

*u*<sup>b</sup>

---

b  
UNIVERSITÄT  
BERN

# *De novo* Genome & Transcriptome Assembly Course

## Lecture 2 – Sequencing Technologies

Rémy Bruggmann – Interfaculty Bioinformatics Unit (IBU)

22.09.21



# Content

- Library preparation principles
- Sequencing technologies
  - Sanger
  - NGS
  - third generation (PacBio & Nanopore)
- Sources of Errors
- Hi-C method & Optical mapping
- Is perfect assembly possible?

## New sequencing methods required

### Requirements | Wish List

- much higher throughput
- much lower costs
- parallel library preparation,  
(i.e., no "cloning" of fragments)
- parallel sequencing
- error free (at least no systematic errors)
- ultra long reads



# Next Generation Sequencing (NGS) Instruments

## Thermo Fisher (LifeTechnologies)

- IonTorrent (PGM)
- IonProton
- Ion GeneStudio S5
- Ion Torrent Genexus



## Illumina

- MiSeq
- NextSeq500
- NextSeq550
- HiSeq2500
- HiSeq X ten
- NextSeq550
- HiSeq3000
- HiSeq4000
- HiSeq X five
- MiniSeq
- NextSeq1000/2000
- NovaSeq6000



# 3<sup>rd</sup> Generation Sequencing Instruments

Pacific Biosciences (PacBio)



PacBio RS II

PacBio Sequel I & II



Oxford Nanopore Technology (ONT)

MinION



MinION Mk1C



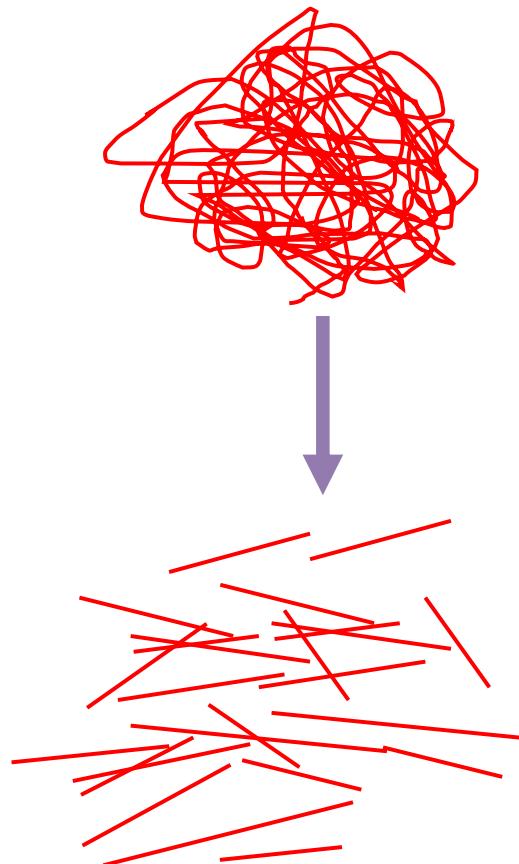
GridION



PromethION

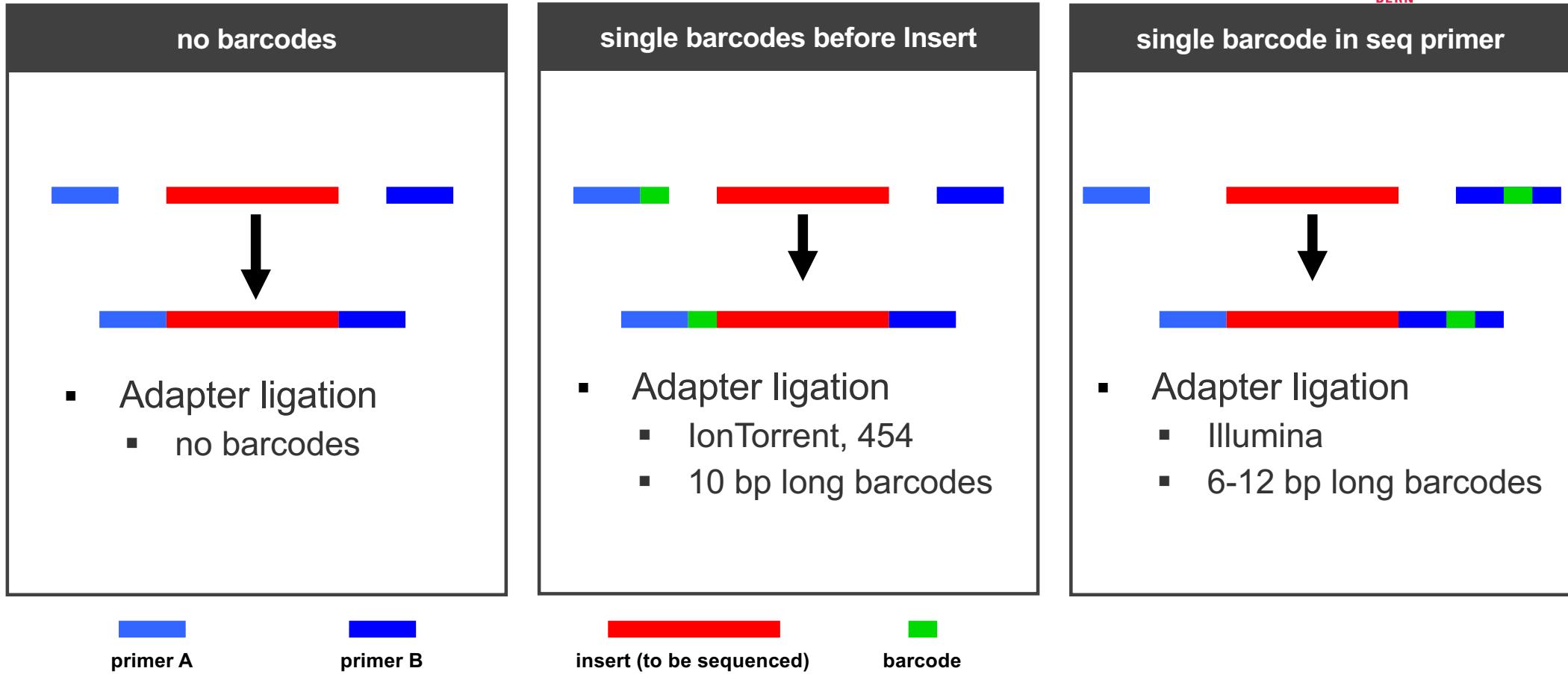


# Library Preparation – (Short Reads)

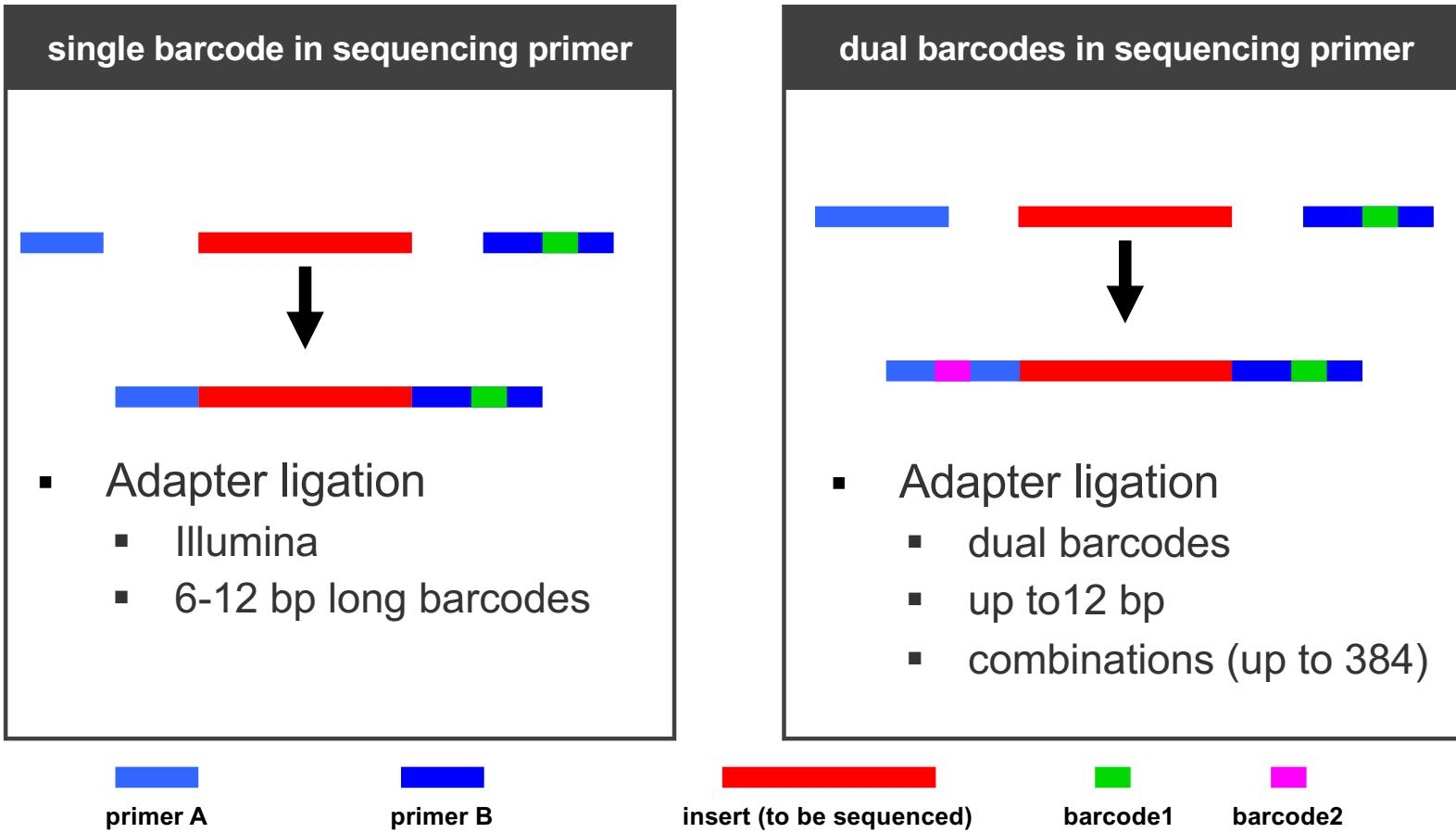


- Input material (DNA)
  - Input material must be broken into smaller pieces suitable for sequencing (and PCR)
  - Not required for amplicon seq
- Fragmentation methods
  - Mechanical
  - Enzymatic
  - Chemical
- Size of fragments follows Poisson distribution

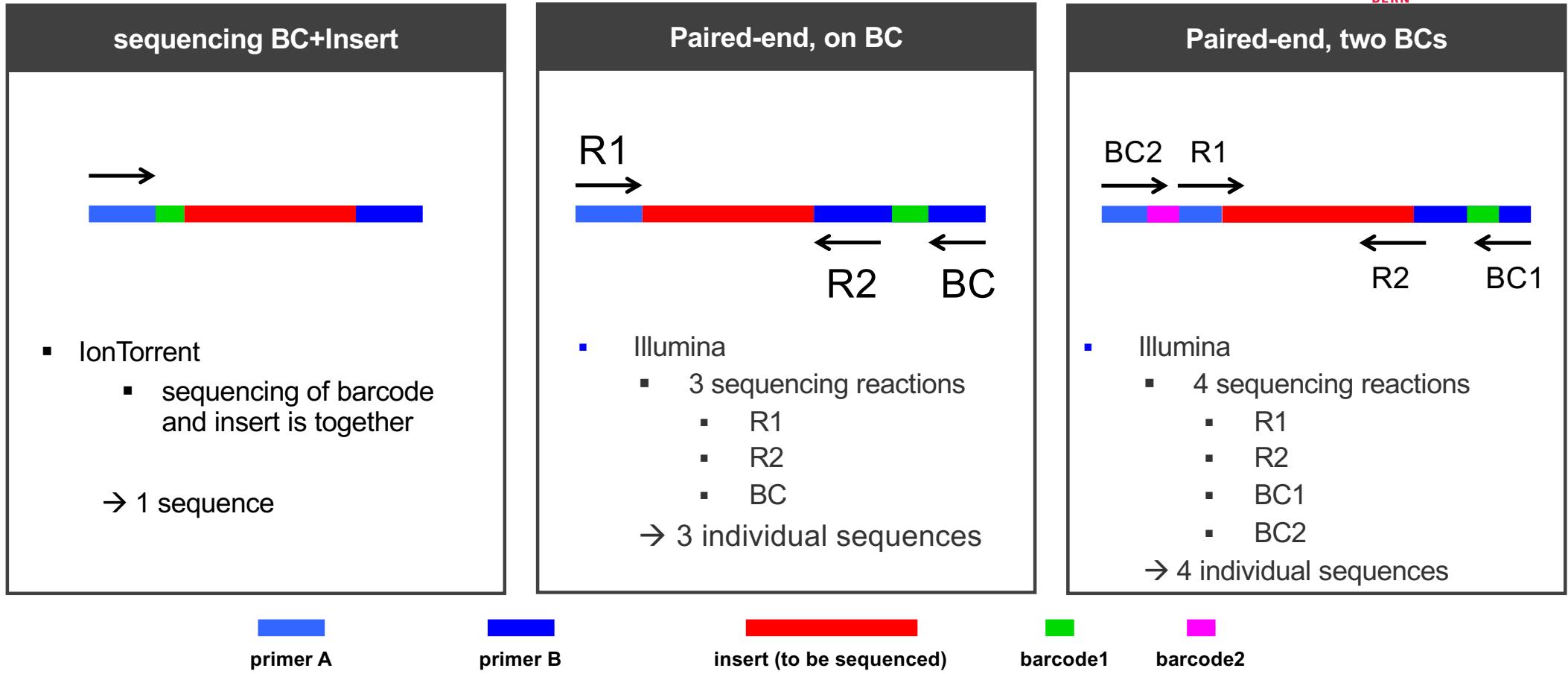
# Library Preparation



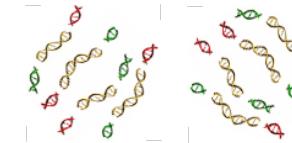
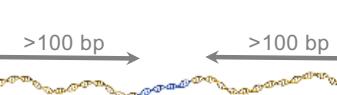
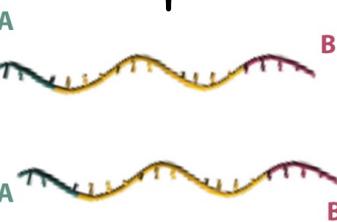
# Library Preparation | Barcoding



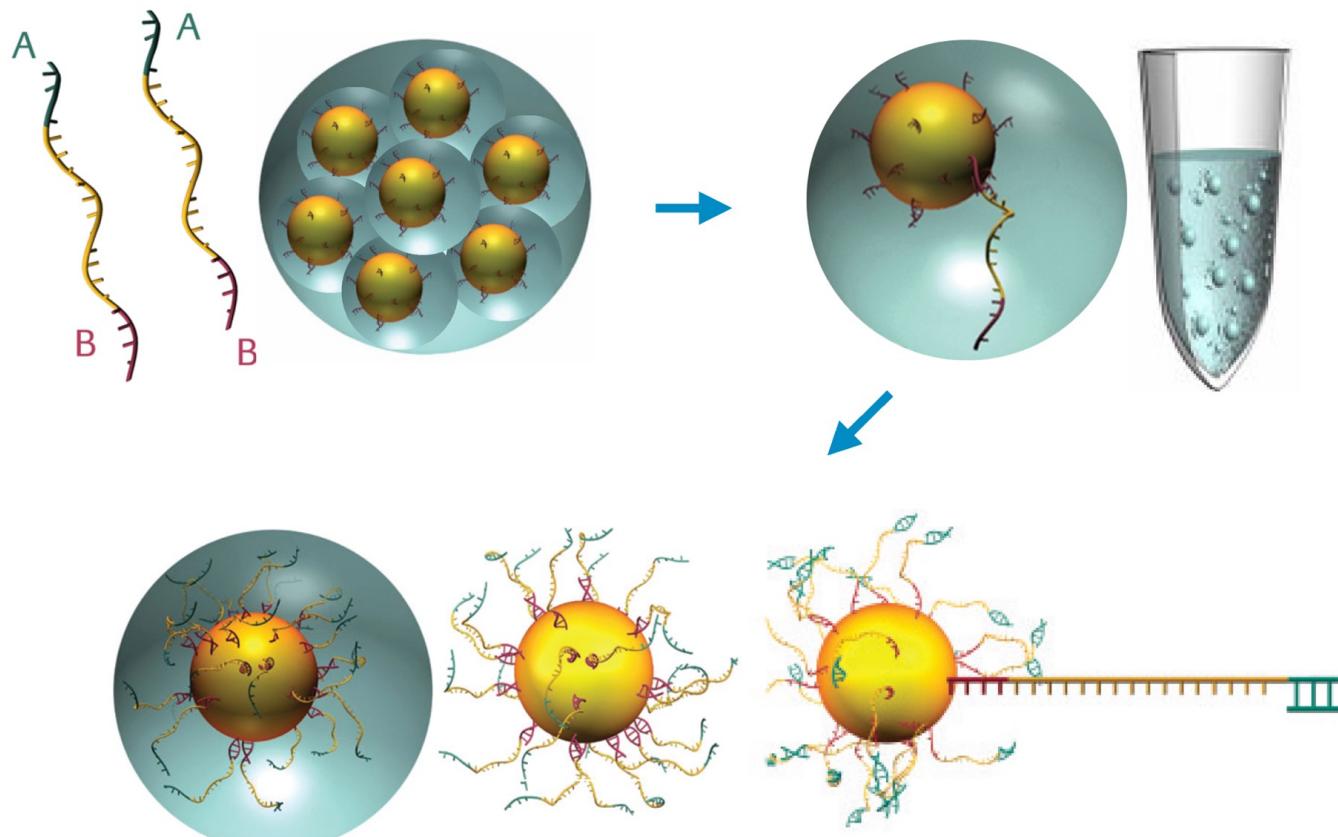
# Sequencing of Libraries



# NGS Library Types

| Shotgun  | Blunt end   | Mate pair  | Amplicons   |
|--|---|--|---|
| <ul style="list-style-type: none"><li>▪ whole genome</li><li>▪ full length transcripts</li><li>▪ BACs</li><li>▪ long range PCR</li></ul> | <ul style="list-style-type: none"><li>▪ ncRNA (e.g. miRNAs)</li><li>▪ ancient DNA</li><li>▪ short DNA fragments</li><li>▪ short PCR fragments</li></ul> | <ul style="list-style-type: none"><li>▪ de novo assembly</li><li>▪ structural variants</li></ul> | <ul style="list-style-type: none"><li>▪ short targeted region</li><li>▪ cancer gene panel</li><li>▪ HIV typing</li><li>▪ metagenomics</li></ul> |
|   |    |               |    |
|    |   |  |   |
| <ul style="list-style-type: none"><li>▪ ssDNA Library</li><li>▪ 20-800 bp insert size</li><li>▪ ready for sequencing</li></ul>           |   |  |   |

