# 3000788 Intro to Comp Molec Biol

**Lecture 3: DNA sequencing applications** 

August 23, 2022

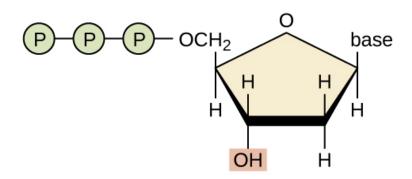


#### Sira Sriswasdi, PhD

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

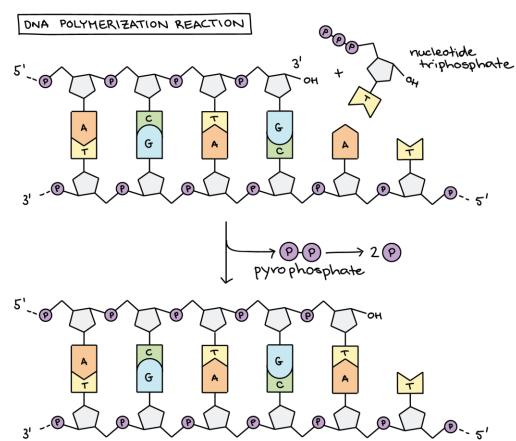
# **Next-Generation Sequencing (NGS)**

#### **DNA** polymerization



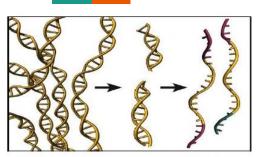
deoxynucleotide (dNTP)

http://www.onlinebiologynotes.com/sangers-method-genesequencing/



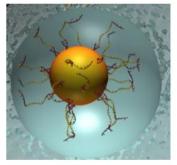
https://www.khanacademy.org

### High throughput from parallel reactions

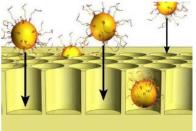


Roche & Ion Torrent wells

1) Adapter-ligated ssDNA library

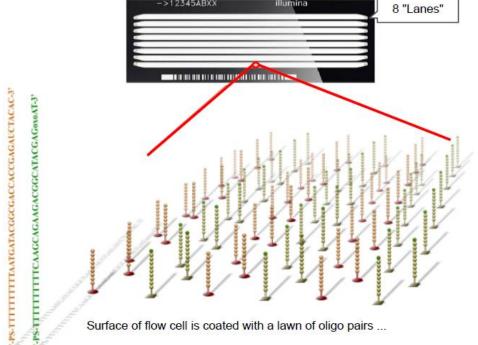


2) Clonal amplification on 28 micron beads ... emulsion PCR



3) Beads deposited on PicoTiterPlate wells

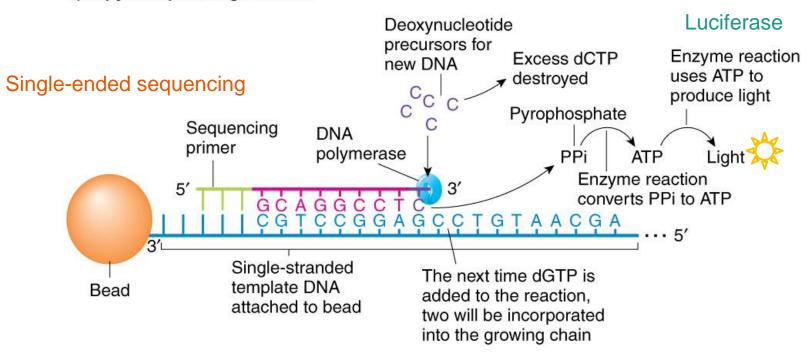
#### Illumina's flow cell



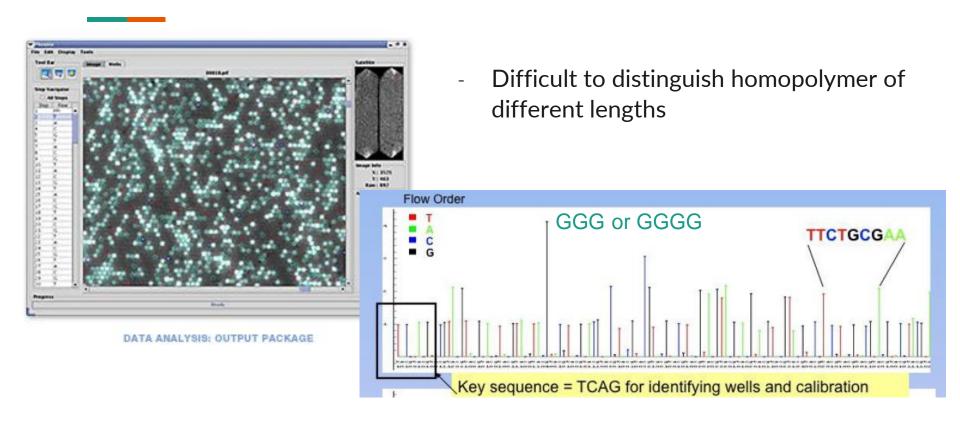
http://training.bioinformatics.ucdavis.edu

