

# INTRODUCTION AU SÉQUENÇAGE À HAUT DÉBIT POUR LA GÉNOMIQUE

Claude THERMES

INSTITUT DE BIOLOGIE INTÉGRATIVE DE LA CELLULE – I2BC

GIF-SUR-YVETTE



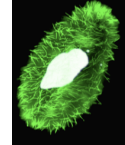
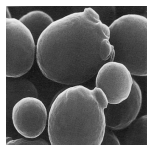
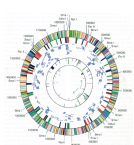
Diplôme Universitaire en Bioinformatique intégrative - DU-Bii - 08/03/2021



# Premiers génomes entièrement séquencés

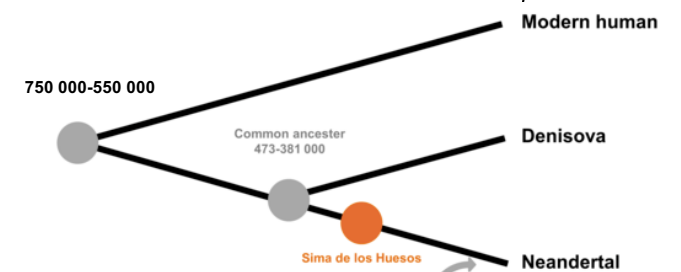
Médecine  
personnalisée  
>> milliers de génomes

2 bactéries    Levure    Ver    Plante    Mouche    Humain    Souris    Poule    Chimpanzé    Paramécie    2 humains



1000  
génomes  
humains

1990 91 92 93 94 95 96 97 98 99 2000 01 02 03 04 05 06 07 08 09 12 ... 21

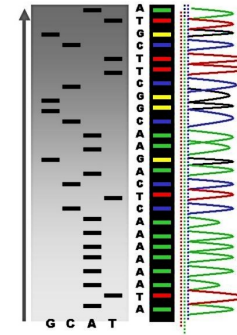


## 1st generation : Sanger sequencing

- Has been the major method up to 2005

### *Limitations*

- Extremely high cost
- Long experimental set up times
- High DNA concentrations needed



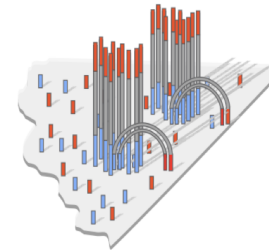
## 2<sup>d</sup> generation

- Single DNA molecules replicated in clusters
- Very high throughput

### *Limitations*

- Maximum read length  $\leq 300\text{bp}$

*Illumina*



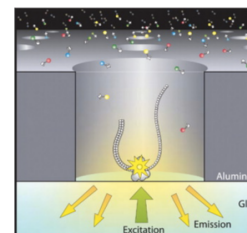
## 3<sup>rd</sup> generation

- Single molecules sequencing
- Very long reads

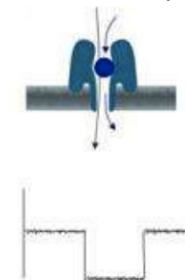
### *Limitations*

- High error rates

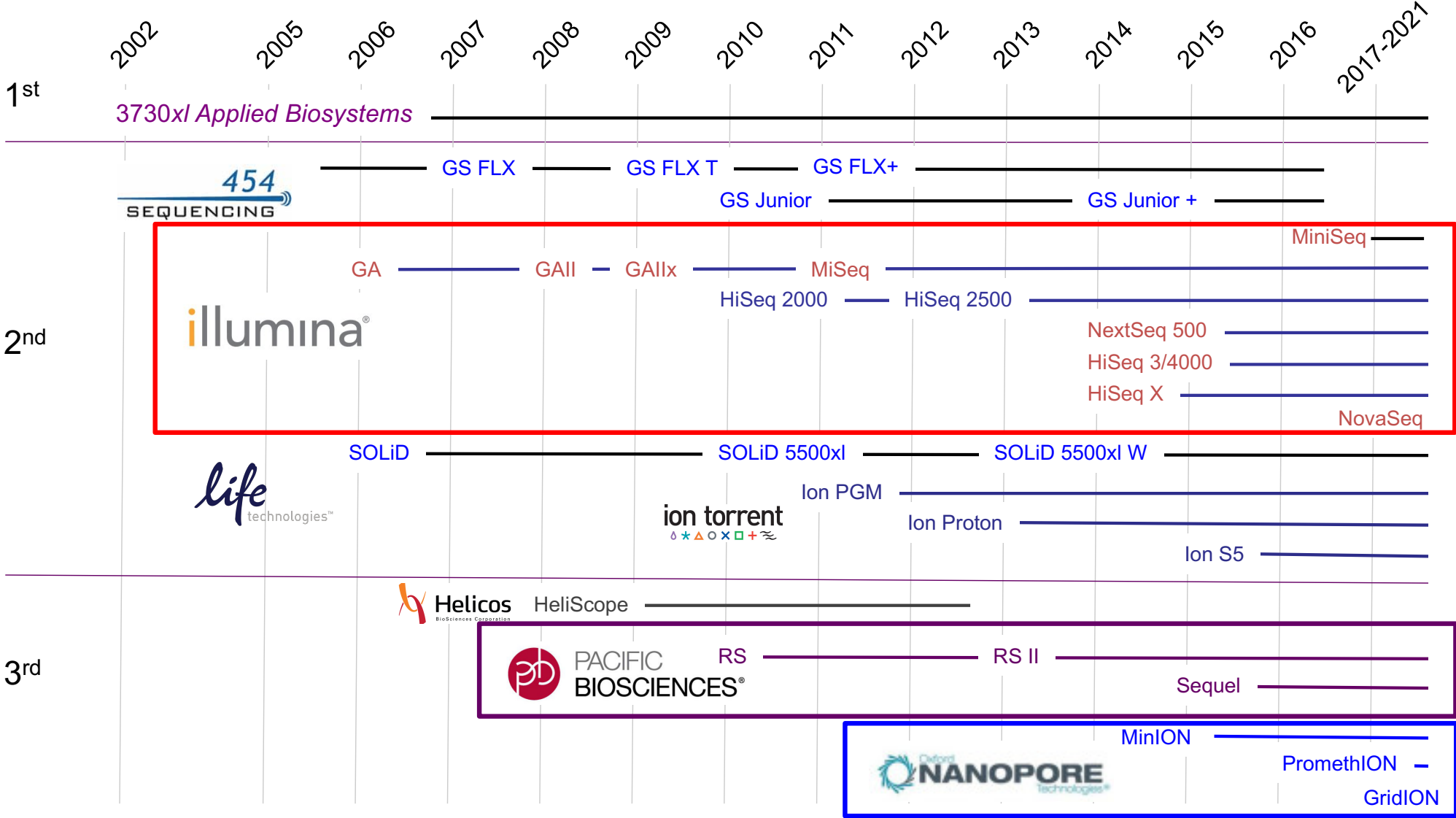
*PacBio*



*Oxford Nanopore*



# Sequencing technologies



# Illumina : the winning technology



**MiniSeq**  
25 million reads



**MiSeq**  
25 millions reads, 2 x 300 bp



**NextSeq**  
400 million reads



**HisSeq 4000**  
5 billion reads

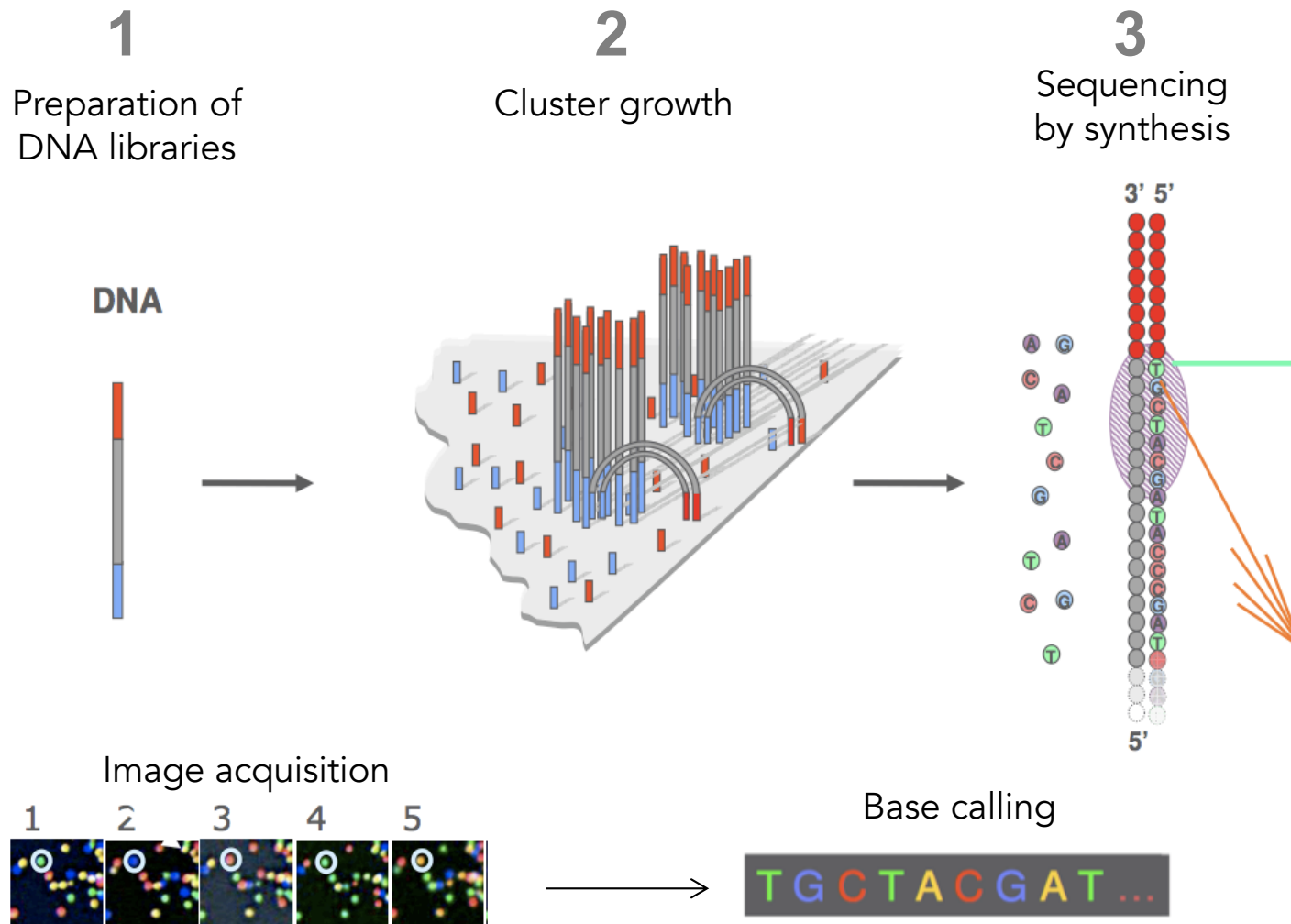


**HisSeq X**  
6 billion reads



**NovaSeq 6000**  
20 billion reads

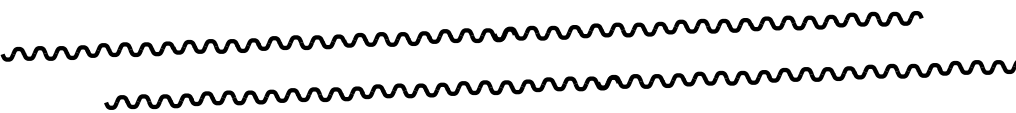
# General scheme of Illumina sequencing



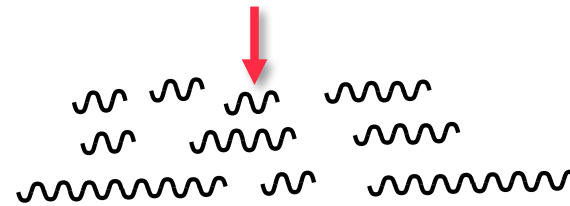
# 1 - Preparation of DNA-seq Libraries

## Illumina TruSeq technology

Genomic DNA



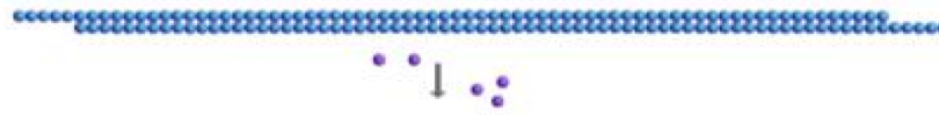
Sonication



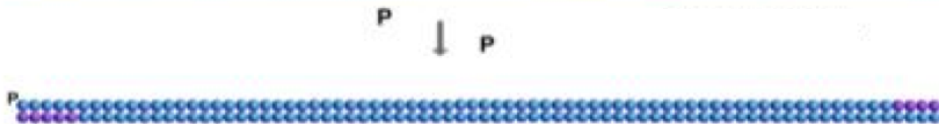
Size selection



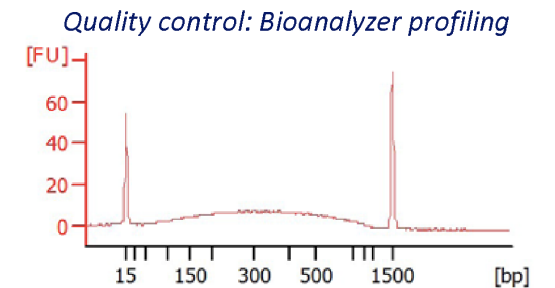
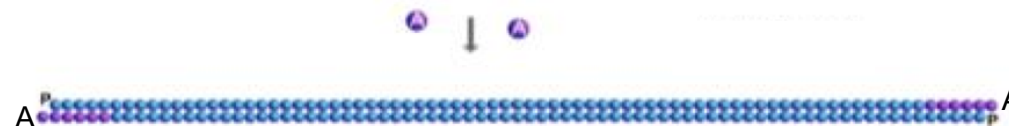
End repair



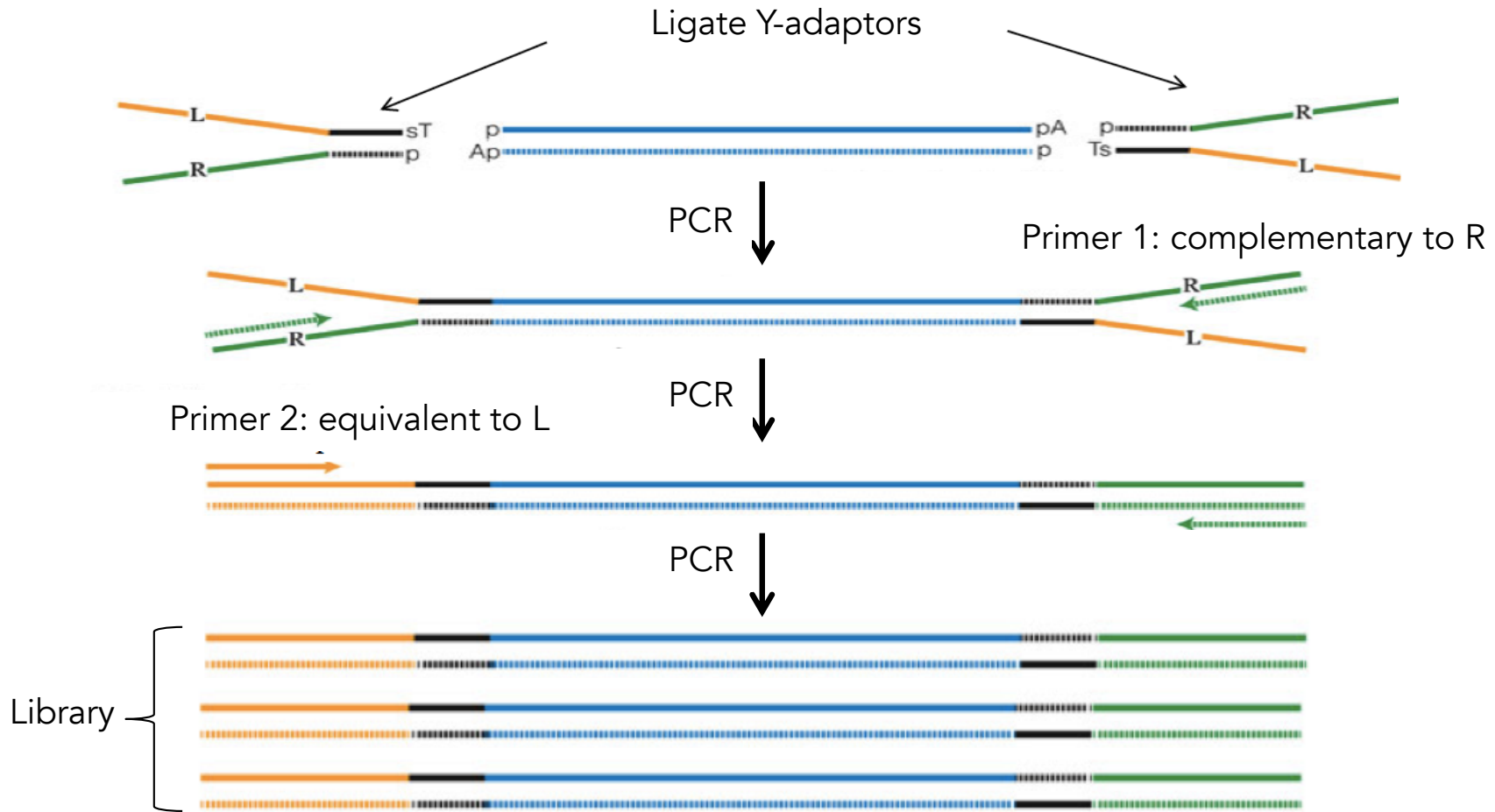
Phosphorylation



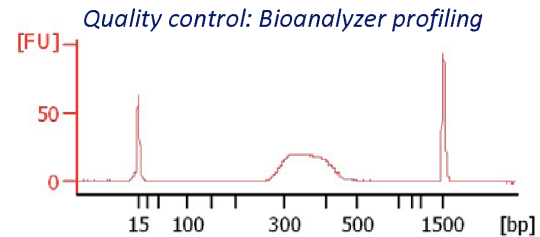
A - overhang



# 1 - Preparation of DNA-seq Libraries



sequencing



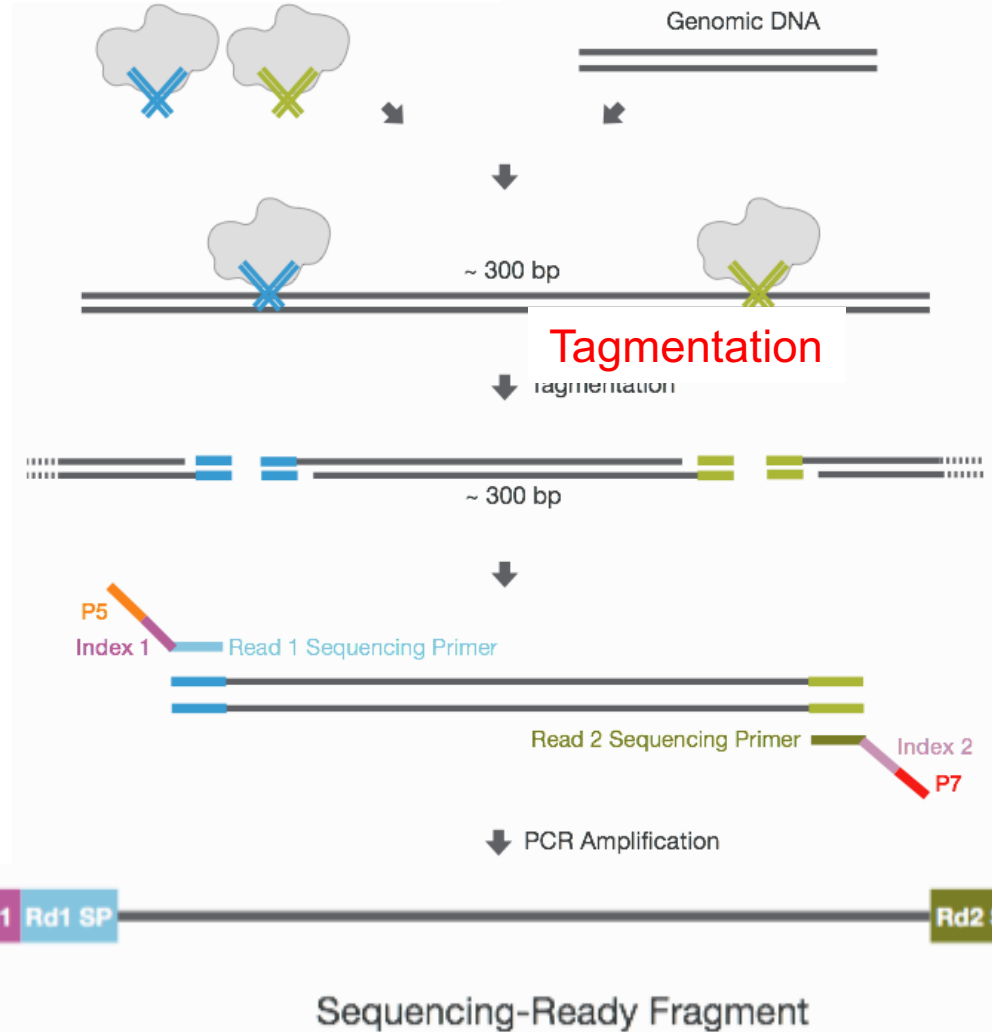


# 1 - Preparation of DNA-seq Libraries

## Nextera "tagmentation"

Tagment Enzyme fragments DNA and attaches junction adapters (blue and green) to both ends of the tagmented molecule

Transposomes / Tagment Enzyme



Dual barcode approach



up to 96 indexed samples

requires small quantities 1ng (bacteria) to 50 ng (human)

# 1 - Preparation of DNA-seq Libraries

## SINGLE READ and PAIRED-END SEQUENCING

- **Single end**: Sequence one physical end of DNA fragment

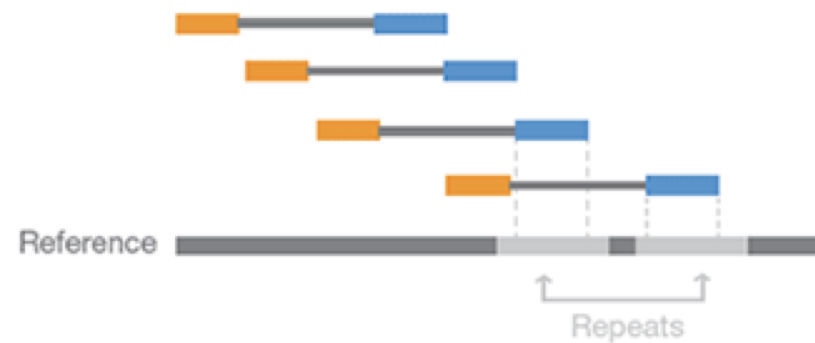


- **Paired End**: Sequence both physical ends of DNA fragment
  - End distance: < 800nt

Paired-End Reads



Alignment to the Reference Sequence



# 1 - Preparation of RNA-seq Libraries

RNA after rRNA depletion

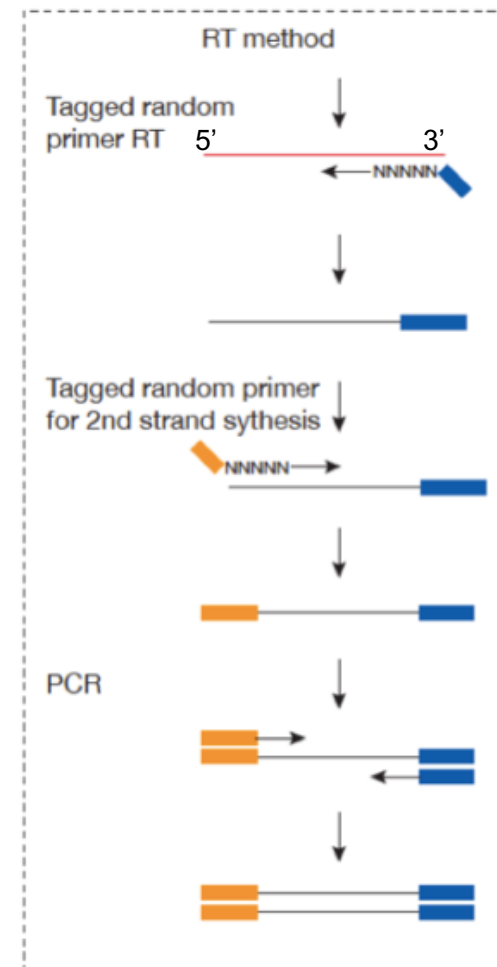
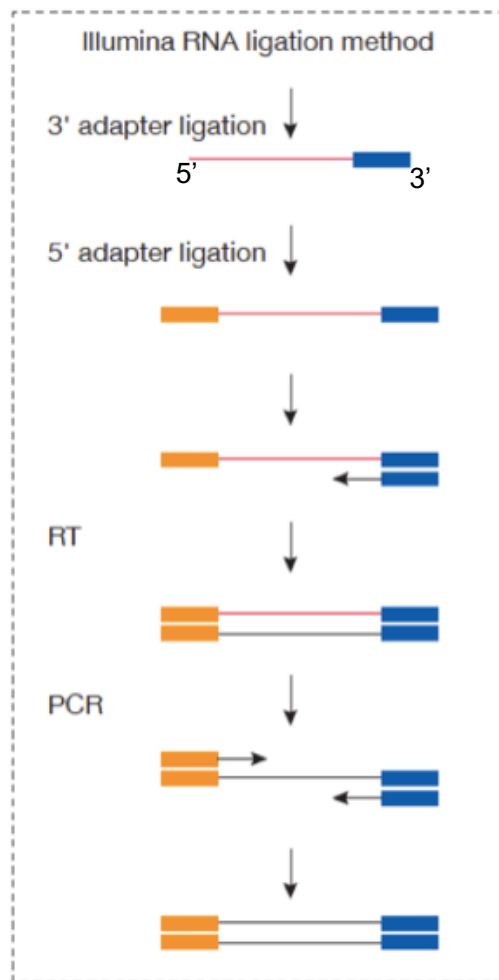


RNA fragmentation



both are directional

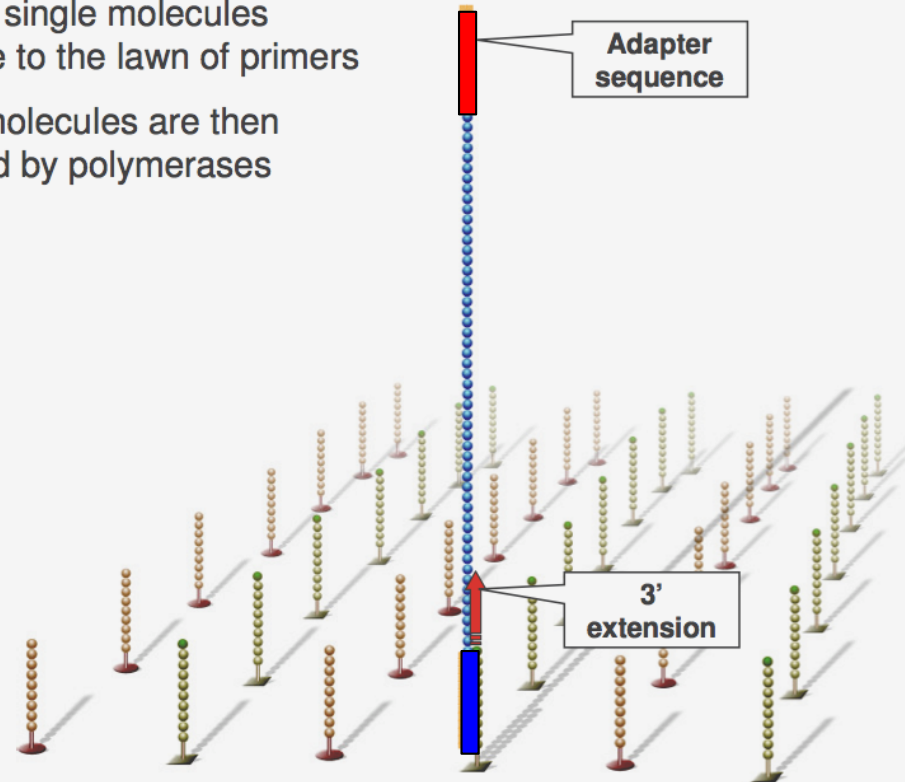
Kit smaRNA illumina



## 2 – Cluster growth

### Cluster Generation: *Hybridize Fragment & Extend*

- ▶ > 100 M single molecules hybridize to the lawn of primers
- ▶ Bound molecules are then extended by polymerases



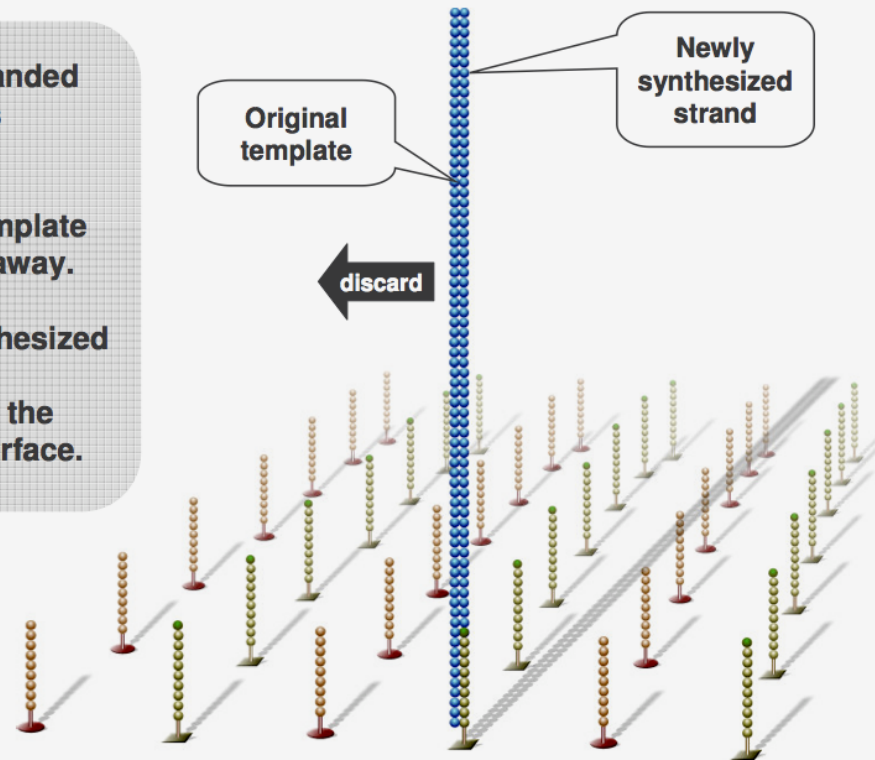
## 2 – Cluster growth

### Cluster generation: *Denature double-stranded DNA*

Double-stranded molecule is denatured.

Original template is washed away.

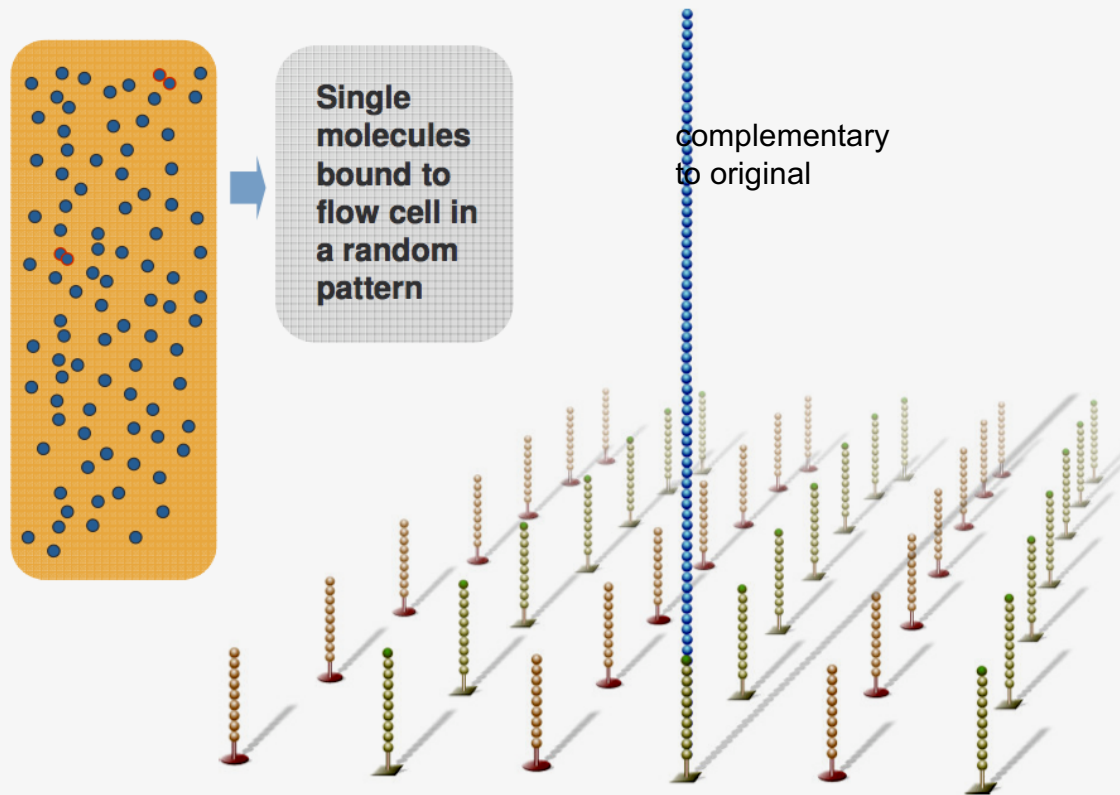
Newly synthesized covalently attached to the flow cell surface.



## 2 – Cluster growth

### Cluster generation:

*Covalently bound spatially separated single molecules*



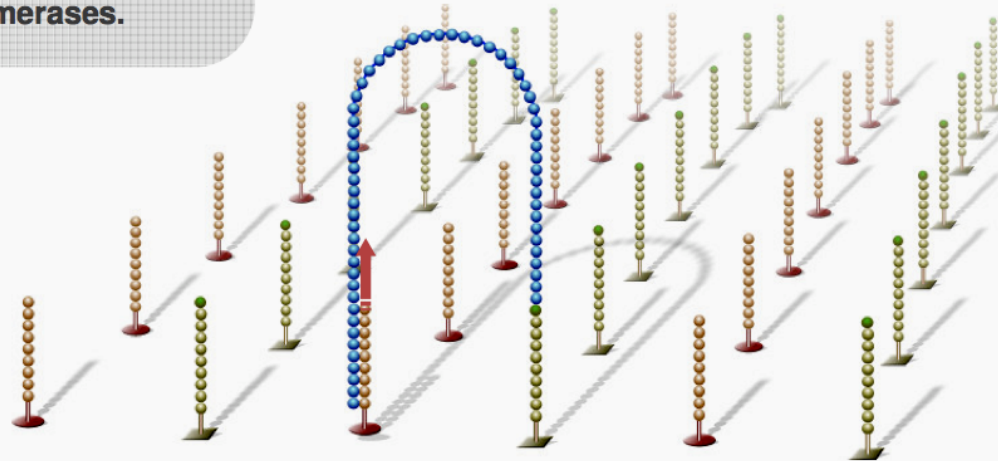
## 2 – Cluster growth

### Cluster generation:

*Bridge amplification*

Single-strand flips over to hybridize to adjacent primers to form a bridge.