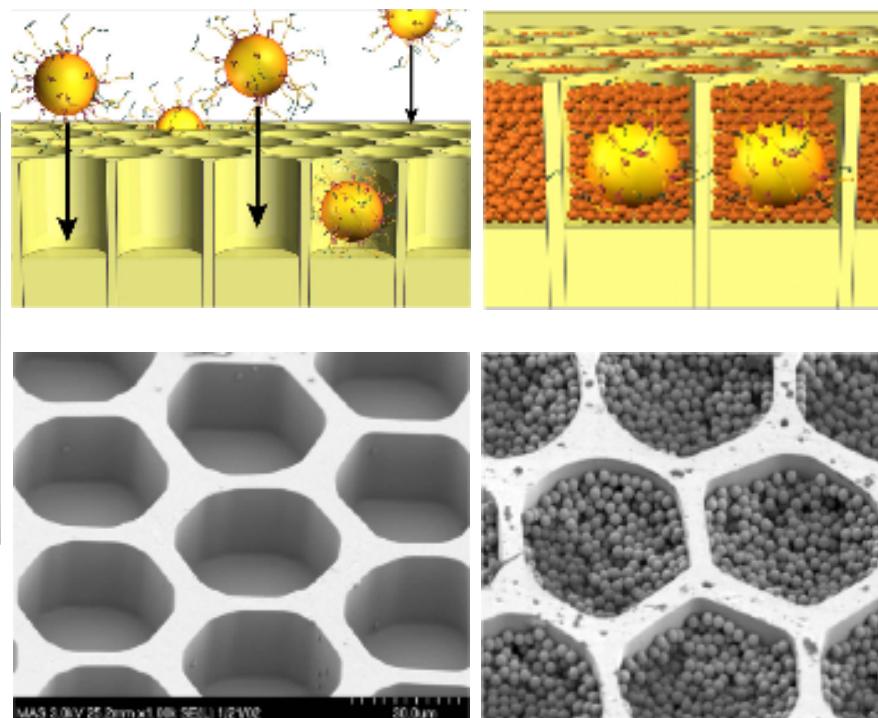
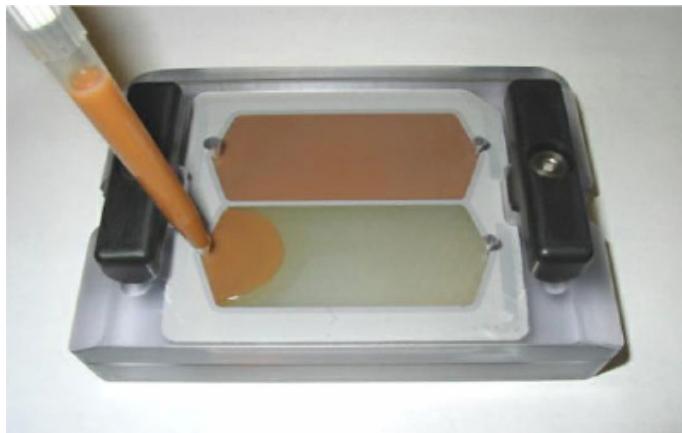
## Preamplification of Libraries



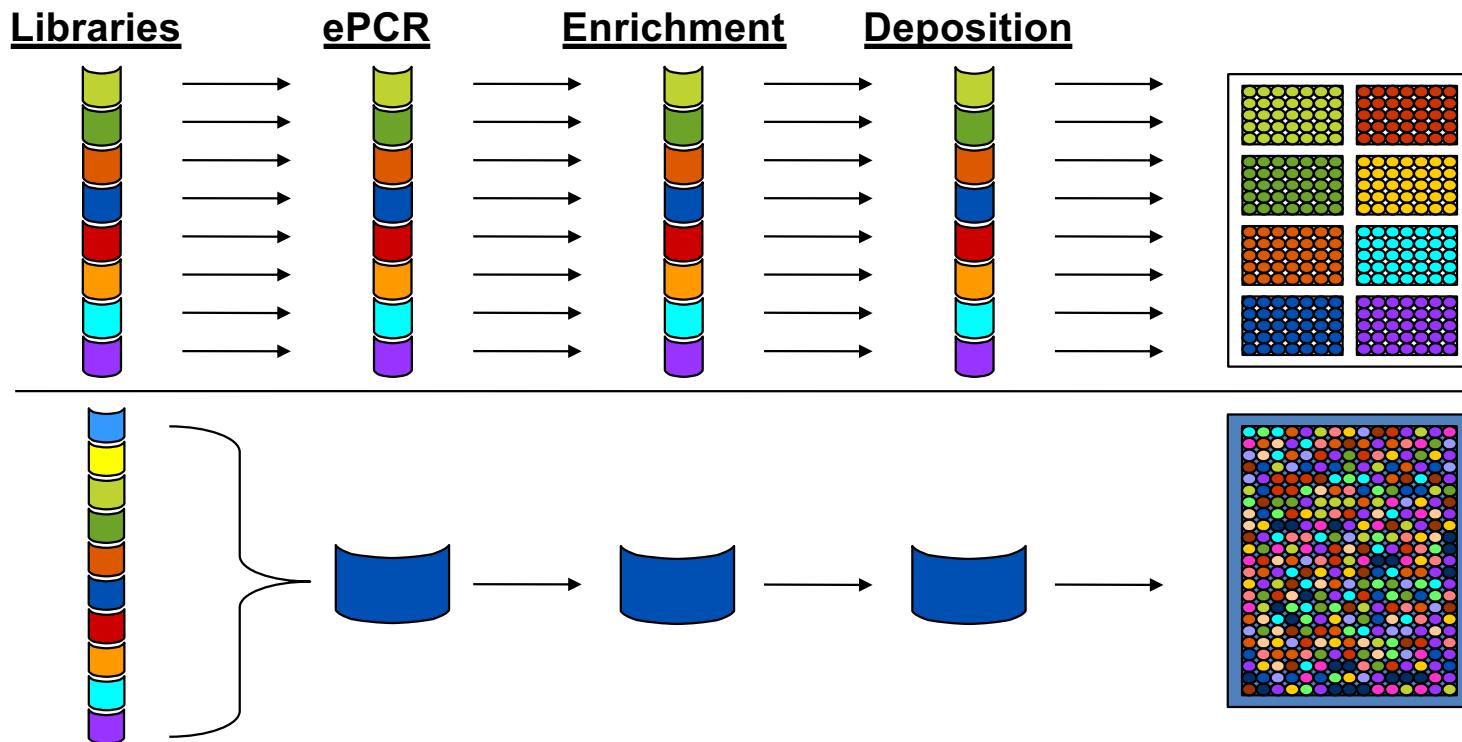
*u*<sup>b</sup>

*b*  
UNIVERSITÄT  
BERN

## Picotiter Plate



## Multiplexing – Barcoding

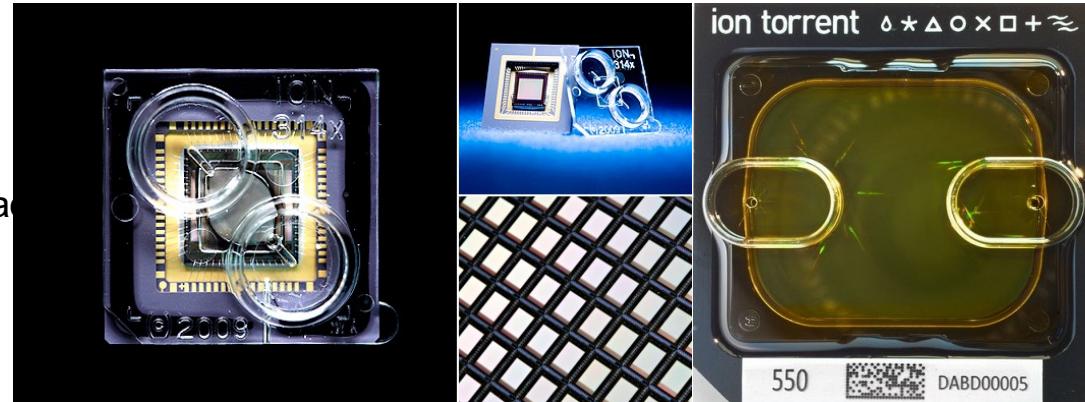


# Ion Torrent Sequencing | Technology

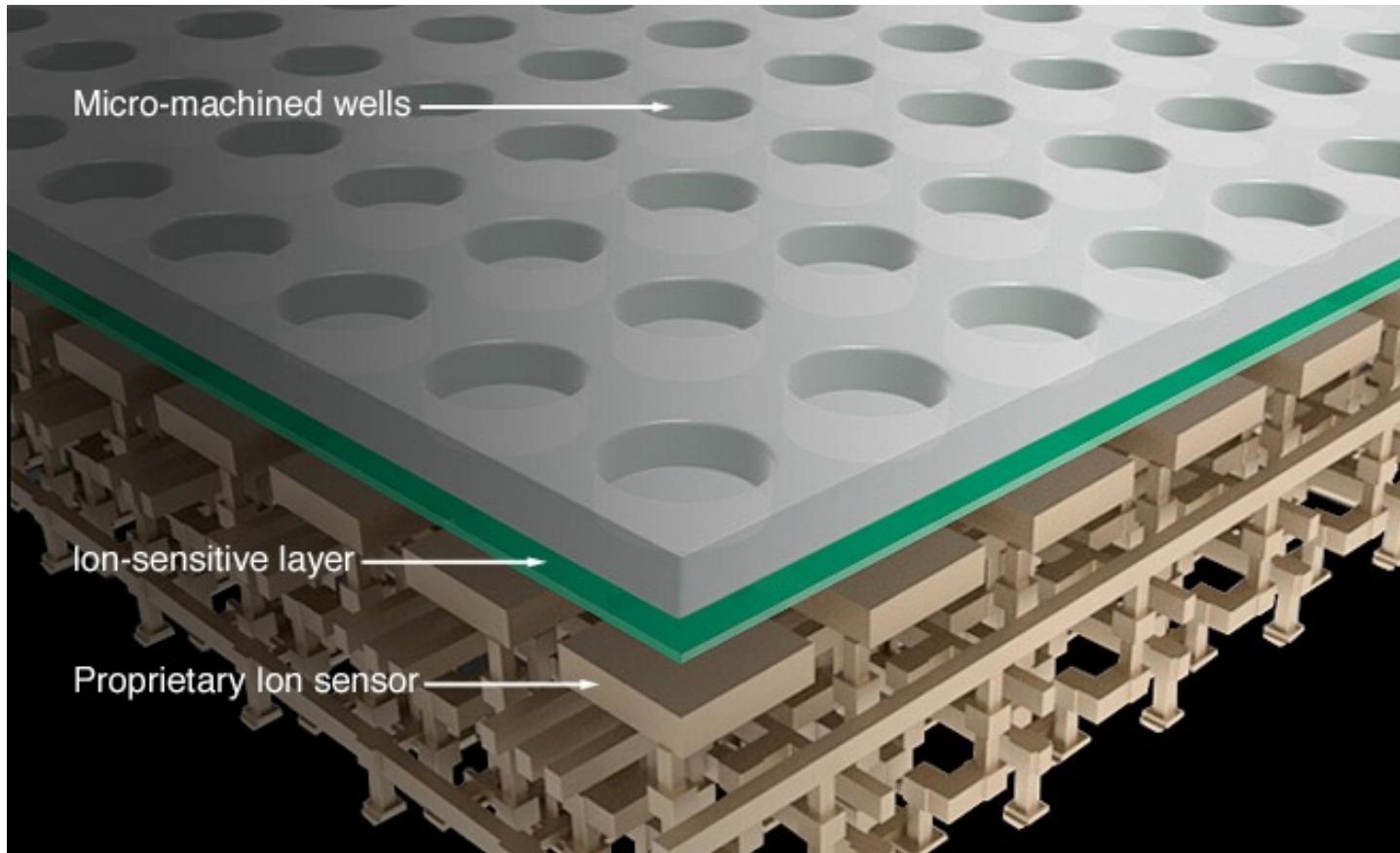
- Short read to long read technology (low to high throughput)
- Clonal amplification by emulsion PCR (ePCR)
- Sequencing by synthesis (SBS) – pyrosequencing
- Label free!