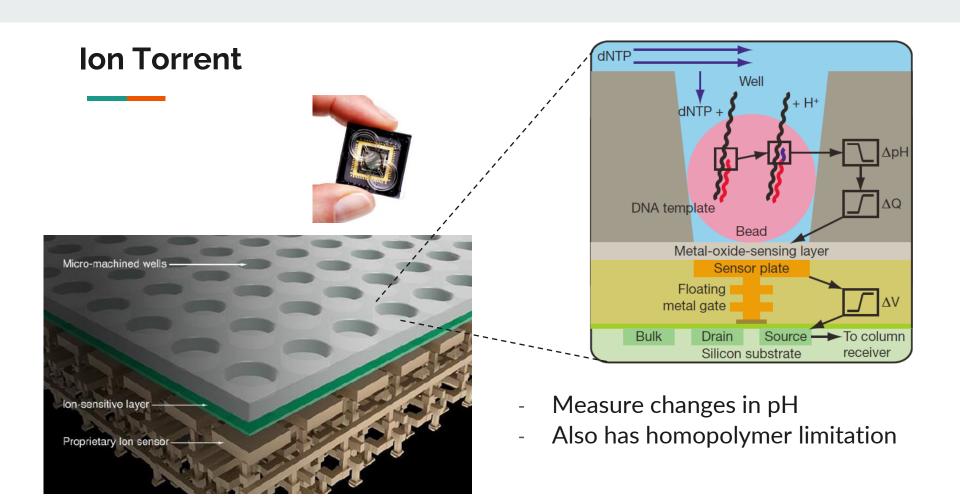
### **Pyrosequencing**

#### a) A pyrosequencing reaction

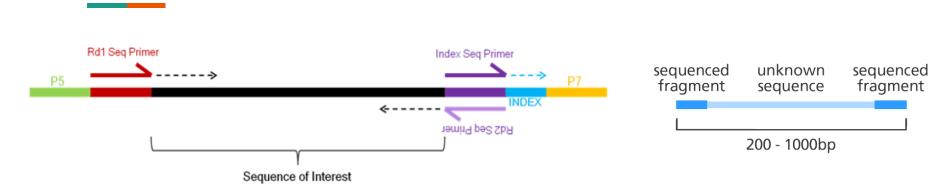


### Limitation of pyrosequencing

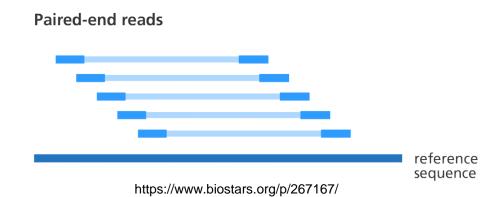


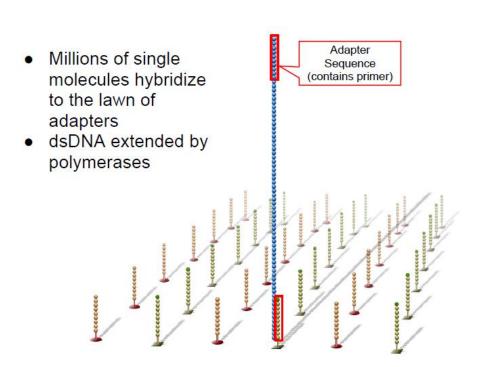


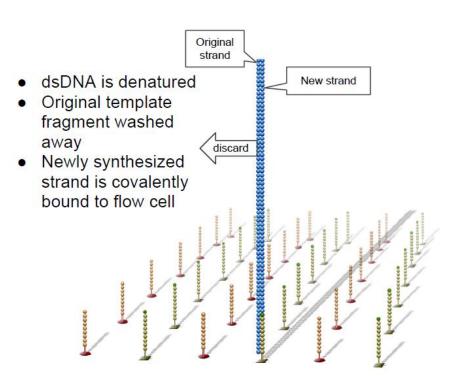
#### Illumina / Solexa



- Enable paired-end sequencing
- Improve mappability
- Identify splice junction
- Identify gene fusion
- Identify translocation

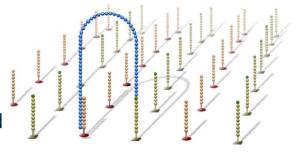


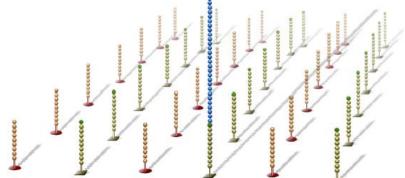


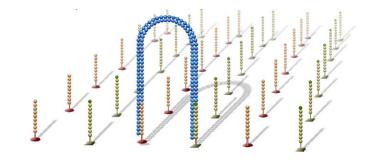


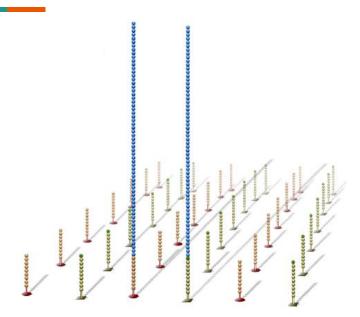
 Resulting covalentlybound DNA fragments are bound to the flow cell surface in a random pattern  Single-strand flops over to hybridize to adjacent adapter, forming a bridge