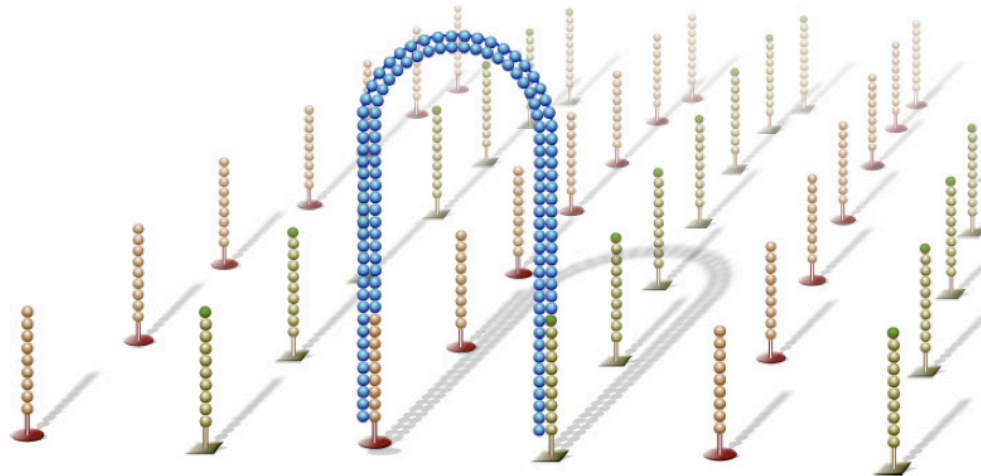
Hybridized primer is extended by polymerases.



## 2 – Cluster growth

### Cluster generation: *Bridge amplification*

→ double-stranded  
bridge is formed.





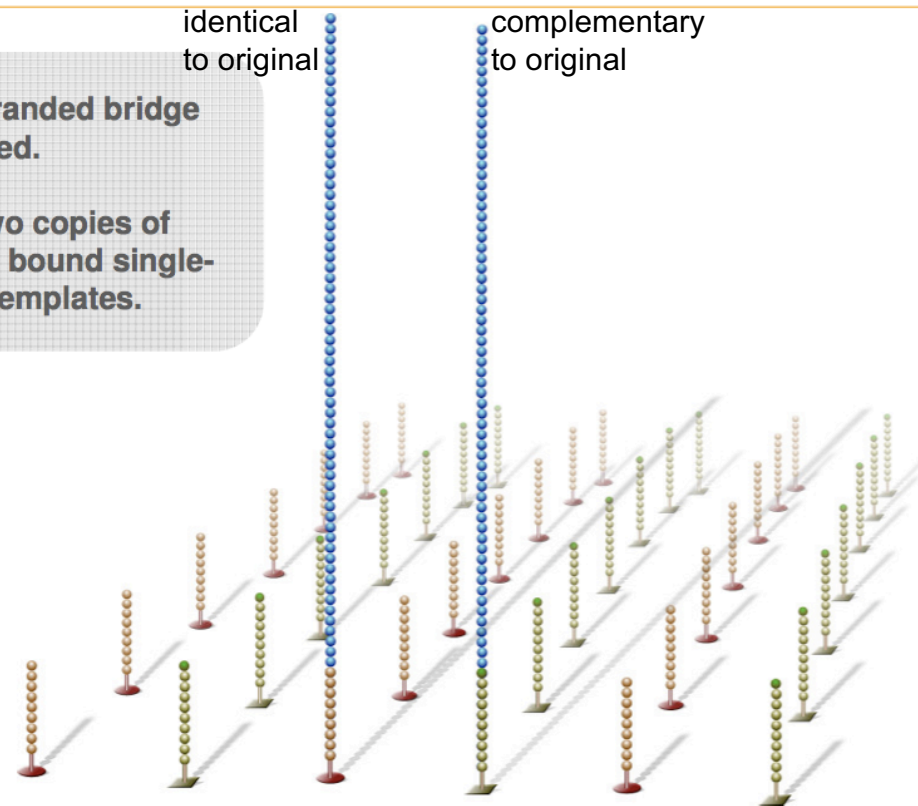
## 2 – Cluster growth

### Cluster generation:

*Bridge amplification*

**Double-stranded bridge is denatured.**

**Result: Two copies of covalently bound single-stranded templates.**

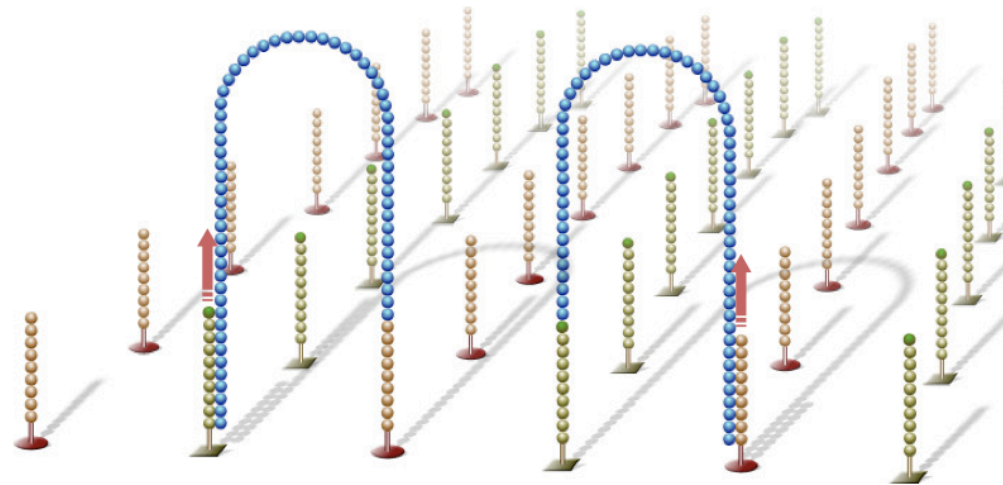


## 2 – Cluster growth

### Cluster generation: *Bridge amplification*

Single-strands flip over to hybridize to adjacent primers to form bridges.

Hybridized primer is extended by polymerase.

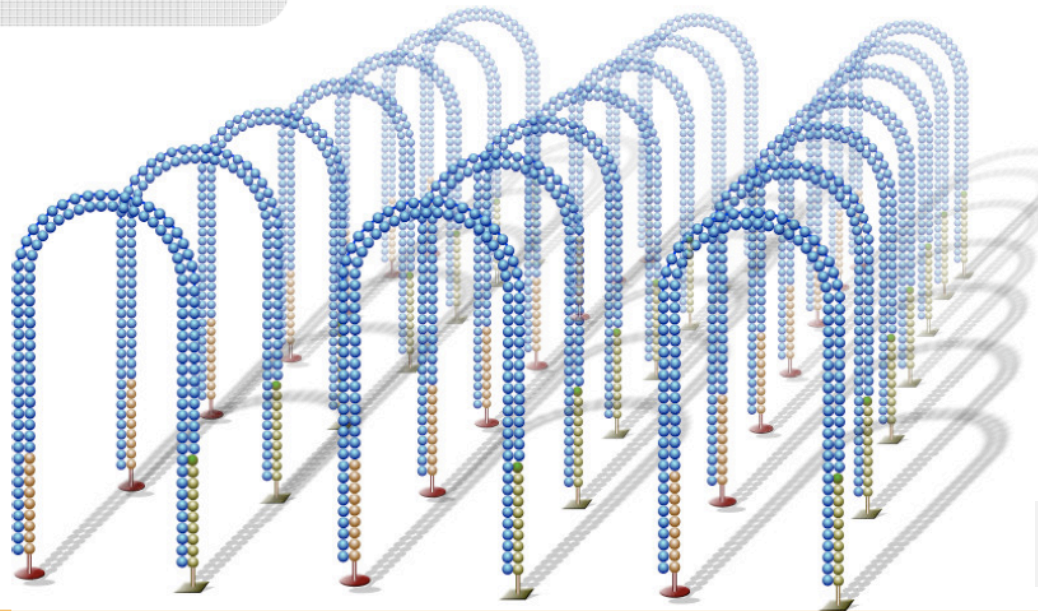


## 2 – Cluster growth

### Cluster generation:

*Bridge amplification*

Bridge amplification cycle repeated till multiple bridges are formed

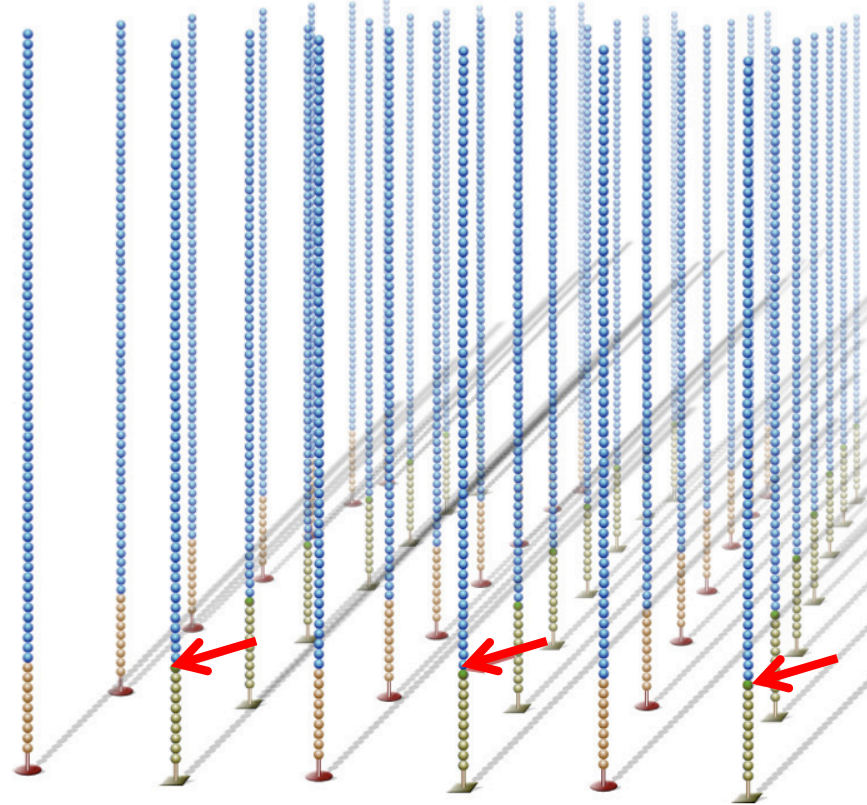


## 2 – Cluster growth

### Cluster generation

dsDNA  
bridges  
denatured.

Reverse  
strands  
cleaved  
and  
washed  
away.

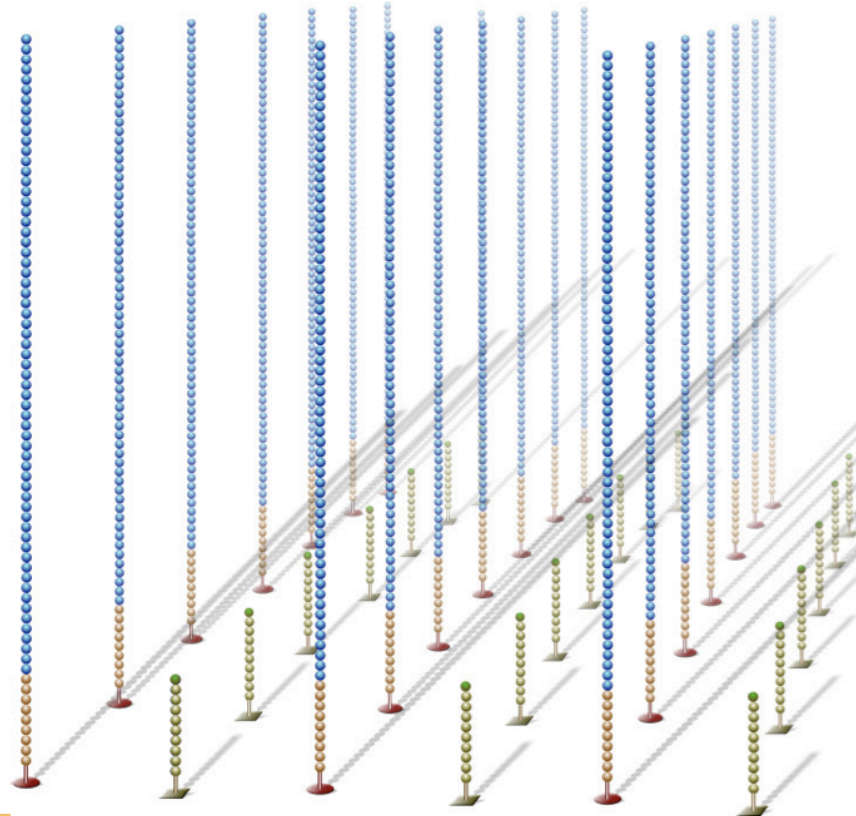


Cleavage of  
a chemically  
modified  
nucleotide

# 2 – Cluster growth

## Cluster generation

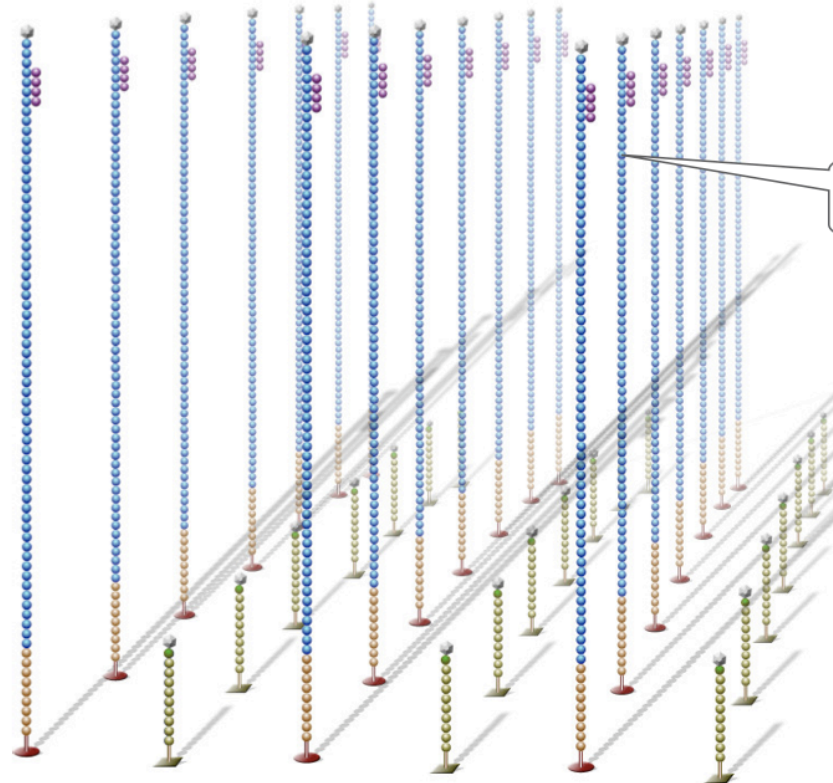
... leaving  
a cluster  
with forward  
strands only.



# 3 – Sequencing



Sequencing primer is hybridized to adapter sequence.



Sequencing primer

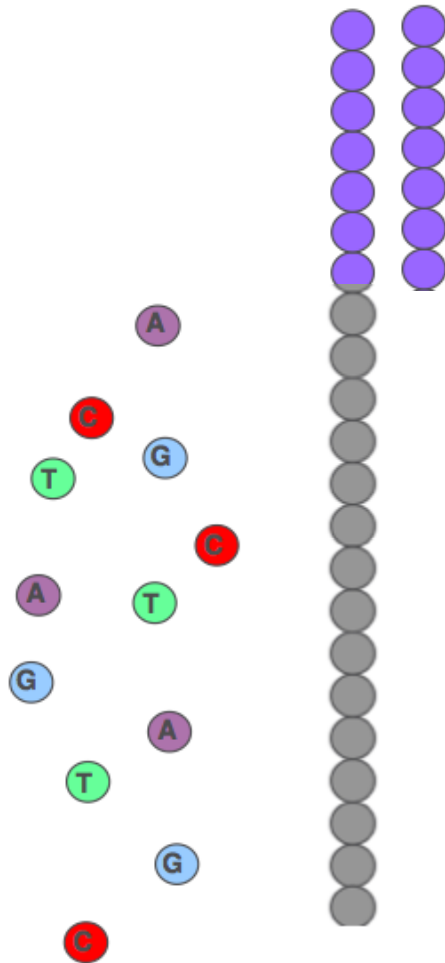
### 3 - Sequencing By Synthesis (SBS)

---



### 3 - Sequencing By Synthesis (SBS)

---

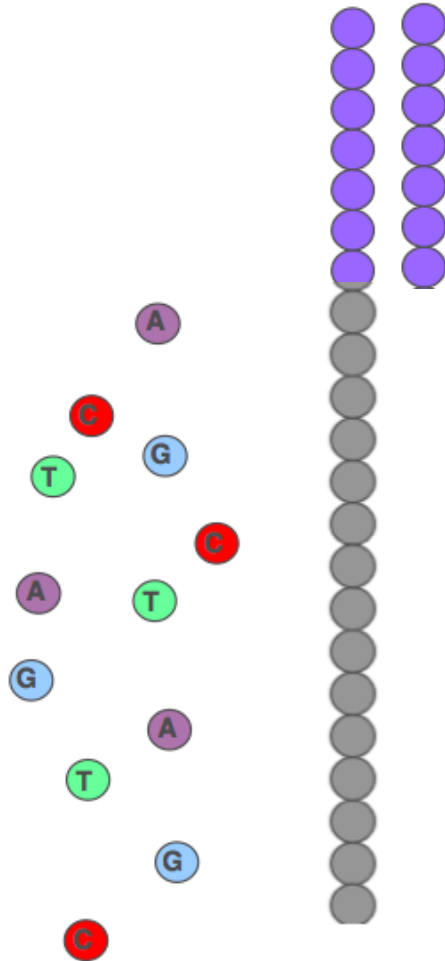


fluorescent-labeled terminator bound to each dNTP



### 3 - Sequencing By Synthesis (SBS)

---

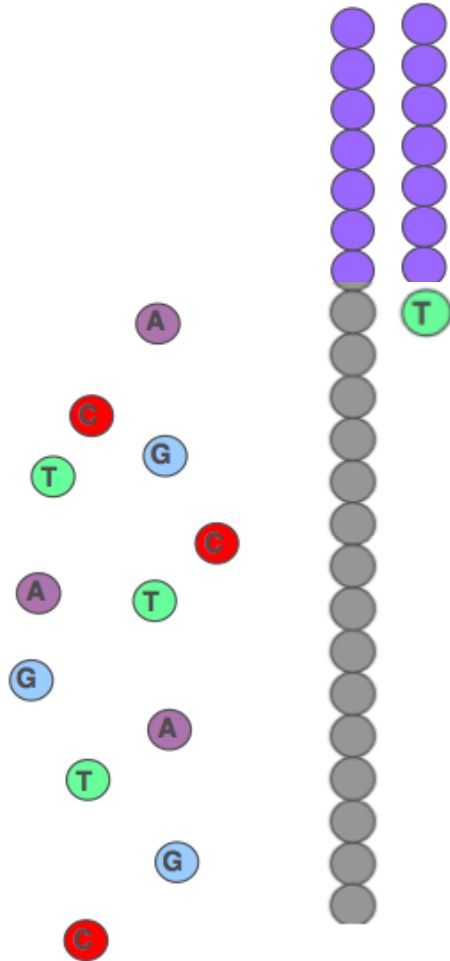


fluorescent-labeled terminator bound to each dNTP

Cycle 1 : add sequencing reagents

### 3 - Sequencing By Synthesis (SBS)

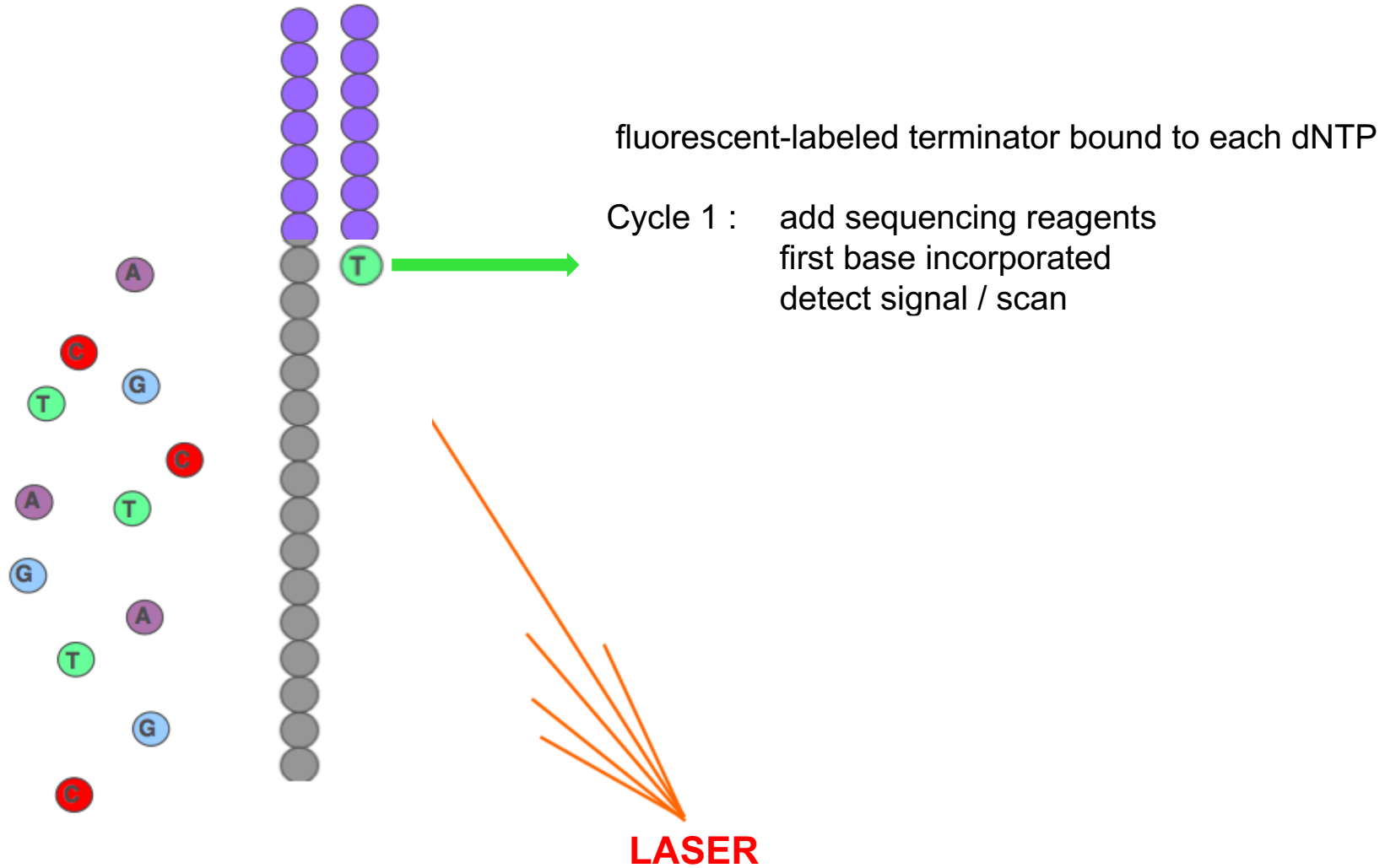
---



fluorescent-labeled terminator bound to each dNTP

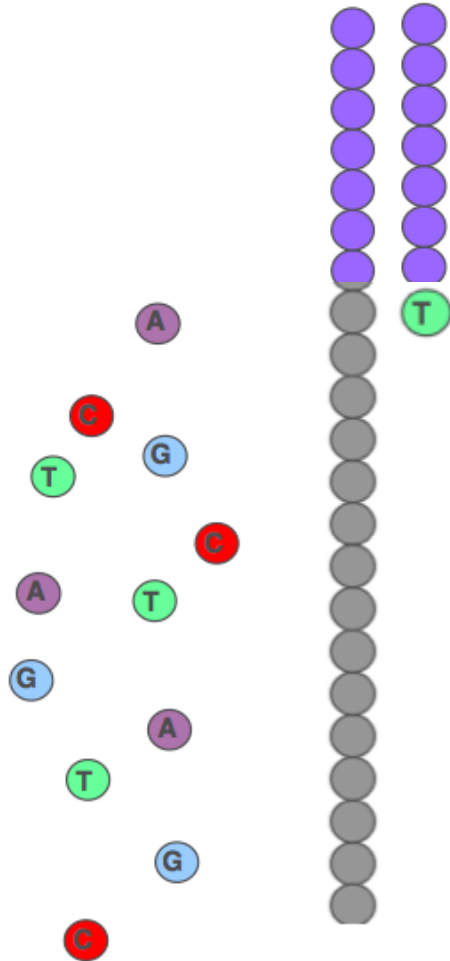
Cycle 1 : add sequencing reagents  
first base incorporated

### 3 - Sequencing By Synthesis (SBS)



### 3 - Sequencing By Synthesis (SBS)

---

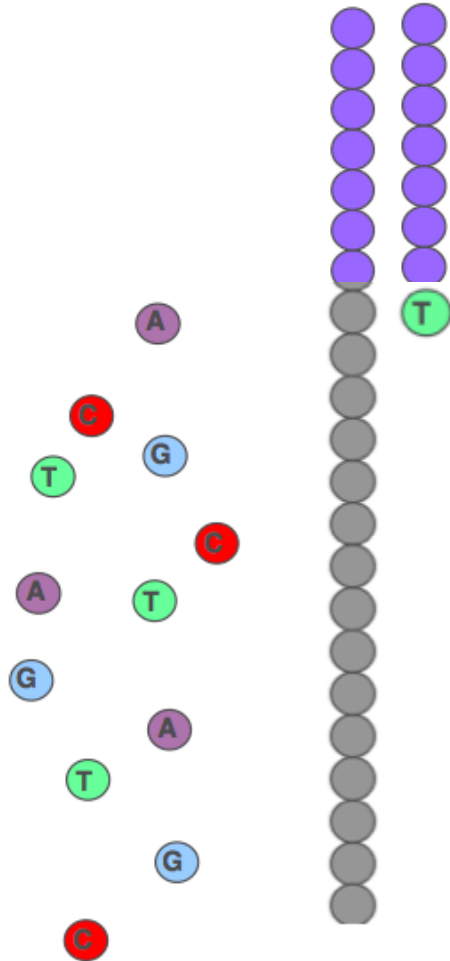


fluorescent-labeled terminator bound to each dNTP

Cycle 1 :    add sequencing reagents  
              first base incorporated  
              detect signal / scan  
              cleave terminator and dye

### 3 - Sequencing By Synthesis (SBS)

---

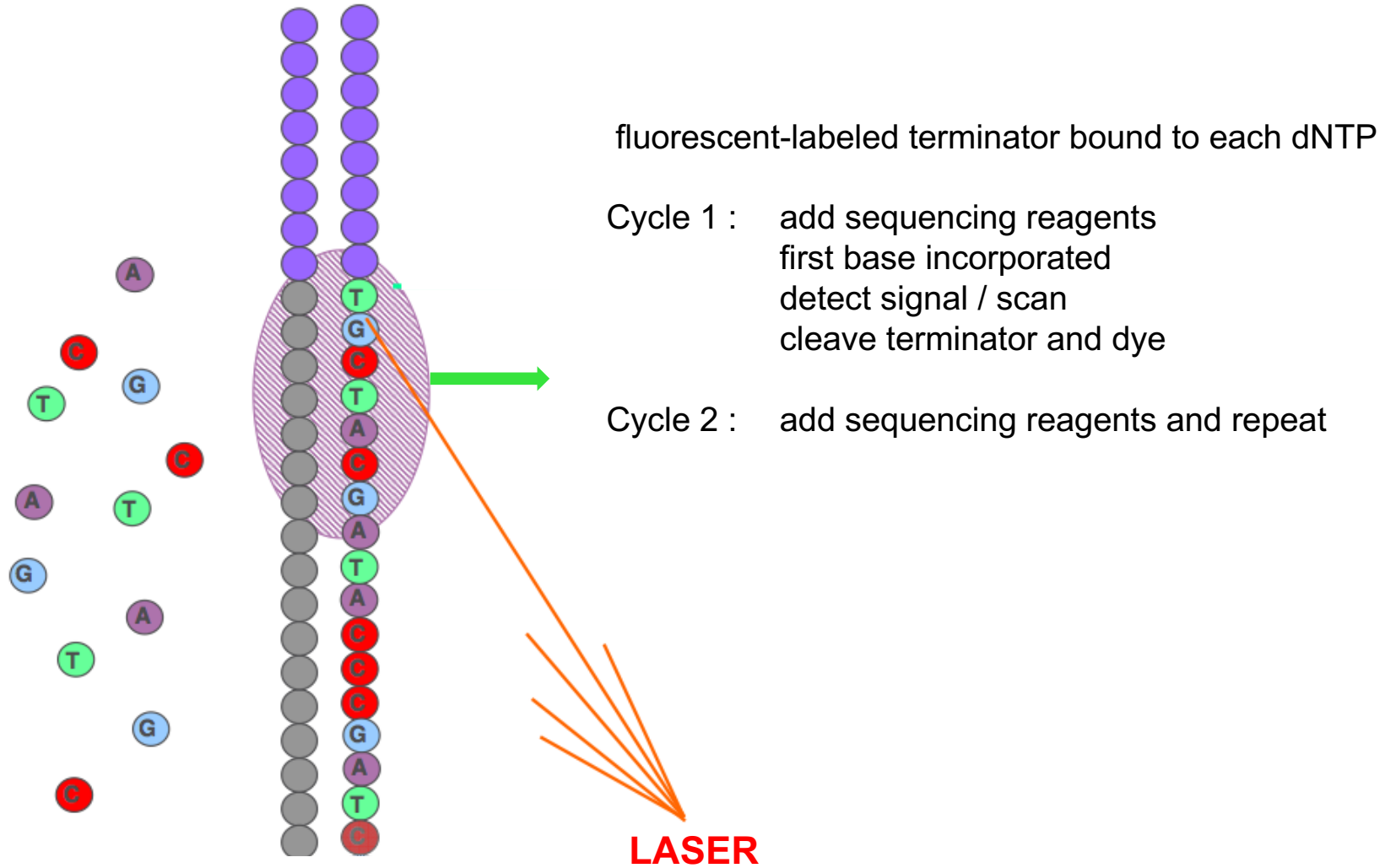


fluorescent-labeled terminator bound to each dNTP

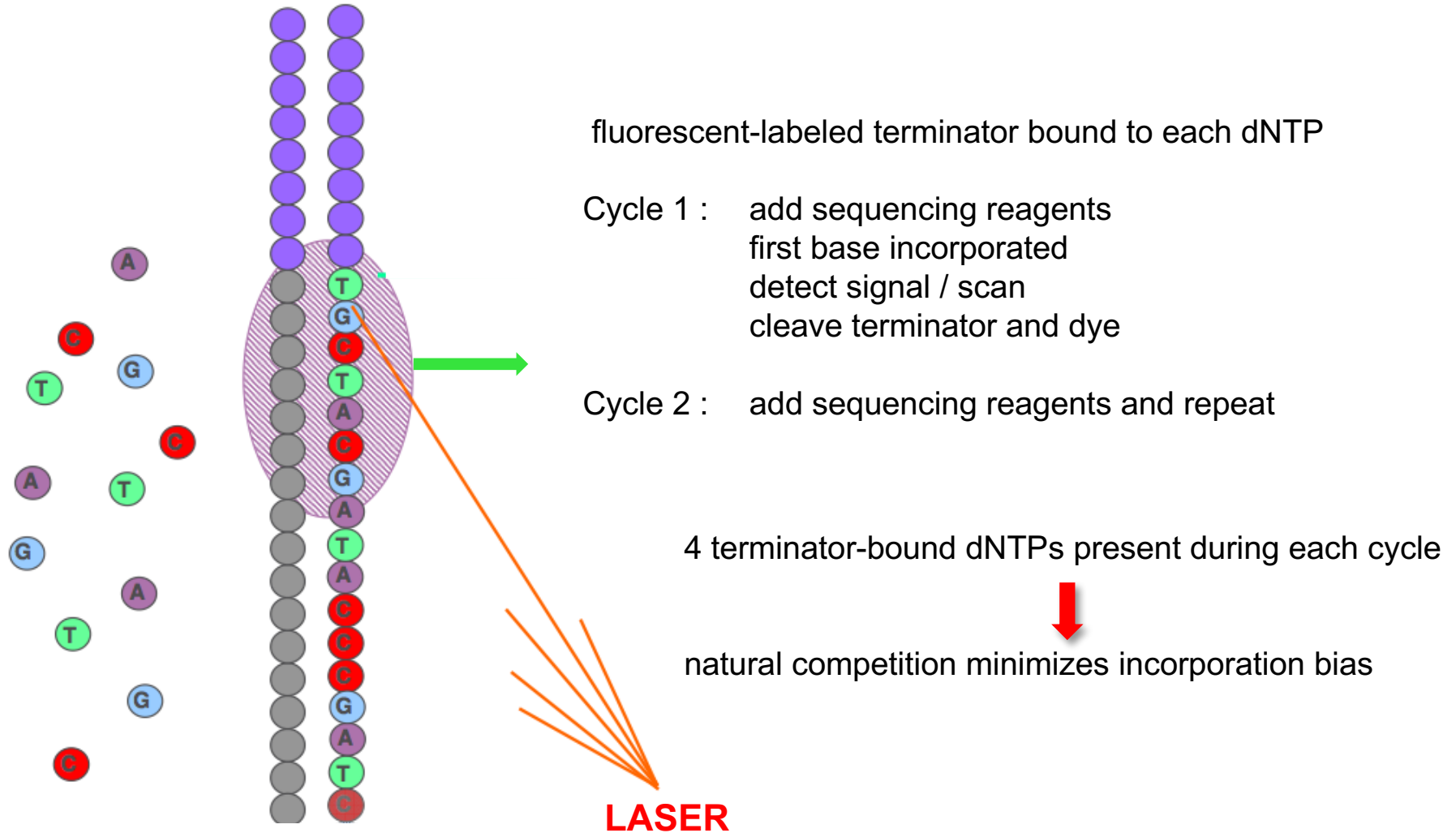
Cycle 1 : add sequencing reagents  
first base incorporated  
detect signal / scan  
cleave terminator and dye

Cycle 2 : add sequencing reagents and repeat

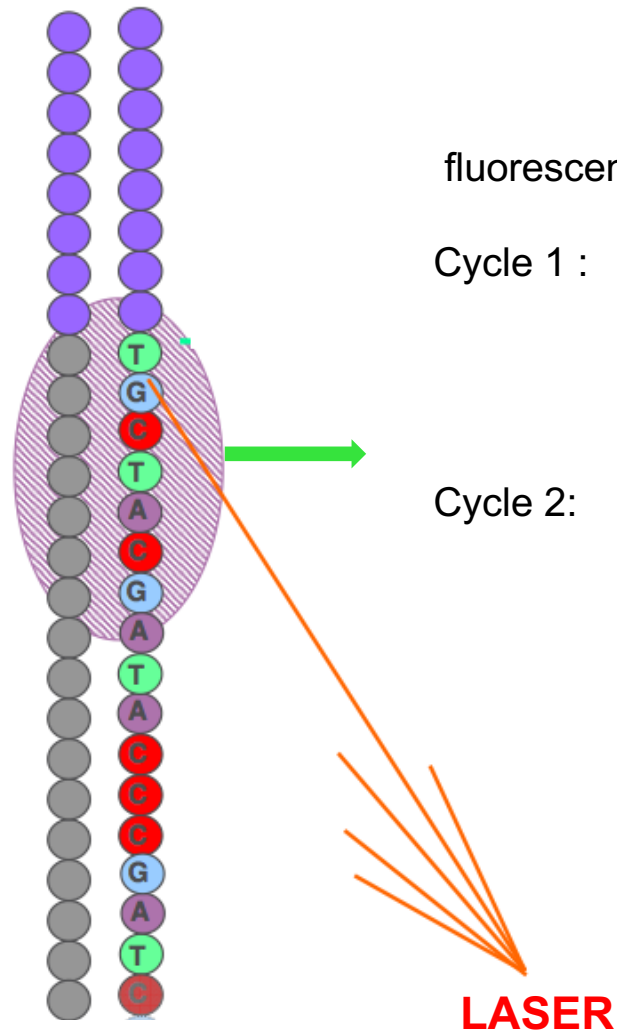
### 3 - Sequencing By Synthesis (SBS)



### 3 - Sequencing By Synthesis (SBS)



### 3 - Sequencing By Synthesis (SBS)



fluorescent-labeled terminator bound to each dNTP

Cycle 1 : add sequencing reagents  
first base incorporated  
detect signal / scan  
cleave terminator and dye

Cycle 2: add sequencing reagents and repeat

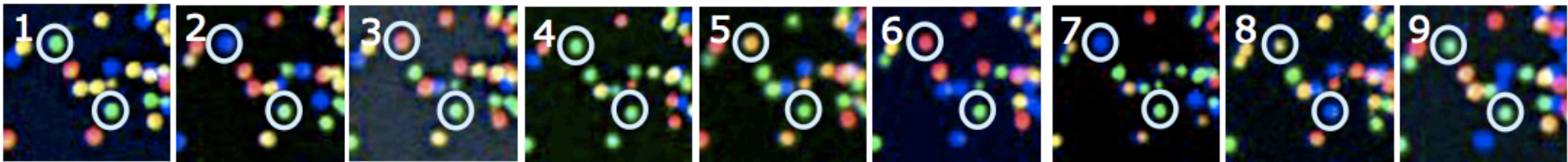
$$(C_{eff})^{FL} = 0.5$$

Cycle Efficiency (%)	Read-length (bases)
90.0	7
95.0	14
98.0	35
99.0	69
99.5	149
99.9	693



## Base calling from raw data

T G C T A C G A T ...