# IonTorrent | Overview

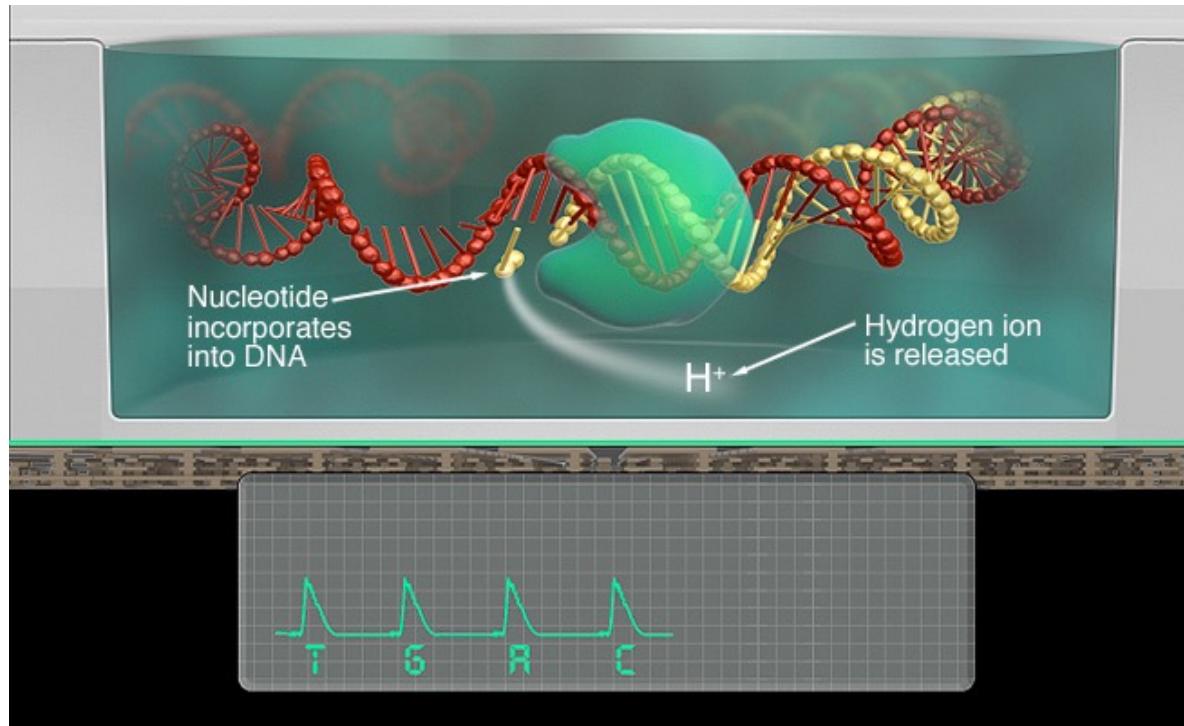
- Ion semiconductor sequencing chips
- The simplest sequencing chemistry
  - natural nucleotides, no enzymatic cascade
- True semiconductor sequencing
- No scanners, no cameras, no light
  - The chip is the machine.
- ***emulsion-PCR still required!***
- Instrument costs:  
1/5 of traditional sequencers (\$125'000)



# World's Smallest pH Meter

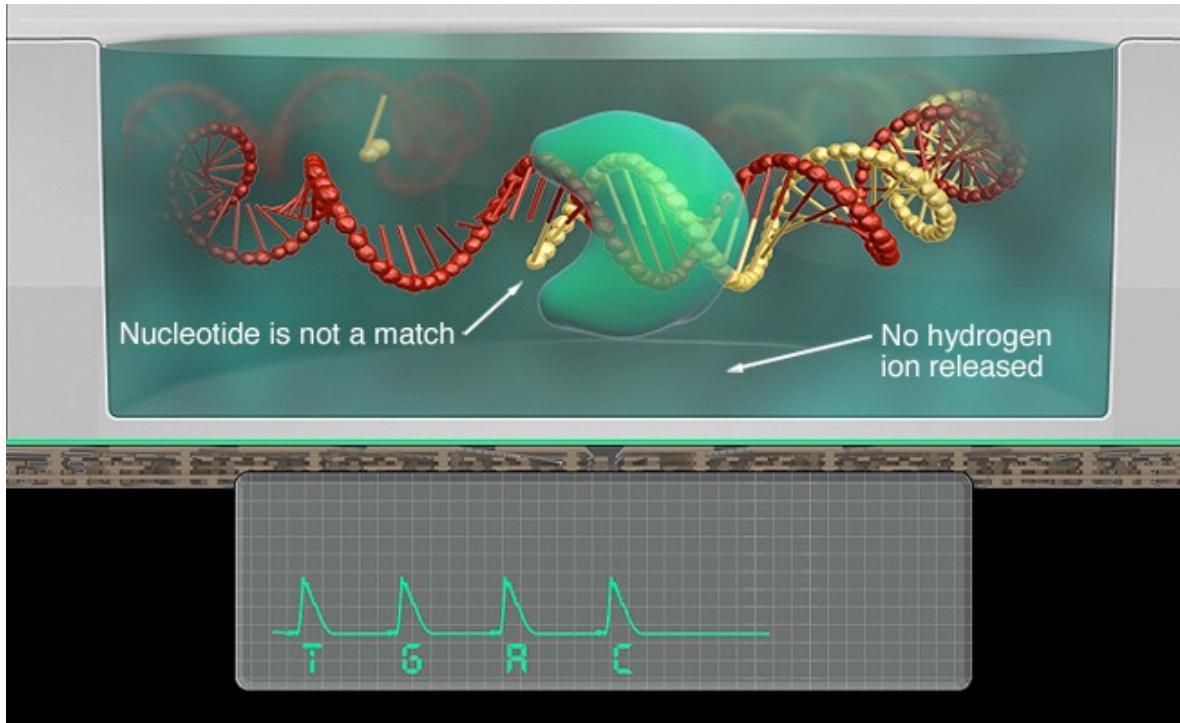


# World's Smallest pH Meter

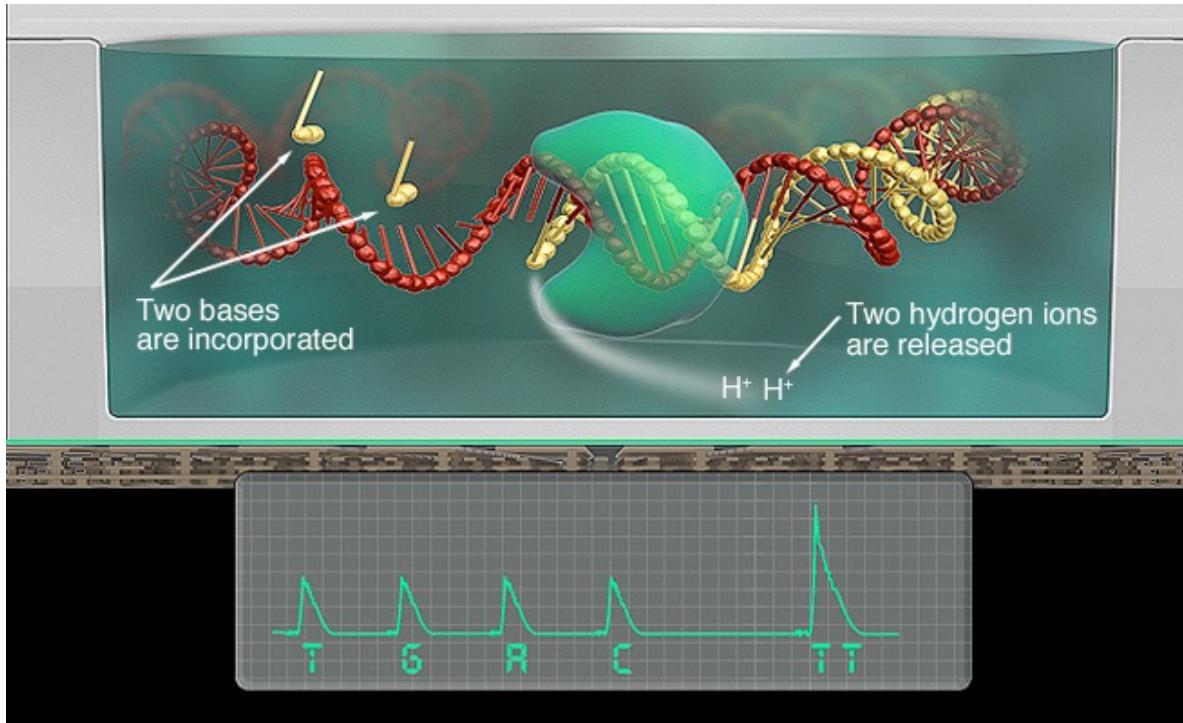


Chemical information is directly transferred into digital information

# World's Smallest pH Meter



# World's Smallest pH Meter



Homopolymers are problematic and the major source of errors

*u*<sup>b</sup>

---

*b*  
**UNIVERSITÄT  
BERN**





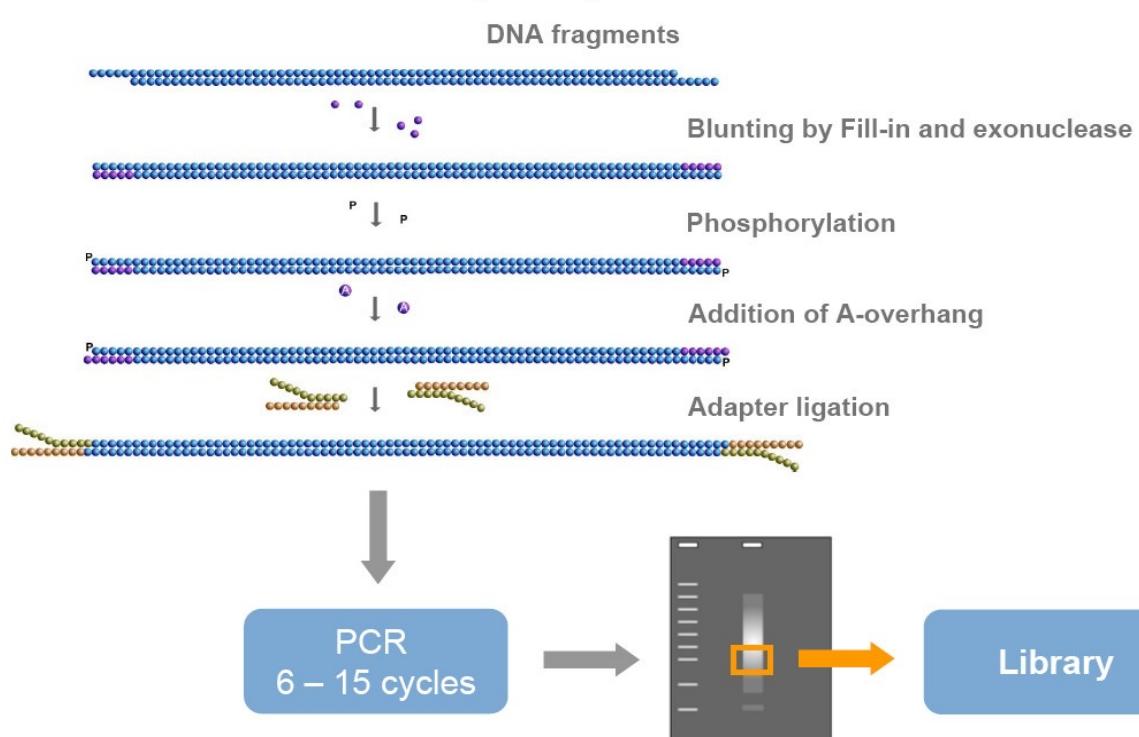


# Illumina Sequencing | Technology

- Short read sequencing technology | low to medium to extremely high throughput
- Clonal amplification by cluster/bridge PCR
- Sequencing by synthesis (SBS)– one base at a time

# Illumina Sequencing | Technology

## Genomic DNA Library Prep



PCR-free Genomic DNA Libraries are now standard for high quality DNA

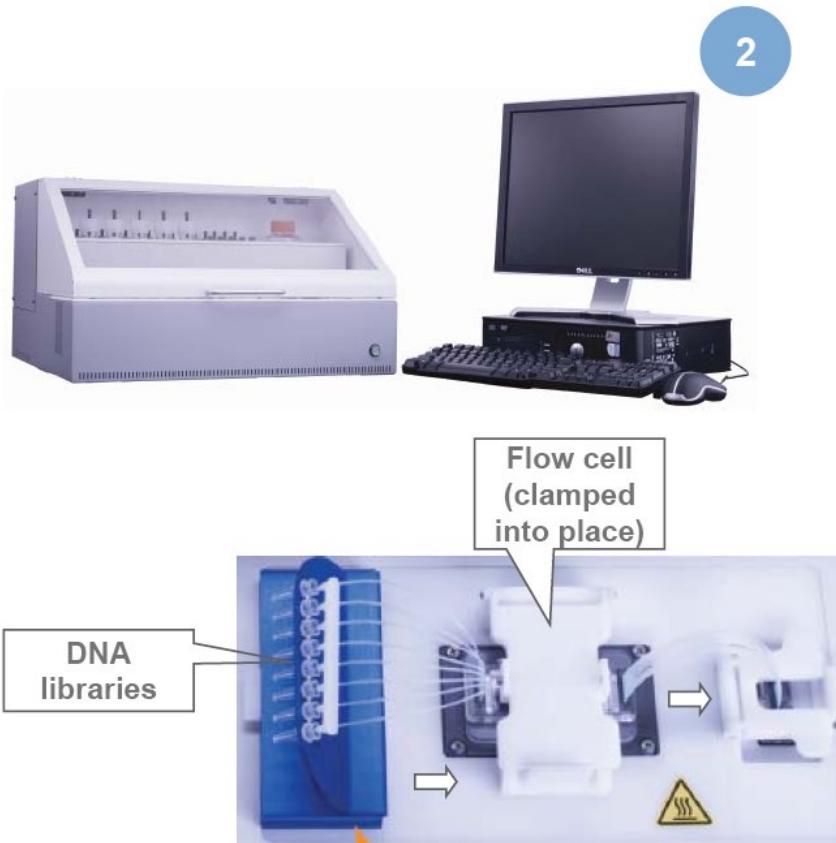
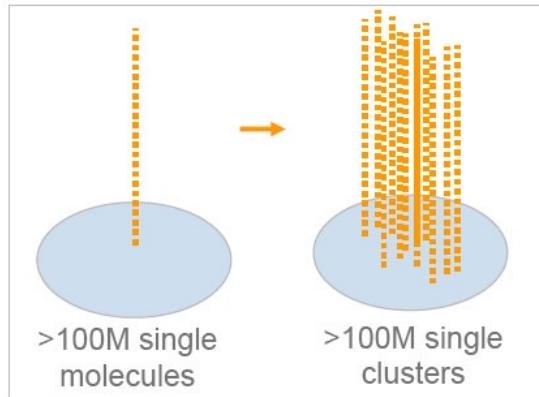


# Illumina Sequencing | Technology

## Cluster Generation

### *Cluster station*

- Aspirates DNA samples into flow cell
- Automated amplified clonal clusters

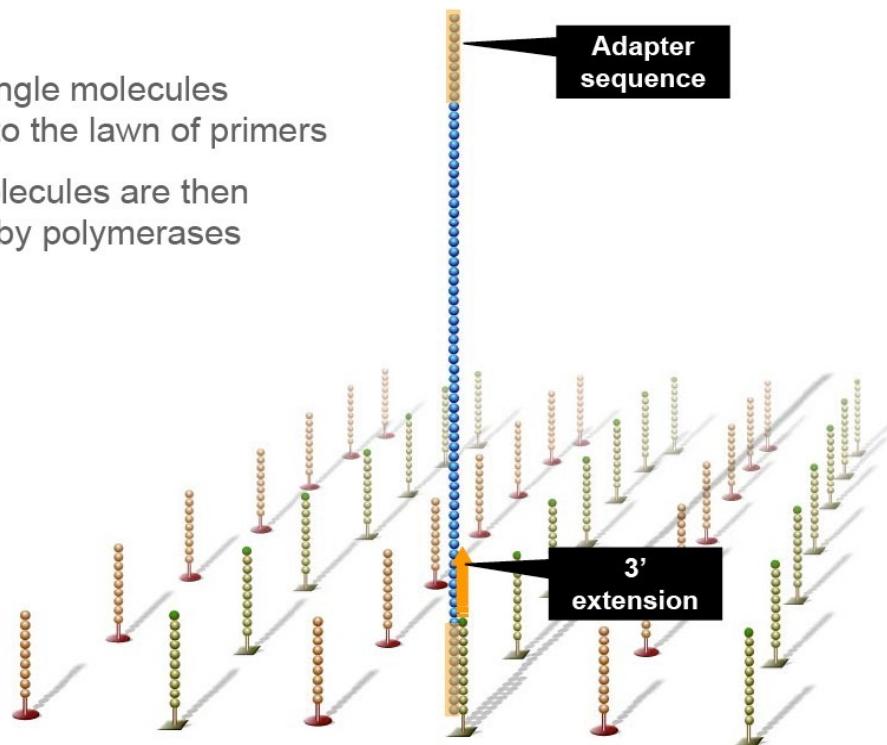


Cluster generation is integrated in the new instruments

# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation *Hybridize Fragment & Extend*

- >150 M single molecules hybridize to the lawn of primers
- Bound molecules are then extended by polymerases