 dsDNA synthesized from primer in hybridized adapter

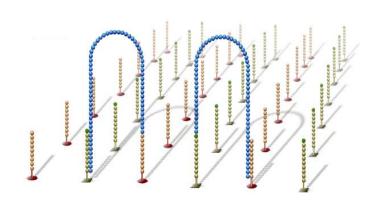






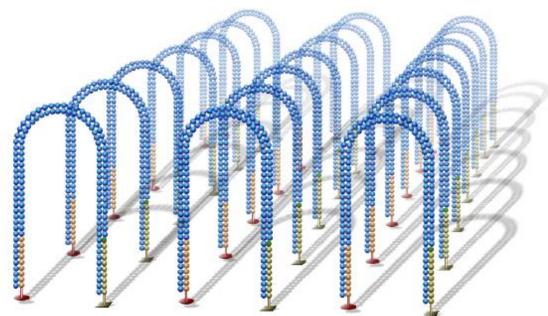


 dsDNA bridge is denatured

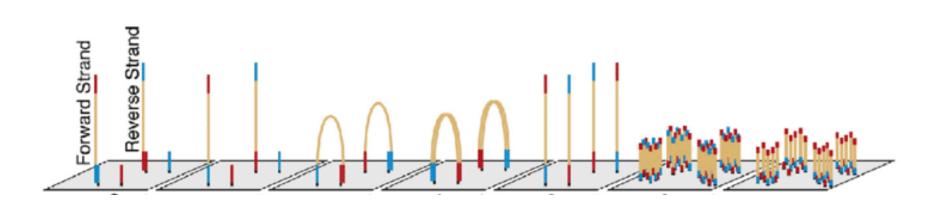


- Single strands flop over to hybridize to adjacent adapters, forming bridges
- dsDNA synthesized by polymerases

 Bridge amplification cycles repeated many times

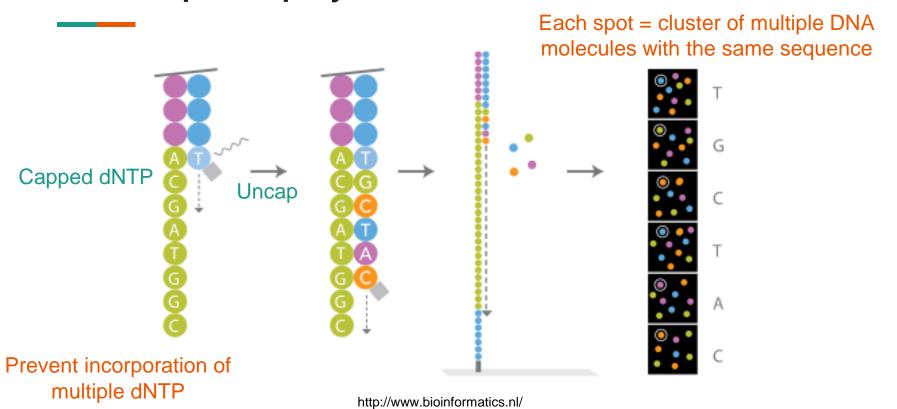


### Illumina / Solexa DNA amplification



 Improve sensitivity by sequencing clusters of the amplified DNA molecules deriving from the same original DNA

#### Multi-step DNA polymerization



### **Pros and cons**

Platform	Read Length	Run Time	Gb/ Run	Advantage	Disadvantage
454 (Pyrosequencing)	400+	1 day	0.7	Long read length	Homopolymer error Single-ended only
Illumina	50-300	10 days	600	Low cost per base	Short reads Long run time
Ion Torrent	200-400	2 hrs	100	Fast run times	Homopolymer error

#### **Tradeoffs**

#### Sanger

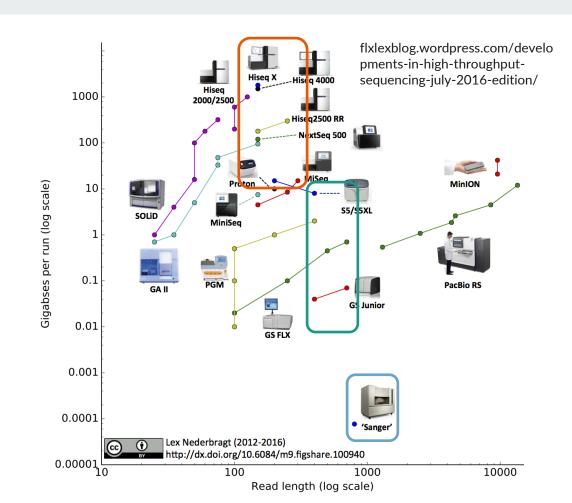
- 1000 bp, low throughput

#### 454 and Ion Torrent

- 400+ bp, medium throughput

#### Illumina

<300 bp, high throughput</p>



#### Use cases

#### Sanger

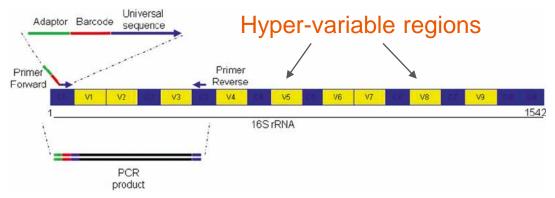
Validate sequences

#### **454 Pyrosequencing**

Metagenomics

#### **Ion Torrent**

Fast turn-around situation



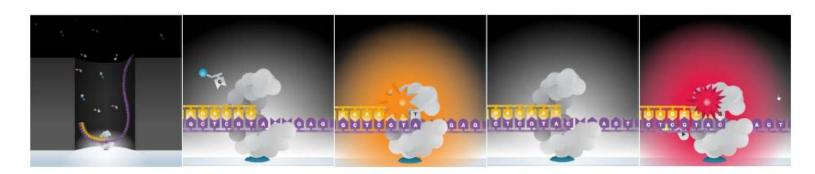
Del Chierico et al. Methods in Molecular Biology 2015

#### Illumina

- Whole-genome
- Practically anything...