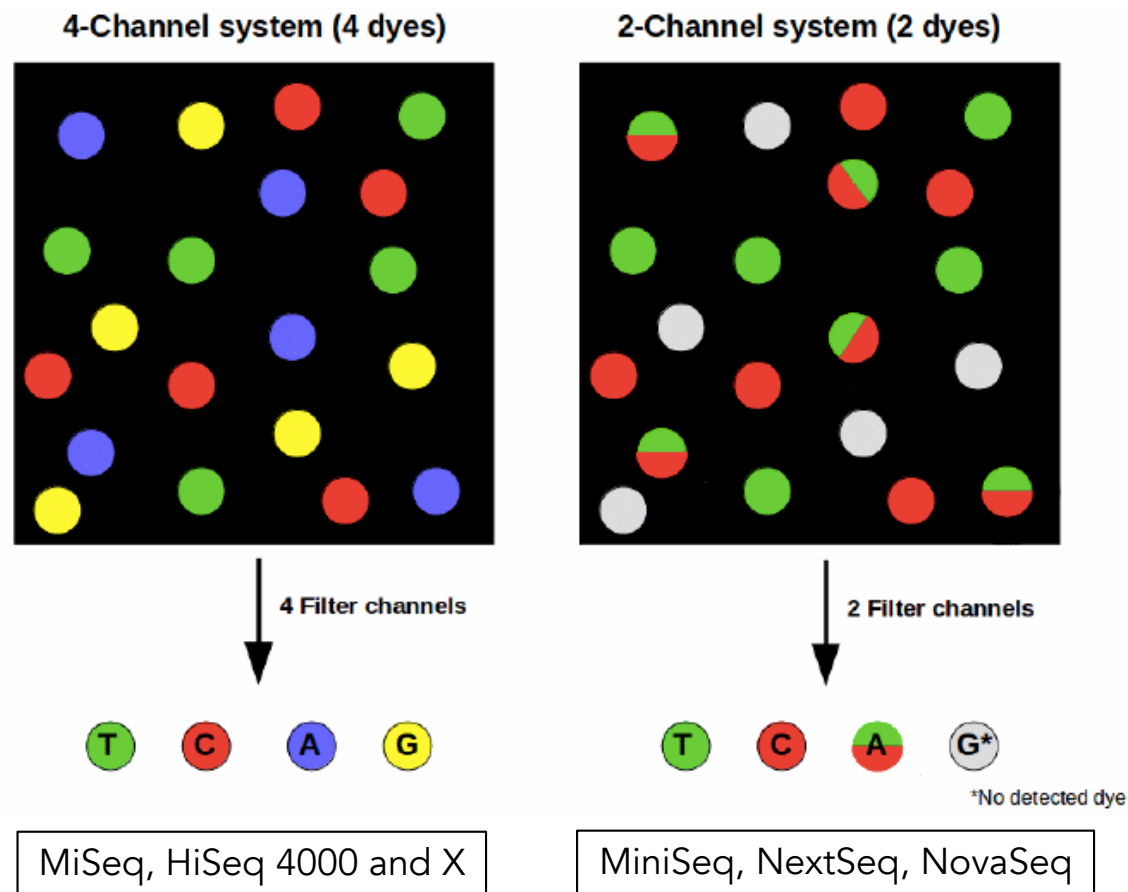


T T T T T T T G T ...

The identity of each base of a cluster is read off from sequential images.

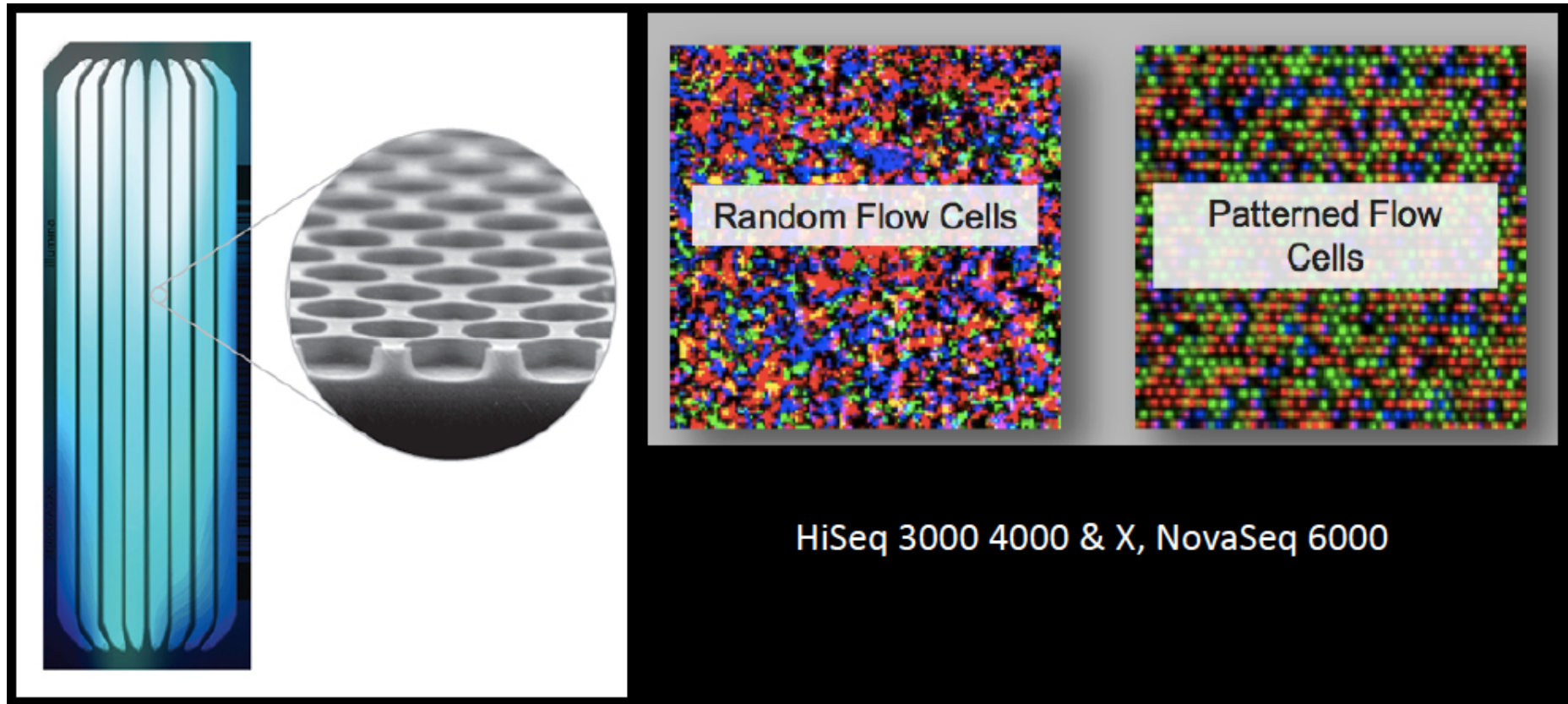
# 4 channel vs 2 channel detection

2 channel detection can cause sequencing errors due to phasing issues and can cause polyG tracts at the end of fragments



# Patterned flow cells

- Improves regularity of densities and qualities
- Reduces analysis time

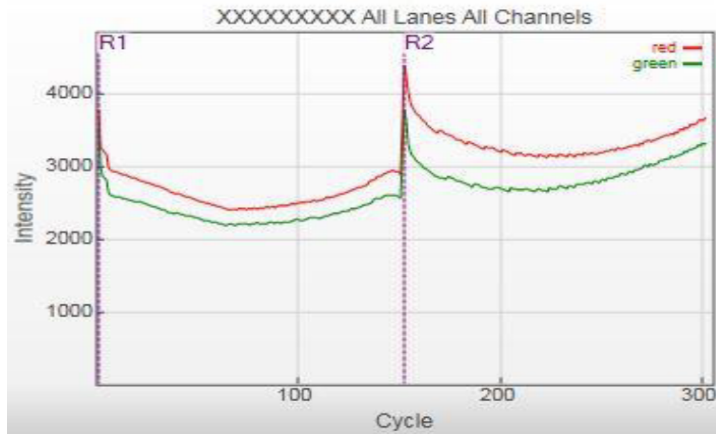


HiSeq 3000 4000 & X, NovaSeq 6000

# Sequencing qualities

## Intensities

Intensity values of the four incorporated bases are measured to determine if a clustering failure has occurred

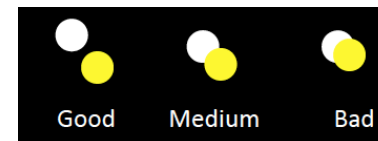
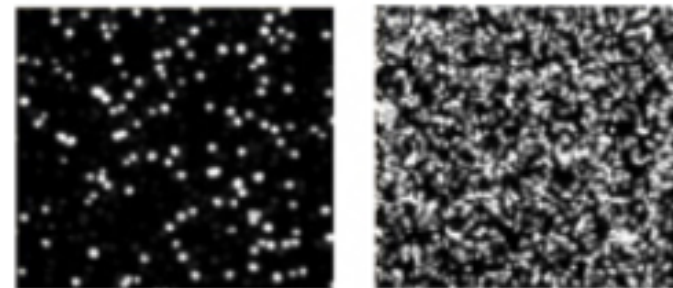


## Percentage of reads passing filter (% PF)

Reflects the purity of the signal from each cluster. Ratio of the brightest intensity divided by the sum of the brightest and second brightest intensities for each cycle. Optimal values for MiSeq and HiSeq 2500 from 80-95%.

## Cluster densities

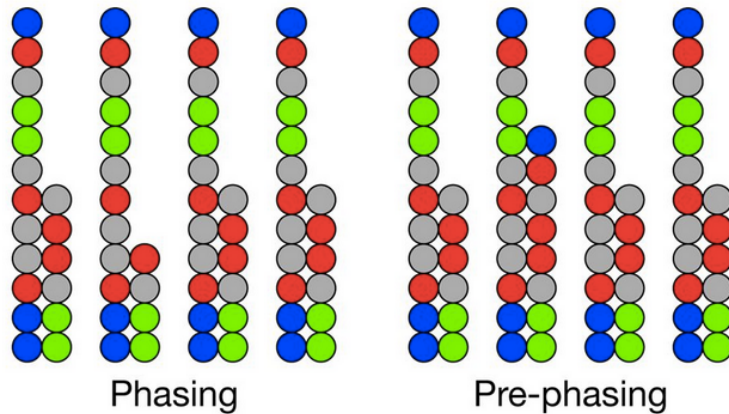
Density of clusters for each tile in thousands/mm<sup>2</sup>. Overloading of the library may cause merging of clusters, lowering of the % PF and reduction of the quality score (Q30)



# Sequencing qualities

## Phasing/Prephasing

Percentages of bases that fell behind or jumped ahead the current cycle within a read. It increases with cycle number, hampering correct base identification for long reads. Optimal values are below 0.5 or 0.2% depending on platform. Issues may arise when read length is above specifications.



## Q-score

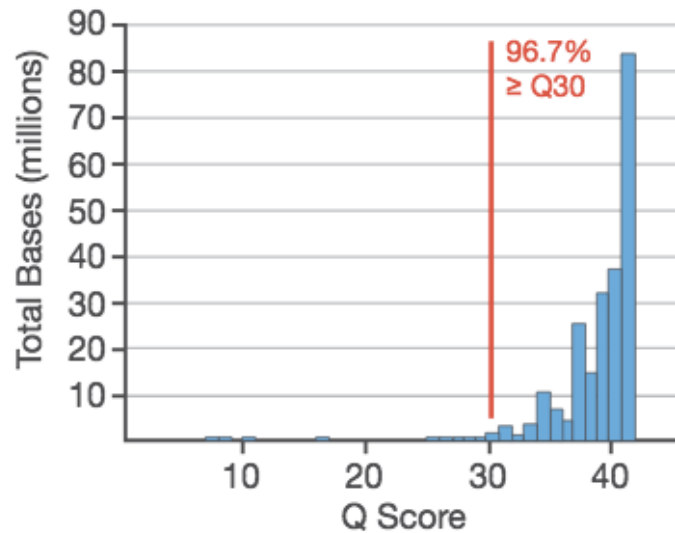
Base calling error probabilities  $P : Q = -10 \log_{10} P$   
% Q-score > 30 : percentage of bases that have a probability of incorrect base calling of 1/1000. Low Q scores can increase false-positive variant calls, which can result in inaccurate conclusions.



# Sequencing qualities

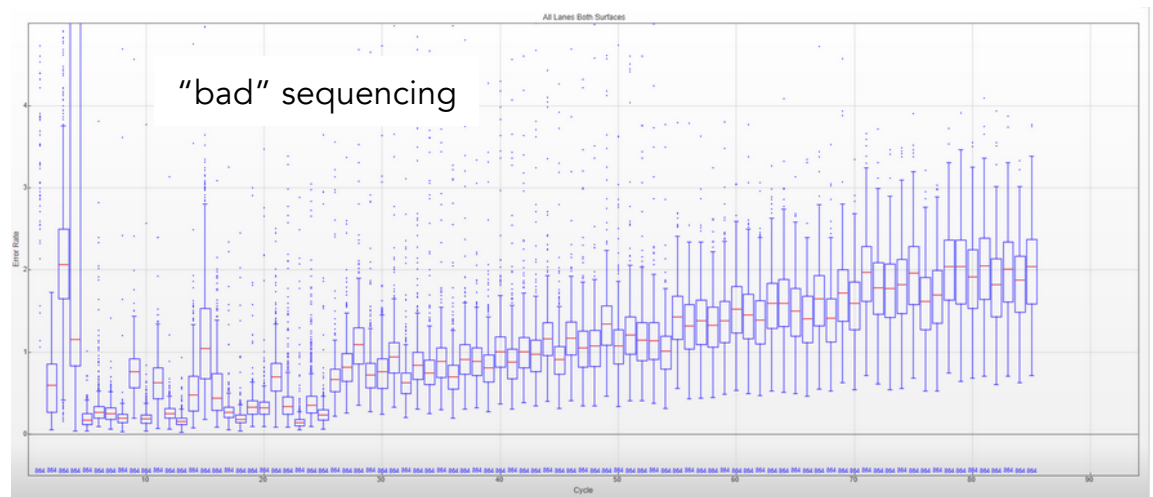
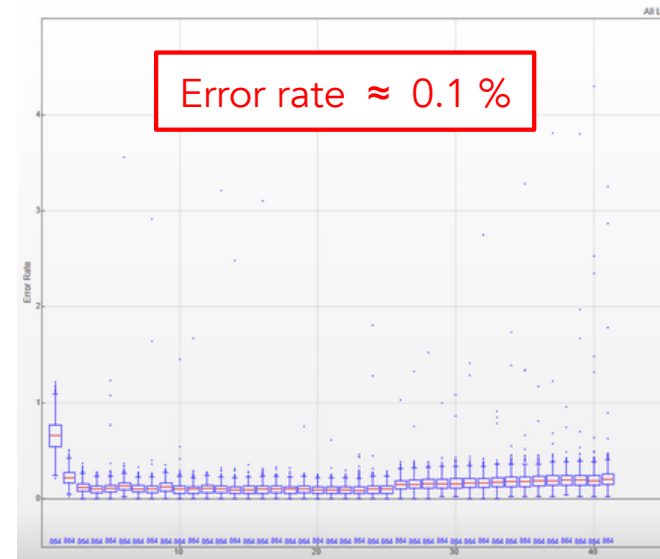
## Percentage of alignment

**% align:** percentage of the reads aligned to PhiX control genome. It is important to estimate the error rate. Indication on the success of the clustering process.



## Error rate

Calculated for each read based on the aligned % of PhiX control sample. Important to determine if the sequencing has proceeded as expected.



# SINGLE CELL TECHNOLOGIES

Single cell transcriptomics allows to study transcriptome heterogeneity, to investigate differences in transcript expression and gene regulation *in individual cells* :

- ❖ Differences in transcript abundance
- ❖ Alternative splicing and differential expression of isoforms

Most widely used device to study single-cell transcriptomics : *Chromium controller (10x Genomics)*

Several applications :

Single cell Gene expression

Measures gene activity on a cell-by-cell basis, characterize cell populations, cell types, ...

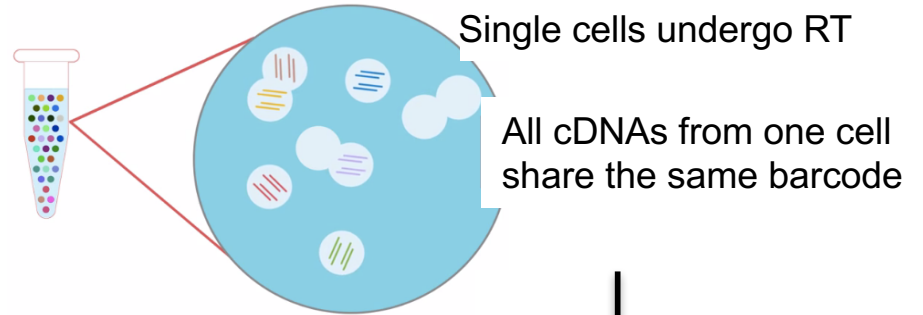
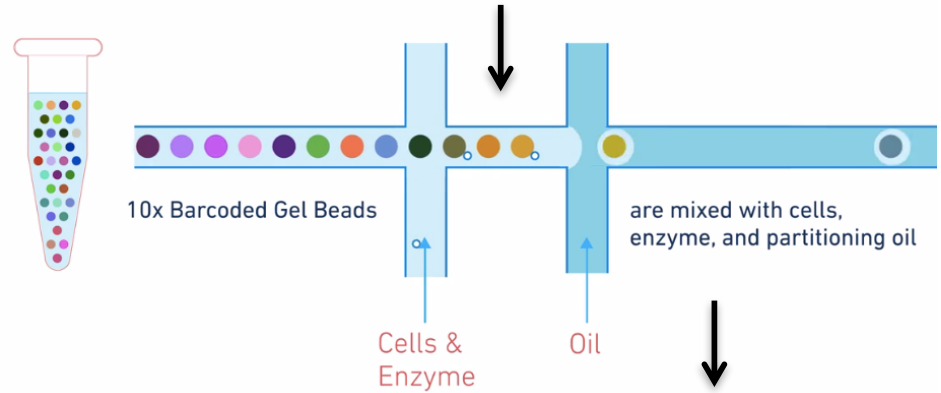
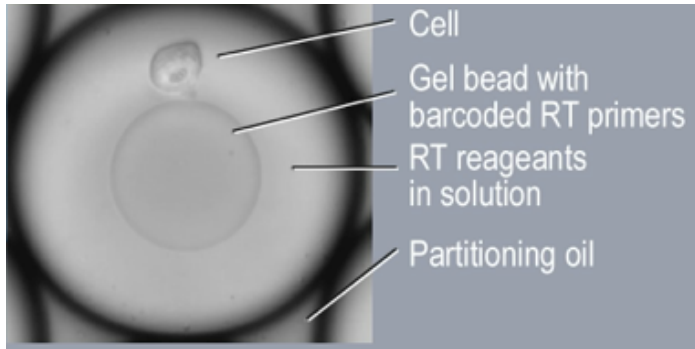
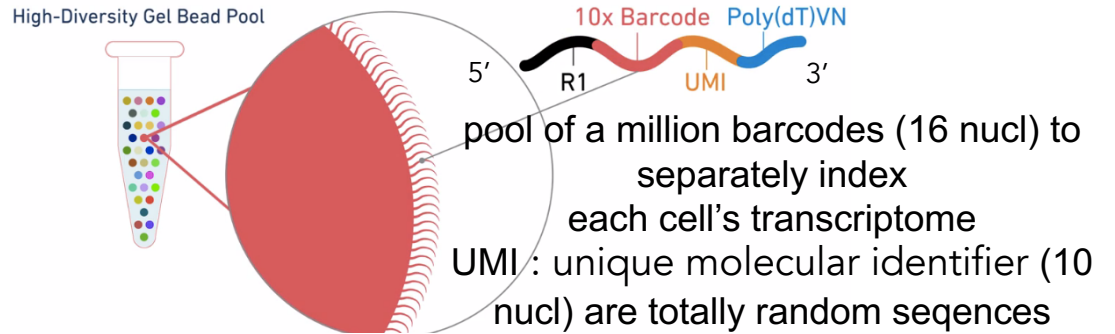
Linked read genomics

Performs diploid de novo assembly, phase haplotypes, genetic variations

Single cell ATAC

Measures epigenetics by detecting open chromatin regions

# CHROMIUM SINGLE-CELL RNA SEQUENCING

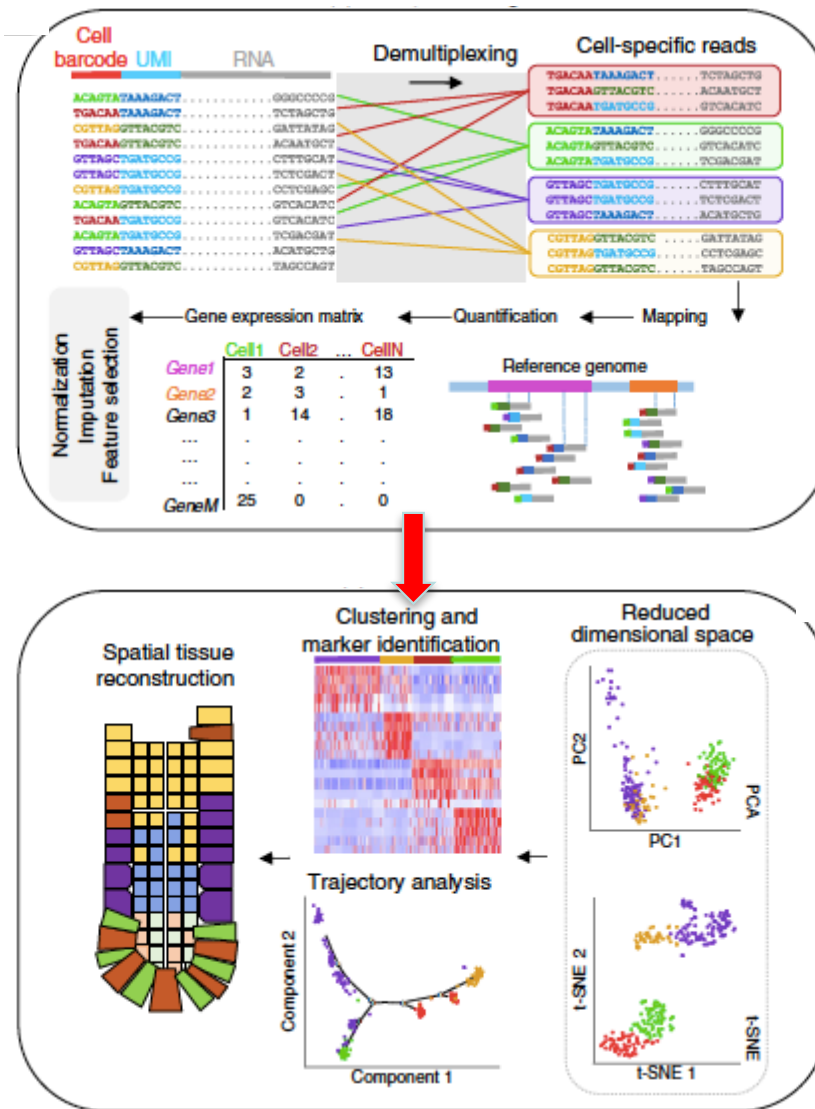


barcoded sequencing library





DATA ANALYSES



CELL RANGER :  
Chromium data analysis pipelines

# Recent advances : Illumina sequencing with Smart-seq3

Single-cell RNA counting at allele- and isoform-resolution using Smart-seq3  
Hagemann-Jensenn et al. *Nature Biotechnology*, 2020

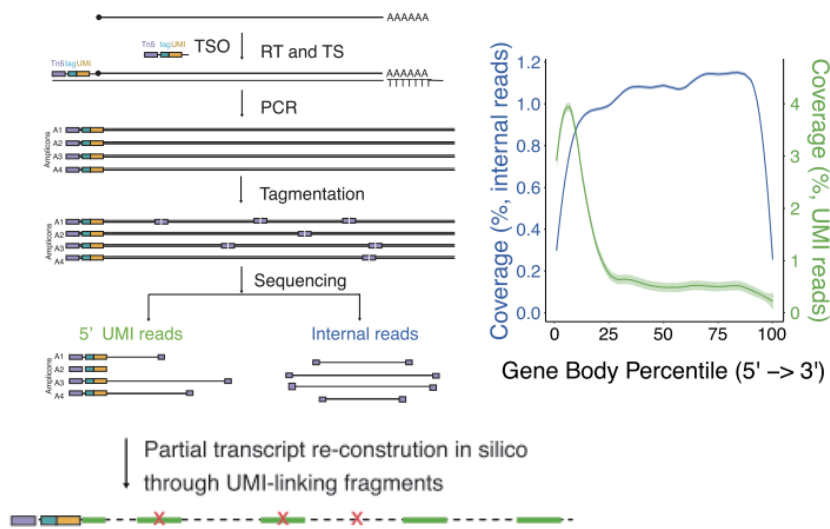
- Single cell methods capture only a small fraction of cellular mRNA (~5–10%)
- They frequently lack splice isoform information

Smart-seq3 :

- full-length transcriptome coverage
- 5' UMI RNA -> in silico reconstruction of thousands of RNA molecules/cell
- 60% assignments to allelic origin
- 30–50% assignments to specific isoforms

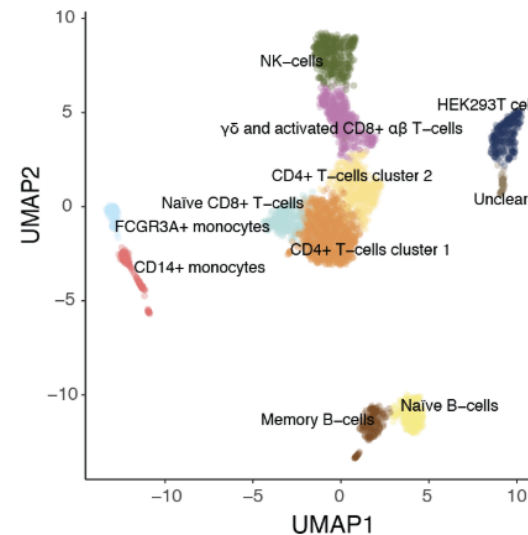
Smart-seq3 greatly increased sensitivity compared to Smart-seq2 detecting thousands more transcripts/cell

TSO : Template-Switching Oligo  
UMI : unique molecular identifier

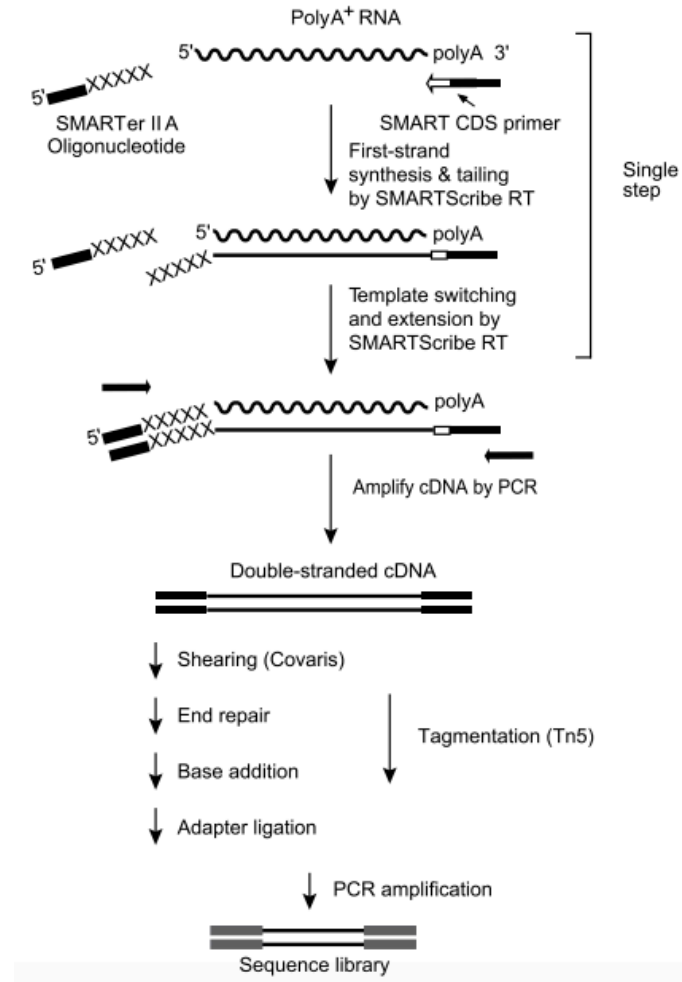


Count RNA molecules at allelic and transcript isoform resolution

Dimensionality reduction of 3,890 human cells colored by annotated cell type



# The Smart-Seq protocol



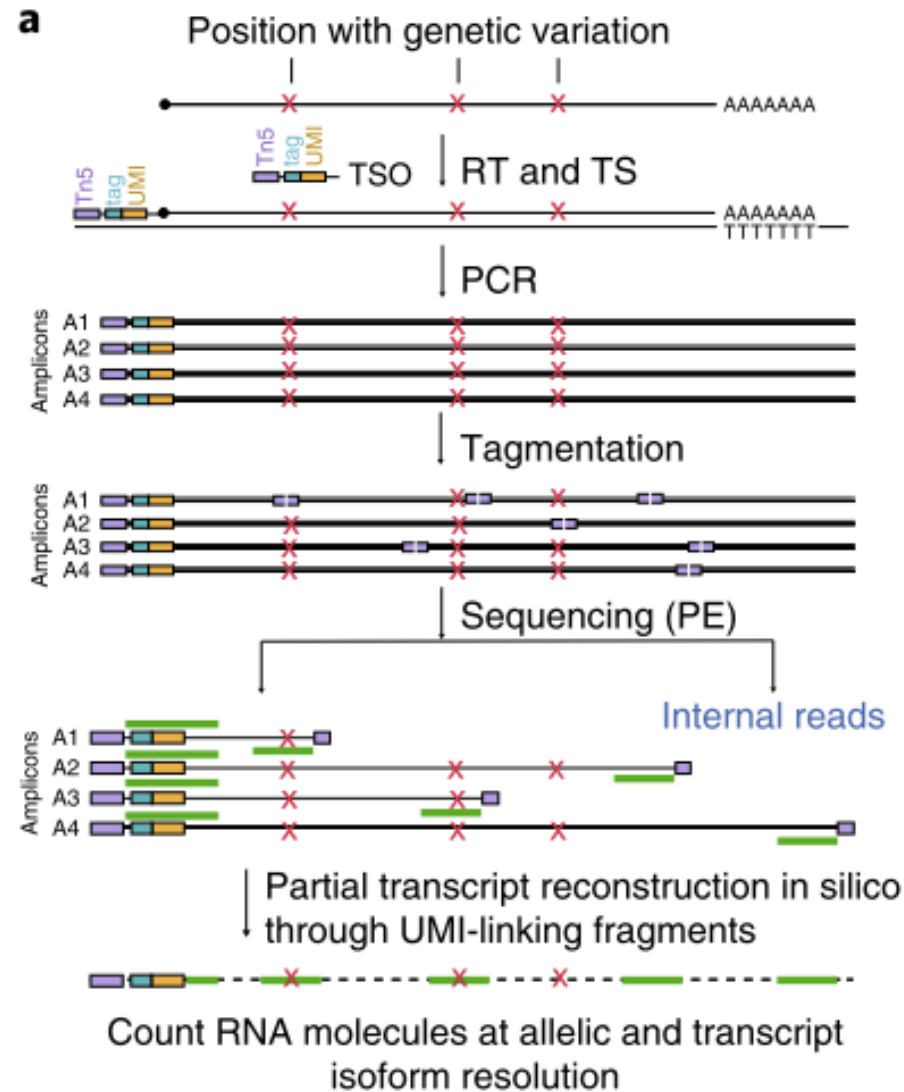
# Overview of single-cell RNA sequencing in Smart-seq3