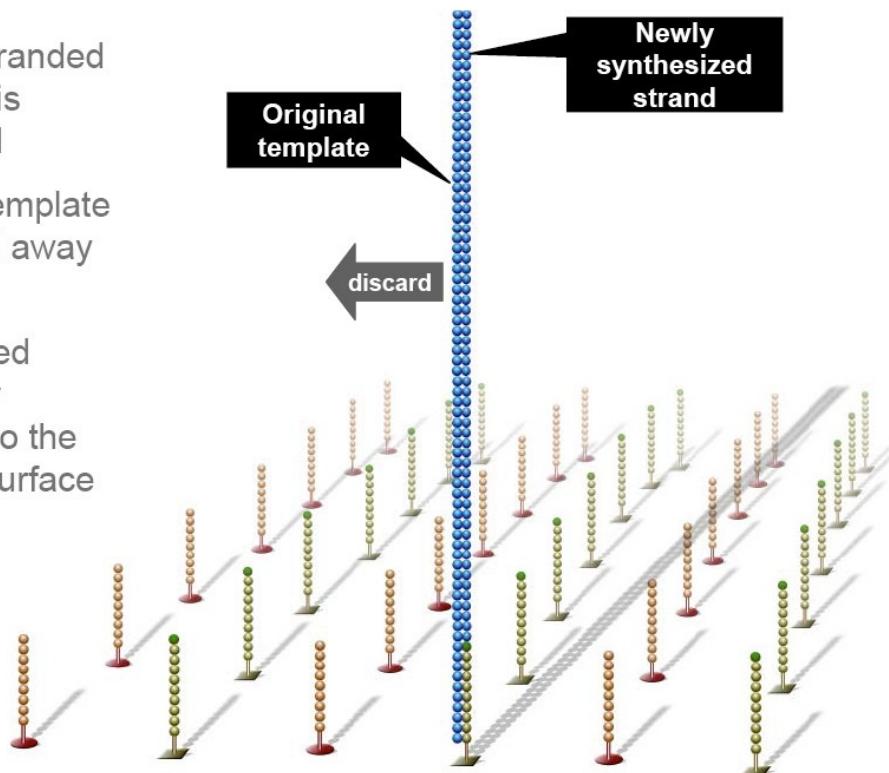


# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation

### *Denature Double-stranded DNA*

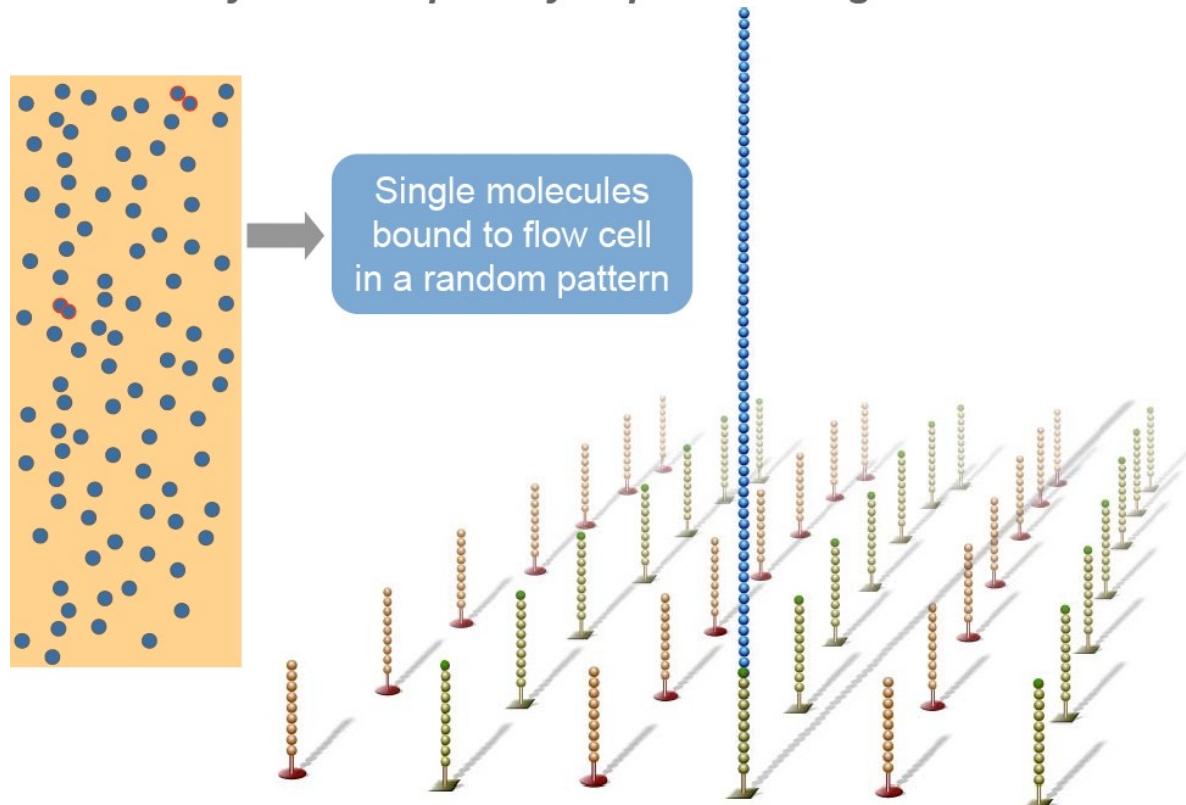
- Double-stranded molecule is denatured
- Original template is washed away
- Newly synthesized covalently attached to the flow cell surface



# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation

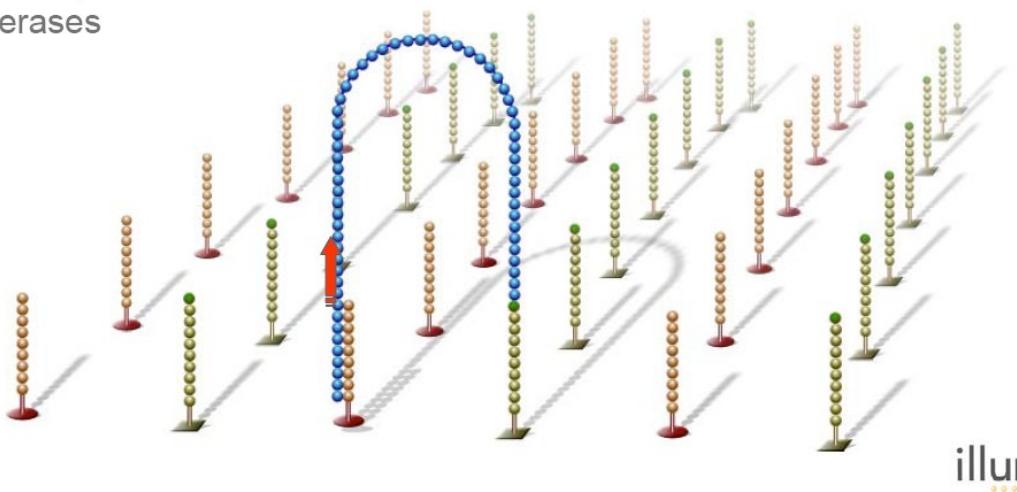
*Covalently-Bound Spatially Separated Single Molecules*



# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation *Bridge Amplification*

- Single-strand flips over to hybridize to adjacent primers to form a bridge
- Hybridized primer is extended by polymerases

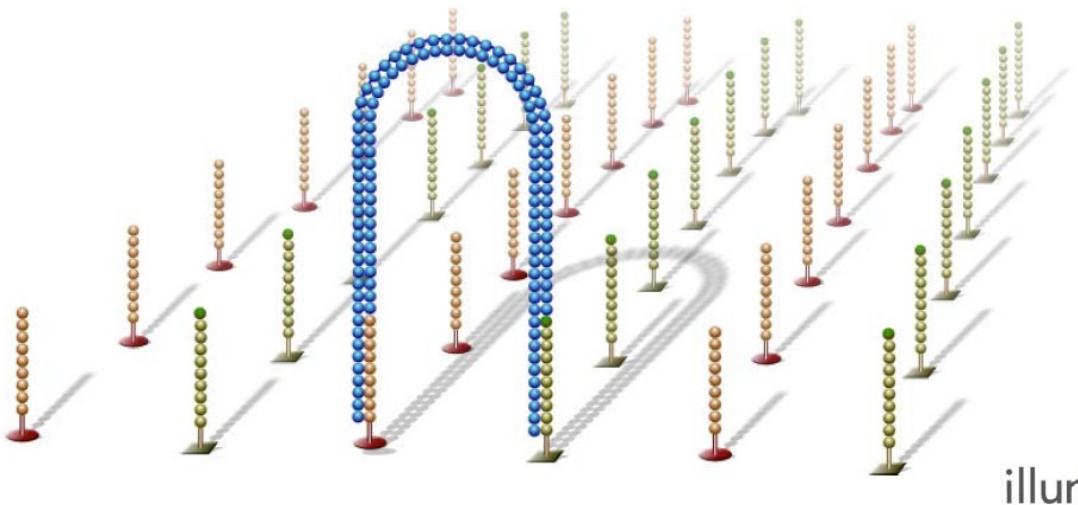


illui

# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation *Bridge Amplification*

- Double-stranded bridge is formed

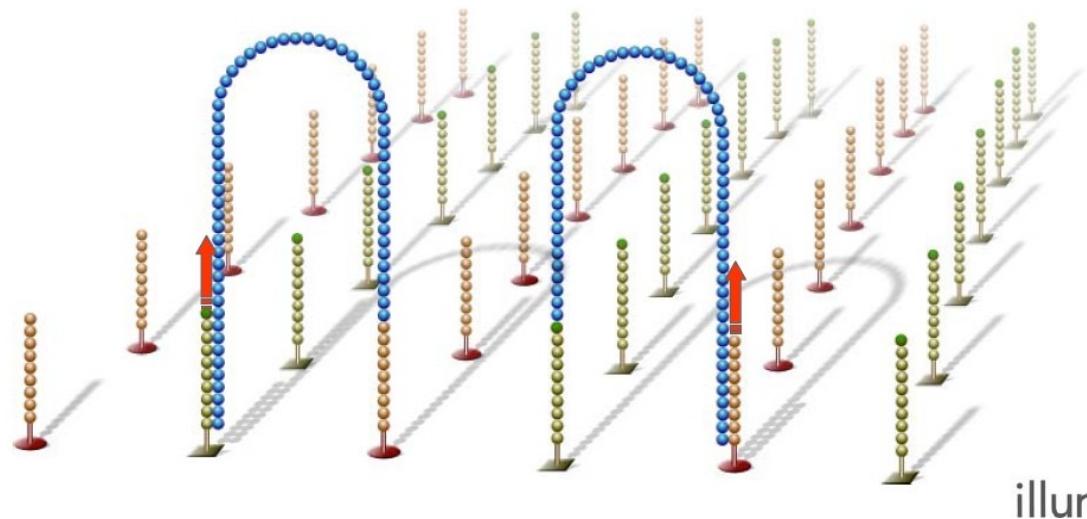


illur

# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation *Bridge Amplification*

- Single-strands flip over to hybridize to adjacent primers to form bridges
- Hybridized primer is extended by polymerase

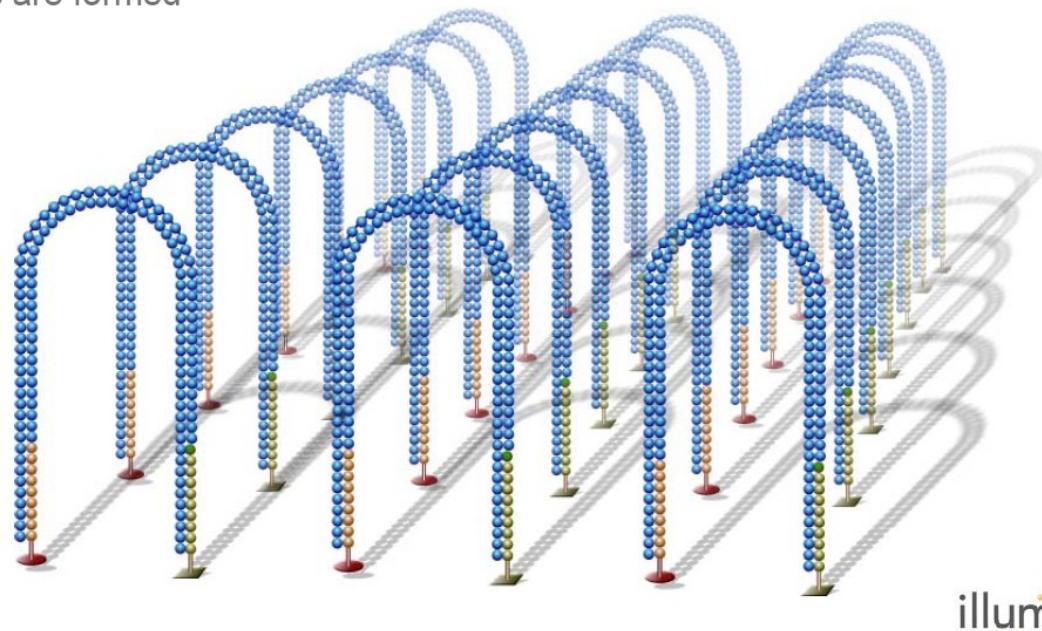


illur

# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation *Bridge Amplification*

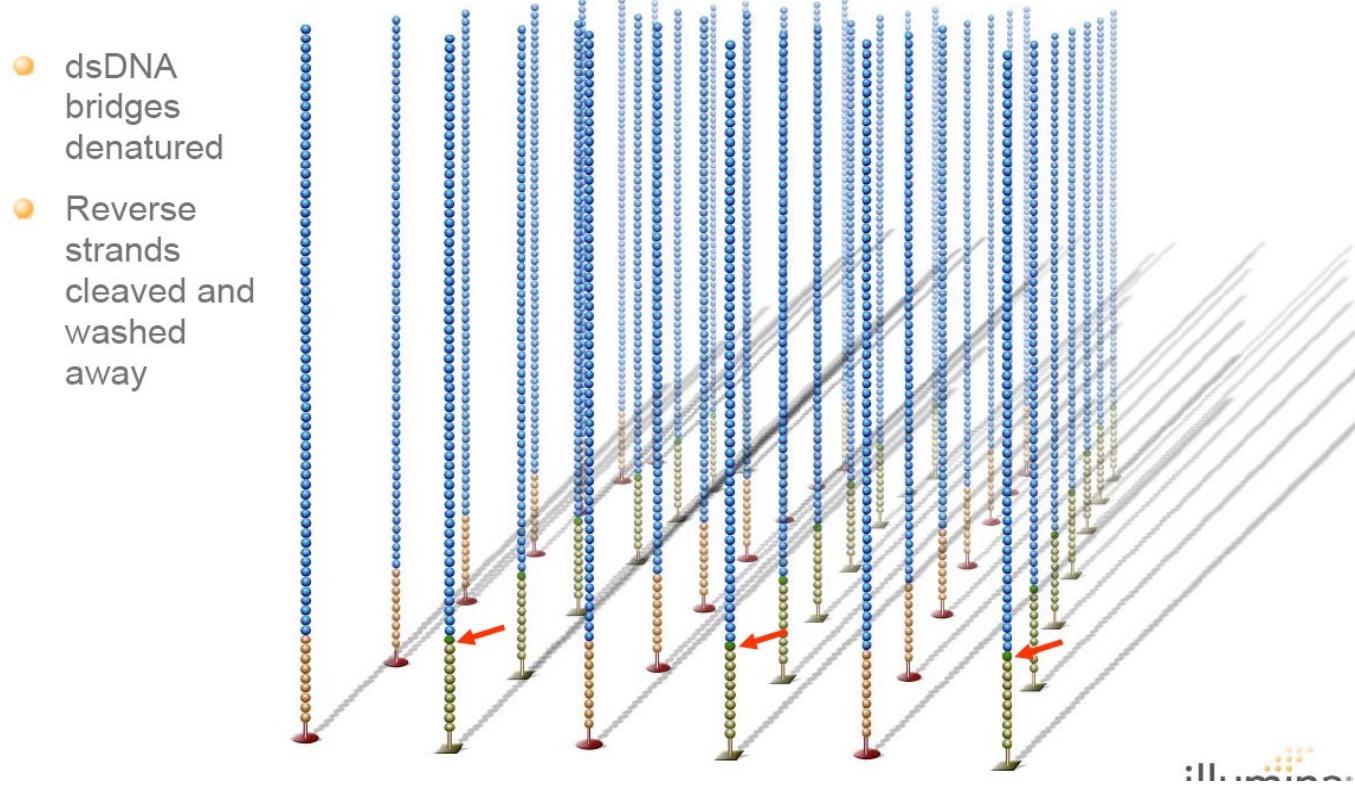
- Bridge amplification cycle repeated until multiple bridges are formed



illum

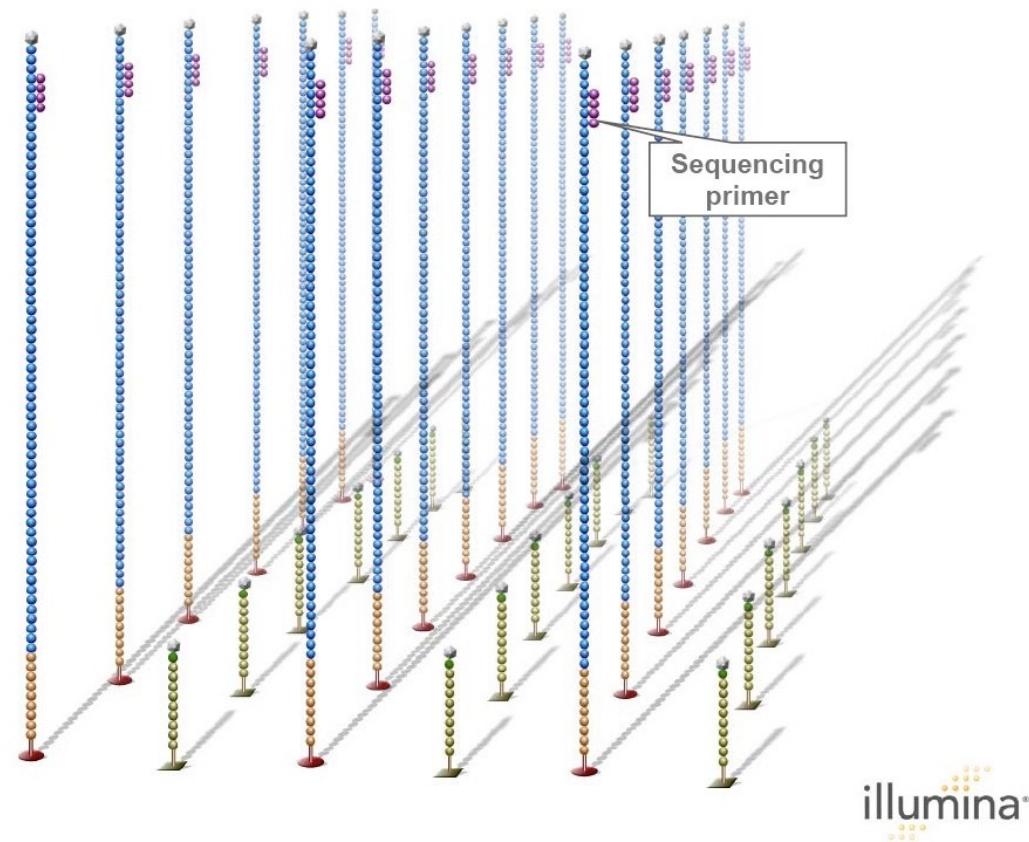
# Illumina Sequencing | Cluster Generation/PCR