## 3<sup>rd</sup> Generation Sequencing (Long-Read)

### Single-Molecule Real-Time (SMRT) sequencing



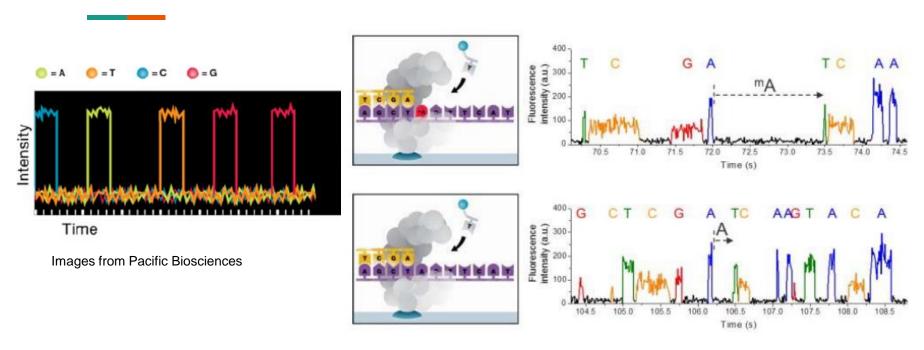
Zero-mode waveguide (ZMW)

Phospholinked nucleotide

Images from Pacific Biosciences

- Faster, more durable DNA polymerase
- Small wells with single DNA molecule
  - Zero-mode waveguide = nanophotonic confinement structure
  - Allow monitoring of fluorescence signal from individual reaction
- No amplification = direct quantification of DNA/RNA abundance

#### Video data



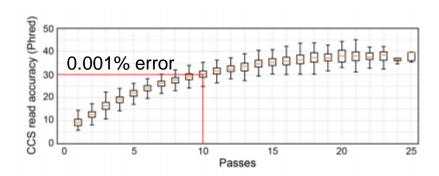
- Compared to image data from Illumina platform
- Video gives time information → identification of modified DNA/RNA

#### High error rate

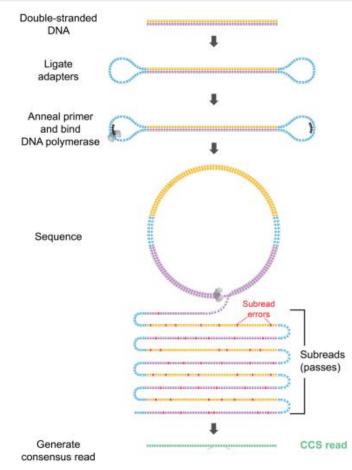
ICCGGAGCGACGCGTACGATTAAAGCACGTACTGCGTATGCGTATCCCTAGCTTGCTAGGCTAGTATGCTAGATTAAAGCTC GTACTGC**A**TATG**T**GTATGCCTAGCTAGCTAGG**A**TAG**C**ATGCTAGATTAAAGCT GTACCGCGTATGCGTATGCCAGGTAGCTAGGCTAGTATGCT PCCGGATCTACGCGTACGATTAAAGCTAGTACTGCGTATGCGTTTTGCCTATGTAGCTA  $ext{PCCGGATCGACGTGTACGATTAGCTCTTACTGCGTATACGTATGCCTAGGTAGCTAGGCTAGTATGCTAGATTAAAGCTCGAAC1}$  $\mathtt{PCT}_{\mathsf{GGATCGACGCGTACGACAGCTCGTACTGTGTATGCGTATGCCTAGCTCGCTACGCTAGTATGCTC}$  ${ t PCC}{ t C}{ t GATCGACGCGT}{ t IC}{ t GATTAAAGCTCGT}{ t C}{ t T}{ t A}{ t T}{ t G}{ t C}{ t C}{ t T}{ t A}{ t G}{ t C}{ t C}{ t T}{ t A}{ t G}{ t C}{ t C}{ t T}{ t A}{ t G}{ t C}{ t C}{ t T}{ t A}{ t G}{ t C}{ t C}{ t T}{ t A}{ t G}{ t C}{ t C}{ t T}{ t A}{ t C}{ t C}{ t C}{ t T}{ t A}{ t C}{ t C}{ t C}{ t C}{ t T}{ t C}{ t C}{ t T}{ t C}{ t C}{ t C}{ t T}{ t C}{ t C}{ t T}{ t C}{ t C}{ t C}{ t T}{ t C}{ t$ PCCGGATCGGCGCGTACGATTAAAGCTCGTACTGCGGATGCGTATGCCTAGCTGGCTAGGCGAGTATGCTAGATGAAAGGTCGTAC1 

- 5-15% error compared to 0.01% of Illumina
- How do we solve this?

