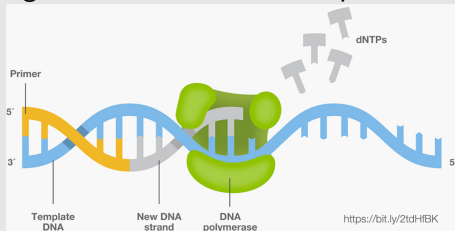
Biases of Illumina-based DNA sequencing

Somewhat **sequencing-machine**-specific problems

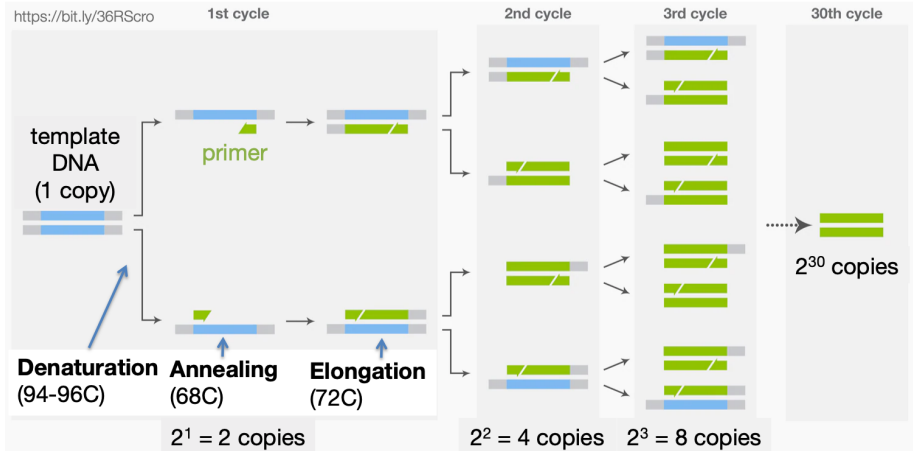
- sequencing errors
- miscalled bases

Sample-specific problems: **PCR artifacts**

- duplicated fragments (low library complexity)
- GC bias: fragments with moderate GC content are preferably amplified
- length bias: fragments between 250-700bp are strongly favored



The most important biochemical assay for NGS: PCR



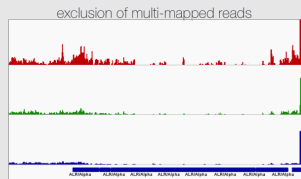
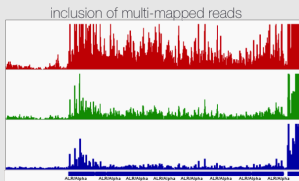
For NGS applications, template DNA fragments vary in size and GC content!
 Exponential nature of the amplification process \Rightarrow small differences in the starting population can lead to strongly skewed final populations.

Always keep the number of PCR cycles to an absolute minimum!

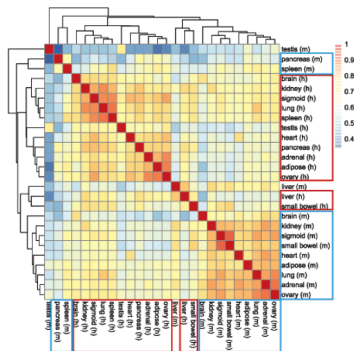
Biases of Illumina-based DNA sequencing

Bioinformatics problems

- **DNA**: long, repetitive elements are difficult to align to with short reads (“mappability” issue)
 - ▶ abundance of (structural) variants may complicate alignments
- **RNA**: great dynamic range (lowly expressed to extremely abundant)
 - ▶ saturation point is hardly reached: number of distinct transcripts depends on the overall make-up of the library
 - ▶ strongly affected by contaminations (DNA, rRNA, ...)
- inappropriate **data processing**, e.g. wrong parameter choices



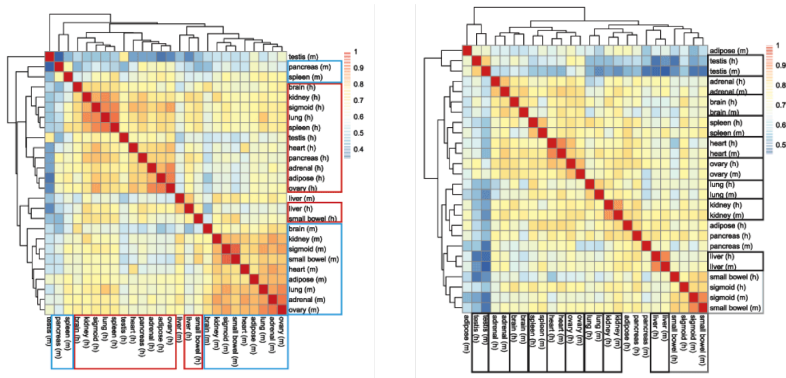
Case study: ENCODE's comparison of mouse and human tissues



“Overall, our results indicate that there is considerable RNA expression diversity between humans and mice, well beyond what was described previously, likely reflecting the fundamental physiological differences between these two organisms.”