## Cluster Generation

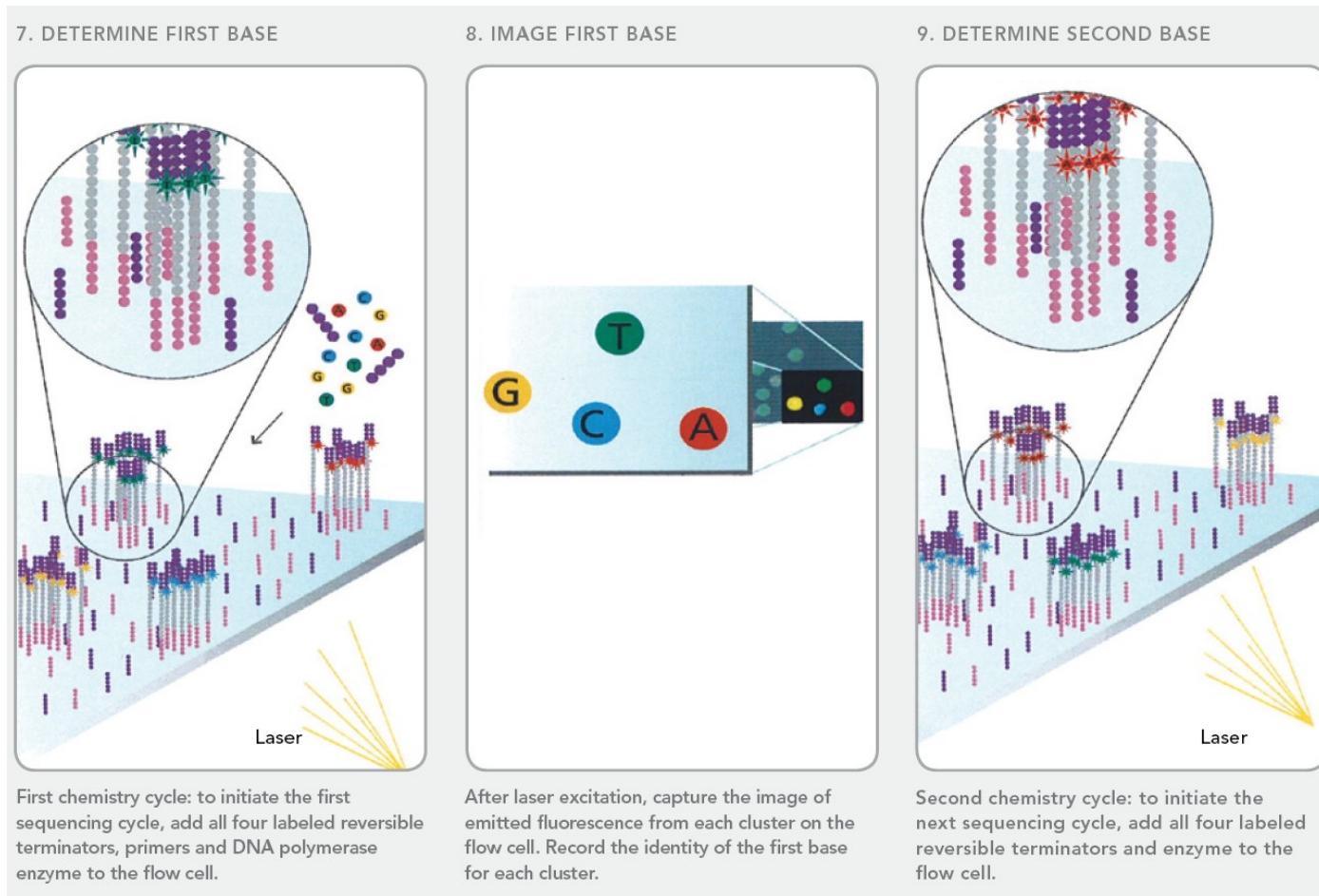


# Illumina Sequencing | Sequencing Reaction

- Sequencing primer is hybridized to adapter sequence

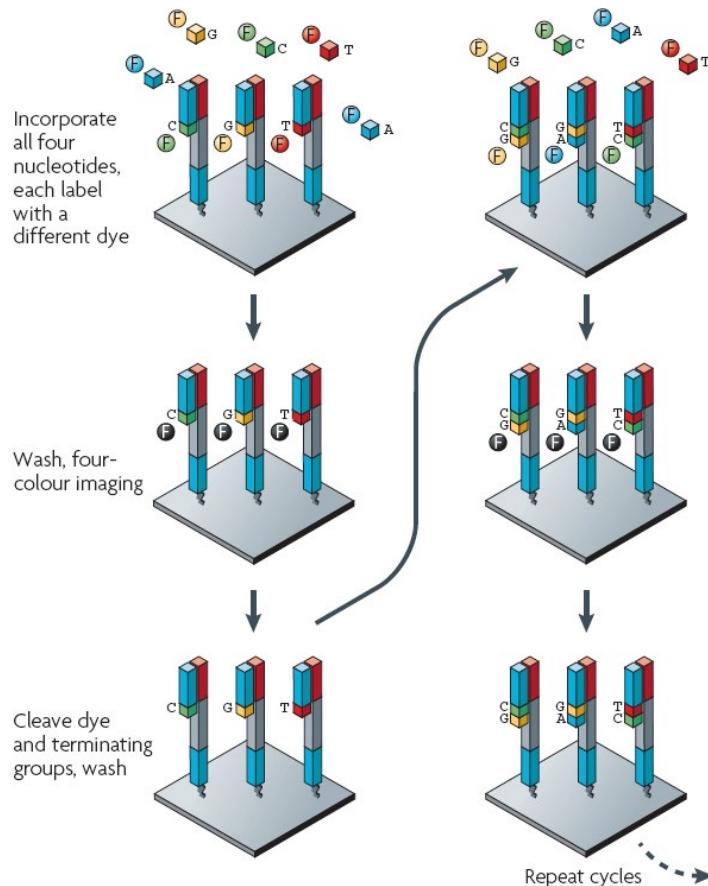


# Illumina Sequencing | Sequencing Reaction

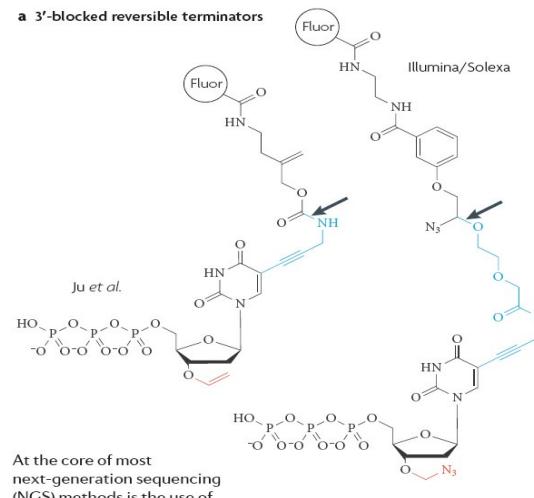


# Illumina Sequencing | Sequencing Reaction

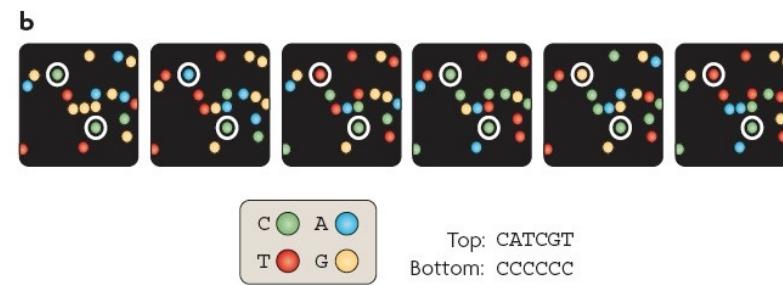
## a Illumina/Solexa — Reversible terminators



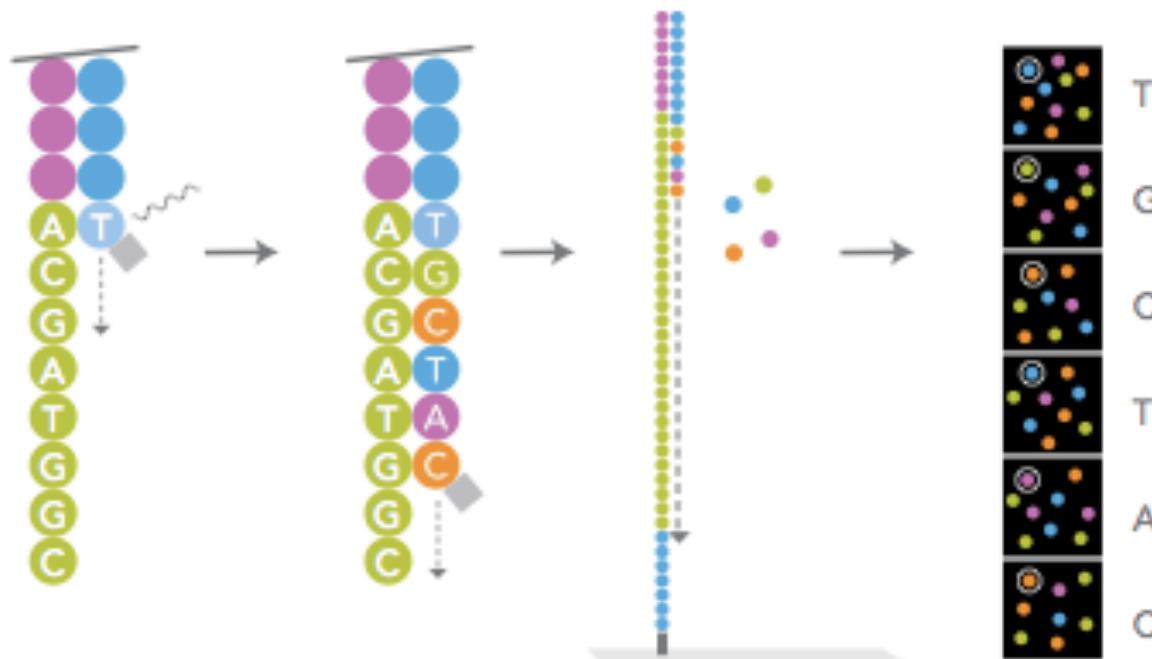
## a 3'-blocked reversible terminators



Illumina uses 3'-O-Azidomethyl as reversible blocking group



# Illumina Sequencing | Sequencing Reaction



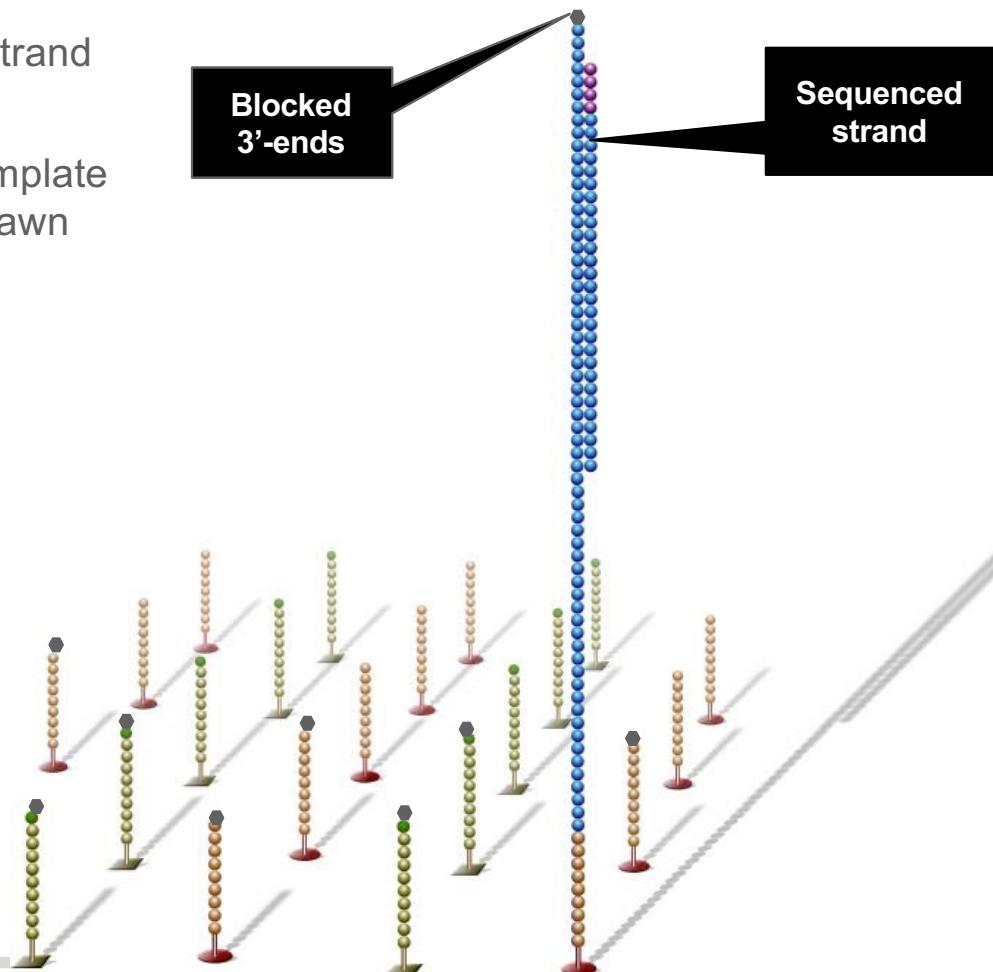
**Sequencing – one base per cycle**

# Illumina Sequencing | Paired End Sequencing

At the moment, we only have sequenced from one end of the library,  
though we want to sequence from both ends

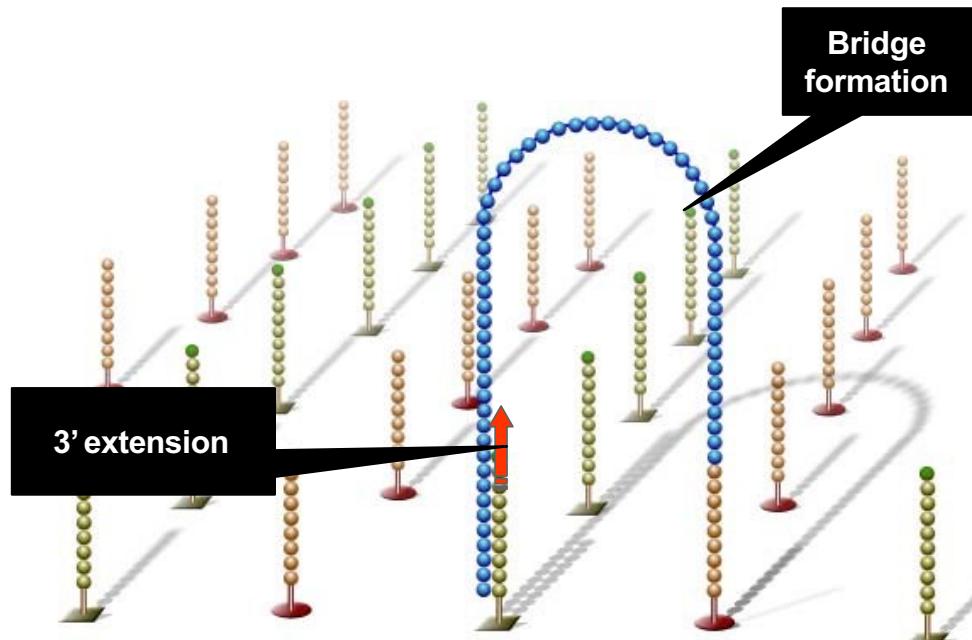
## Illumina Sequencing | Paired End Sequencing