#### Circular consensus sequencing

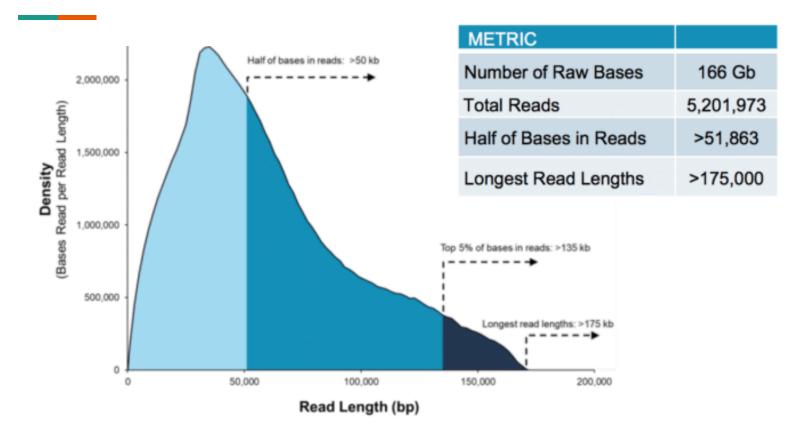


- Circular extension of each DNA molecule
- Read the extended molecules = multiple resequencing of the original sequence
- Take the consensus (majority vote)
- P(correct base in >k of N passes)  $\sim$  Binomial

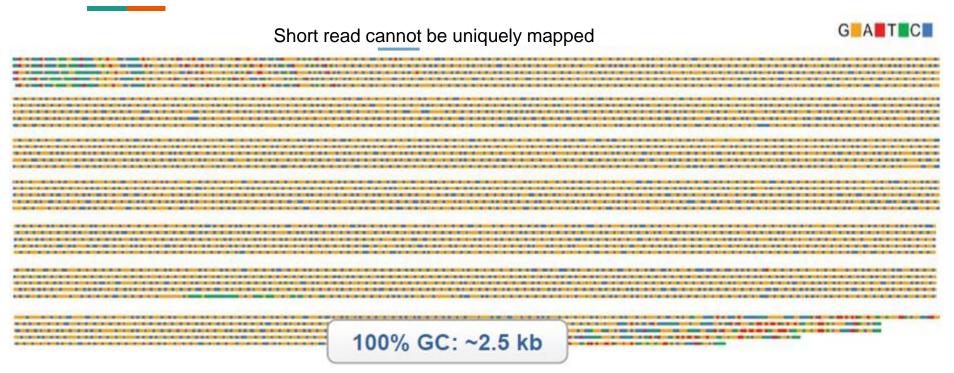


Images from Pacific Biosciences

#### Read length >> 10kb



### Resolve repetitive region

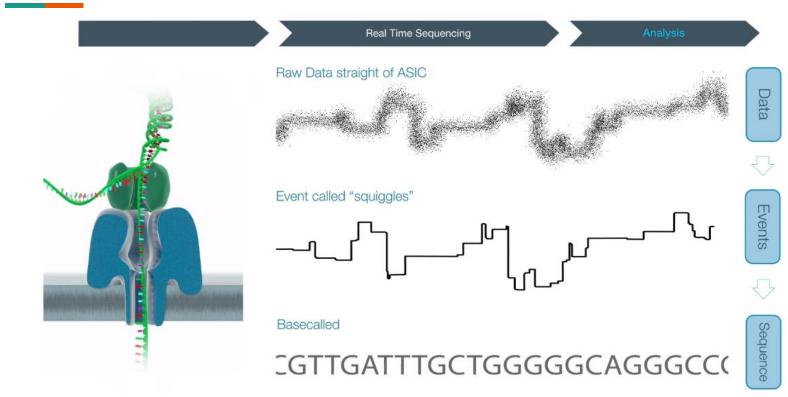


~20% of human genome

#### Resolve haplotype

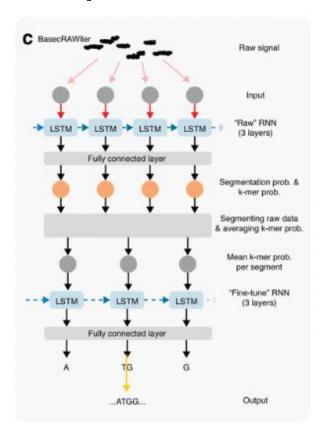


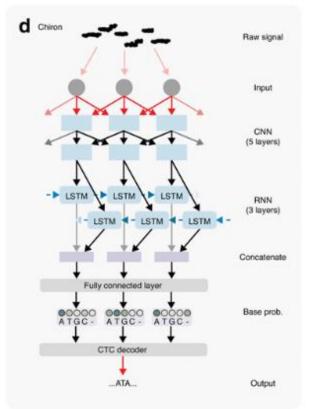
#### **Nanopore**



#### Basecalling with deep neural networks

- Trained using data from synthetic DNA
- 14% base error
- Improved to 3-5% using bioinformatics and machine learning