Template-Switching Oligo

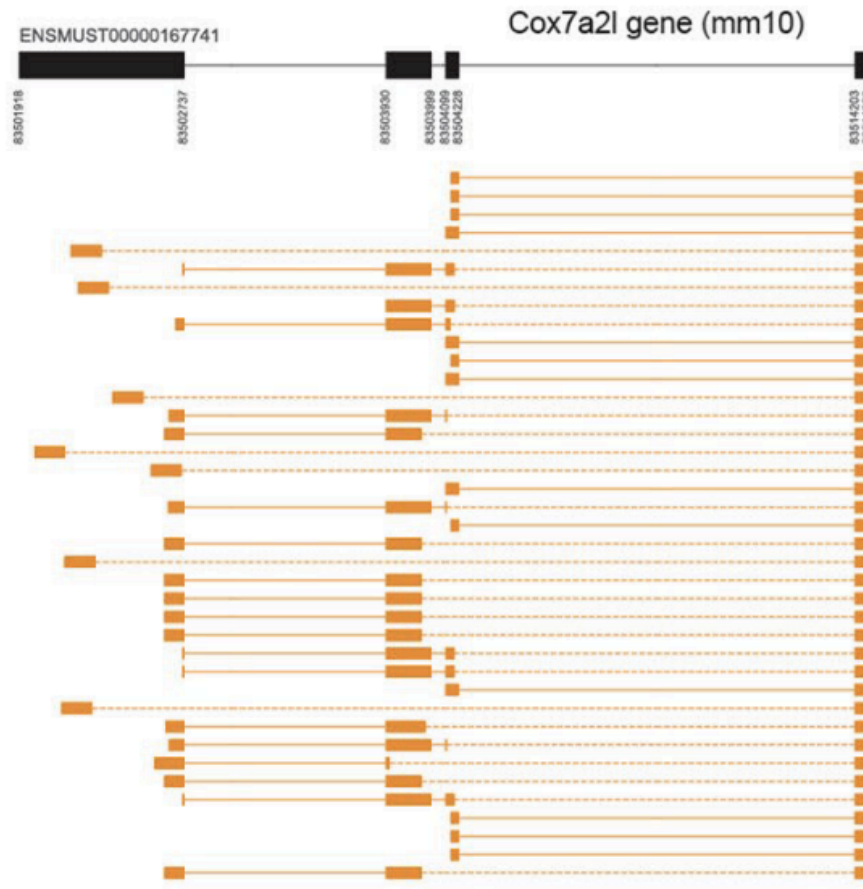
□ TSO consists of :

- a partial Tn5 motif
- a 11-bp tag sequence
- a 8-bp UMI sequence and 3 riboguanosines



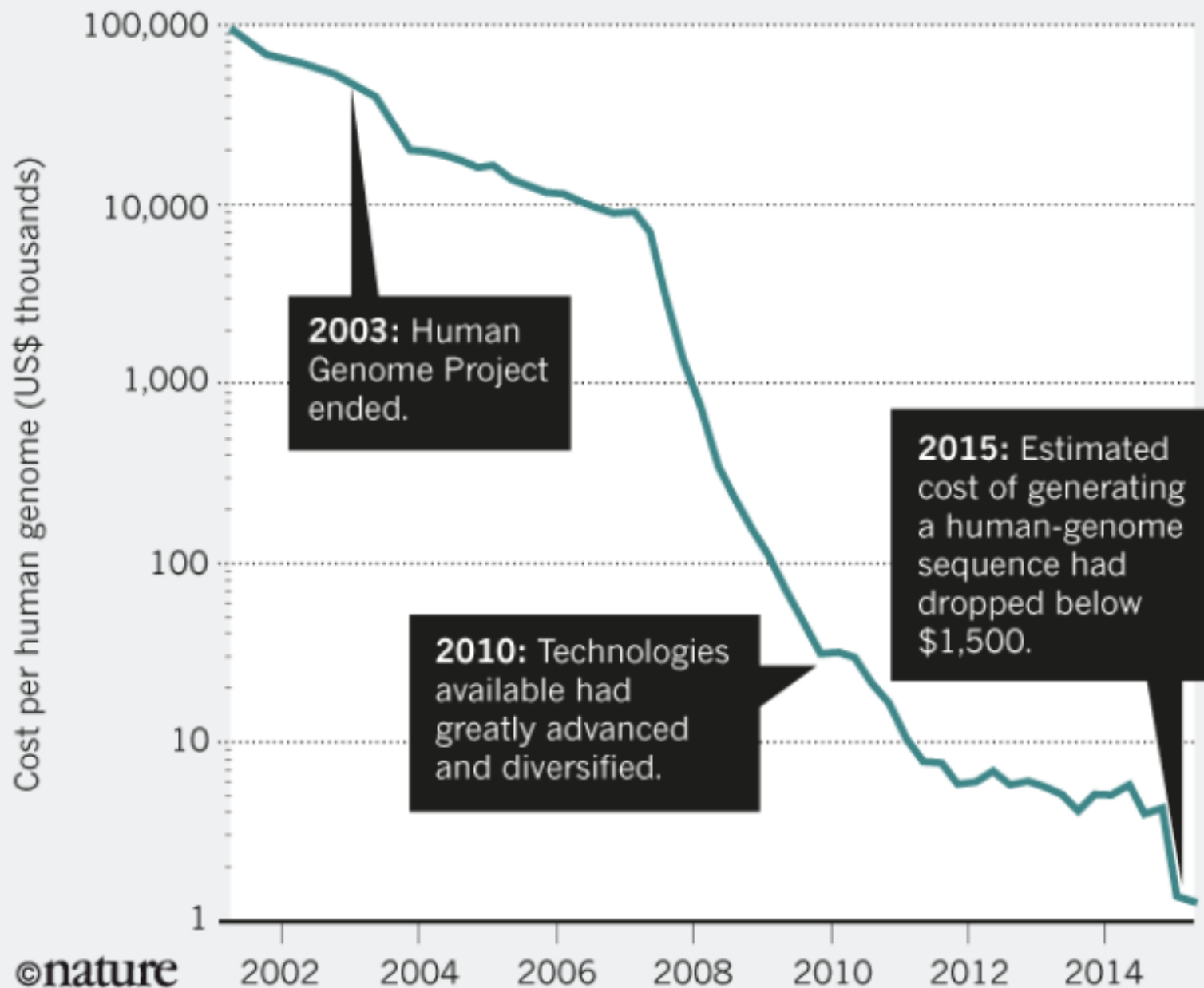
# Example of read pairs from a single transcribed molecule

Each row shows a unique read pair



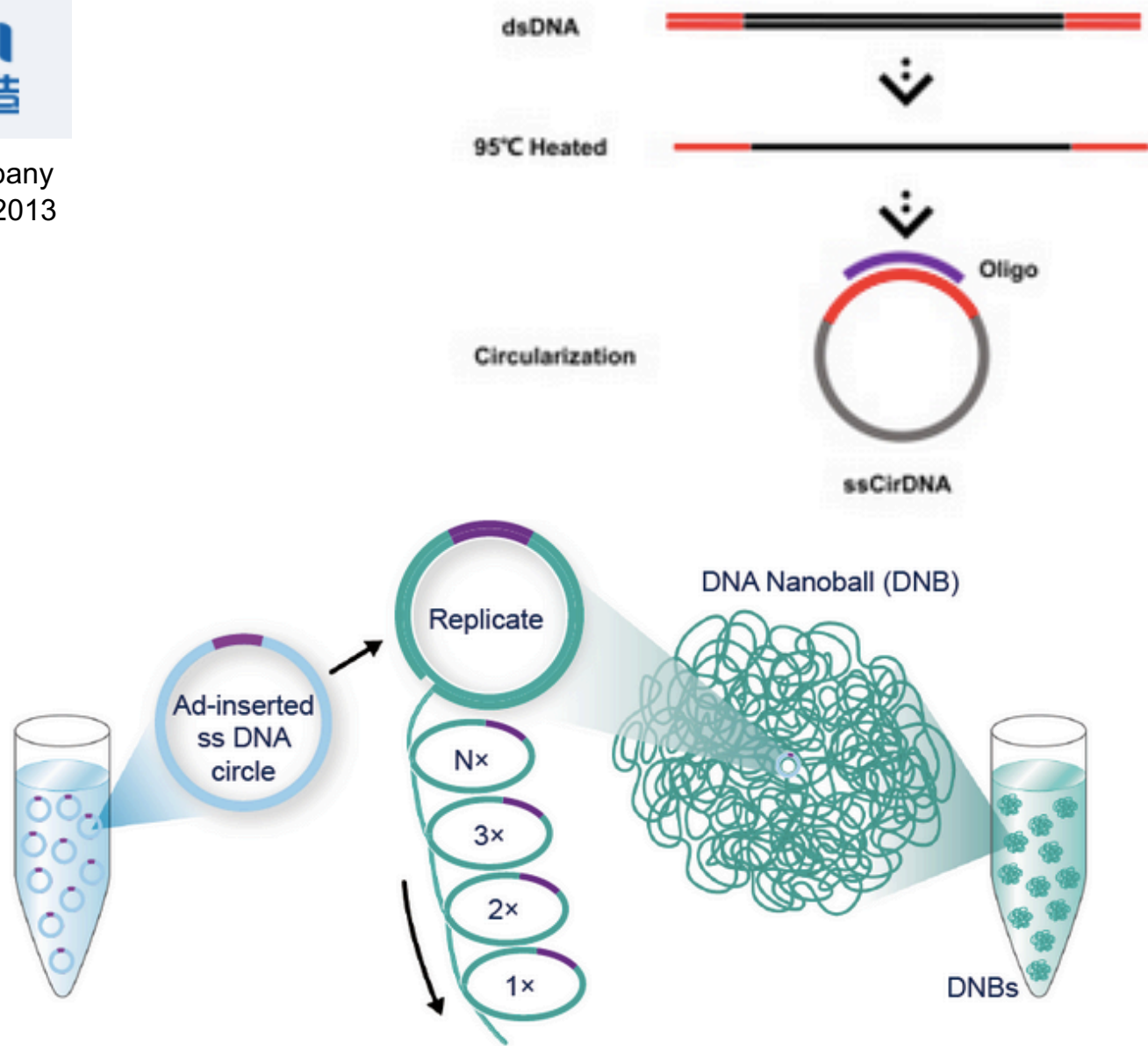
# BETTER, CHEAPER, FASTER

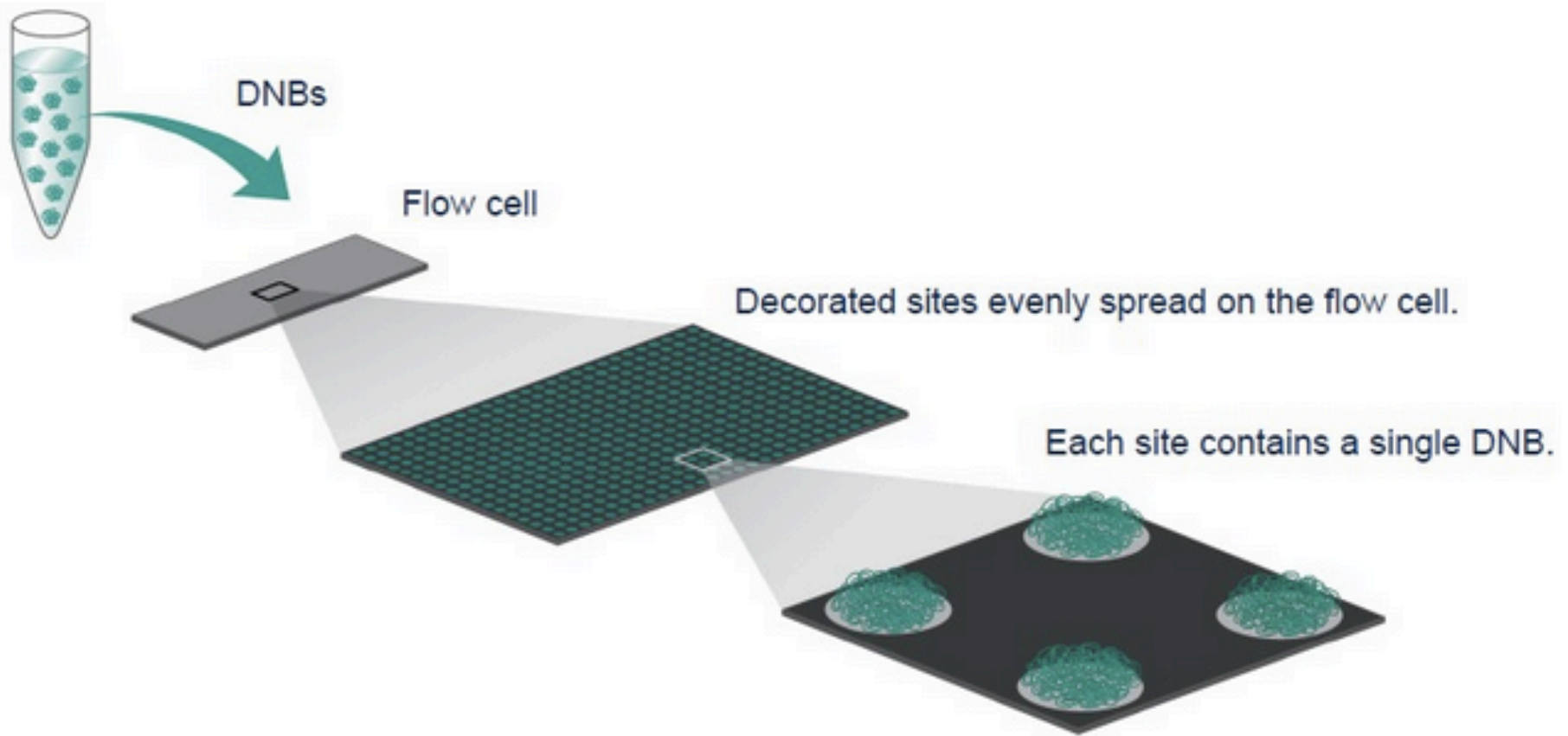
The cost of DNA sequencing has dropped dramatically over the past decade, enabling many more applications.





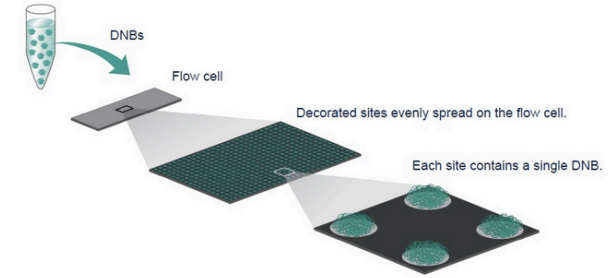
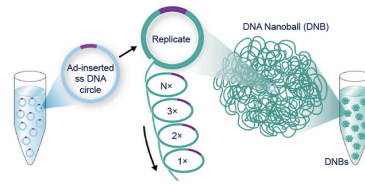
acquisition of U.S. company Complete Genomics in 2013







# NEW SEQUENCING TECHNOLOGY



## Tests préliminaires (Genoscope, CNRGH)

- Qualité : Q30 moyen élevé (même en 2x200)
- Taux d'erreur MGI < taux d'erreur NovaSeq
- % reads mappés MGI > % reads mappés Illumina

Mais :

- Runs plus longs (+ lavage 6h)
- Preparation of nanoballs « délicate »
- Régions riches A/T et G/C moins couvertes

	G400	NovaSeq S1	NovaSeq S4
Run time	66h	24h	48h
Output	500 Gb	500 Gb	1600-2000 Gb
Cost/Gb (€)	3.8 - 4.5	9.9 - 12.4	6.6 - 8.3

## PART 2

### 3<sup>rd</sup> GENERATION SEQUENCING

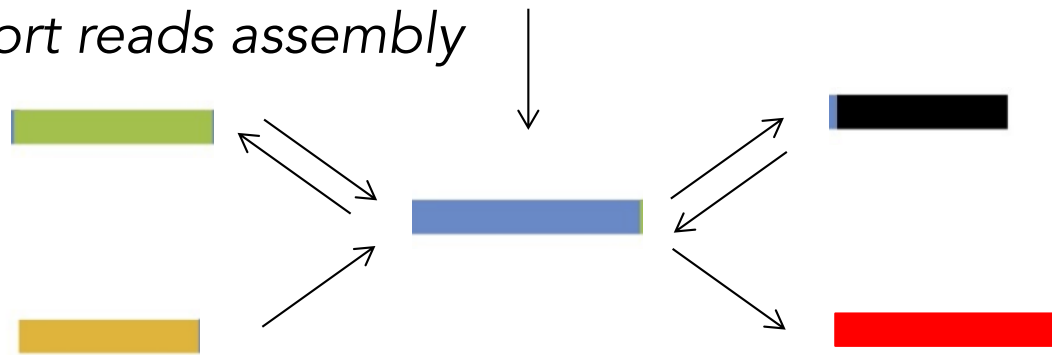
#### LONG READS

# LONG-READS VERSUS SHORT-READS

Assembly of DNA fragments with repeated sequences



*NGS short reads assembly*



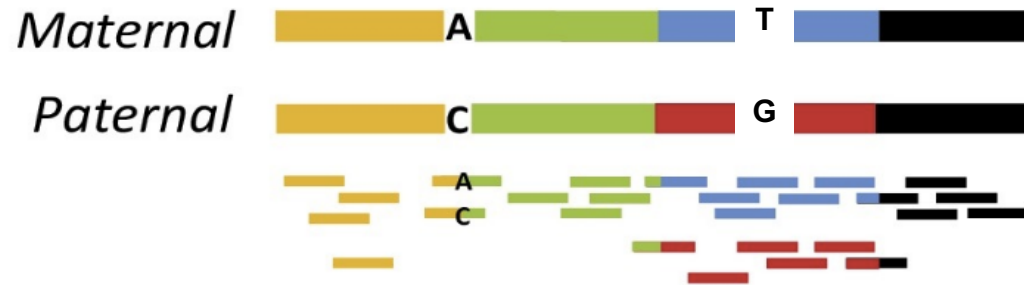
Several contigs → incomplete assembly, underestimation of repeats

*Long reads assembly*

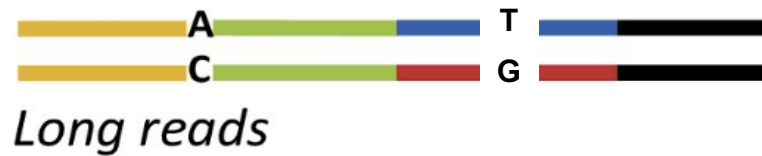
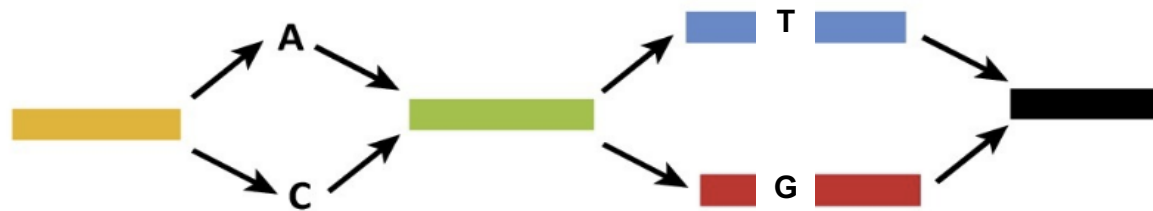


# LONG-READS VERSUS SHORT-READS

## Haplotype phasing

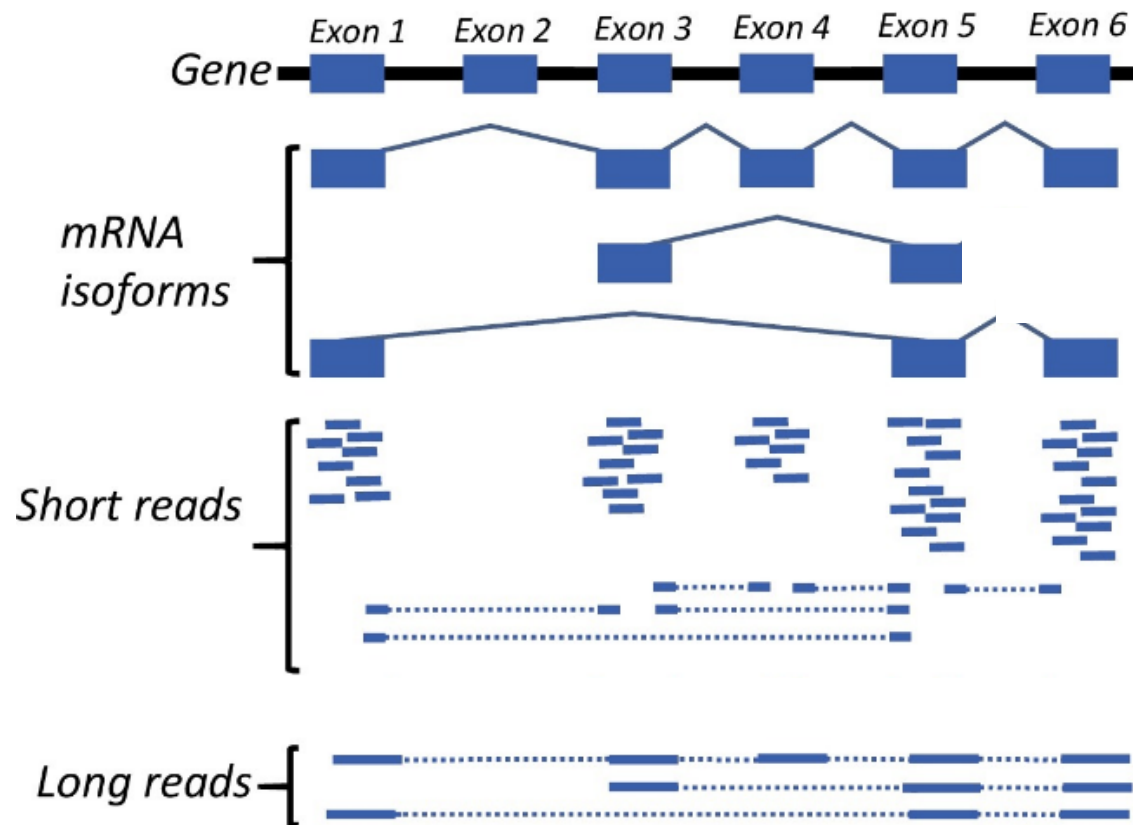


*NGS short reads  
assembly*



# LONG-READS VERSUS SHORT-READS

## Detection of splicing isoforms



# The 3rd generation winning technologies



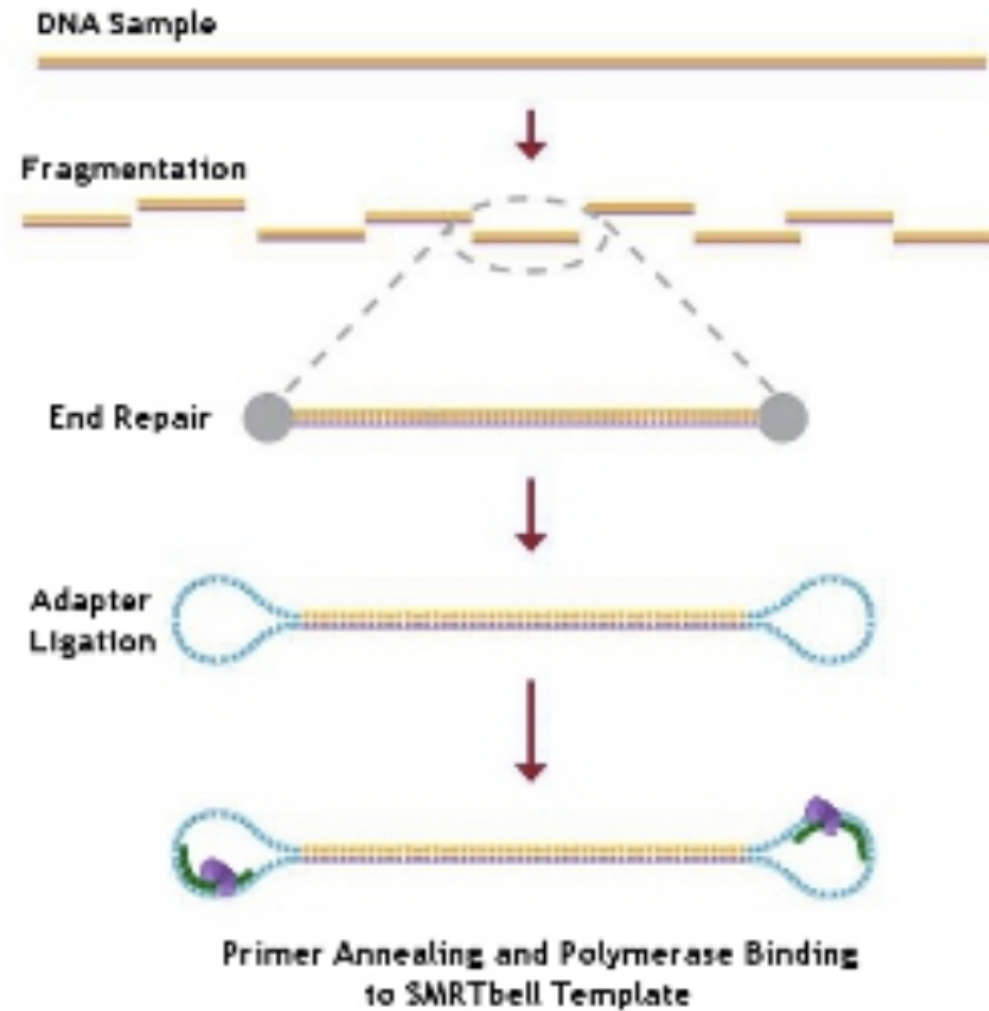
Sequel - Pacific Biosciences  
Single molecules  
Up to 80,000 bp long  
Error rate  $\approx$  10-15 % - CCS: <1%  
Compensated by coverage



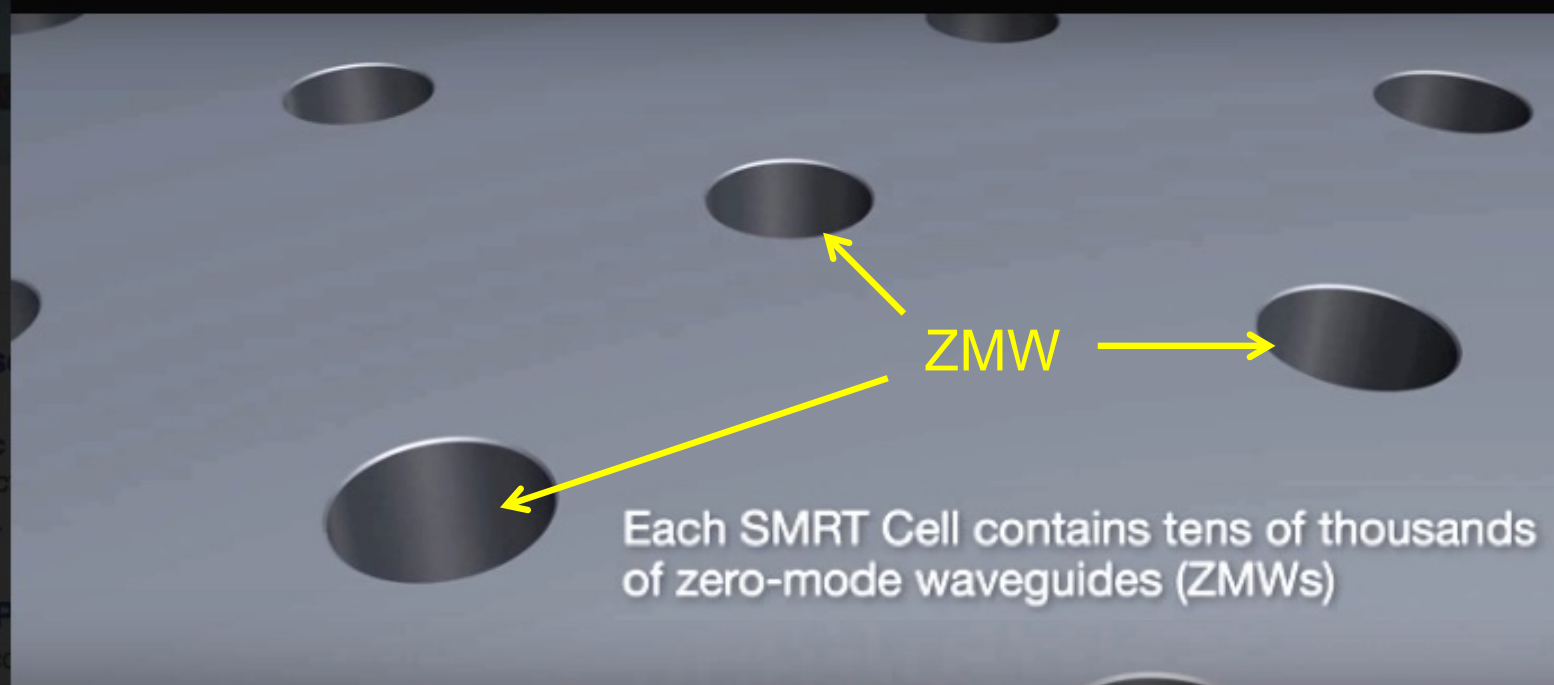
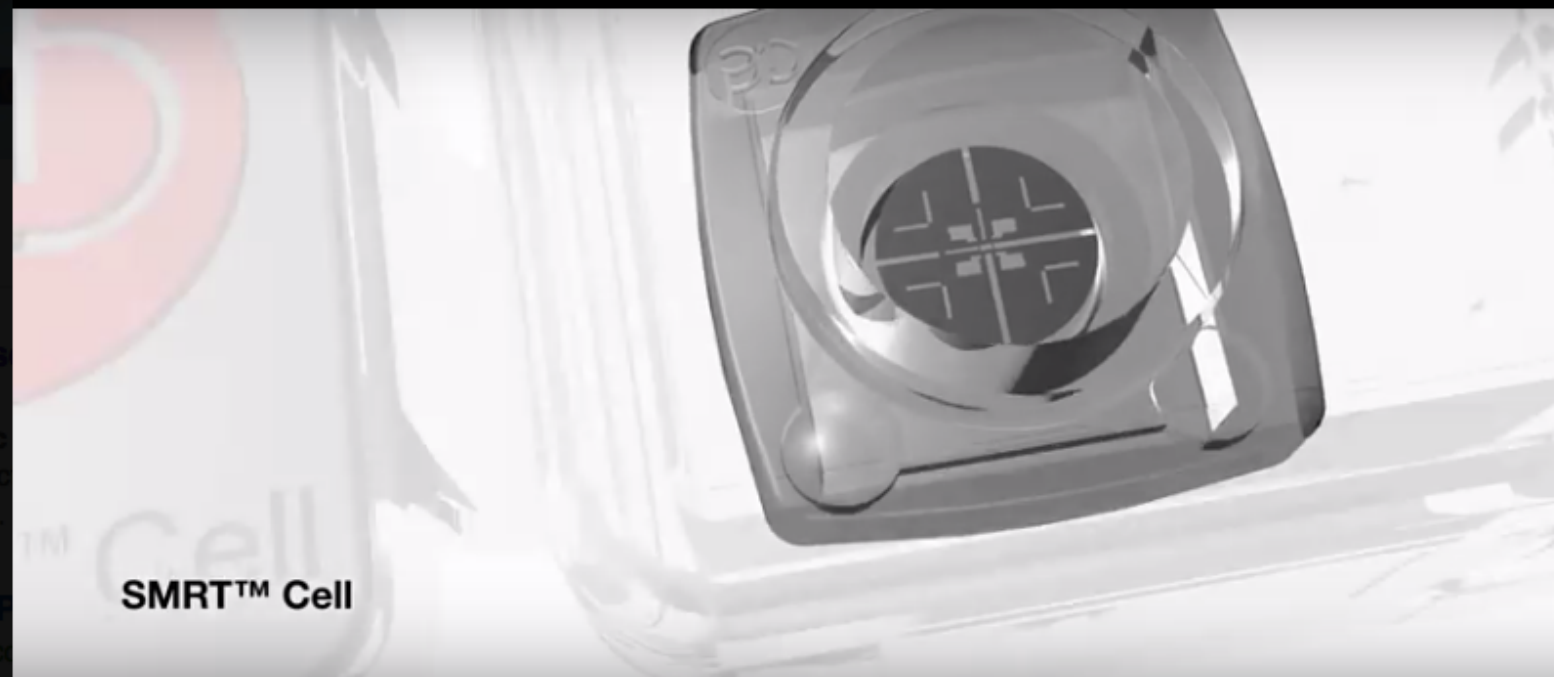
MinION - Oxford Nanopore  
Single molecules  
> 200 000 bp long  
Error rate  $\approx$  10-15 %  
Compensated by coverage

# PacBio : Single Molecule Real Time (SMRT) sequencing

## PacBio DNA-seq library




# PACIFIC BIOSCIENCES





# PACIFIC BIOSCIENCES

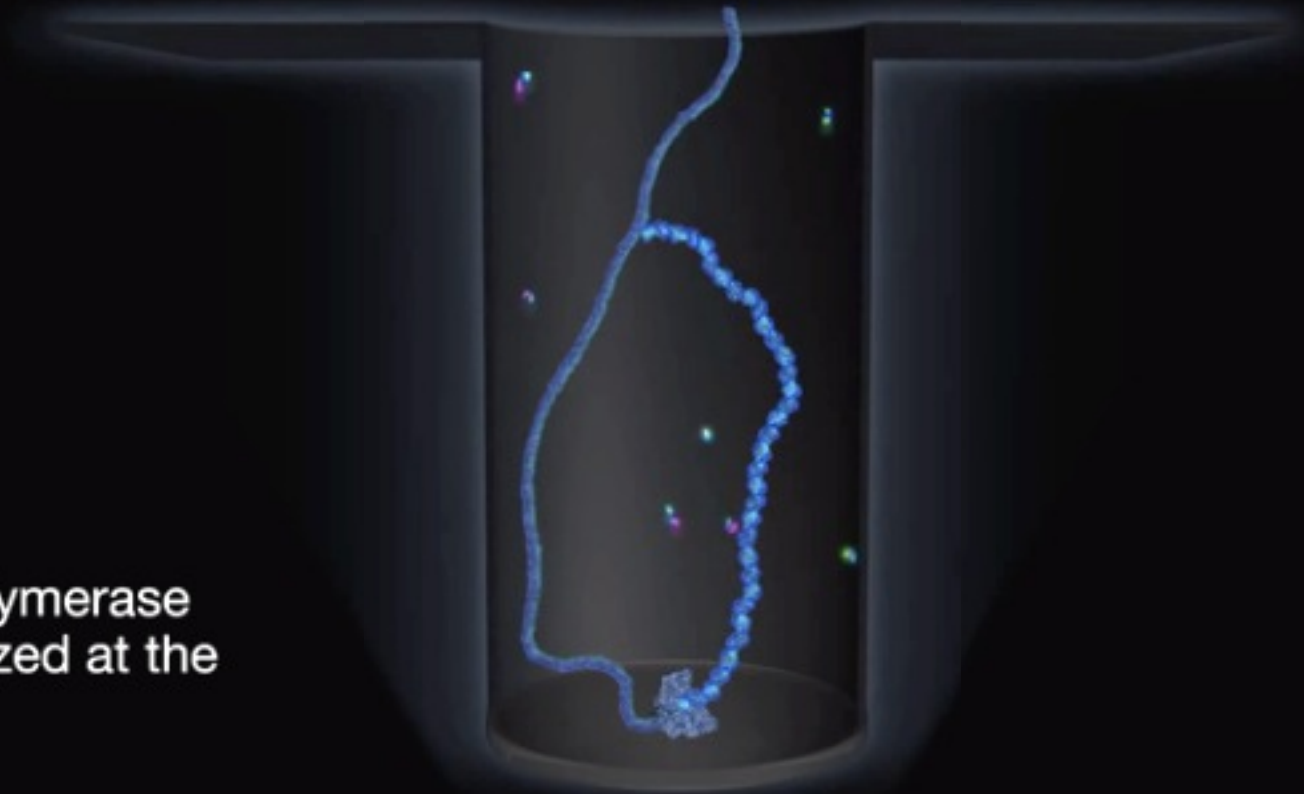


ZMW : optical waveguide that guides light energy into a volume that is small compared to the wavelength of the light