- Sequenced strand is tripped off
- 3'-ends of template strands and lawn primers are unblocked



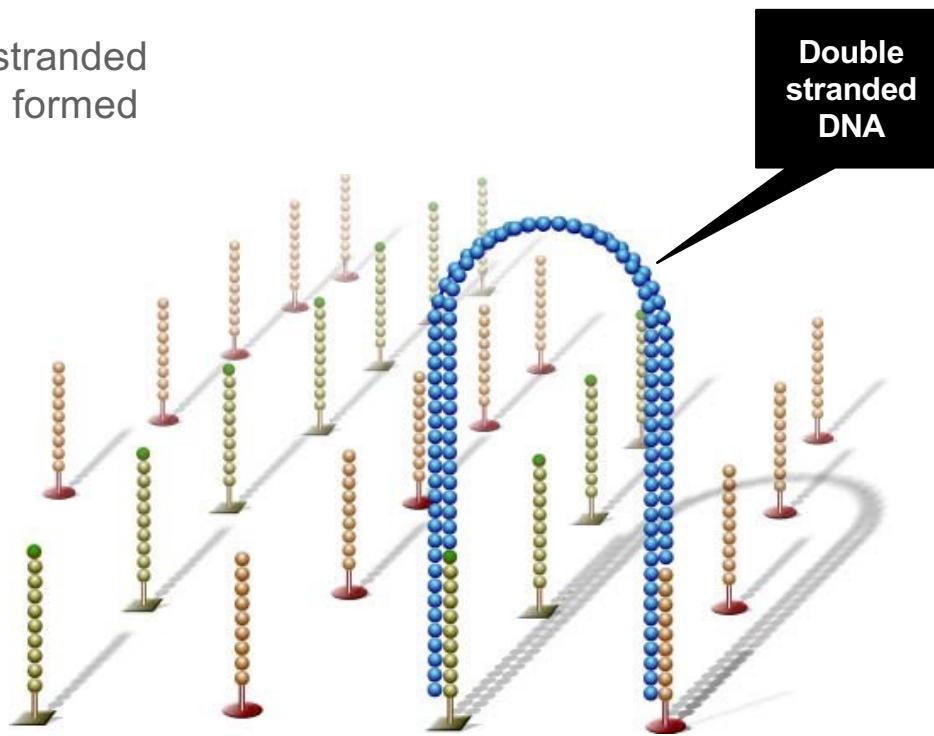
# Illumina Sequencing | Paired End Sequencing

- Single-stranded template loops over to form a bridge by hybridizing with a lawn primer
- 3'-ends of lawn primer is extended



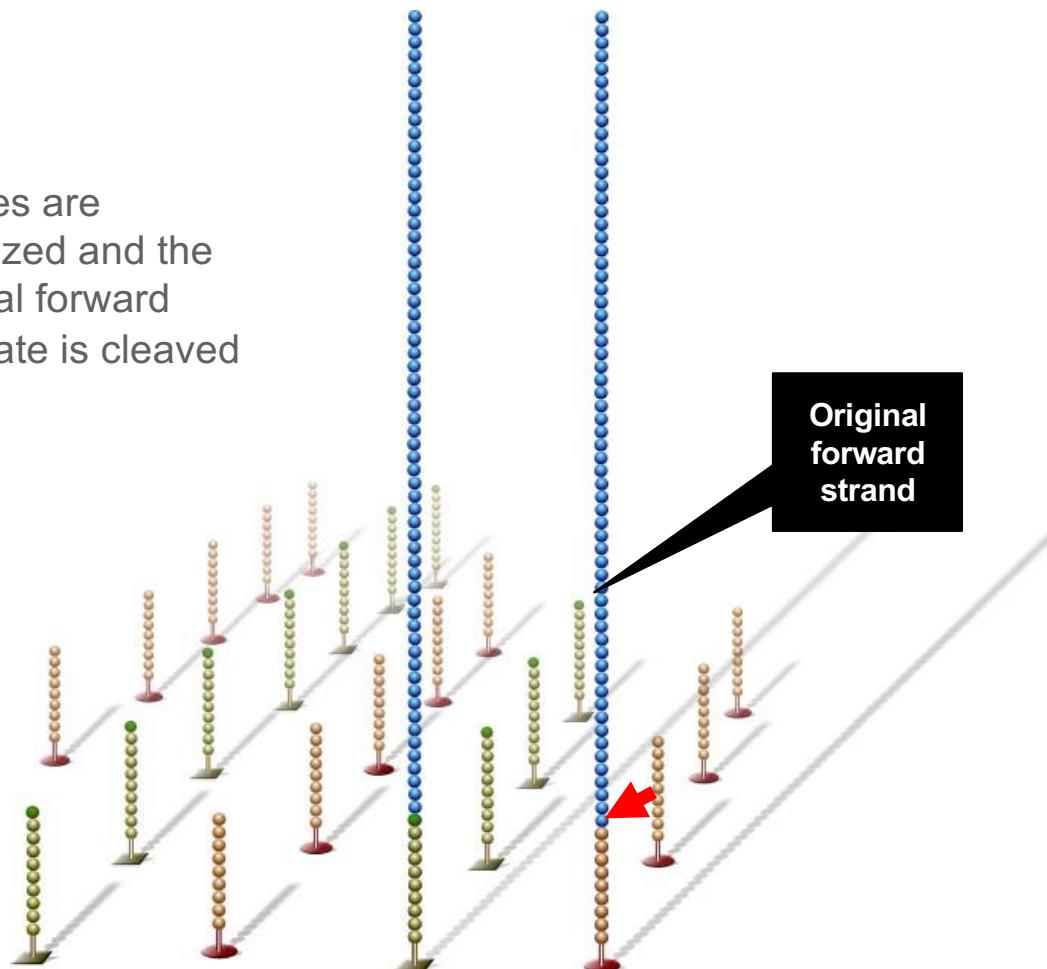
# Illumina Sequencing | Paired End Sequencing

- Double-stranded bridge is formed



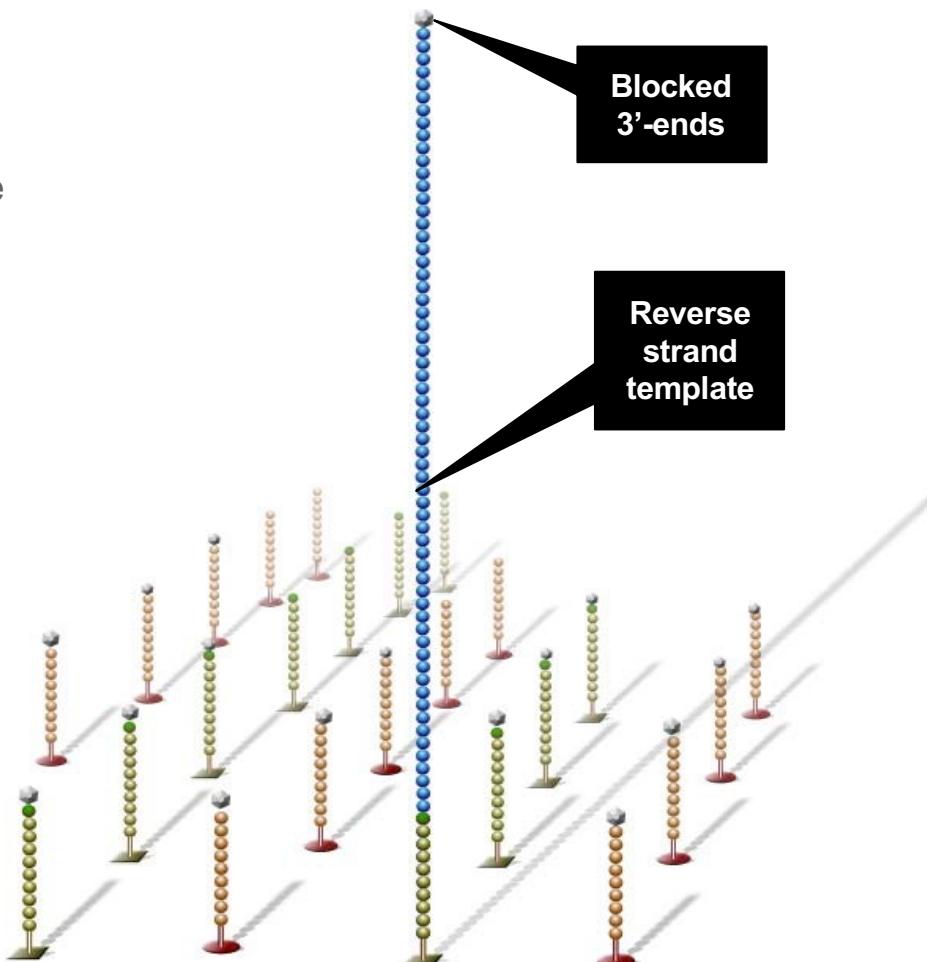
# Illumina Sequencing | Paired End Sequencing

- Bridges are linearized and the original forward template is cleaved off

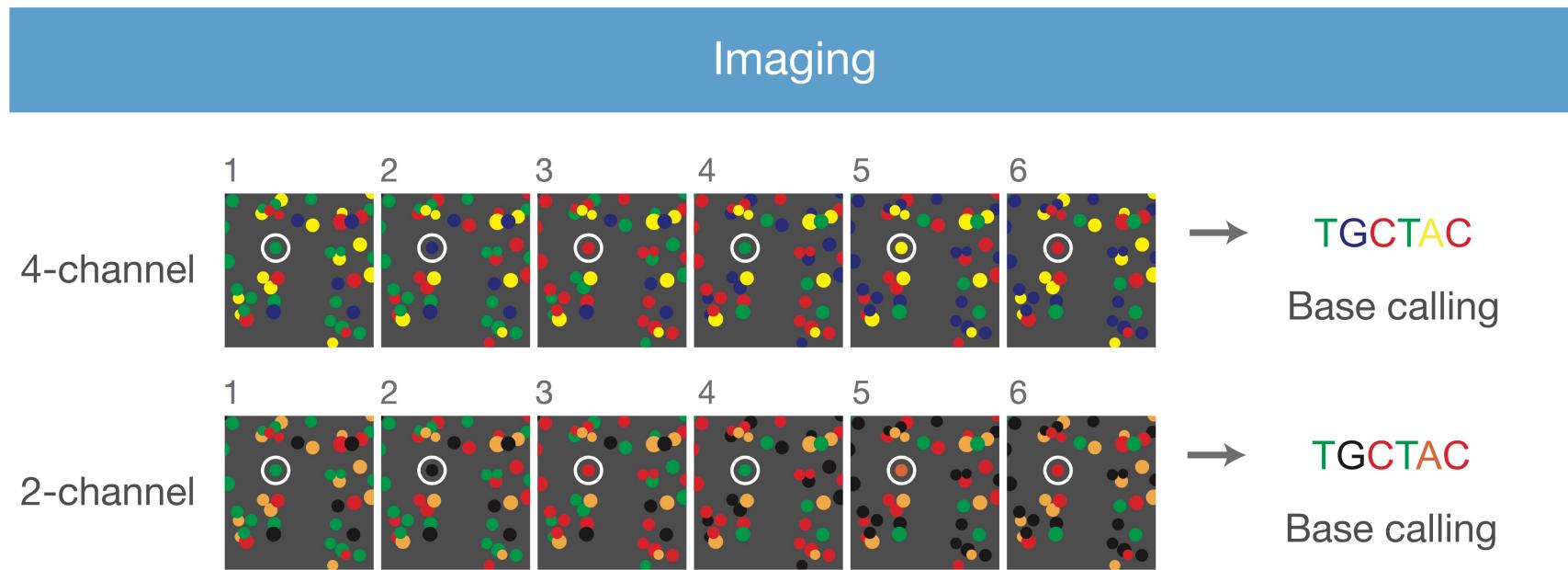


## Illumina Sequencing | Paired End Sequencing

- Free 3' ends of the reverse template and lawn primers are blocked to prevent unwanted DNA priming



## Illumina Sequencing | 4-channel vs. 2-channel



## Advantages/disadvantages 4-channel vs. 2-channel

### 4-channel

- 4 images taken per cycle
- Twice the data
- Proven accuracy/reliability
- Each base has its own color/signal

### 2-channel

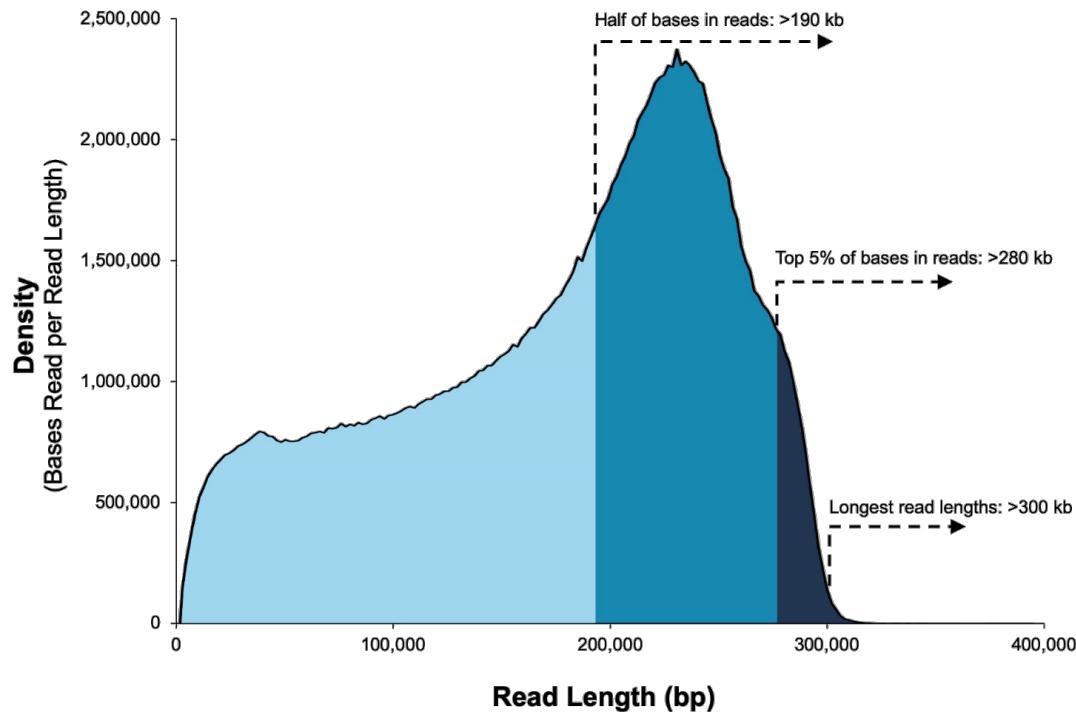
- Only 2 images taken per cycle (2-channel)
- Less storage and computation
- Faster
- Cheaper (camera)
- No-color: problematic → same as background

all new Instruments use 2-channel

# Pacific Bioscience (PacBio)

- 3<sup>rd</sup> generation sequencing technology
- **single molecule sequencing!** (relatively large amounts of DNA required as input!)
- measures polymerase activity in real time  
**(Single Molecule Real Time: SMRT)**
- can directly detect modified (methylated) nucleotides
- high error rate (random!)
- Very long reads

# HIFI Sequencing (20 KB INSERT) – Raw Data

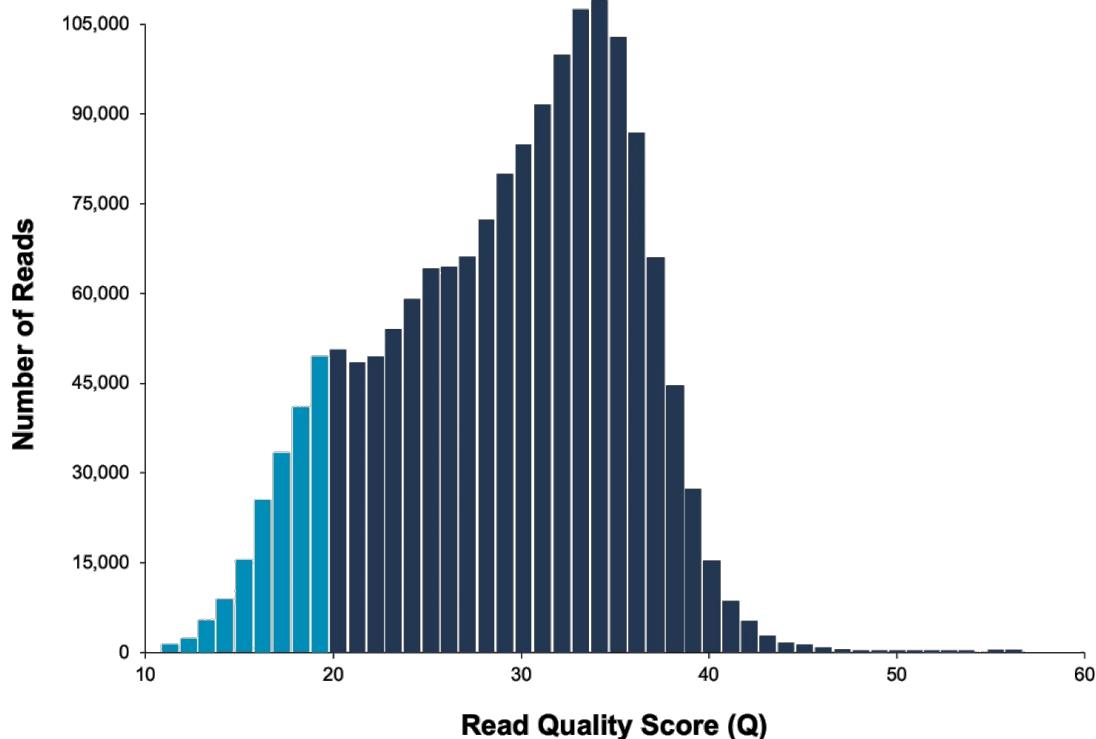


| Up to 500 Gb<br>Customer average: 325 Gb |           |
|--|-----------|
| Metrics                                  |           |
| Insert Size                              | 20 kb     |
| Number of Raw Bases (Gb)                 | 377       |
| Total Reads                              | 4,266,403 |
| Half of Bases in Reads                   | >193,403  |
| Longest read lengths                     | >300,000  |