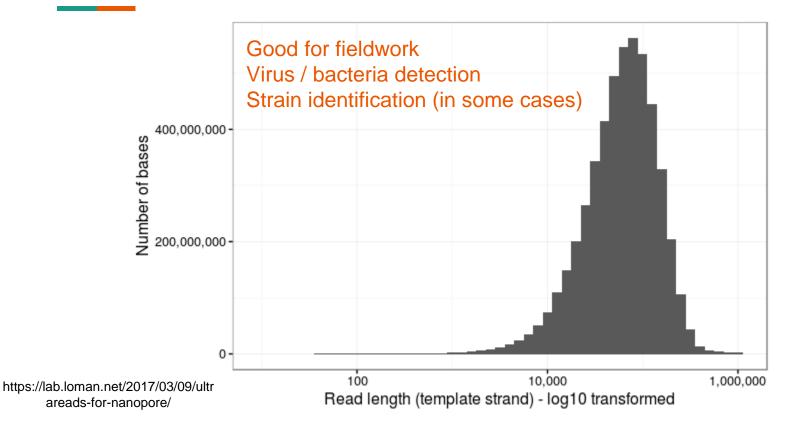


Rang, Kloosterman, and de Ridder. Genome Biology 2018

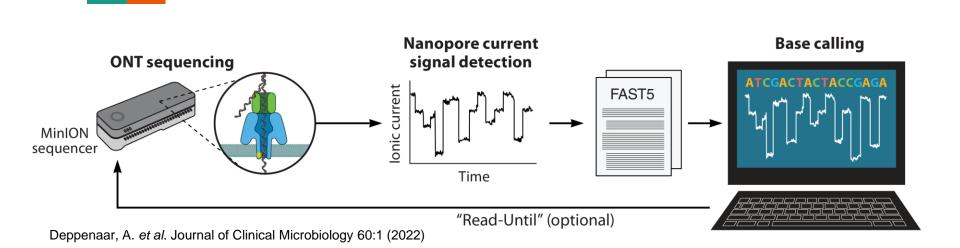
### Portability & fast turn-around time

	Flongle	MinION	GridION (5 flow cells)	PromethION (48 flow cells)
		www.	OF THE STATE OF	O O
Maximum run time	16 hours	72 hours	72 hours	64 hours
Theoretical 1D maximum yield	Up to 3.3 Gb	Up to 40 Gb	Up to 200 Gb	Up to 15 Tb
Current 1D maximum yield	Up to 2 Gb	Up to 30 Gb	Up to 150 Gb	Up to 8.6 Tb
Available channels	Up to 126	Up to 512	Up to 2,560	Up to 144,000

### Read length up to Mb

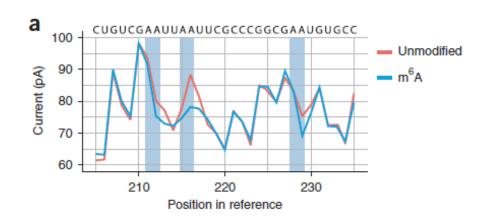


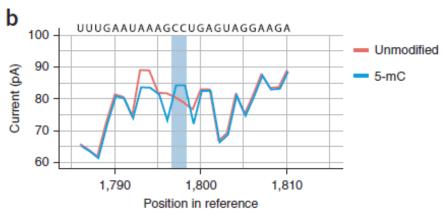
#### Real-time data



- Real time ionic flow signals
- Ability to manipulate individual pore and terminate unwanted reads
- Rapid decision making (no need to wait for the full 16-72hr run)

#### **Detection of modified nucleotides**

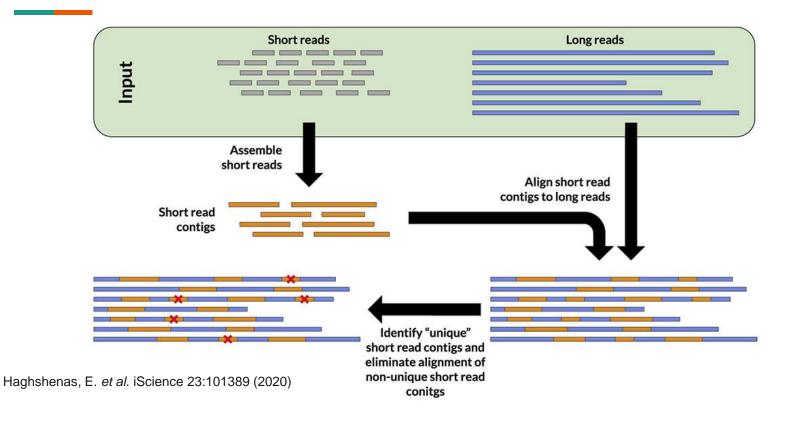




Geralde et al. Nature Methods 15, 201-206 (2017)

- Modified nucleotides = different 3D structure = different change in ionic flow
- Trained using synthetic nucleotides

### Combining short and long read data



## Applications of DNA/RNA sequencing

#### **Sequencing scope**

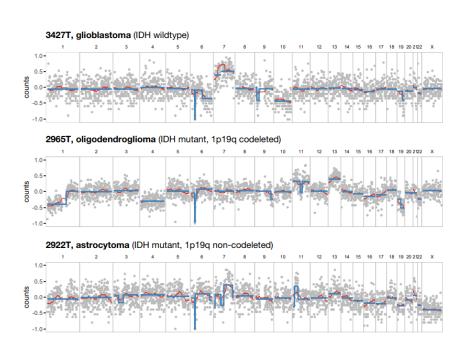
Cost = Base Pair = Scope x Depth

#### Reduced scope

- Exome sequencing = exons only
- Amplicon sequencing = selected loci
  - 16S rRNA, RDRP gene
  - (Cancer) gene panels