Lin, Lin, and Snyder (2014). PNAS 111:48

Case study: ENCODE's comparison of mouse and human tissues



“Overall, our results indicate that there is considerable RNA expression diversity between humans and mice, well beyond what was described previously, likely reflecting the fundamental physiological differences between these two organisms.”

Lin, Lin, and Snyder (2014). PNAS 111:48

“Once we accounted for the batch effect (...), the comparative gene expression data no longer clustered by species, and instead, we observed a clear tendency for clustering by tissue.”

Gilad & Mizrahi-Man (2015). F1000Research 4:121

Suboptimal study design

Most human samples were sequenced separately from the mouse samples:

D87PMJN1 (run 253, flow cell D2GUAACXX, lane 7)	D87PMJN1 (run 253, flow cell D2GUAACXX , lane 8)	D4LHBFN1 (run 276, flow cell C2HKJACXX , lane 4)	MONK (run 312, flow cell C2GR3ACXX , lane 6)	HWI-ST373 (run 375, flow cell C3172ACXX , lane 7)
heart	adipose	adipose	heart	brain
kidney	adrenal	adrenal	kidney	pancreas
liver	sigmoid colon	sigmoid colon	liver	brain
small bowel	lung	lung	small bowel	spleen
spleen	ovary	ovary	testis	● Human
testis		pancreas		● Mouse

Many tissues were not sex-matched

Tissue	Human	Mouse
adipose	FEMALE	MALE
adrenal	MALE	FEMALE
brain	FEMALE	MALE
heart	FEMALE	FEMALE
kidney	MALE	FEMALE
liver	MALE	FEMALE
lung	FEMALE	FEMALE
ovary	FEMALE	FEMALE
pancreas	FEMALE	FEMALE
sigmoid colo	MALE	FEMALE
small bowel	FEMALE	FEMALE
spleen	FEMALE	MALE
testis	MALE	MALE

- human data: deceased organ donors
- mouse data: 10-week-old littermates

Not all variables can be controlled for! Know the limitations of your study before making bold claims! Recommended reading:
<https://f1000research.com/articles/4-121/v1>

Avoiding bias by relying on randomization

Completely randomized design

STRESS	A	B	A	A	B	A	B	A	A	B	B	B
DIET	1	2	1	2	2	1	1	2	2	1	2	1

Restricted randomized design

GENOTYPE	A	A	A	A	A	A	B	B	B	B	B	B
DIET	1	2	1	2	2	1	1	2	1	1	2	2

Blocked & randomized design

GENOTYPE	A	A	B	B	A	A	B	B	A	A	B	B
DIET	1	2	1	2	1	2	1	2	1	2	1	2
WEIGHT	•	•	•	•	•	•	•	•	●	●	●	●



Block what you can,
randomize what you cannot.

*What factors are of **interest**? Which ones might introduce noise? Which nuisance factors do you absolutely need to account for?*

Experimental design considerations

- How many replicates?
- How to avoid batch effects?
- **How many reads?**

How deep is deep enough?

lower limit should usually be whatever ENCODE says:

<https://www.encodeproject.org/about/experiment-guidelines/>

Application	Recommended seq. depth
differential gene expression	20 - 50 mio SR, 75 bp
variant calling	30-200x coverage
whole-genome bisulfite sequencing	30x coverage
ChIA-PET	200 mio PE

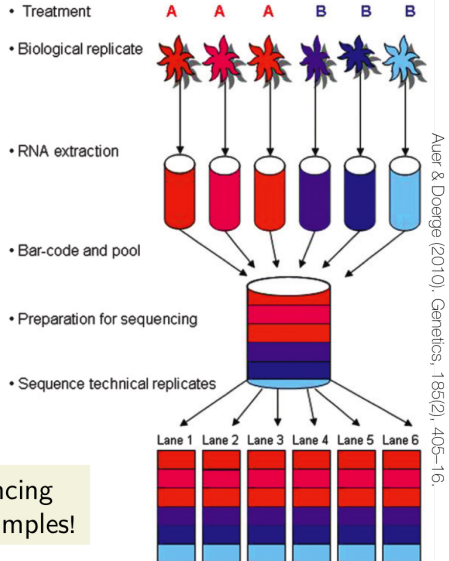
- you may need more, longer, and possibly paired-end reads
 - ▶ novel transcript identification
 - ▶ alternative splicing
 - ▶ ChIP-seq for broad histone marks
 - ▶ 3D chromatin structure assessment assays

Oftentimes the addition of replicates is more meaningful than increased sequencing depth! [Rapaport et al., 2013]

Typical experimental setup

- keep the **technical nuisance** factors (harvest date, RNA extraction kit, sequencing date...) to a minimum
- cover only as much of the **biological variation** as needed (but keep possible limitations for the final conclusions in mind)

Make sure the sequencing core **multiplexes** all samples!