As each ZMW is illuminated from below, the wavelength of the light is too large to allow it to pass through the waveguide

# PACIFIC BIOSCIENCES

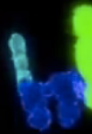
A DNA template-polymerase complex is immobilized at the bottom of the ZMW



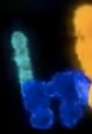
# PACIFIC BIOSCIENCES

Phospholinked Nucleotides

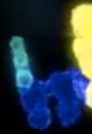
A



C



G



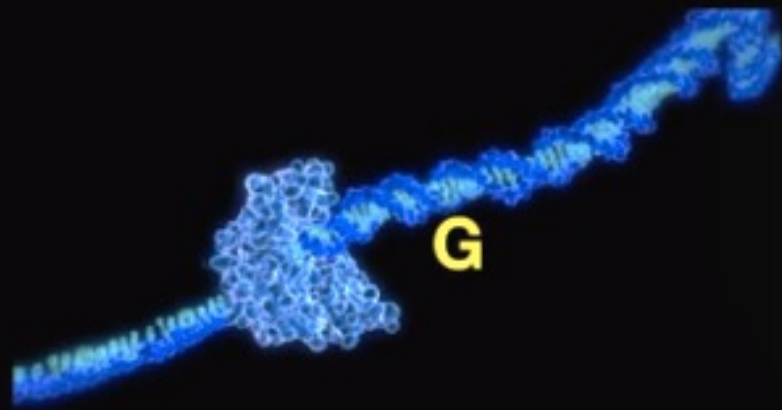
T



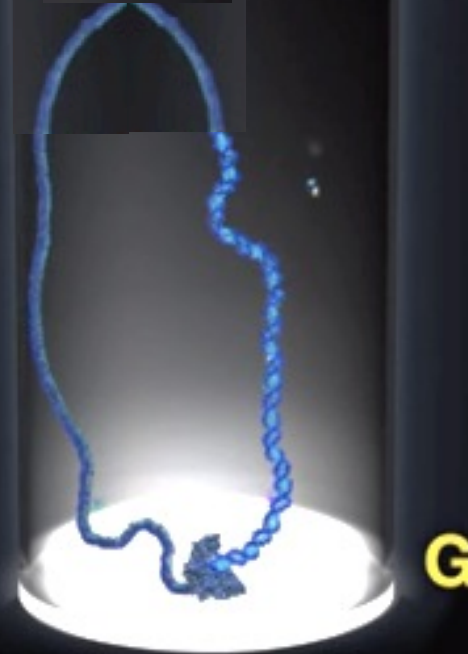
Phospholinked nucleotides are introduced into the ZMW chamber



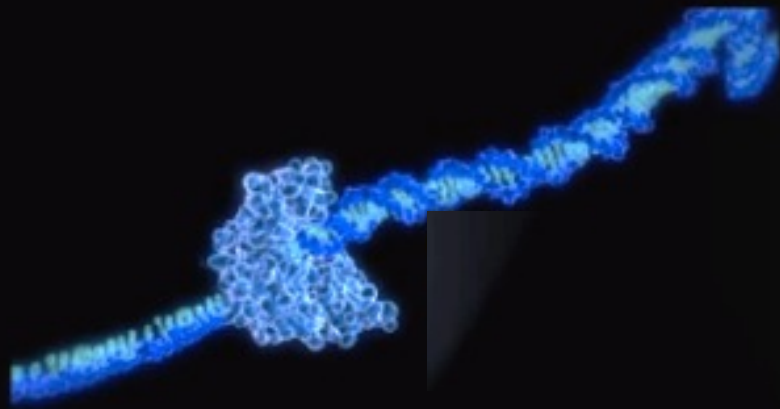
# PACIFIC BIOSCIENCES



As a base is held in the detection volume, a light pulse is produced



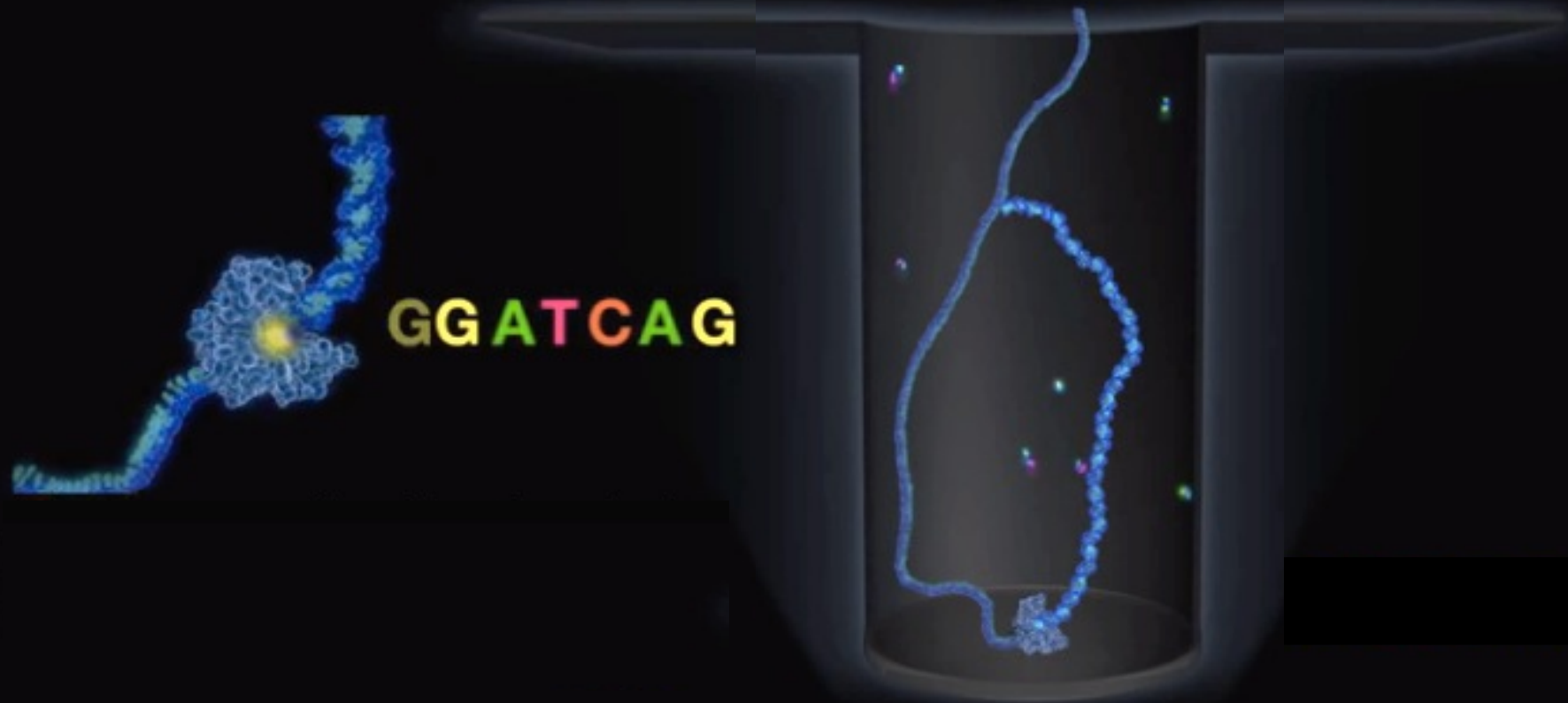
# PACIFIC BIOSCIENCES

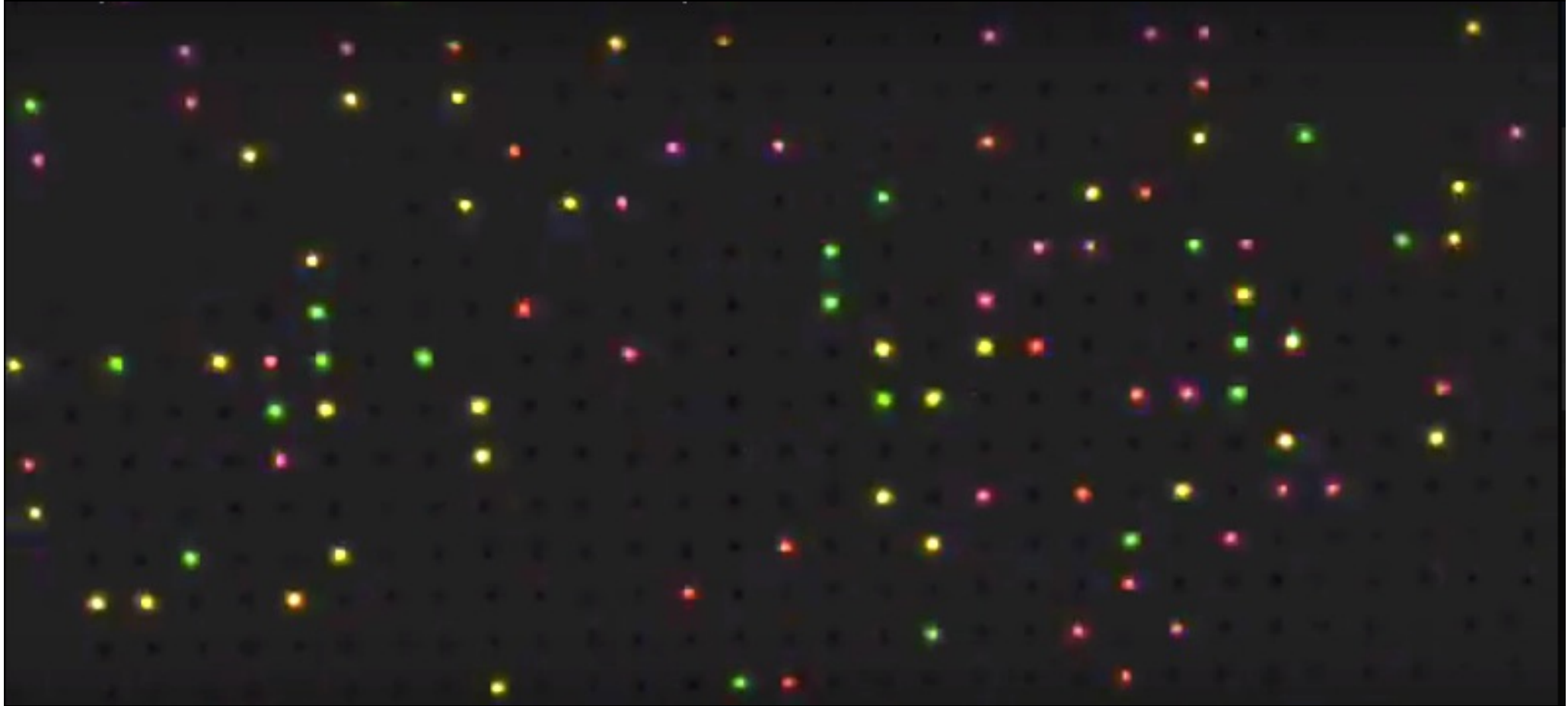


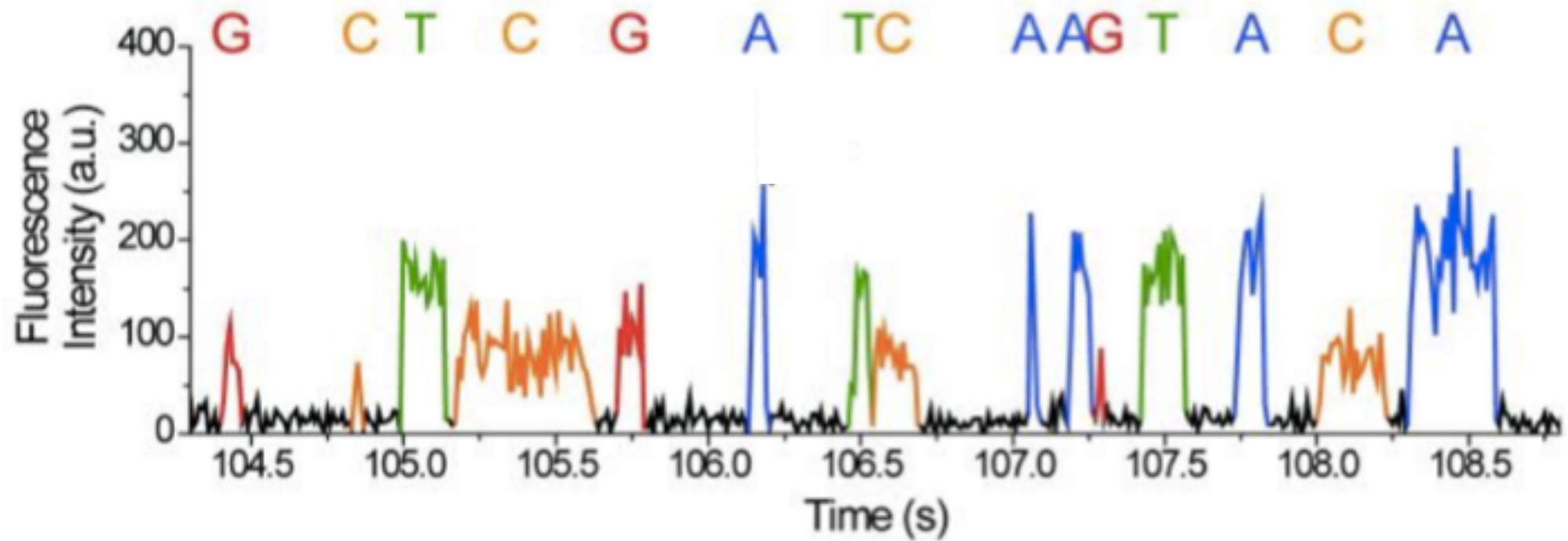
After incorporation the phosphate chain is cleaved, releasing the attached fluorophore



# PACIFIC BIOSCIENCES

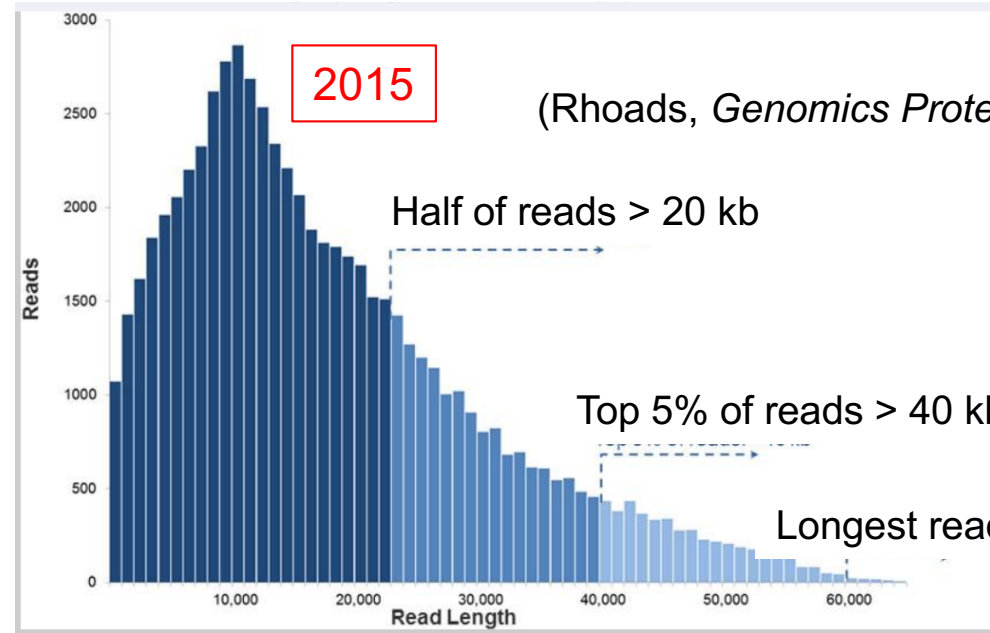




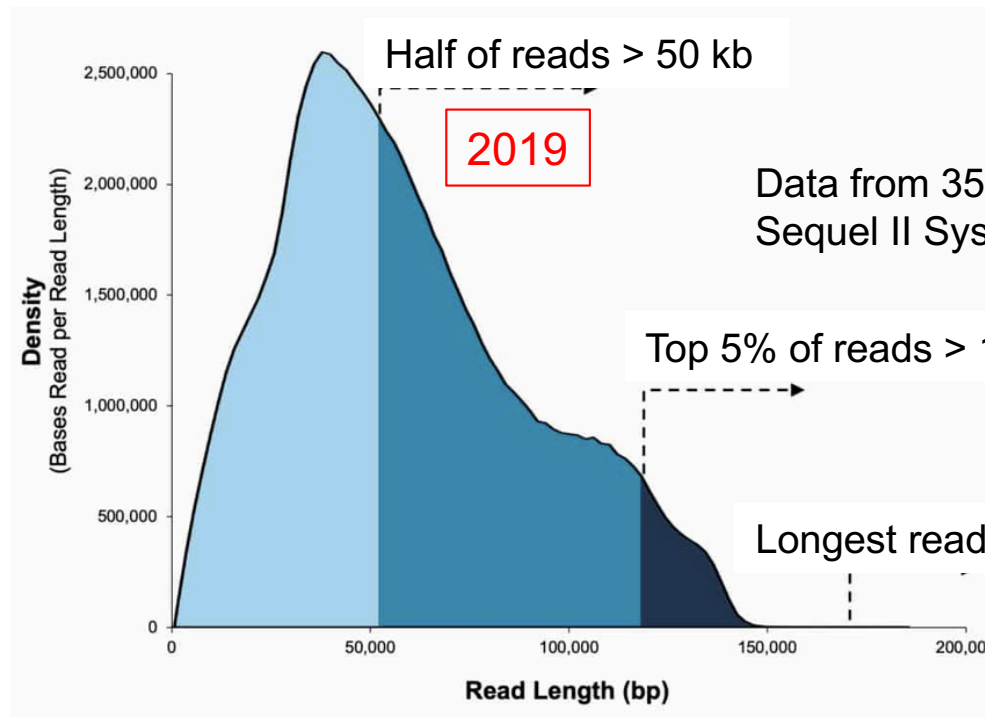




# Length of PacBio reads

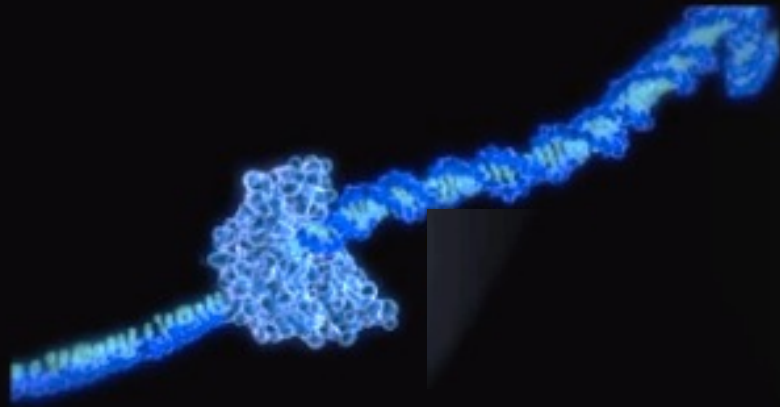


(Rhoads, *Genomics Proteomics Bioinformatics*, 2015)



Data from 35 kb size-selected *E. coli* library  
Sequel II System ; 2.0 Chemistry

## RECENT IMPROVEMENT WITH NEW CHEMISTRY

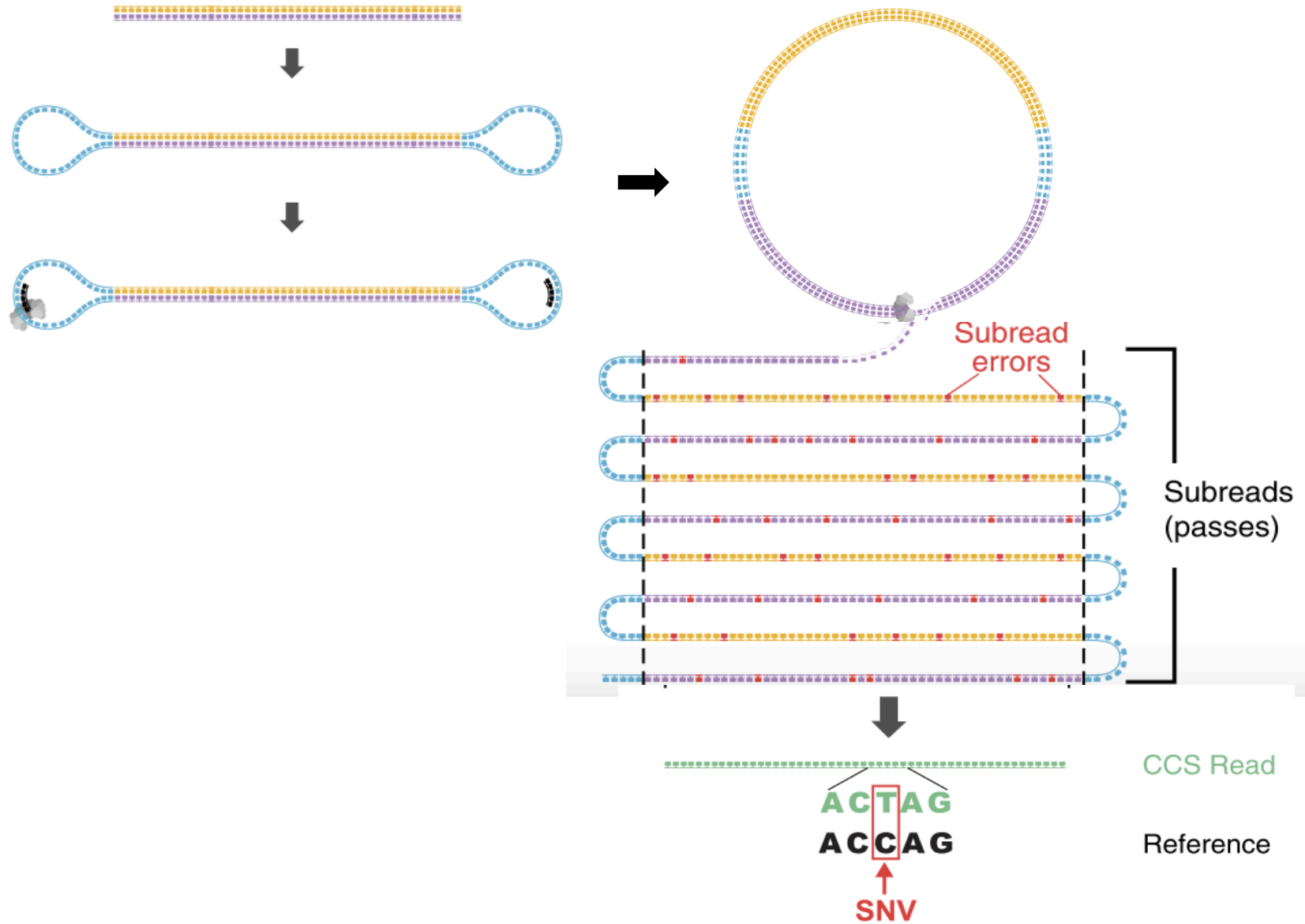


Circular consensus sequencing (CCS) reads are obtained when the SMRT bell template is replicated several times by the polymerase

This allows a highly accurate sequencing by correction of random errors



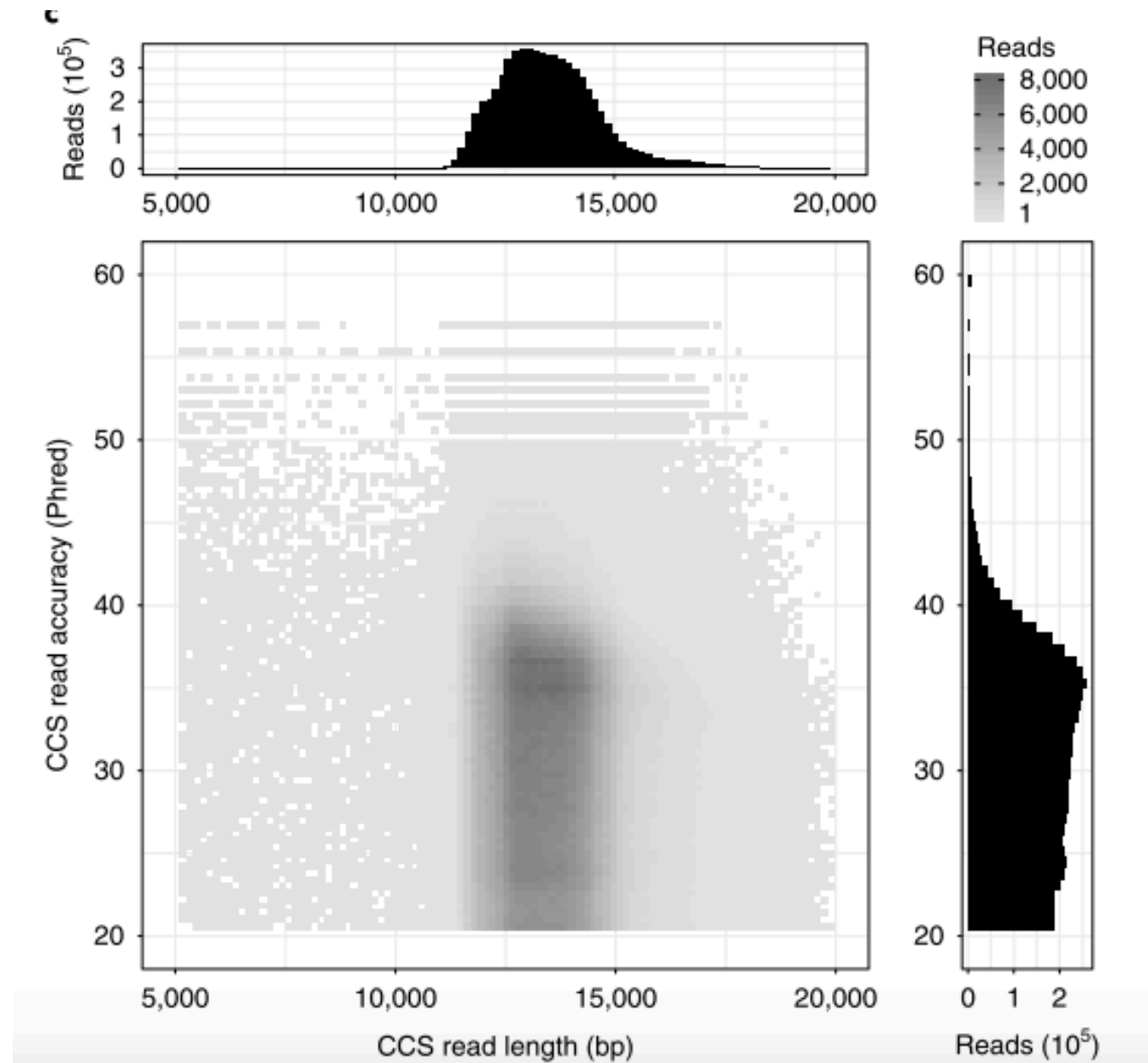
# Circular Consensus Sequences (CCS): HIFI READS



# RECENT IMPROVEMENT: GENOME ASSEMBLY WITH CCS

Circular consensus assembly of a human genome

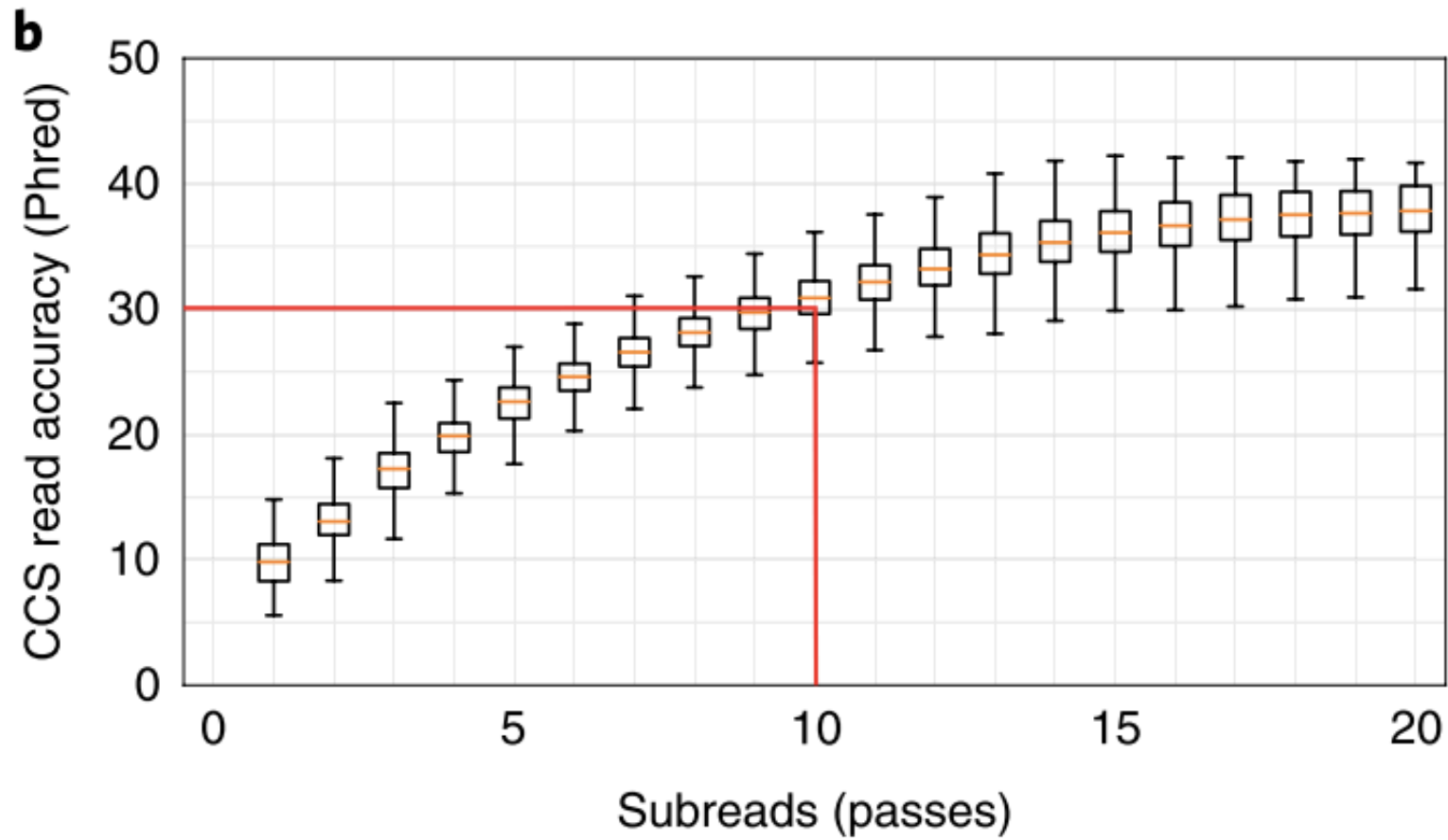
Wenger et al. *Nat. Biotechnol.* oct. 2019



# RECENT IMPROVEMENT: GENOME ASSEMBLY WITH CCS

Circular consensus assembly of a human genome

Wenger et al. *Nat. Biotechnol.* oct. 2019



# Genome assembly with CCS

Circular consensus assembly of a human genome  
Wenger et al. *Nat. Biotechnol.* oct. 2019

CCS reads alone : high quality contiguous genome : concordance of 99.997%

Assembler	Total size (Gb)	Contigs	N50 (Mb)	Ensembl genes (%)
Canu	3.42	18,006	22.78	93.2
FALCON	2.91	2,541	28.95	97.6
wtdbg2	2.79	1,554	15.43	96.1

## Canu assembly

- genome size > expected haploid genome because it resolves some heterozygous alleles into separate contigs

## Majority of CCS read discordances

- 3.4% mismatches → 1 mismatch every 13,048 bp
- 4.6% indels in non homopolymers. → 1 non-homopolymer indel every 9,669 bp
- 92.0% indels in homopolymers → 1 homopolymer indel every 477 bp

## Comparison with NovaSeq

- CCS mismatch rate is 17× lower than reads from NovaSeq
- CCS indel rate is 181× higher than reads from NovaSeq

# SEQUENCING cDNA USING CIRCULAR CONSENSUS SEQUENCES

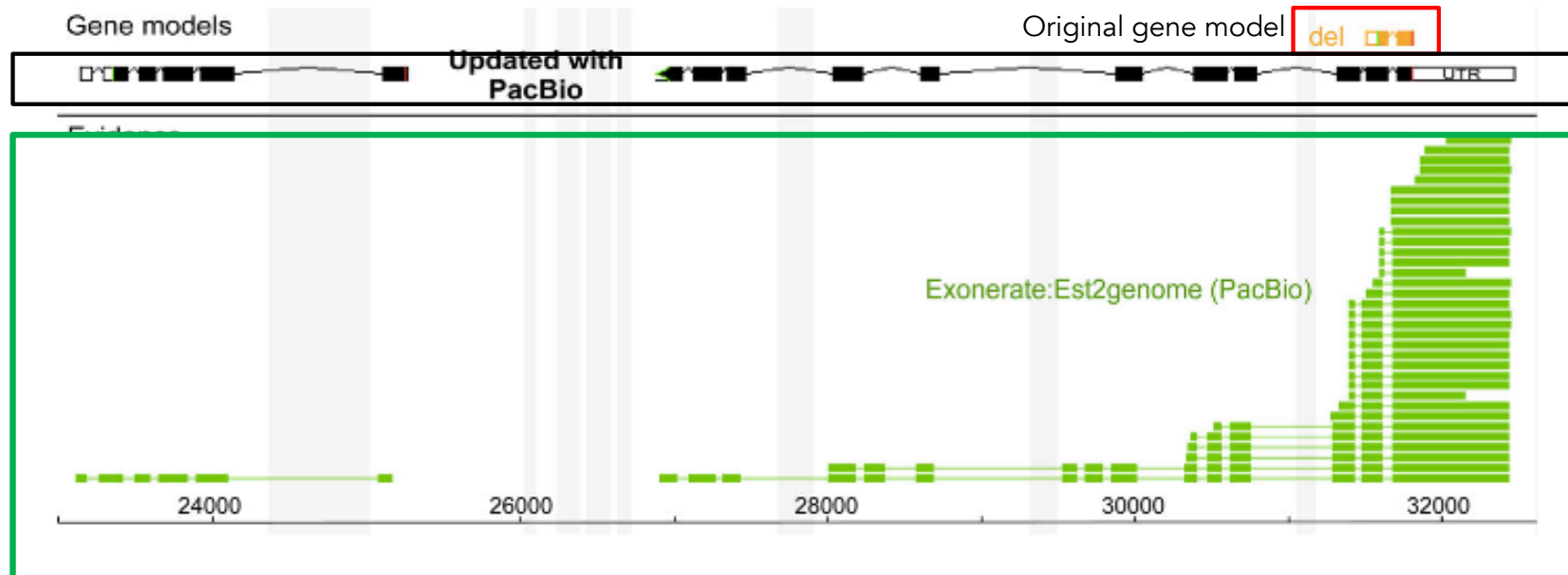
Genome annotation of the parasitic hookworm *Ancylostoma ceylanicum*  
using single molecule mRNA sequencing