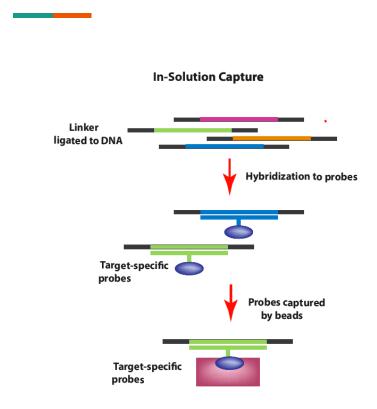
#### Reduced depth

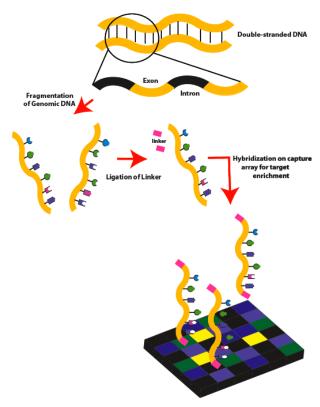
- Ultra-low pass
  - Detect chromosomal copy alternation
  - Estimate tumor fraction



Euskirchen, P. et al. Acta Neuropathol 134:691-703 (2017)

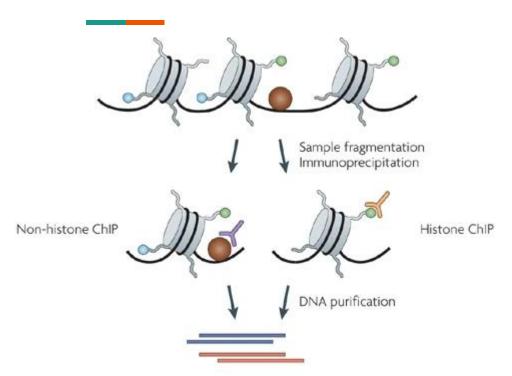
#### **Enrichment for targeted sequencing**

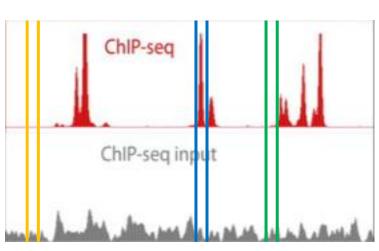




https://en.wikipedia.org/wiki/Exome\_sequencing

### **Chromatin immunoprecipitation**

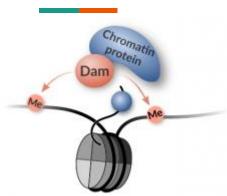




Park et al. Nat Rev Genet 10:669-680 (2009)

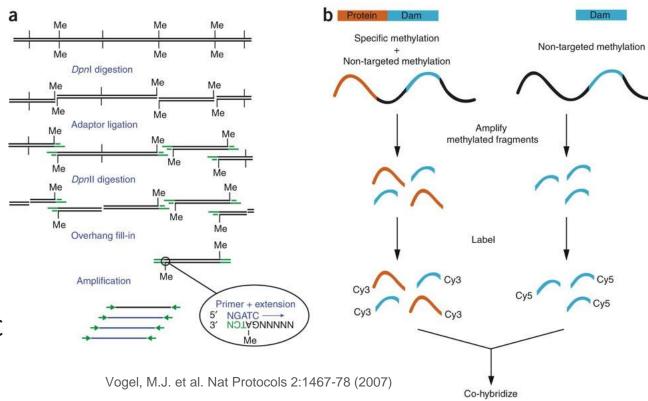
DNA-bound protein / histone modification

### DNA adenine methylatransferase (DamID)

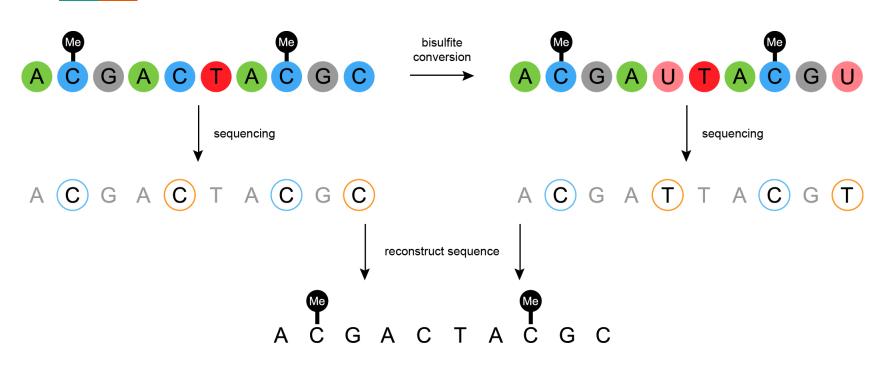


https://marshall-lab.org/damid/

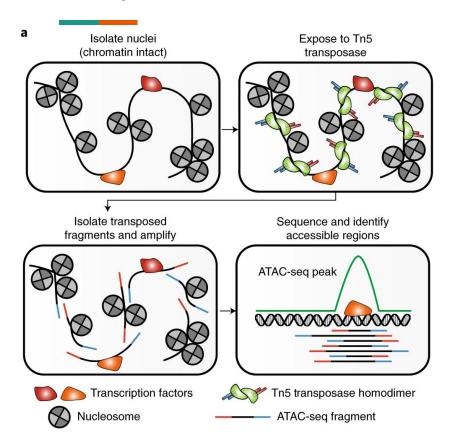
- Dam attached to protein of interest
- Methylation of GATC
- DpnI/DpnII enzymes