Data shown above from a 20 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

from PacBio

# HIFI Sequencing – Post-CSS processing



Q20 (99%) HiFi accuracy

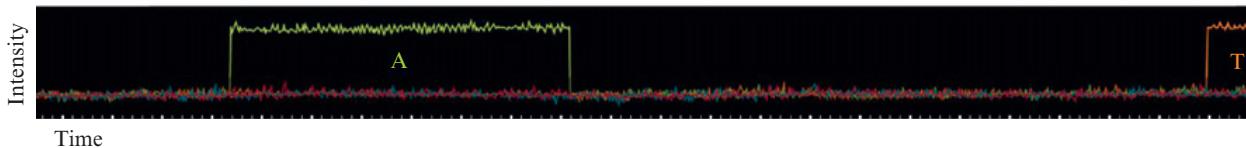
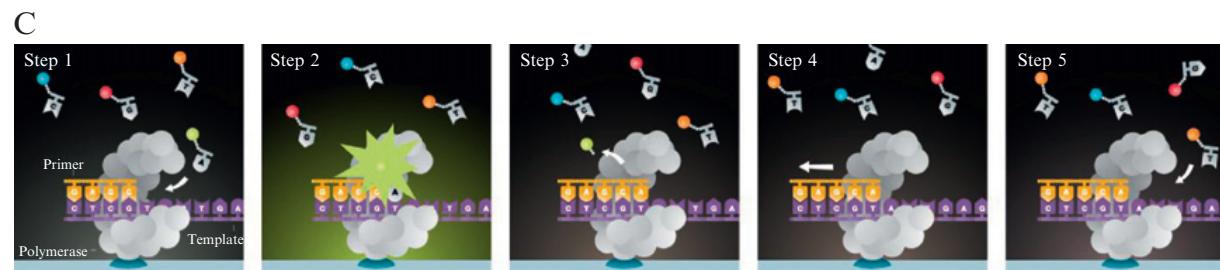
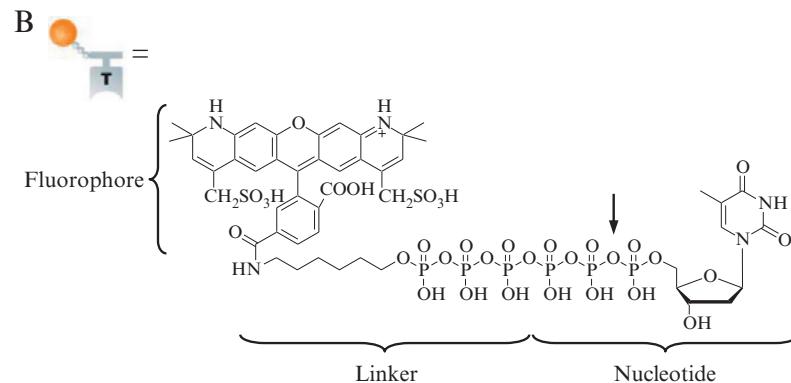
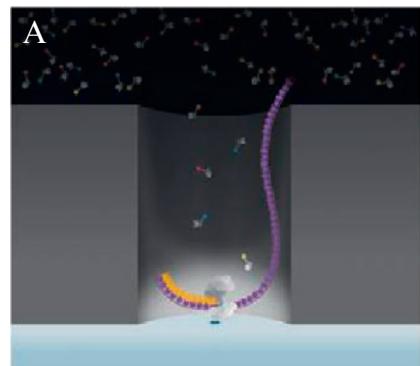
| Metrics                  |               |
|--------------------------|---------------|
| Insert Size              | 20 kb         |
| Number of CCS Bases (Gb) | 30 Gb         |
| Number of >Q20 Bases     | 26 Gb         |
| Number of >Q20 Reads     | 1,423,277     |
| <b>Accuracy (Mean)</b>   | <b>99.92%</b> |

Customer average:  
25-30 Gb

Data shown above from a 20 kb size-selected human library using the SMRTbell Template Prep Kit on a Sequel II System (2.0 Chemistry, Sequel II System Software v8.0, 30-hour movie). Read lengths, reads/data per SMRT Cell 8M and other sequencing performance results vary based on sample quality/type and insert size.

from PacBio

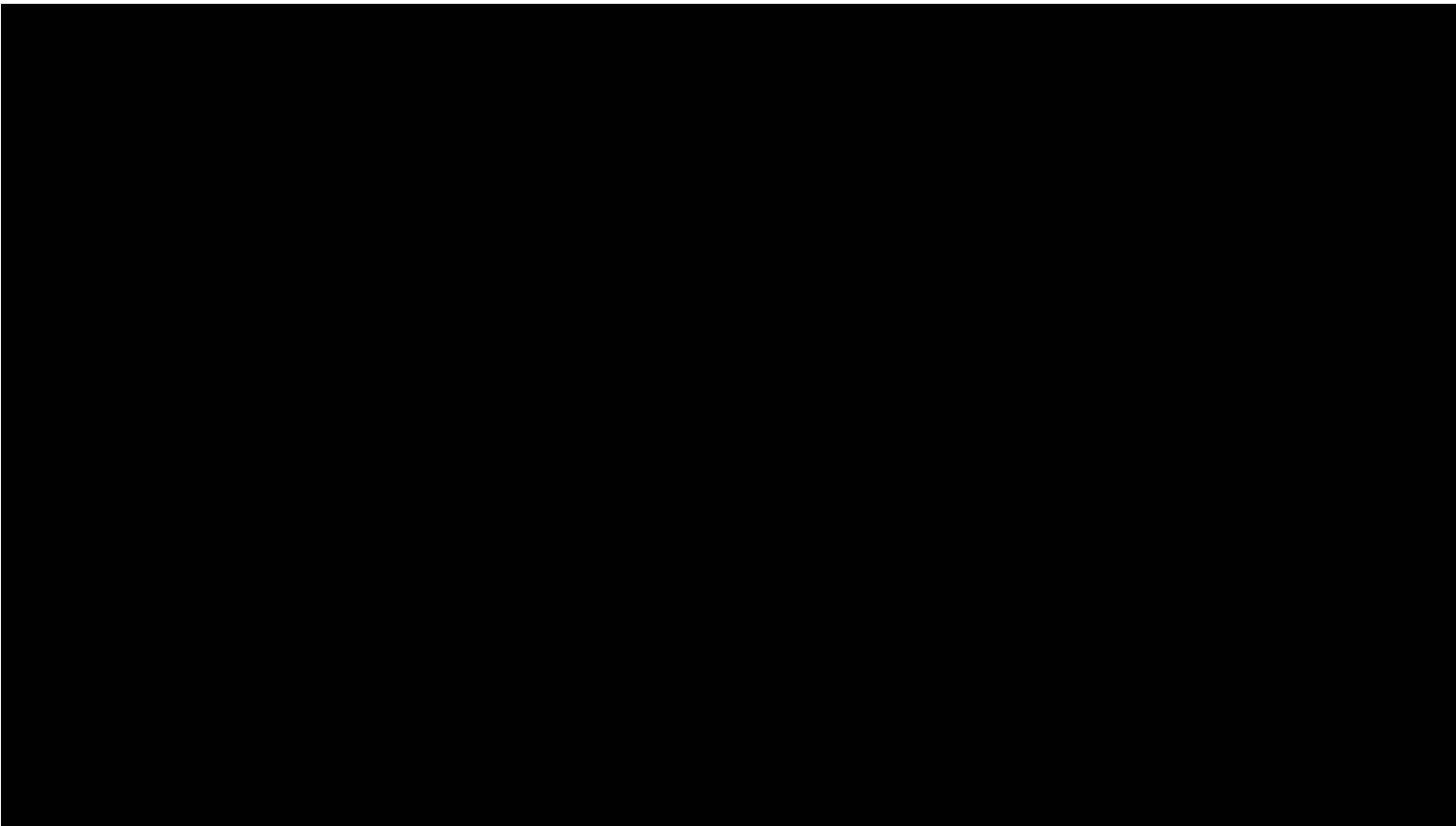
# PacBio Sequencing Principle



# PacBio Sequencing Principle

*u*<sup>b</sup>

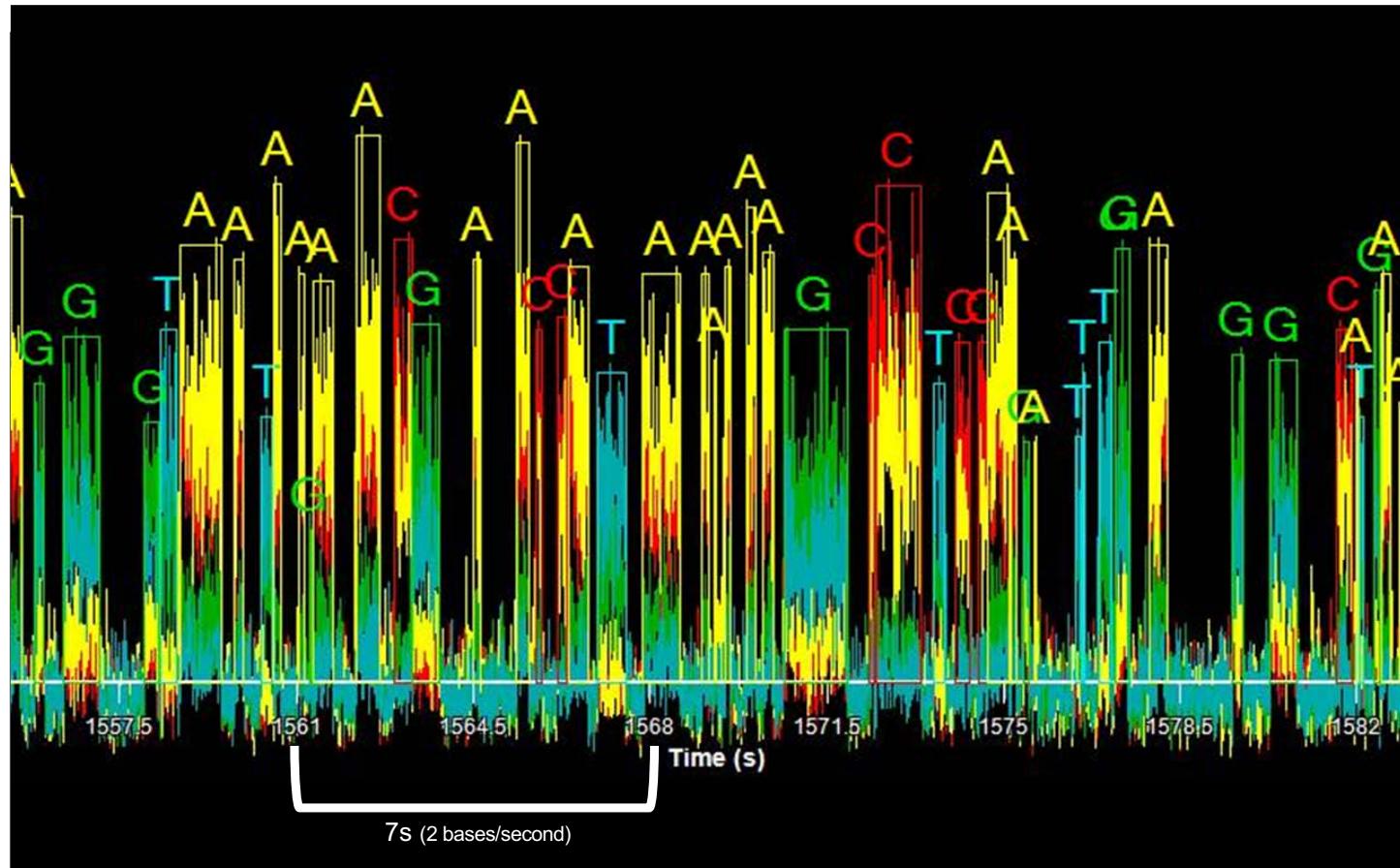
*b*  
**UNIVERSITÄT  
BERN**



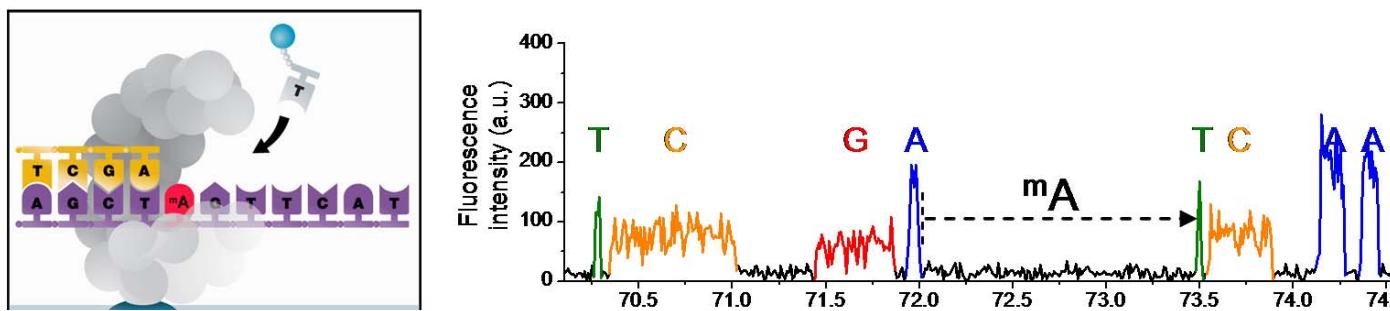
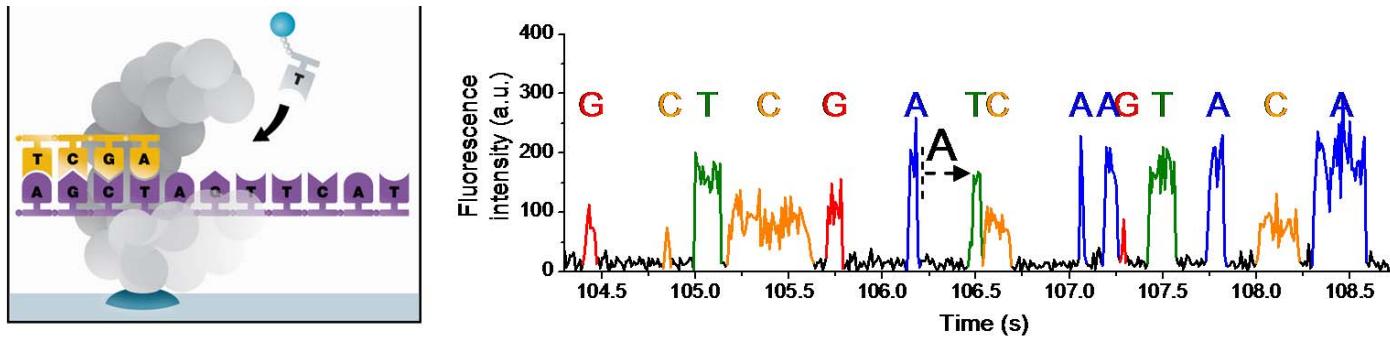
$u^b$