Magrini et al. *BMC Genomics*, 2018

RNA → cDNA  $\xrightarrow{\text{PacBio sequencing}}$  193 000 CCS



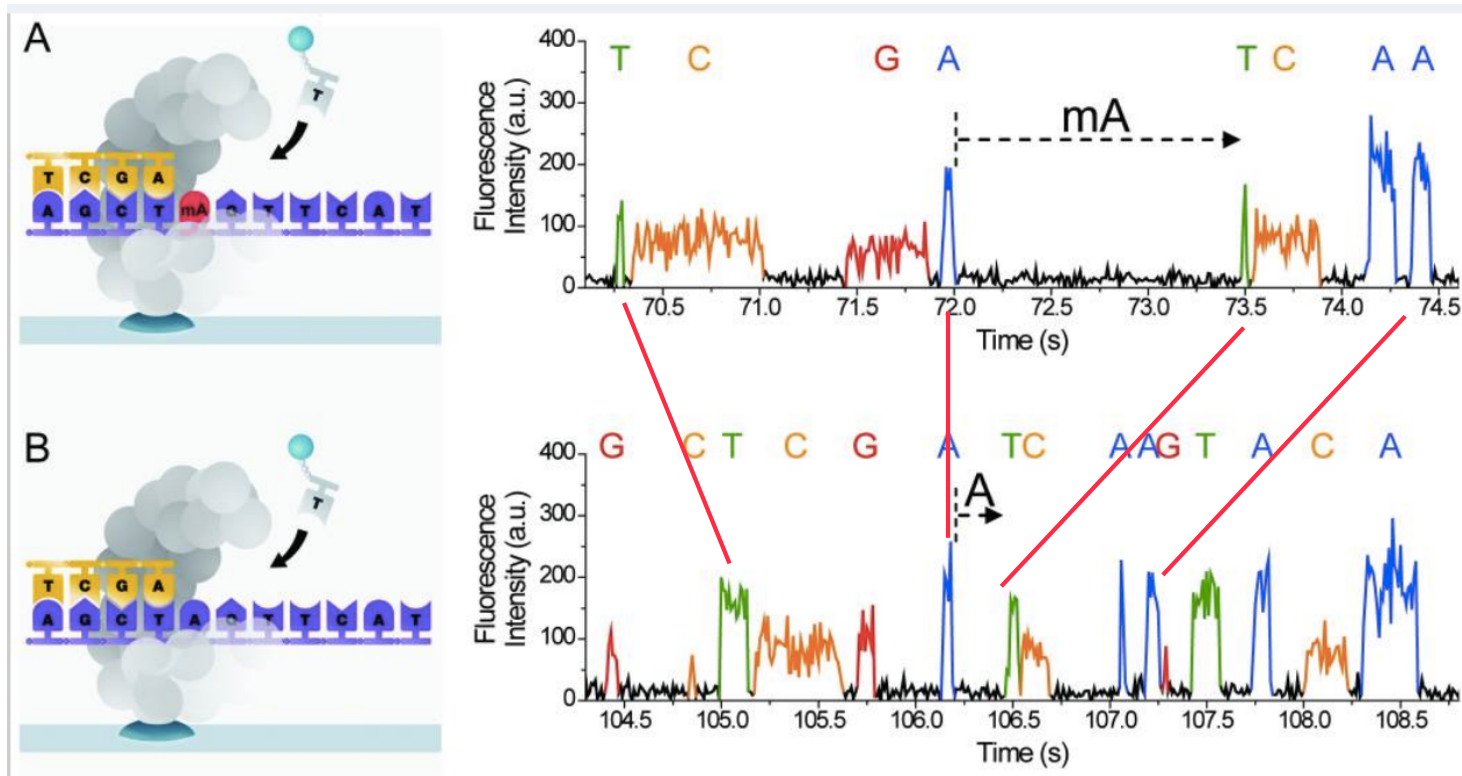
- Increased the total genomic exon length by 1.9 Mb (12.4%)
- 1609 (9.2%) new genes







# DETECTION OF MODIFIED BASES



from Fusberg et al. *Nature Methods* (2010)

Detection of 5mA with strong influence of sequence contexts: requires high coverage

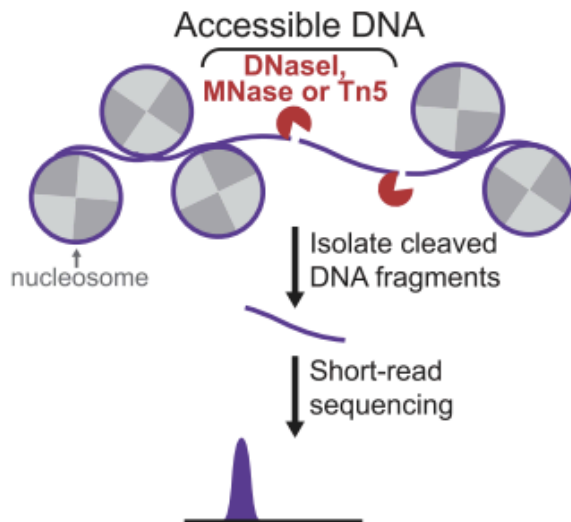
Feng et al. *PLOS Comput Biol* (2013)

# Detection of m6A with CCS

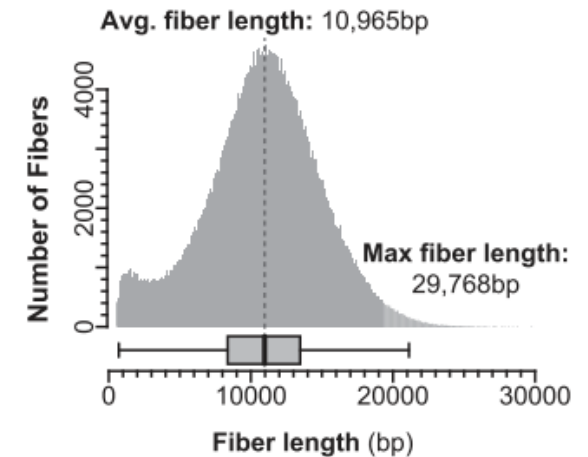
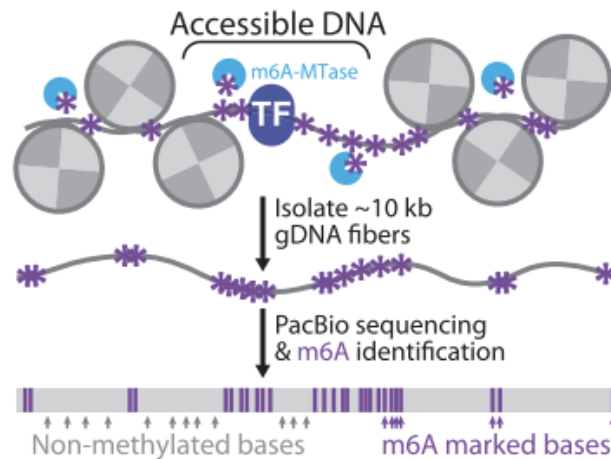
Single-molecule regulatory architectures captured by chromatin fiber sequencing  
Stergachis et al. *Science* (2020)

## DnaseI-seq.

### Cleavage-based assay:



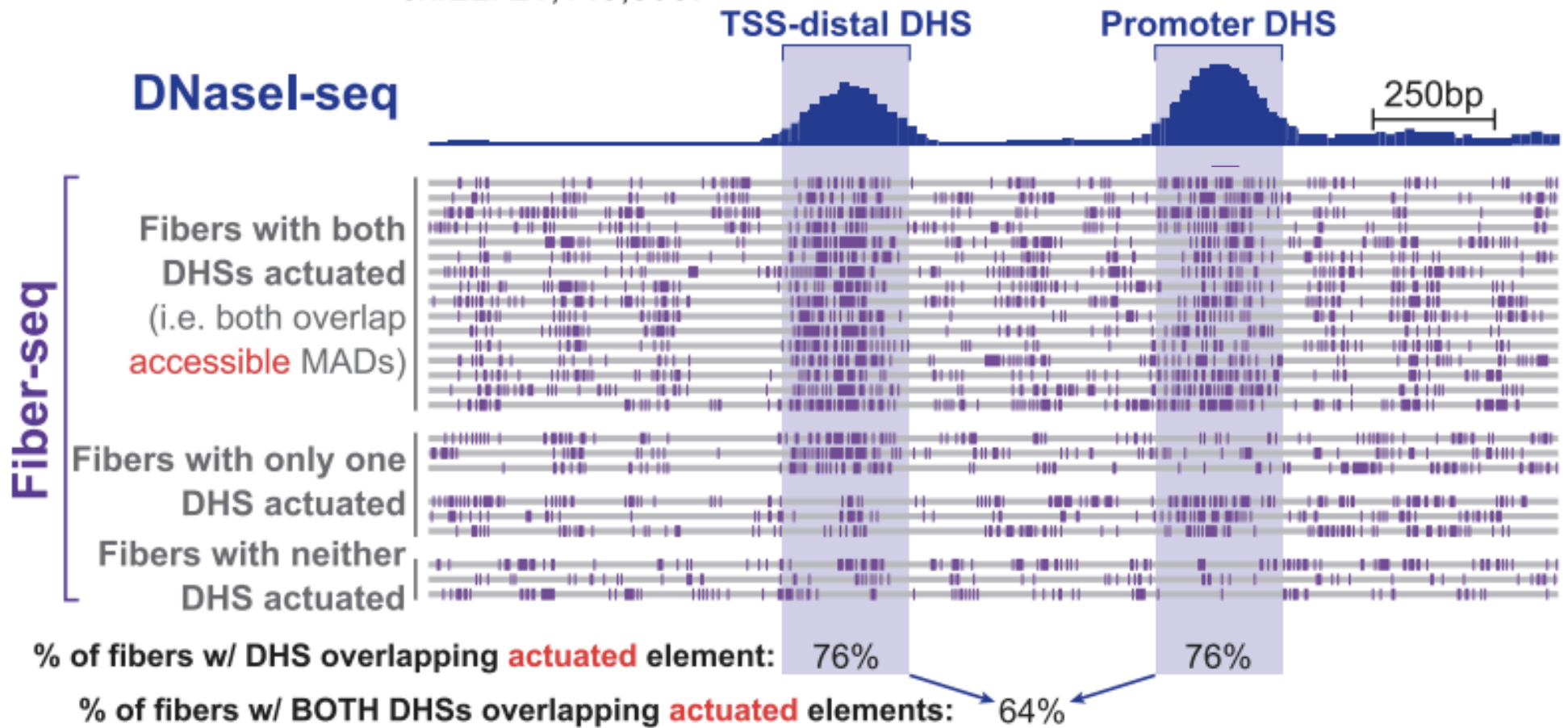
## Fiber-seq.



# Detection of m6A with CCS

Single-molecule regulatory architectures captured by chromatin fiber sequencing  
Stergachis et al. *Science* (2020)

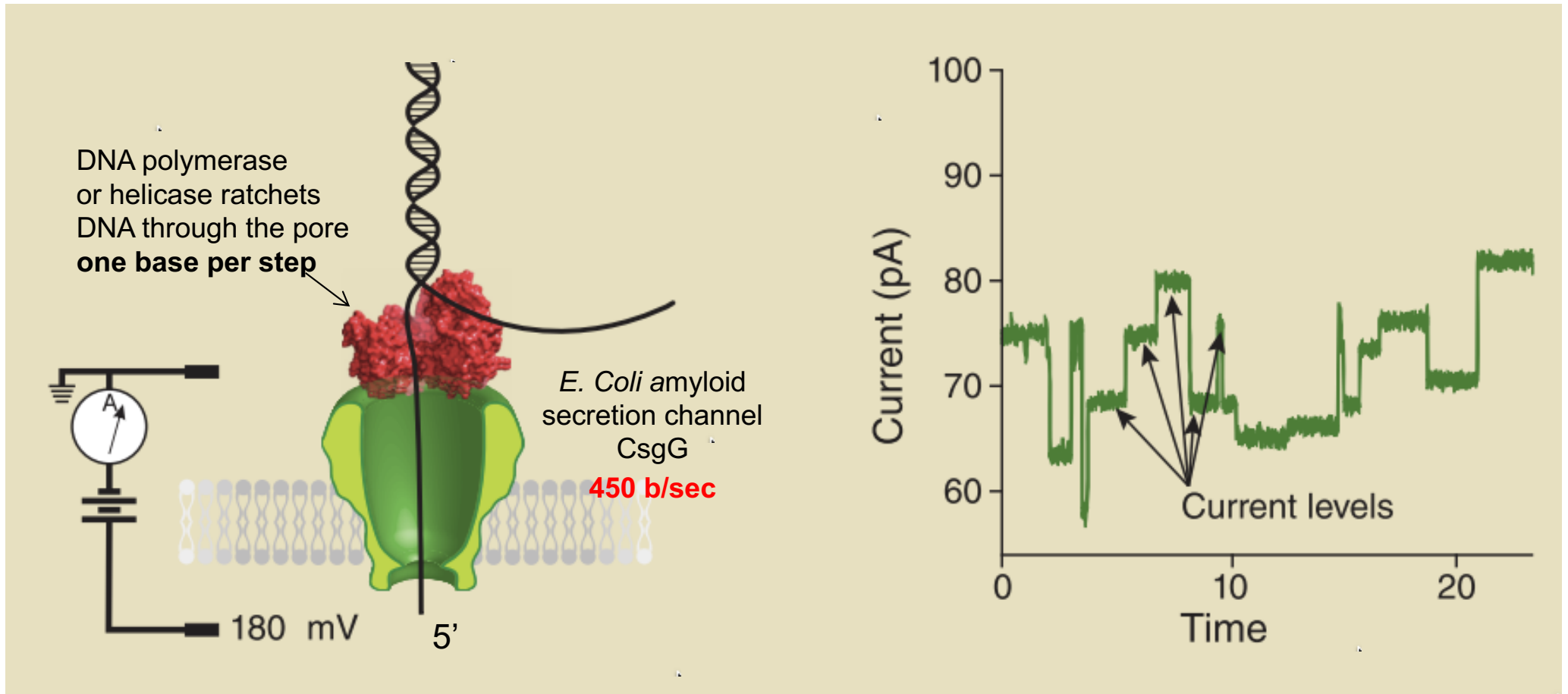
DHS : DNaseI Hypersensitive Site



# Next Generation Sequencing

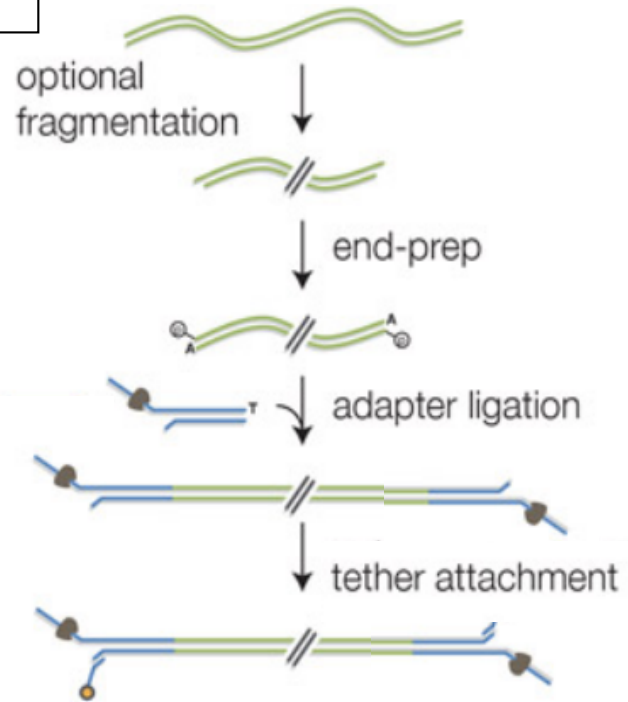


# BASIC CONCEPTS

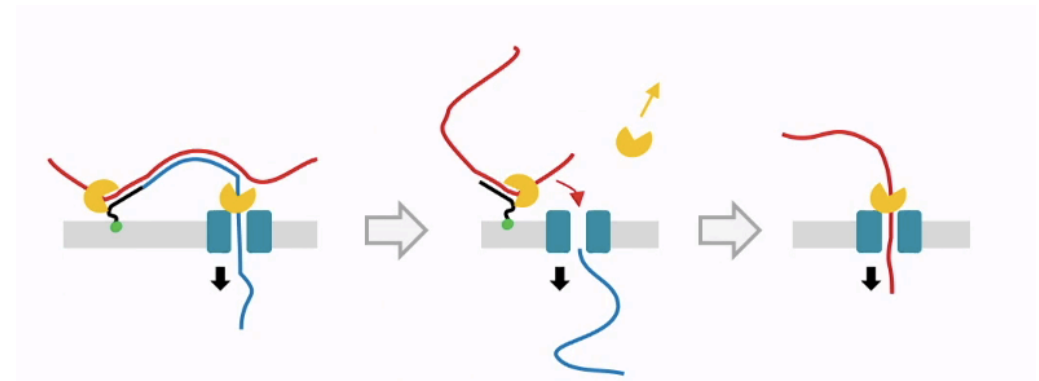
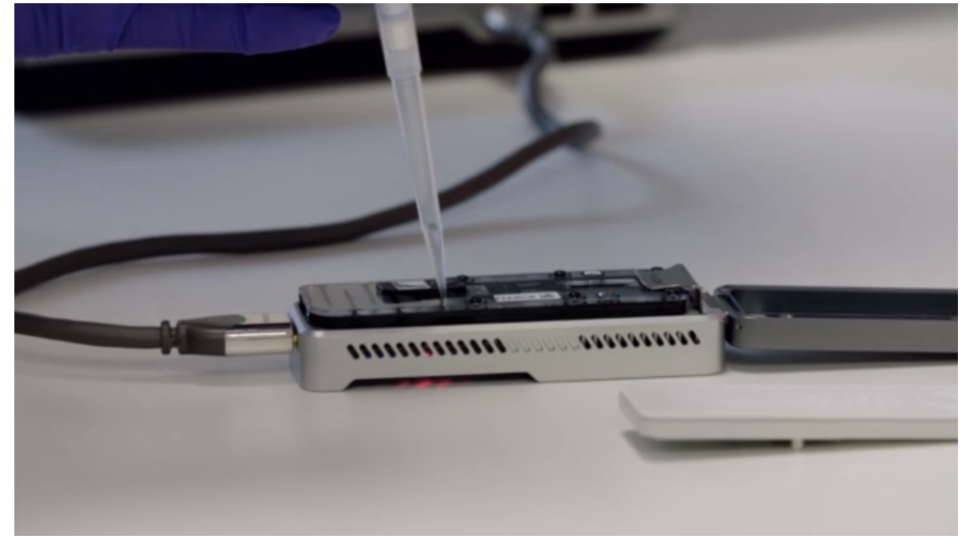


# SEQUENCING PROCESS

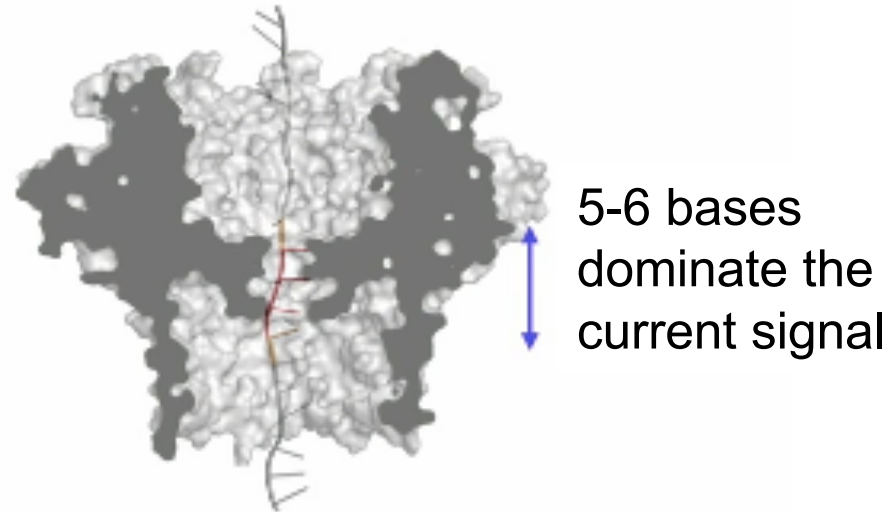
1D<sup>2</sup> Library  
(2017)



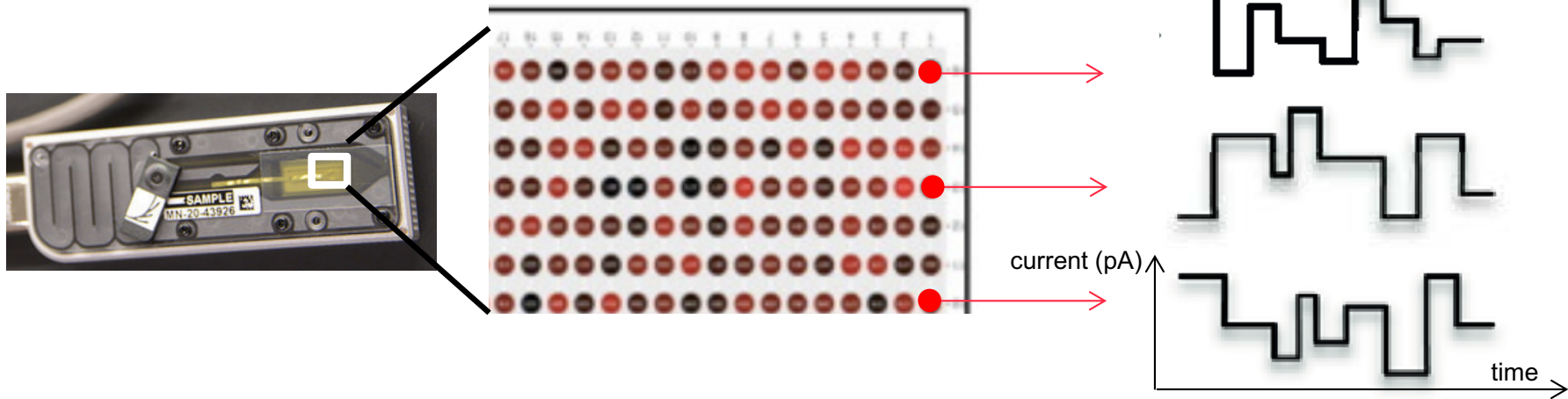
SEQUENCING



# SEQUENCING PROCESS : MinION FLOW CELL

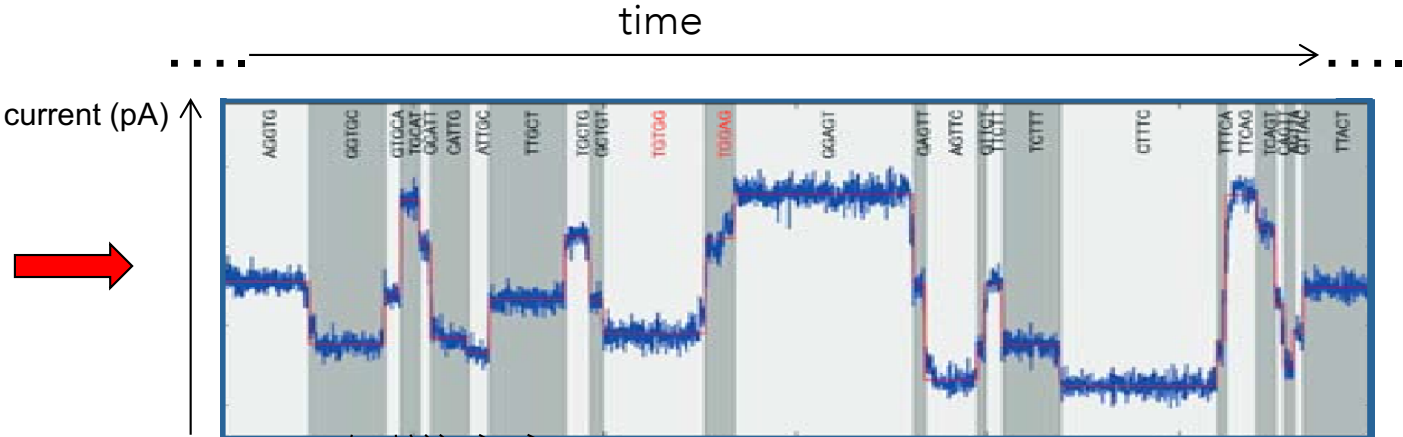
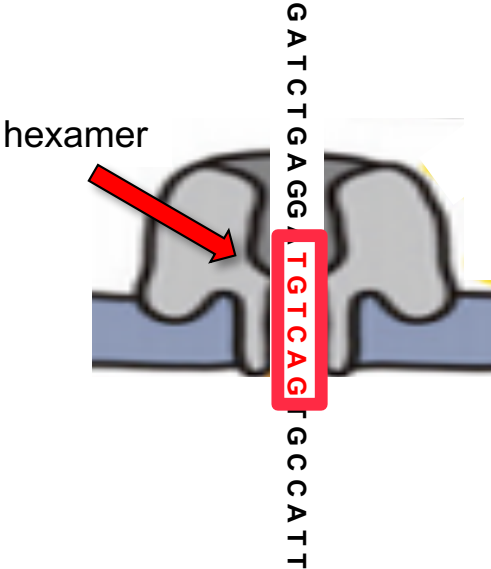


MinION : 512 pores

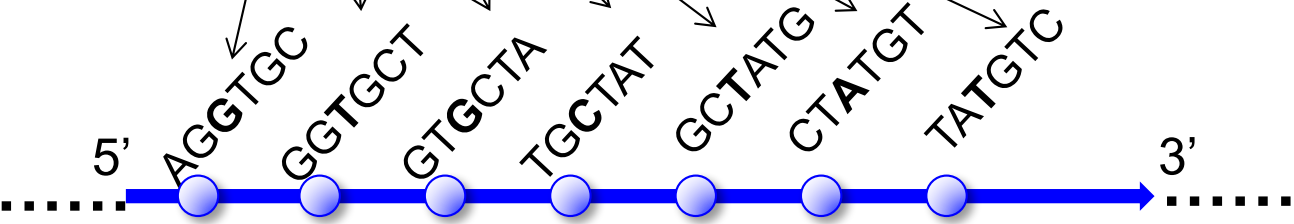


PromethION : 144000 pores (48 x 3000)

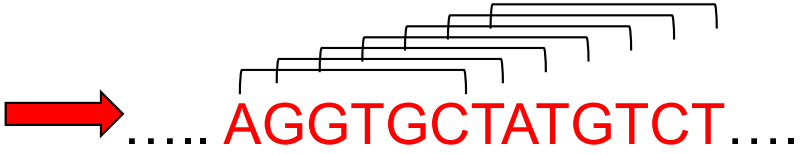
# BASECALLING



3000 values/s

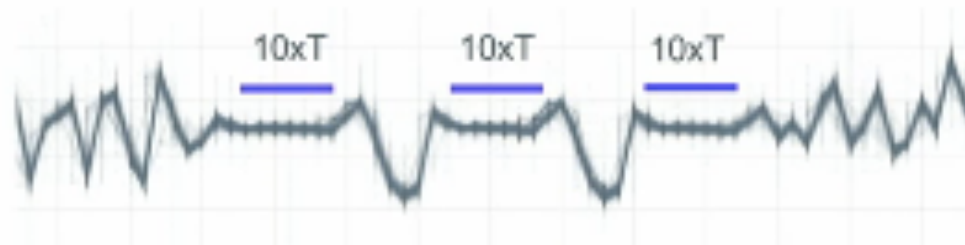
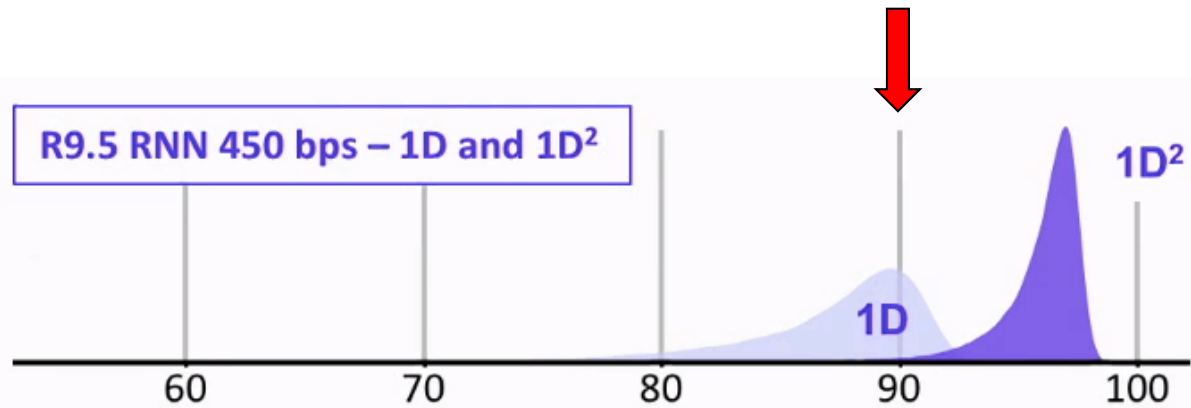
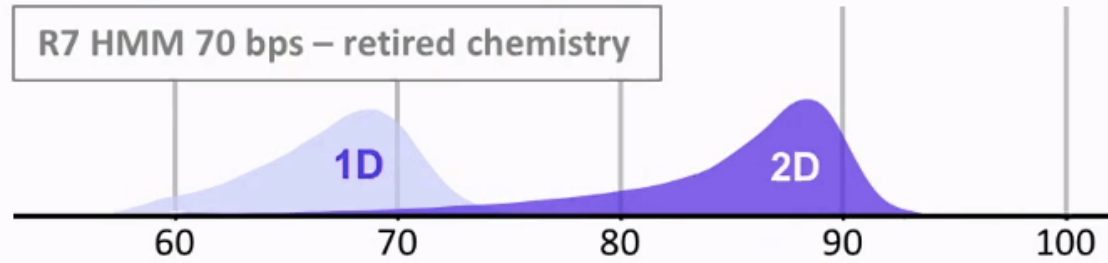


Basecalling : finding the optimal path of successive 6-mers





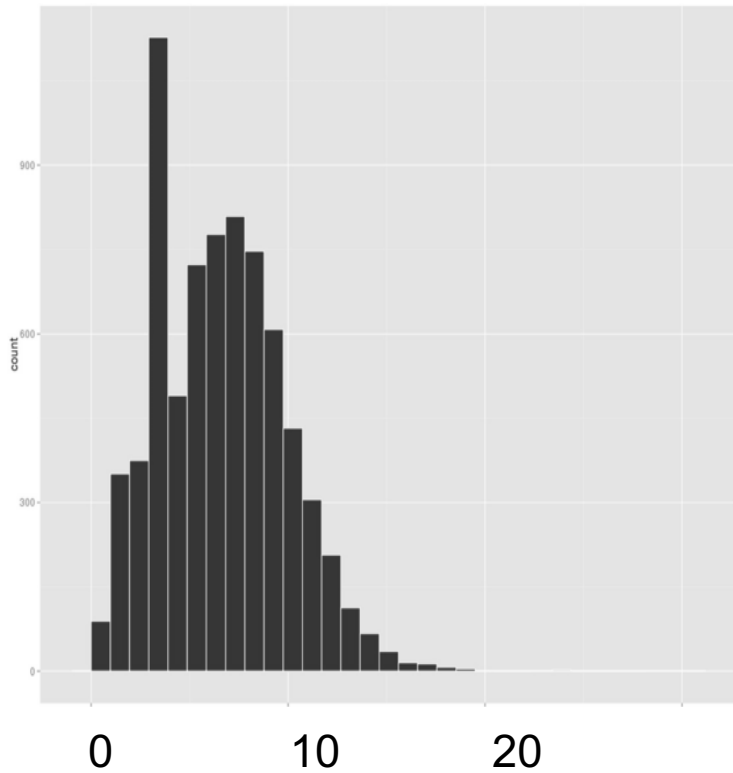
# QUALITY



Homopolymers difficult to sequence

# SIZE OF SEQUENCED DNA FRAGMENTS

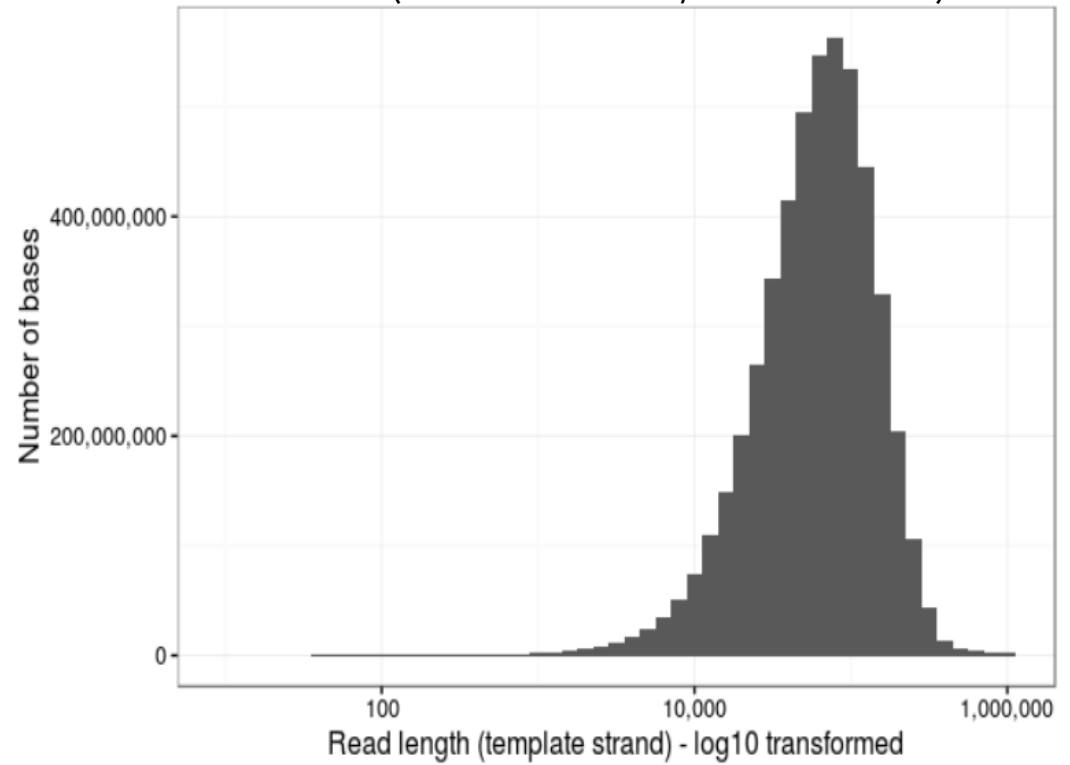
Typical profile of fragment size



Length of sequenced fragments (kb)