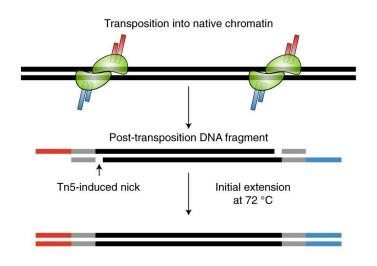


#### Bisulfite sequencing



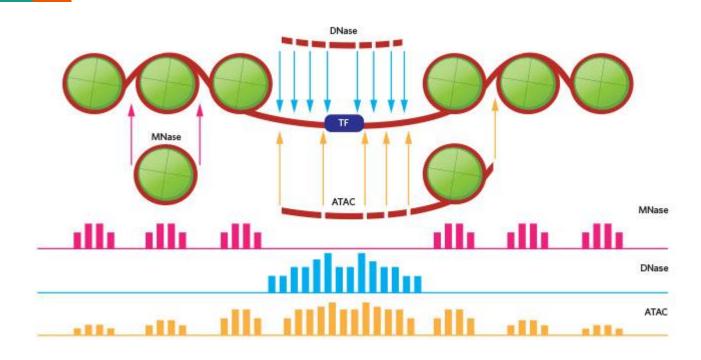
#### Assay for transposase-accessible chromatin (ATAC)



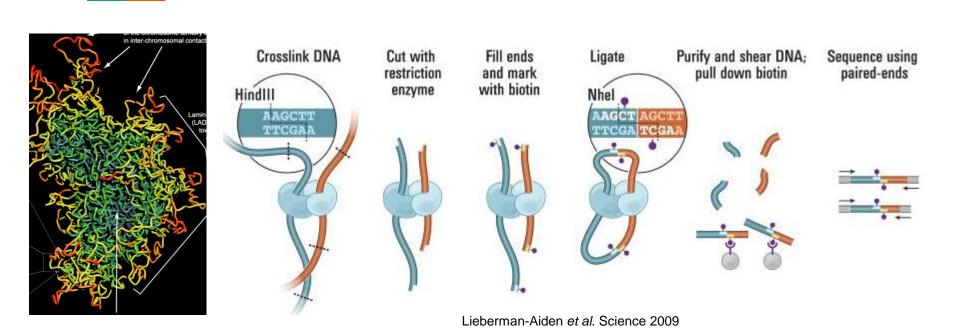


Transposase with sequencing adapter insertion into open chromatin

### Targetting bound or unbound chromatin

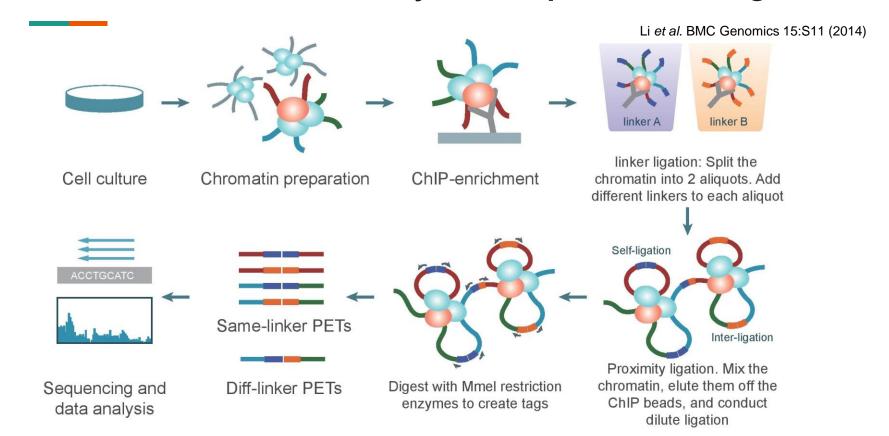


### Chromatin conformation capture

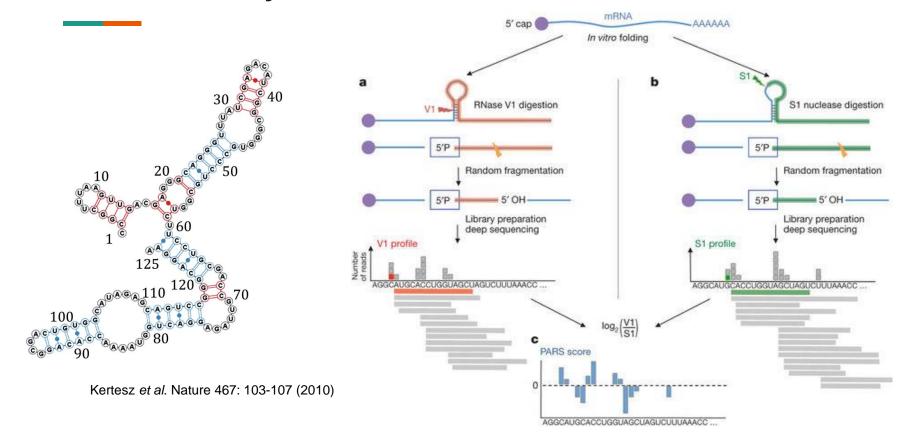


Cross-link proximal DNA  $\rightarrow$  join ends from different regions  $\rightarrow$  sequencing

### Chromatin interaction analysis with paired-end tag



### RNA secondary structure



### Sequencing application key points

- Scope of sequencing = enrich target DNA
- DNA-binding protein = antibody pull-down
- Detection of DNA modification
- Targeting bound / unbound DNA
- Enzyme specificity

### Any question?

- See you again on August 25<sup>th</sup> 1-2:30pm