References

Figures taken from the following publications:

[Auer and Doerge, 2010, Gilad and Mizrahi-Man, 2015, Levy and Myers, 2016, Lin et al., 2014, Park, 2009, Pieterman et al., 2014, Reuter et al., 2015, Yang et al., 2015]

- Paul L Auer and R W Doerge. Statistical design and analysis of RNA sequencing data. *Genetics*, 185(2):405–16, Jun 2010. ISSN 1943-2631. doi: 10.1534/genetics.110.114983.
- Dakota Z. Derryberry, Matthew C. Cowperthwaite, and Claus O. Wilke. Reproducibility of SNV-calling in multiple sequencing runs from single tumors. *PeerJ*, 2016. doi: 10.7717/peerj.1508.
- Yoav Gilad and Orna Mizrahi-Man. A reanalysis of mouse ENCODE comparative gene expression data. *F1000Research*, 2015. doi: 10.12688/f1000research.6536.1.
- Sara Goodwin, John D Mcpherson, and W Richard Mccombie. Coming of age : ten years of next- generation sequencing technologies. *Nature Genetics*, 17(6):333–351, 2016. doi: 10.1038/nrg.2016.49.
- Martin Krzywinski and Naomi Altman. Points of significance: Importance of being uncertain. *Nature Methods*, 10(9):809–810, 2013. ISSN 1548-7091. doi: 10.1038/nmeth.2613. URL <http://www.nature.com/doifinder/10.1038/nmeth.2613>.

- Shawn E. Levy and Richard M. Myers. Advancements in Next-Generation Sequencing. *Annual Review of Genomics and Human Genetics*, 2016. doi: 10.1146/annurev-genom-083115-022413.
- Shin Lin, Yiing Lin, Joseph R. Nery, Mark A. Urich, Alessandra Breschi, Carrie A. Davis, Alexander Dobin, Christopher Zaleski, Michael A. Beer, William C. Chapman, Thomas R. Gingeras, Joseph R. Ecker, and Michael P. Snyder. Comparison of the transcriptional landscapes between human and mouse tissues. *Proceedings of the National Academy of Sciences*, 2014. doi: 10.1073/pnas.1413624111.
- Peter J Park. ChIP-seq: advantages and challenges of a maturing technology. *Nature Reviews Genetics*, 10(10):669–80, Oct 2009. doi: 10.1038/nrg2641.
- C. R.C. Pieterman, E. B. Conemans, K. M.A. Dreijerink, J. M. De Laat, H. Th M. Timmers, M. R. Vriens, and G. D. Valk. Thoracic and duodenopancreatic neuroendocrine tumors in multiple endocrine neoplasia type 1: Natural history and function of menin in tumorigenesis. *Endocrine-Related Cancer*, 2014. doi: 10.1530/ERC-13-0482.

- Franck Rapaport, Raya Khanin, Yupu Liang, Mono Pirun, Azra Krek, Paul Zumbo, Christopher E Mason, Nicholas D Socci, and Doron Betel. Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data. *Genome Biology*, 14(9):R95, January 2013. doi: 10.1186/gb-2013-14-9-r95.
- Jason A. Reuter, Damek V. Spacek, and Michael P. Snyder. High-Throughput Sequencing Technologies. *Molecular Cell*, 58(4): 586–597, May 2015. doi: 10.1016/j.molcel.2015.05.004.
- Kimberly Robasky, Nathan E. Lewis, and George M. Church. The role of replicates for error mitigation in next-generation sequencing. *Nature Reviews Genetics*, 2014. doi: 10.1038/nrg3655.
- Nicholas J Schurch, Pietá Schofield, Marek Gierliński, Christian Cole, Alexander Sherstnev, Vijender Singh, Nicola Wrobel, Karim Gharbi, Gordon G Simpson, Tom Owen-Hughes, Mark Blaxter, and Geoffrey J Barton. How many biological replicates are needed in an RNA-seq experiment and which differential expression tool should you use? *RNA*, pages 1–13, 2016. ISSN 1469-9001. doi: 10.1261/rna.053959.115.

Hyun Jin Yang, Rinki Ratnapriya, Tiziana Cogliati, Jung Woong Kim, and Anand Swaroop. Vision from next generation sequencing: Multi-dimensional genome-wide analysis for producing gene regulatory networks underlying retinal development, aging and disease. *Progress in Retinal and Eye Research*, 2015. doi: 10.1016/j.preteyeres.2015.01.005.