(Risse et al. *GigaScience*, 2015)

“Ultra long” reads  
(lab.loman.net, March 2017)

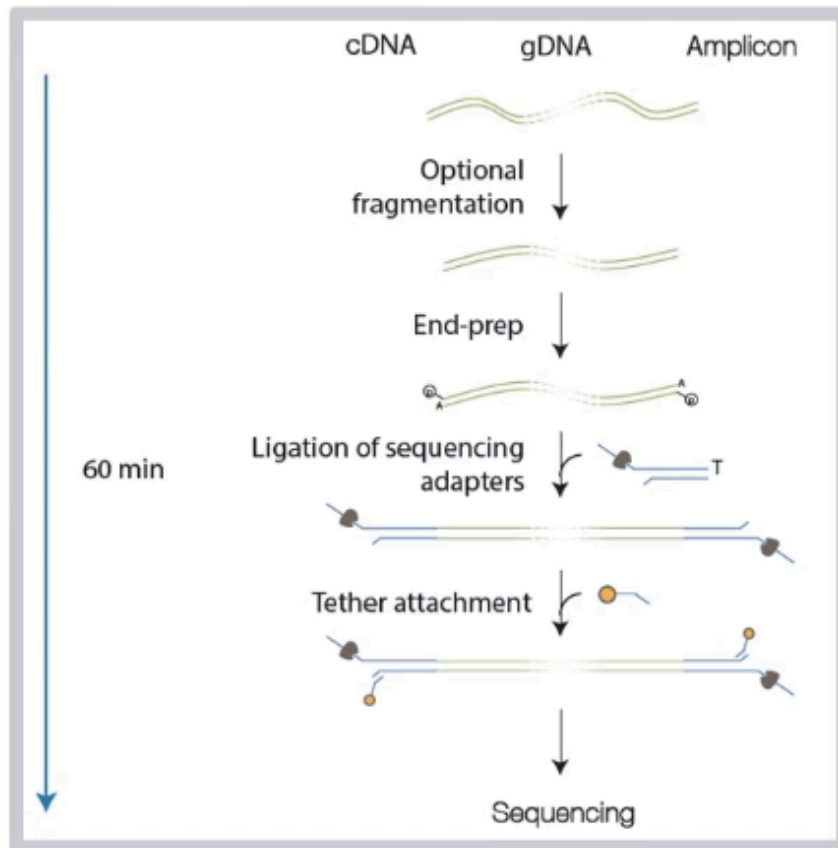


Size of the longest read : 778 kb

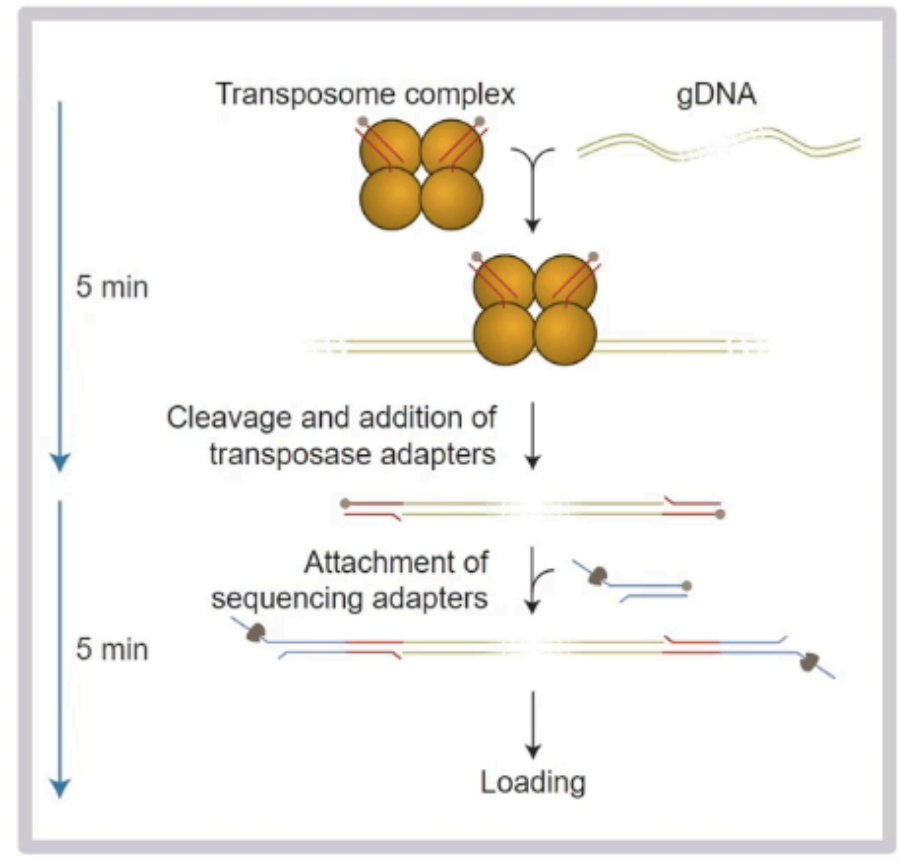
1 contig of the 4.6Mb chromosome of *E. coli*  
obtained with just the 7 longest reads

# SIZE OF SEQUENCED DNA FRAGMENTS

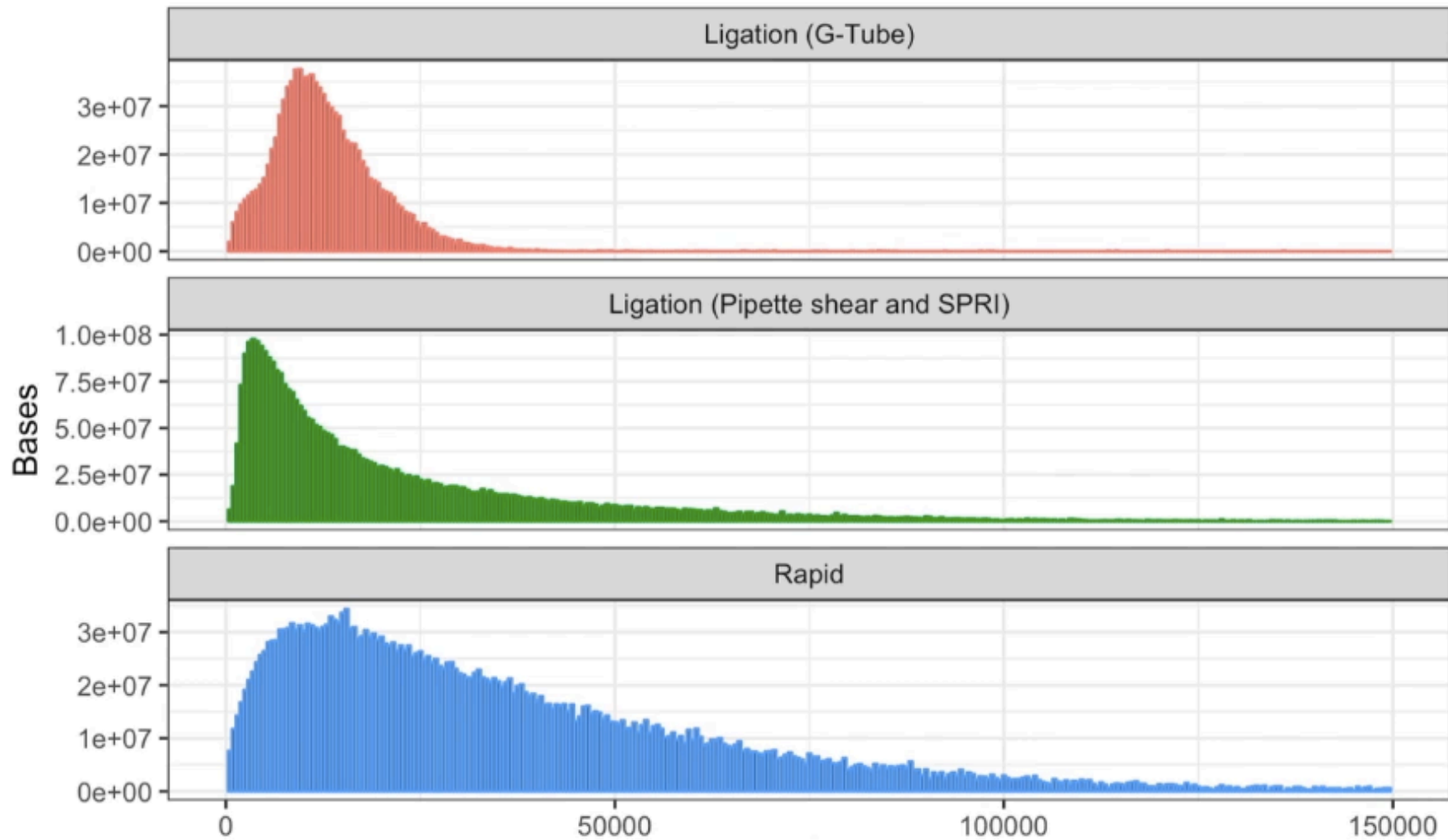
## Ligation method



## Transposase method



# SIZE OF SEQUENCED DNA FRAGMENTS



Josh Quick, Nick Loman

see John Tyson's video (ONT website)

# HYBRID GENOME ASSEMBLY : NANOPORE AND ILLUMINA DATA

## *Acinetobacter baylyi* (data from Oxford Nanopore)

Assemblies	Illumina only	Illumina + MinION
Input Coverage	50X	13X
# contigs	20	1
Assembly size (Mb)	3.59	3.62
N90 size (Kb)	326	3 621
NA75 size (Kb)	194	1 002
Genome fraction (%)	99.73	99.997
# misassemblies	4	2
# local misassemblies	3	4
# mismatches per 100 Kb	6.49	3.11
# indels per 100 Kb	0.33	0.14

**nature  
biotechnology**

Jan. 2018

OPEN

## Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain<sup>1,13</sup>, Sergey Koren<sup>2,13</sup>, Karen H Miga<sup>1,13</sup>, Josh Quick<sup>3,13</sup>, Arthur C Rand<sup>1,13</sup>, Thomas A Sasani<sup>4,5,13</sup>, John R Tyson<sup>6,13</sup>, Andrew D Beggs<sup>7</sup>, Alexander T Dilthey<sup>2</sup>, Ian T Fiddes<sup>1</sup>, Sunir Malla<sup>8</sup>, Hannah Marriott<sup>8</sup>, Tom Nieto<sup>7</sup>, Justin O'Grady<sup>9</sup>, Hugh E Olsen<sup>1</sup>, Brent S Pedersen<sup>4,5</sup>, Arang Rhie<sup>2</sup>, Hollian Richardson<sup>9</sup>, Aaron R Quinlan<sup>4,5,10</sup>, Terrance P Snutch<sup>6</sup>, Louise Tee<sup>7</sup>, Benedict Paten<sup>1</sup>, Adam M Phillippy<sup>2</sup>, Jared T Simpson<sup>11,12</sup>, Nicholas J Loman<sup>3</sup> & Matthew Loose<sup>8</sup>

eserved.

Using nanopore reads alone assembly of a human genome :

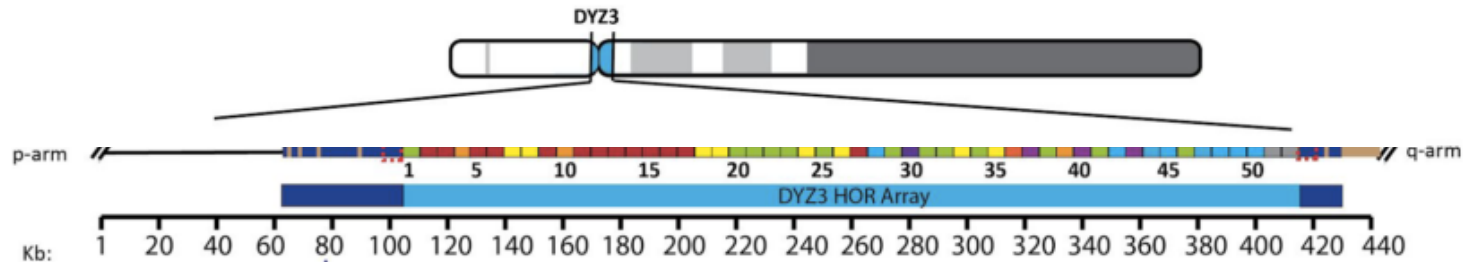
- NG50 contig size of ~6.4 Mb
- covers >85% of the reference
- 99.88% accuracy
- MHC locus on a single contig, phased over its full length
- closure of 12 large (>50 kb) gaps in the reference human genome

# ASSEMBLY OF A HUMAN Y CENTROMERE

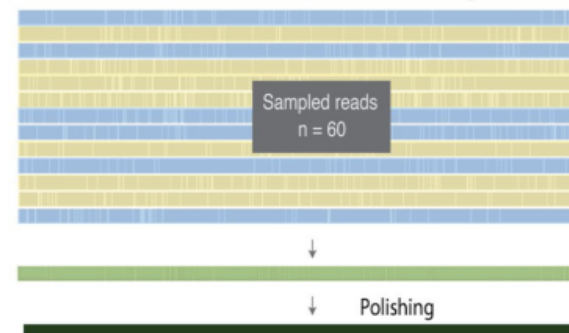
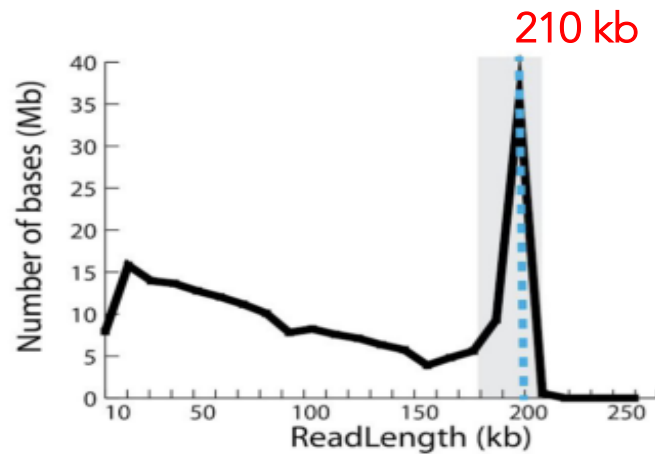
(Jain et al., *bioRxiv*, 2017)

300 kb array of 5.8 kb sequence repeated in an uninterrupted head-to-tail orientation

To date, no technology has been capable of sequencing centromeres due to **requirement for extremely high-quality long reads**



9 BACs  
100kb to 210kb



Final high quality consensus BAC sequence

FIRST COMPLETE SEQUENCE OF A HUMAN CENTROMERE

# GENOME SEQUENCING

Long-read sequencing for non-small-cell lung cancer genomes  
Sakamoto et al. *Genome Research*, Sept. 2020

Nanopore sequencing of **cancer cell lines** (Promethion)

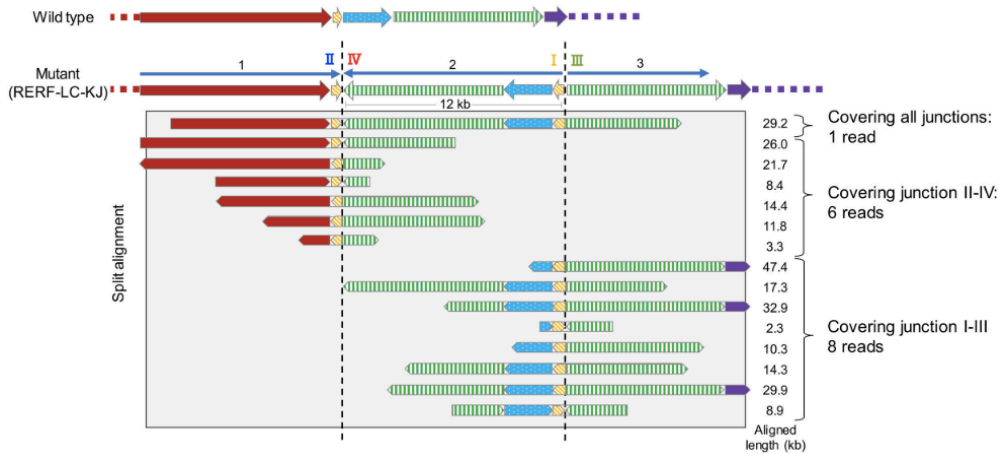
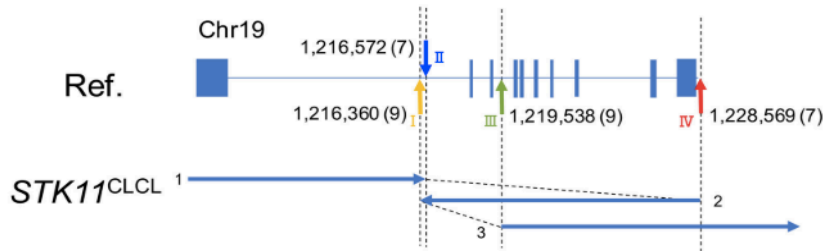
- Maximum length : 0.99 Mb
- N50 : 32 kb
- Average mapped reads : 14 kb

Biological relevance of SV further revealed by :

- epigenome,
- transcriptome,
- protein analyses

Sequencing of clinical tumor samples

- Structural aberrations also found in **clinical lung adenocarcinoma specimens**



Structural variants : comparison with PacBio sequencing

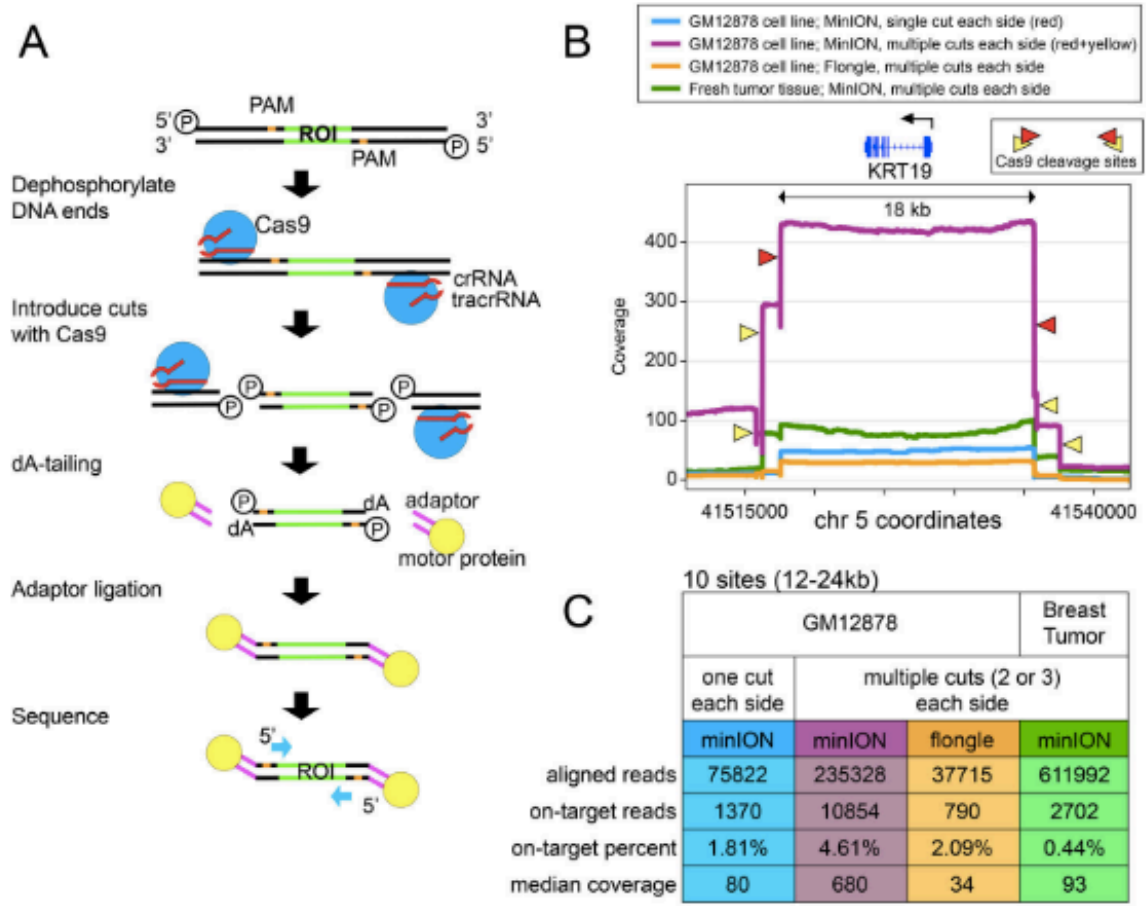


“These results indicated that neither the PacBio nor the PromethION platform is currently perfect; therefore, they should be used to complement each other.”



# GENOME SEQUENCING : TARGETED SEQUENCING

Targeted nanopore sequencing with Cas9-guided adaptor ligation  
 Gilpatrick et al. *Nature Biotechnology* April 2020



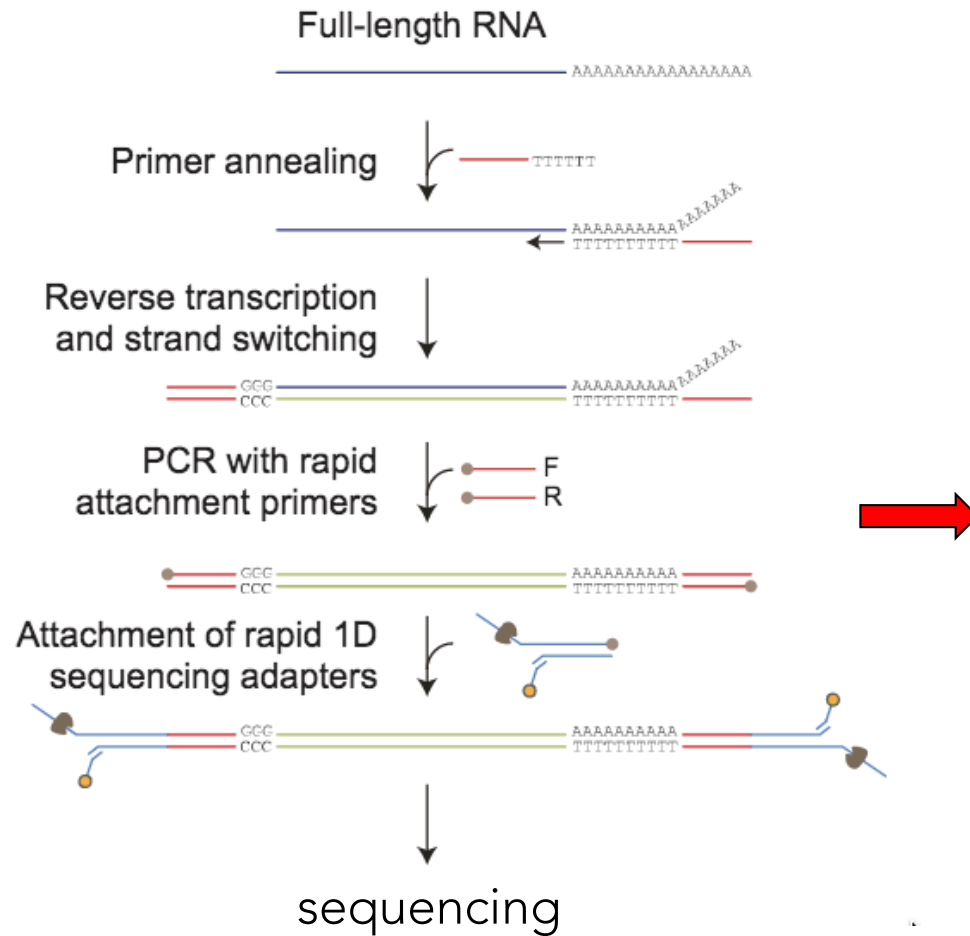
nCATS = nanopore Cas9-targeted sequencing : enrichment strategy using targeted cleavage of DNA to ligate adaptors for nanopore

nCATS can simultaneously assess :

- haplotype-resolved single-nucleotide variants (SNVs)
- structural variations (SVs)
- CpG methylation...
- **Best median sequencing coverage : 680 X**
- nCATS uses only ~3 µg of genomic DNA + can target a large number of loci in a single reaction.

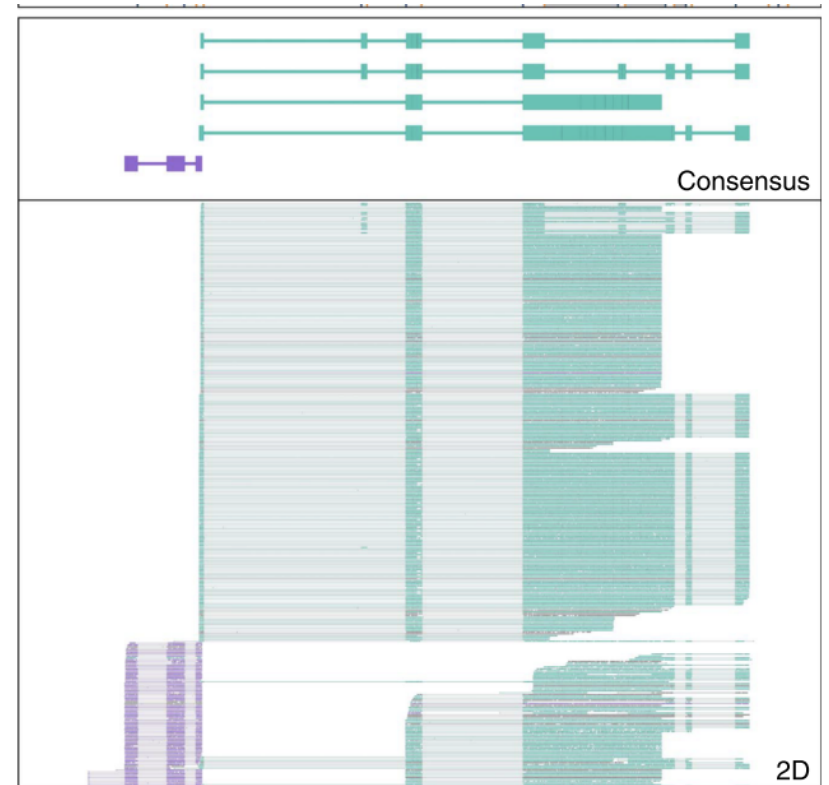
# cDNA SEQUENCING

## Library preparation



## Detection of splice variants in surface receptor of B cells

(Byrne et al. *Nat. Comm.* 2017)

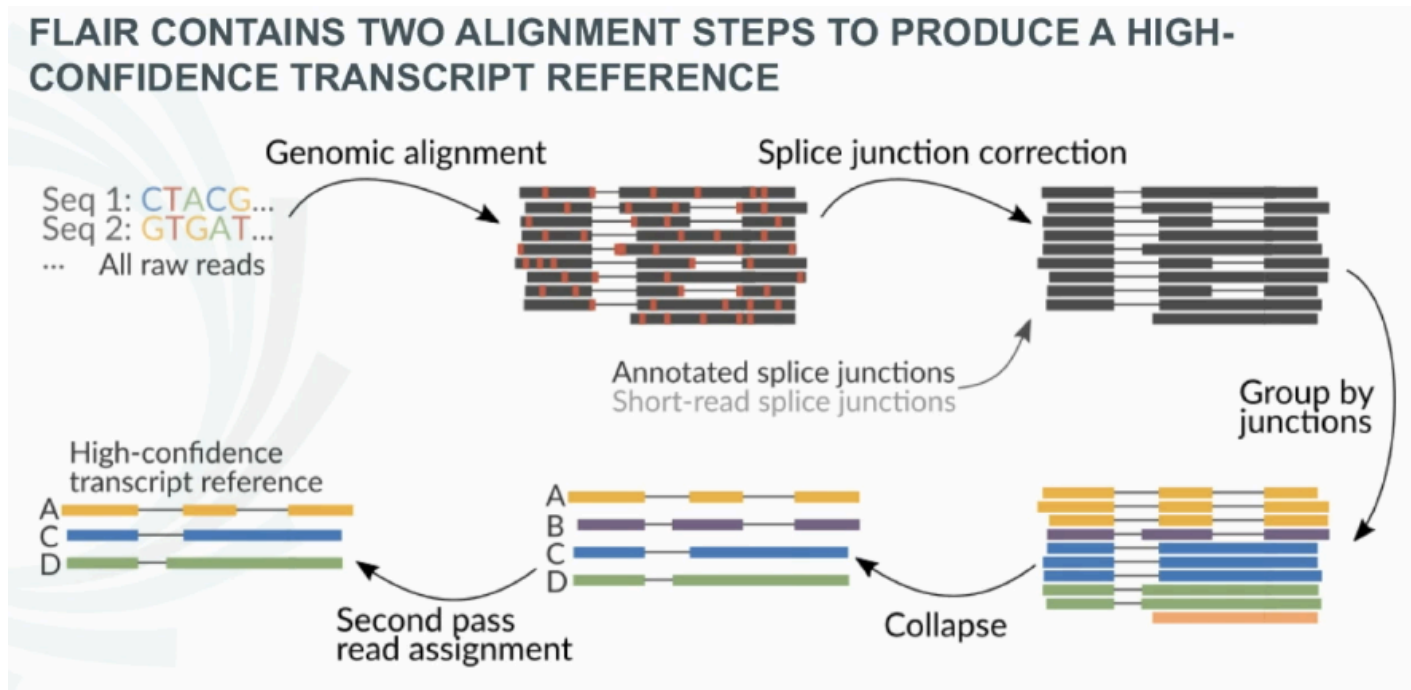


- Splice alignment uneasy due to high (5-10%) error rate
- Reads are frequently truncated from 5' end

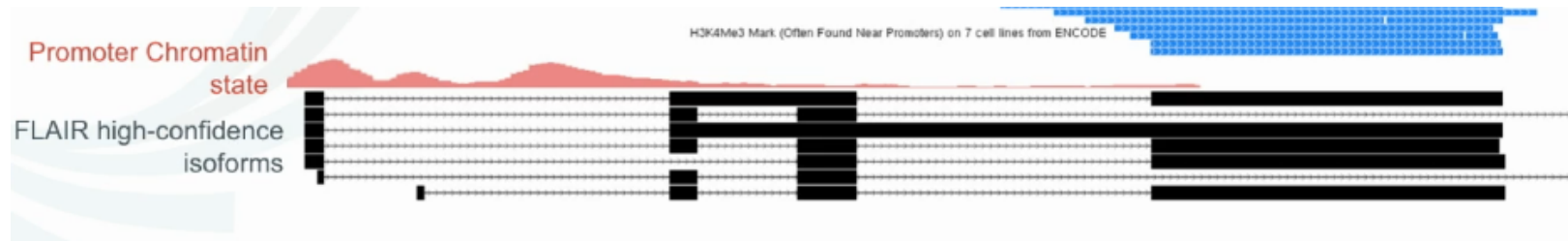
# CHALLENGES OF NANOPORE TRANSCRIPTOME ANALYSIS

FLAIR : a pipeline for splicing isoform determination

Tang et al. *bioRxiv* 2018



FLAIR incorporates promoter chromatin states to distinguish 5' truncations from true novel start sites

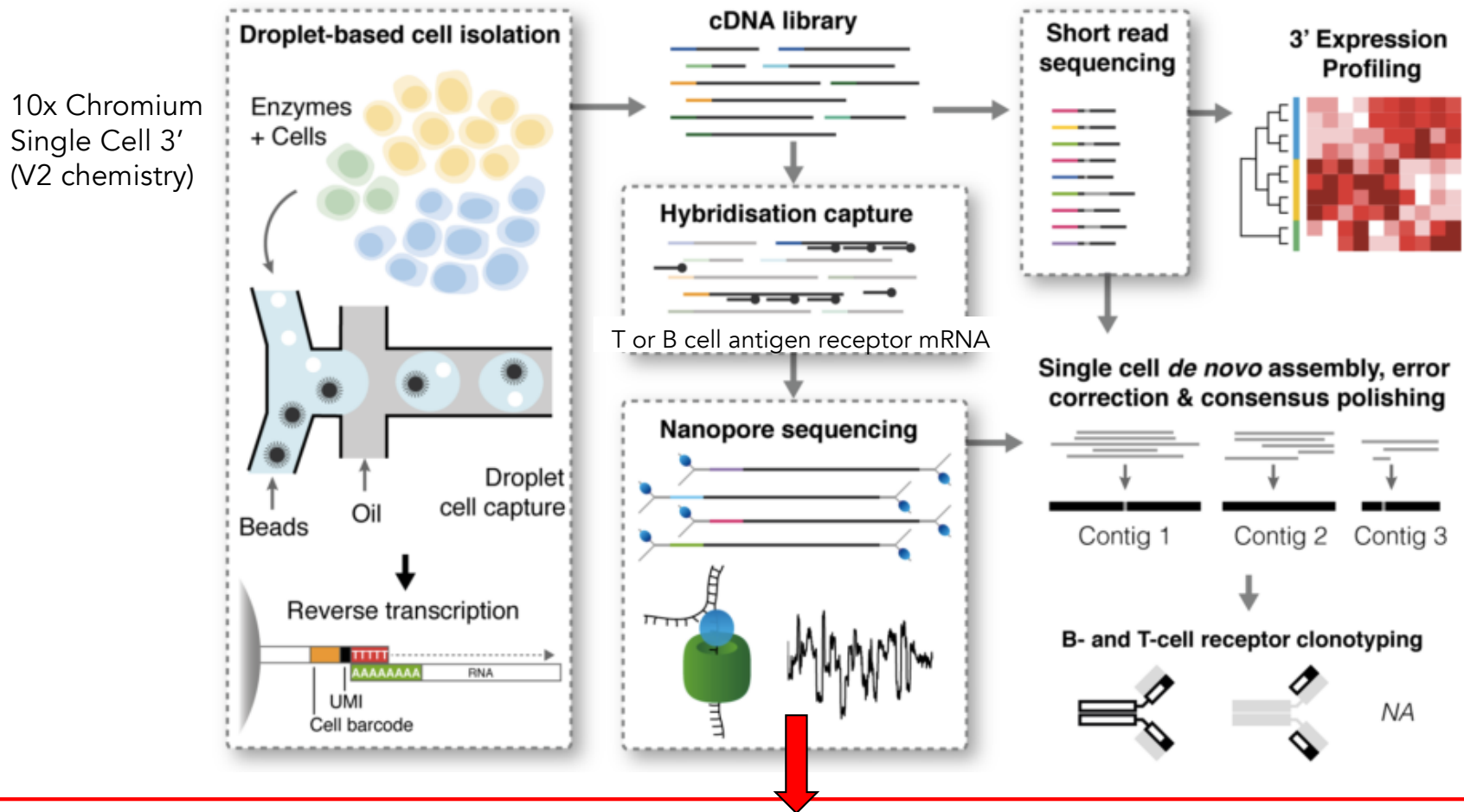


cDNA from chronic lymphocytic leukemia (CLL)

# NANOPORE and SINGLE CELL cDNA SEQUENCING

High-throughput targeted long-read single cell sequencing reveals the clonal and transcriptional landscape of lymphocytes  
Singh et al., *bioRxiv*, 2018

RAGE-seq (Repertoire And Gene Expression sequencing) : combines targeted long-read sequencing with short-read transcriptome of barcoded single cell libraries



Tracking of somatic mutation, alternate splicing and clonal evolution of T and B lymphocytes

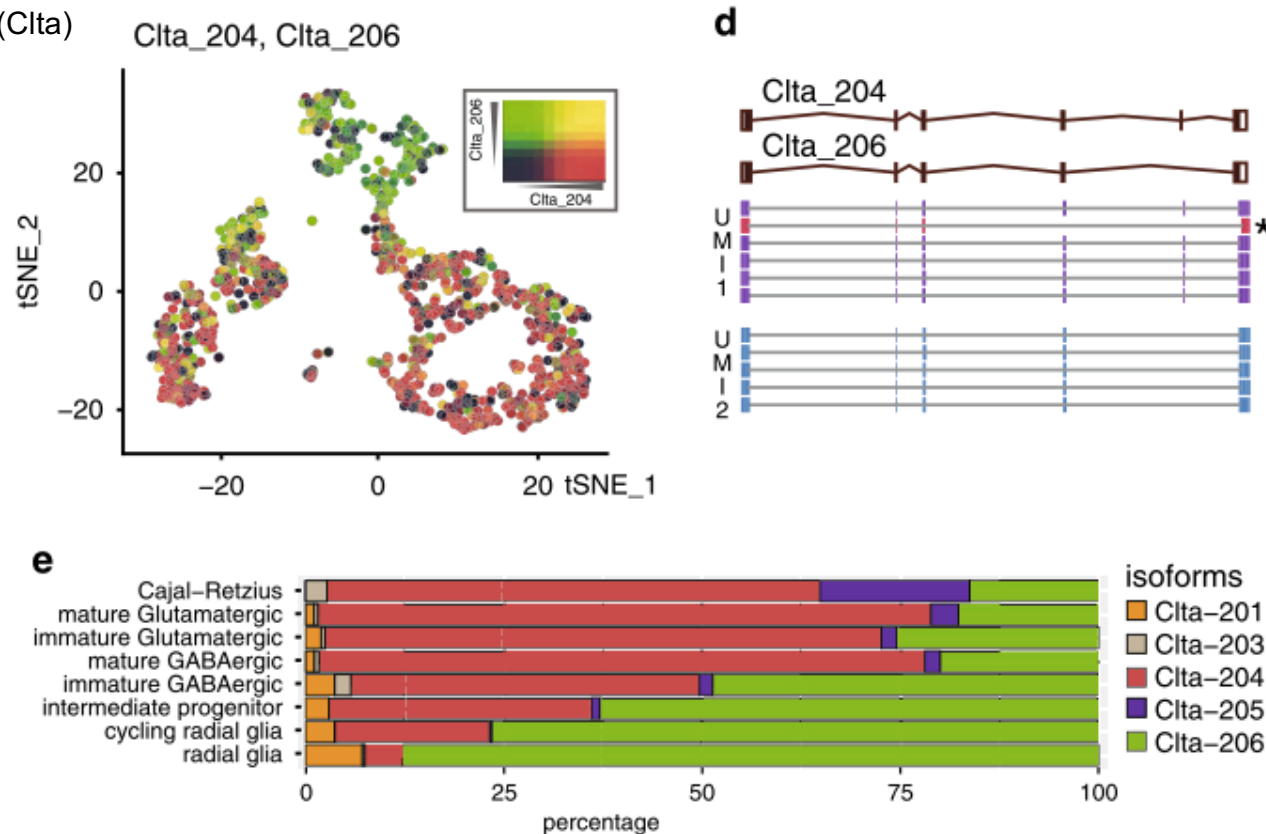
# NANOPORE and SINGLE CELL cDNA SEQUENCING

High-throughput targeted long-read single cell sequencing reveals the clonal and transcriptional landscape of lymphocytes  
Lebrigand et al., *Nature Communications*, 2020

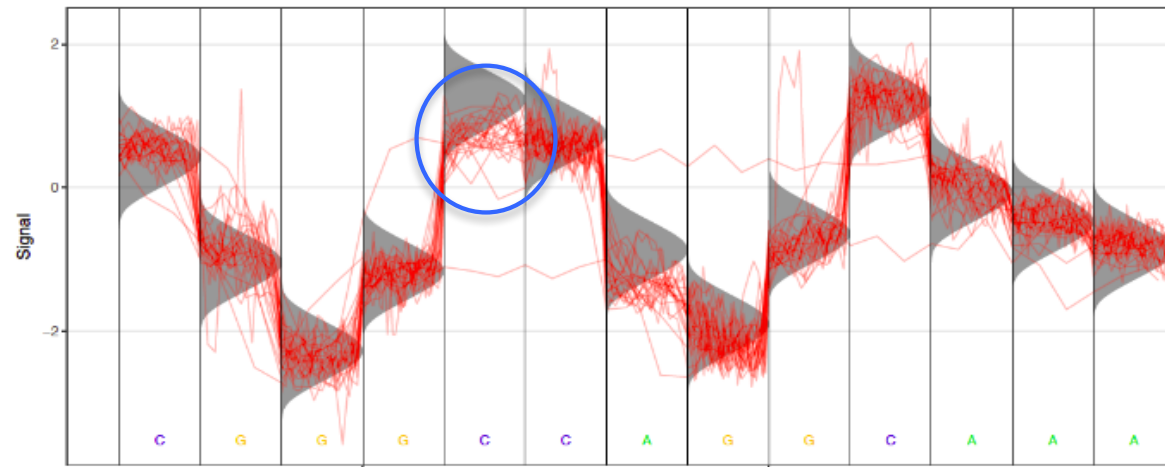
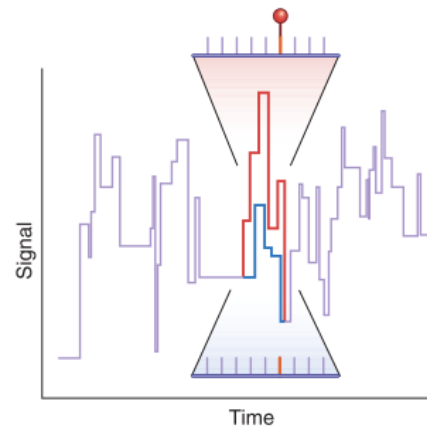
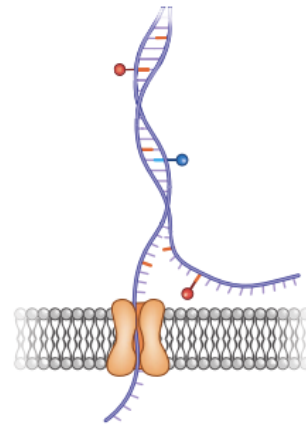
**ScNaUmi-seq** : Single-cell Nanopore sequencing with UMIs (10x Genomics Chromium system)

- High accuracy cell BC and UMI assignment
- Analysis of splicing and sequence variation at the single-cell level

Clathrin light chain A (Clta)



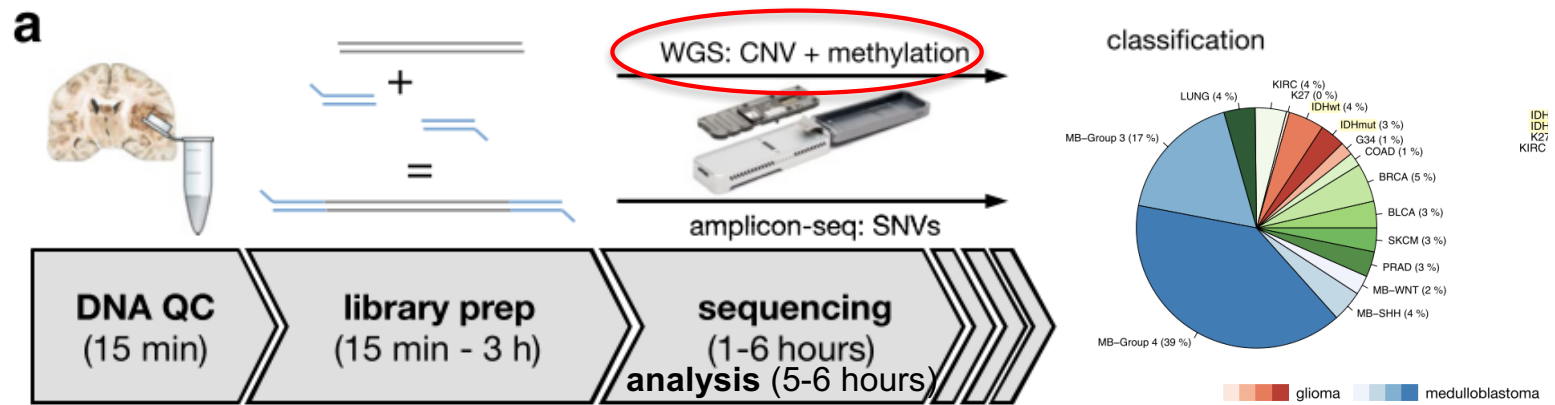
# DETECTION OF MODIFIED BASES



— Electric signal  
▶ Canonical base distribution

# DETECTION OF MODIFIED BASES IN CANCER GENOMES

Same-day genomic and epigenomic diagnosis of brain tumors (gliomas, medulloblastomas) with nanopore sequencing  
Euskirchen et al., *Acta Neuropathol.* (2017)



Same-day detection of :

- structural variants
- point mutations
- methylation profiling

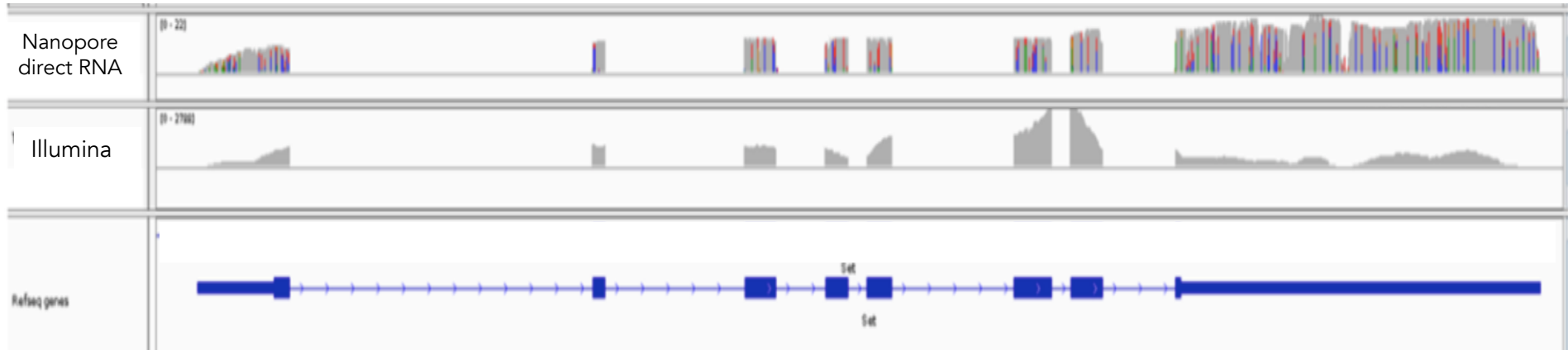
Single device with negligible capital cost :

- outperforms hybridization-based and current sequencing technologies
- makes precision medicine possible for every cancer patient

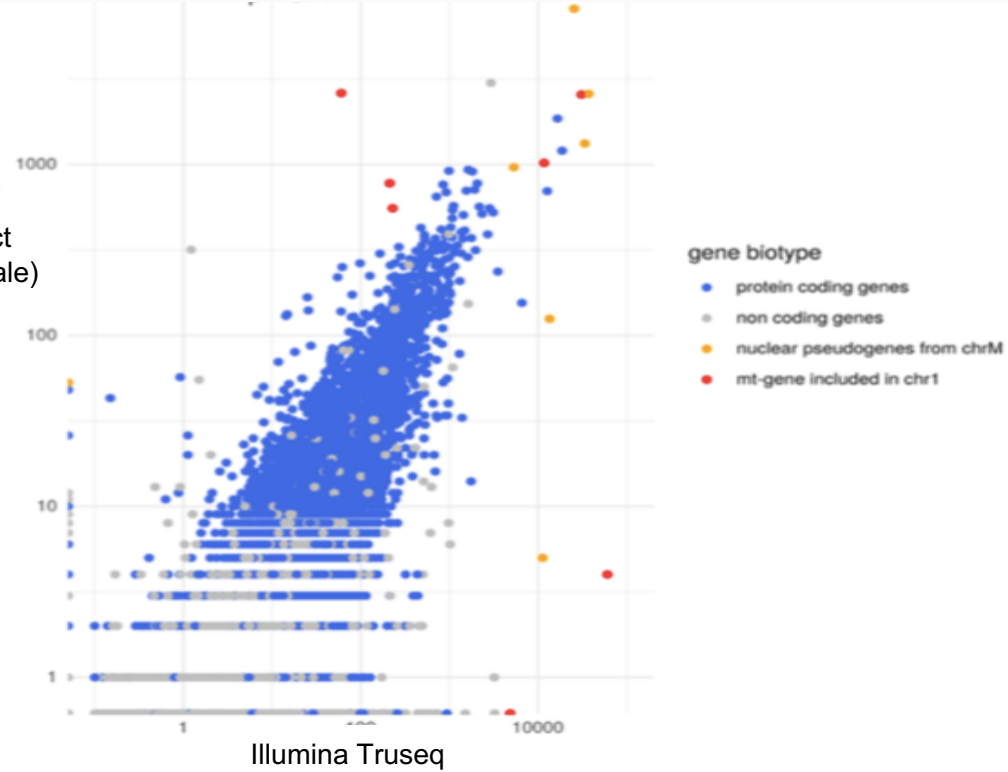




# DIRECT RNA SEQUENCING vs ILLUMINA



Nanopore RNA direct  
(read number ; Log scale)

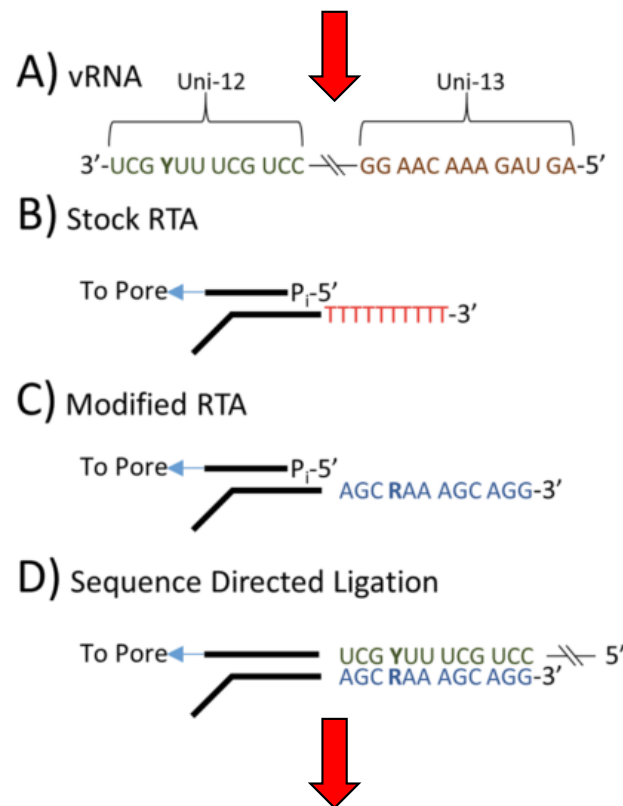


# DIRECT RNA SEQUENCING: INFLUENZA VIRUS GENOME

Direct RNA Sequencing of the complete Influenza A Virus Genome  
Keller et al. *Scientific Reports*, Sept. 2018

For the first time a complete genome of an RNA virus sequenced in its original form

Influenza A viruses are negative-sense segmented RNA viruses (8 segments)



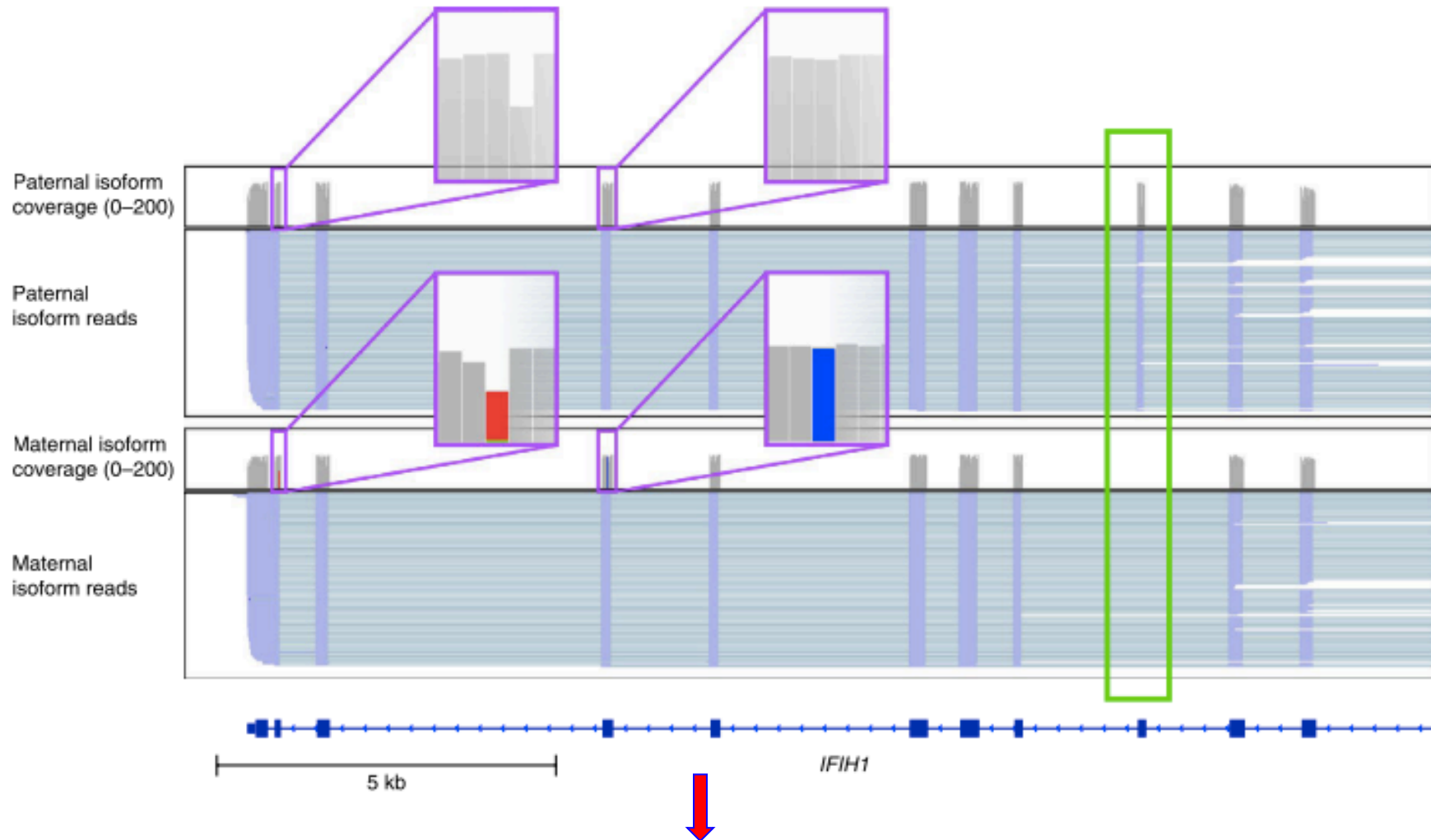
sequencing of complete genome with 100% nucleotide coverage, 99% consensus identity

Potential to identify and quantify splice variants, base modifications  
not practically measurable with current methods

# DIRECT RNA SEQUENCING: TRANSCRIPT HAPLOTYPE

## Nanopore native RNA sequencing of a human transcriptome

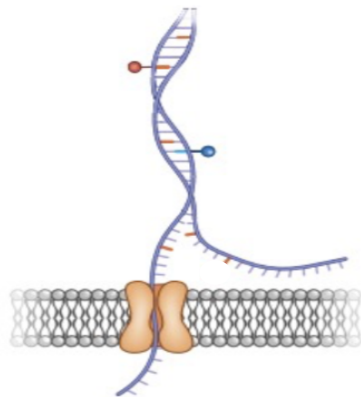
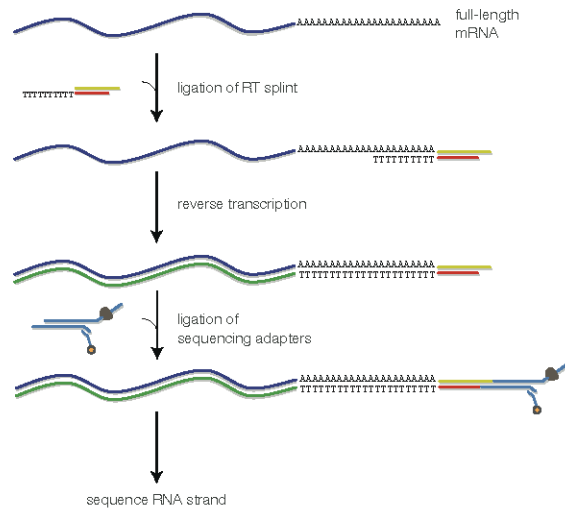
d



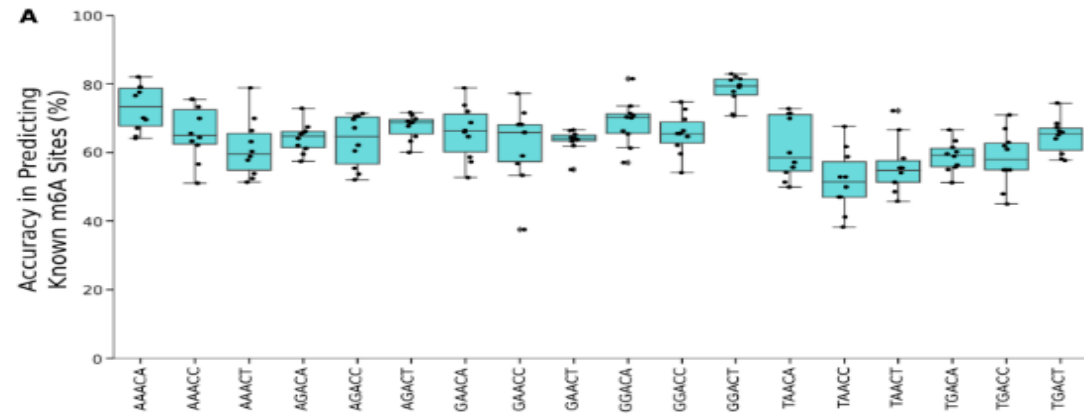
34 genes with discordant allele specificity in two isoforms

# DIRECT RNA SEQUENCING: DETECTION OF m6A

## Library preparation



Lorenz et al. *RNA* 2019



Detection of m6A with Nanopolish :

Different detection efficiency in different sites: 45% to 82%

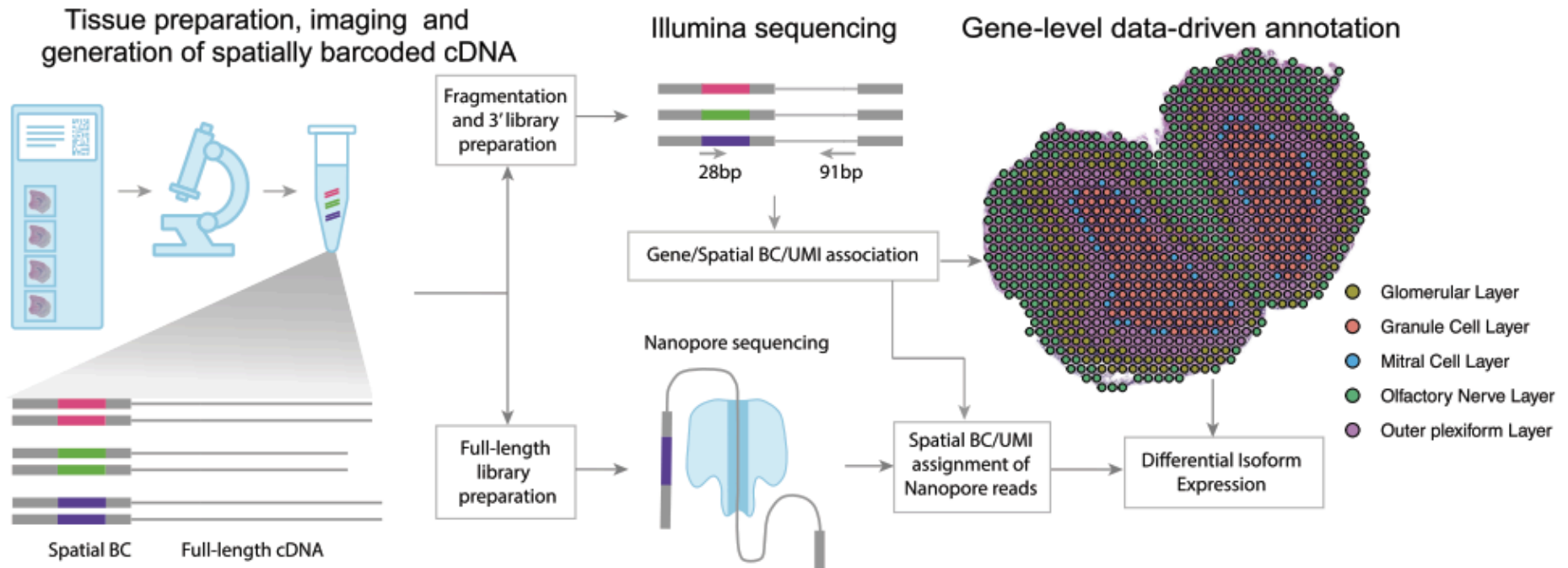
Context dependent detection efficiency

# Recent advances : Nanopore and 10x Genomics Visium

The spatial landscape of gene expression isoforms in tissue sections  
Lebrigand et al., *bioRxiv*, 2020

Spatial Isoform Transcriptomics (SiT) : Genome-wide approach to explore and discover in a tissue context :

- Isoform expression (bi-allelic expression)
- Sequence heterogeneity (SNP expression)



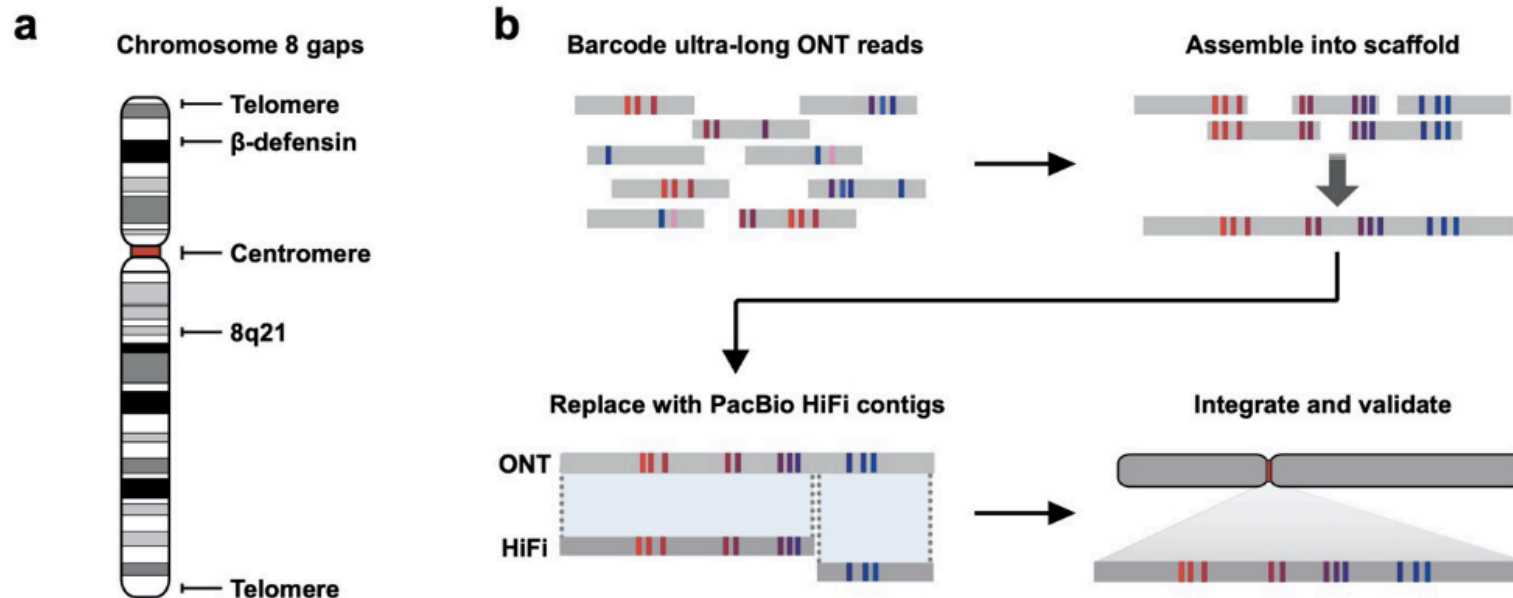
# The structure, function, and evolution of a complete human chromosome 8

Logsdon et al., *bioRxiv*, Sept 2020

First complete linear assembly of a human autosomal chromosome.

It resolves the sequence of five previously long-standing gaps :

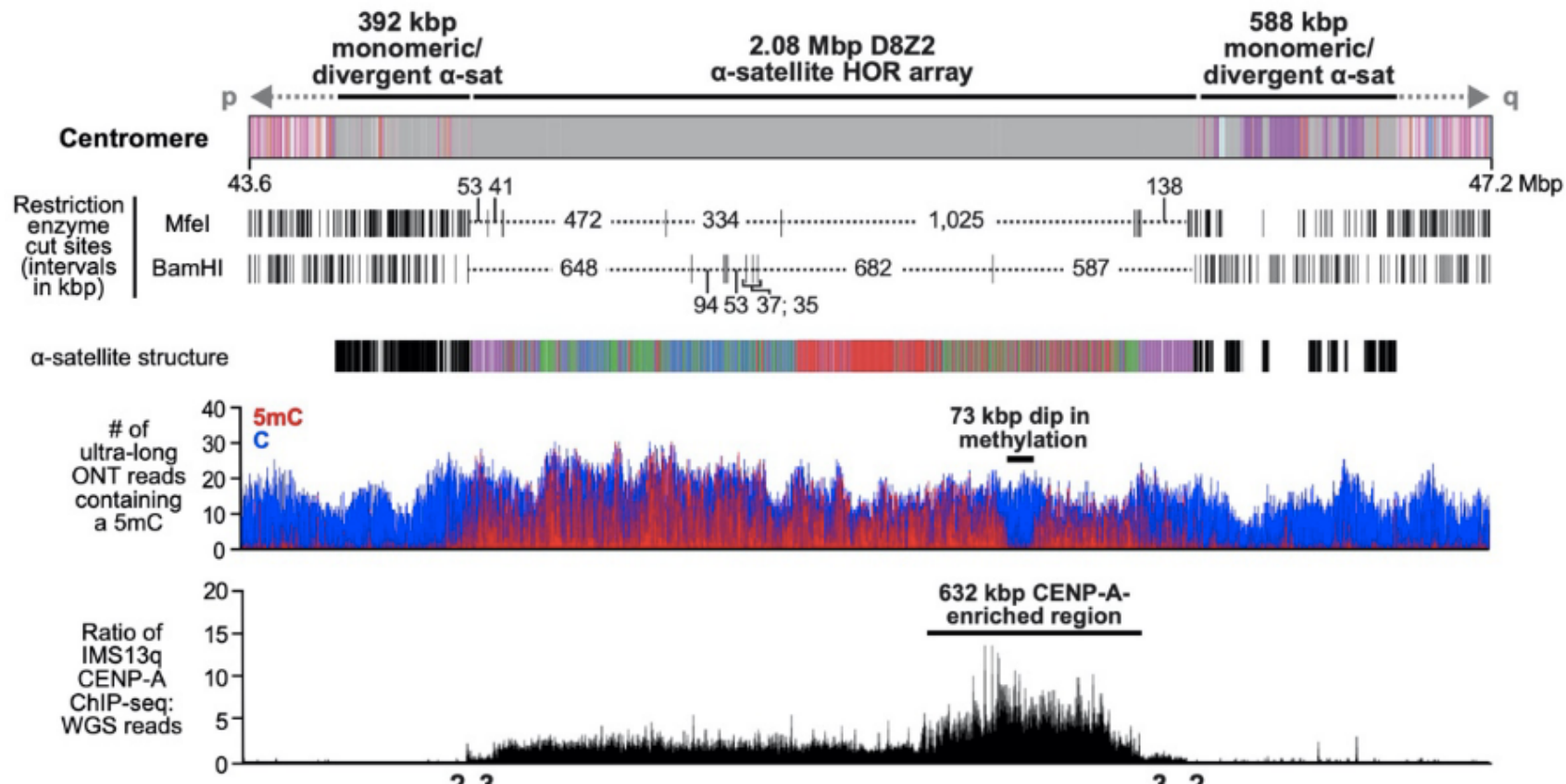
- 2.08 Mbp centromeric  $\alpha$ -satellite array
  - 644 kbp defensin copy number polymorphism
  - 863 kbp variable number tandem repeat at chromosome 8q21.2 (neocentromere)
  - Etc..
- 
- Barcoded **Ultra-long Nanopore reads** assembled into a scaffold
  - Regions within the scaffold with high sequence identity with **PacBio HiFi** contigs are replaced, thereby improving the base accuracy to >99.99%.



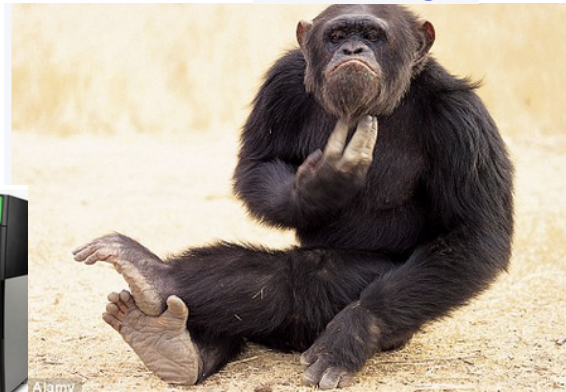
# The structure, function, and evolution of a complete human chromosome 8

Logsdon et al., *bioRxiv*, Sept 2020

## Epigenetic map of the chromosome 8 centromeric region



???



PacBio



Nanopore



Illumina