b  
UNIVERSITÄT  
BERN

## Natural Process Monitored



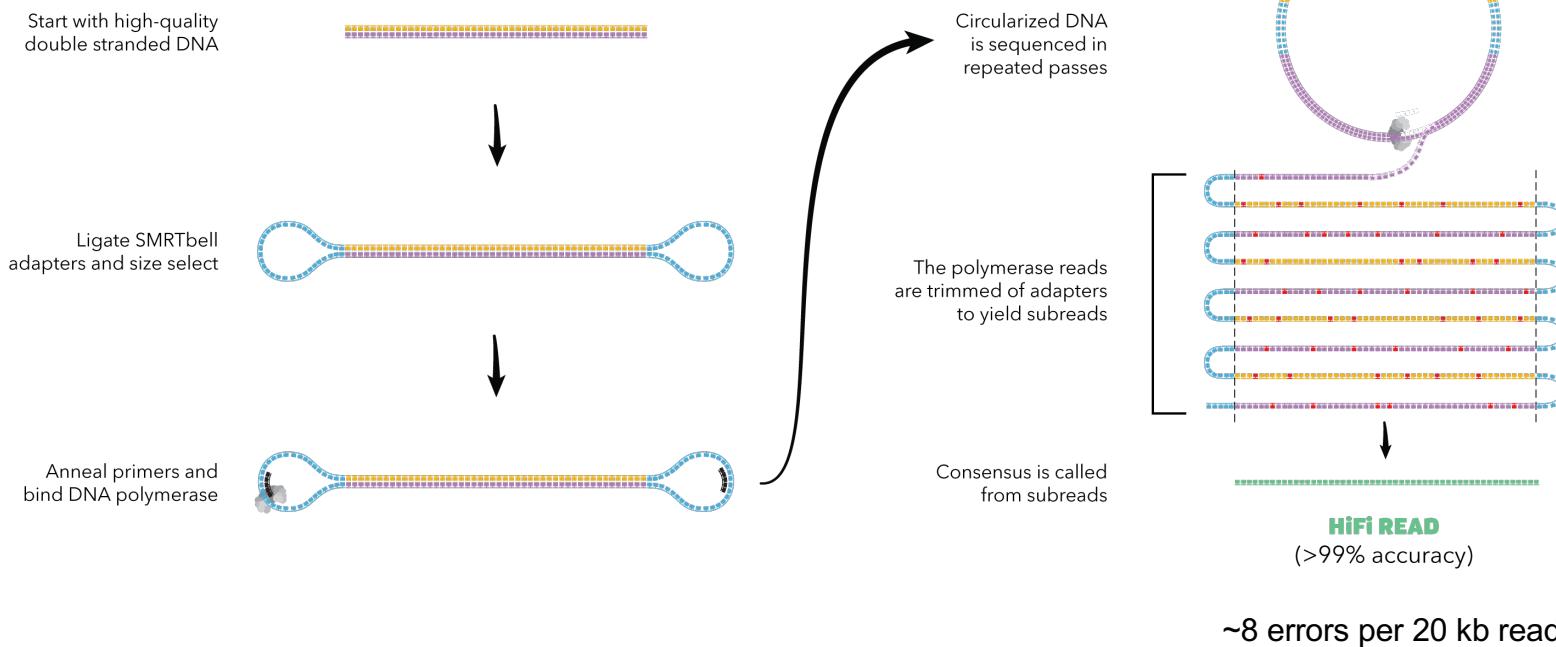
# Kinetic Detection Enables Study of Modified Bases



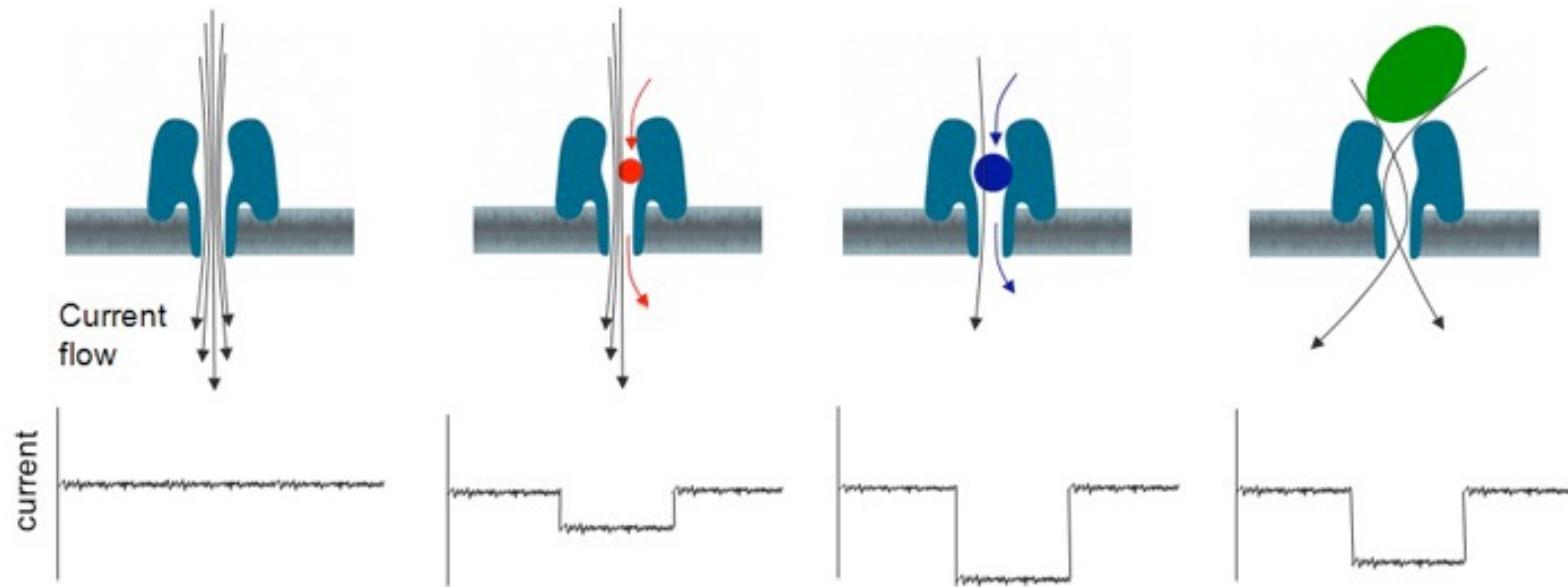
# HIFI – SMRT-Bell



# Pacbio HiFi Reads



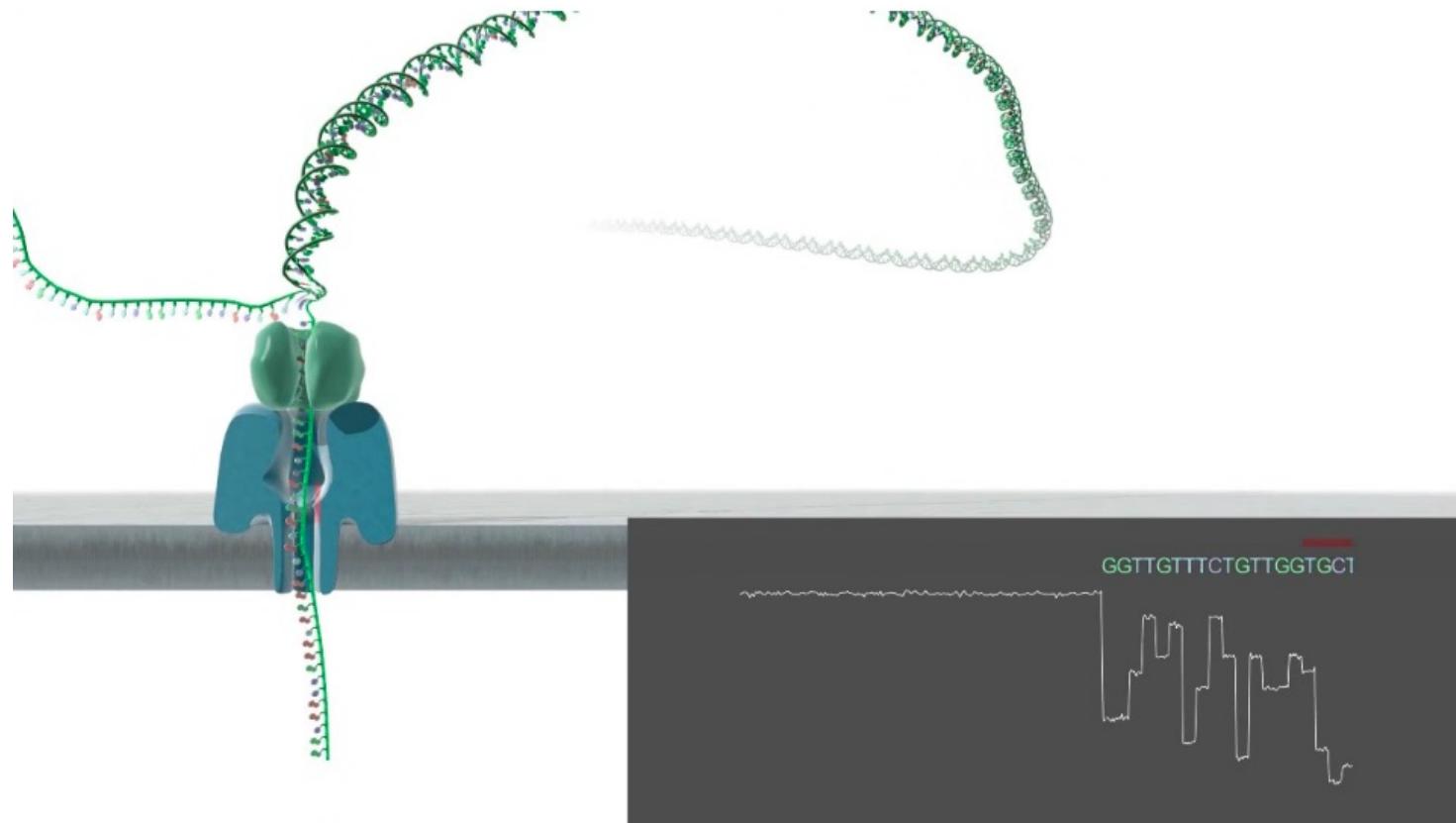
# Oxford Nanopore Technology | ONT



$u^b$

b  
UNIVERSITÄT  
BERN

# Oxford Nanopore Technology | ONT



# ONT Sequencing Principle

*u*<sup>b</sup>

*b*  
**UNIVERSITÄT  
BERN**

<https://www.youtube.com/watch?v=E9-Rm5AoZGw>

*u*<sup>b</sup>

*b*  
UNIVERSITÄT  
BERN

# ONT Sequencing Principle

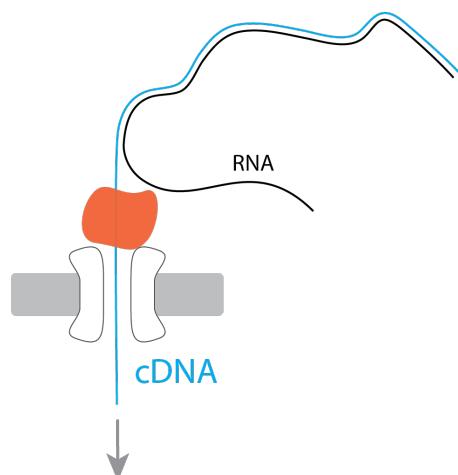


<https://vimeo.com/337258910>

# ONT | RNA Sequencing

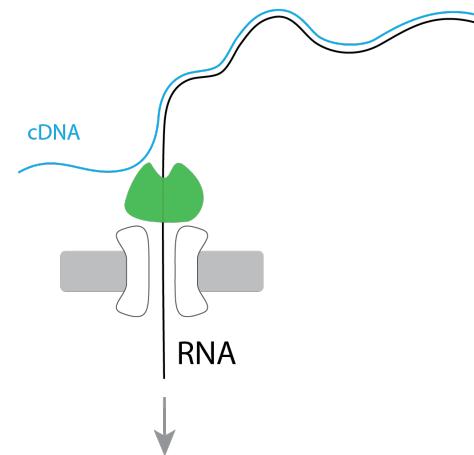
## cDNA

- Isoform identification
- Compatible with cDNA workflows
- PCR free
- amplification for low input



## direct RNA

- **sequence native RNA**
- no PCR bias
- **detection of modified bases**
- suitable for transcripts that fail RT



# ONT | Features

- Real time sequencing/data analysis
  - first sequences appear minutes after the run start
  - provides live feedback on how the experiment is running
  - interesting for e.g. diagnostics applications where time matters
- no fixed run time
  - once the result/objective is reached, one can stop and wash/reuse the flowcell
  - The system can run until sufficient data has been collected, analysing the data in real time

# ONT | Features

|                          |   |
|--------------------------|---|
| Read length              | Nanopores read the entire length of the fragment of DNA/RNA presented to them.<br>Longest read so far: > 4Mb. Accuracy is maintained throughout the fragment. |
| Single Molecule Accuracy | R9 modal >98.3%, New chemistry >99% (coming soon)   |
| Consensus Accuracy       | R9.4.1: Current best Q45 (>99.99%)<br>R10: Current Best Q50 (99.999%)   |
| Sample Preparation time  | Rapid Kit: 10 minutes, Ligation Kit: 60 minutes, other protocols and timings also available   |
| Modified base detection  | Yes - Base modification information available in raw signal   |
| Time to 1st usable data  | 2 minutes   |

# Error Sources

*u*<sup>b</sup>

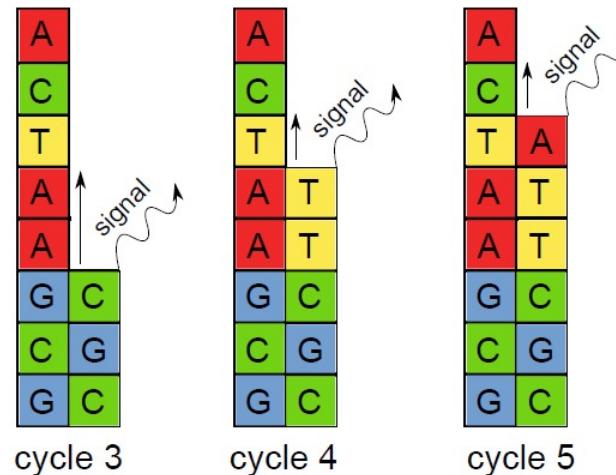
---

*b*  
**UNIVERSITÄT  
BERN**

# Sequencing Errors

- Cloning artifacts
- Errors in PCR amplification
- GC-Bias, enrichment bias
- De-phasing
- Base misincorporation
- Instrument specific errors
- ... lots we still don't know ...

De-phasing: loss of cycle information



cycle 4: two T's incorporated, one signal

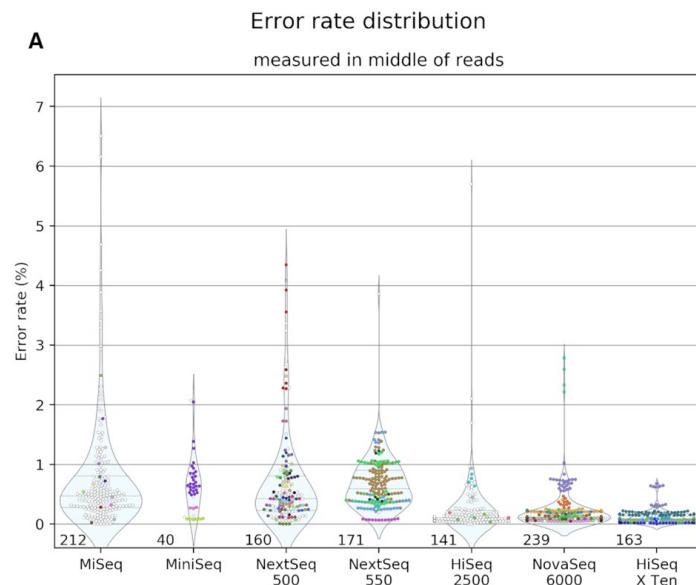
|               |          |
|---------------|----------|
| Sequence read | CGCT-AGT |
| True sequence | CGCTTAGT |

$u^b$

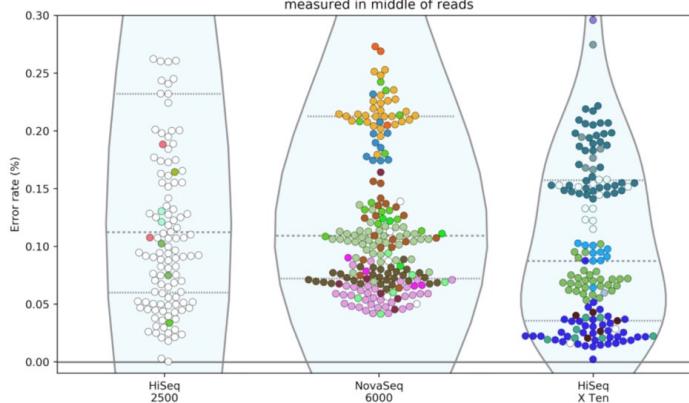
b  
UNIVERSITÄT  
BERN

# Error Rate of Illumina Instruments

A



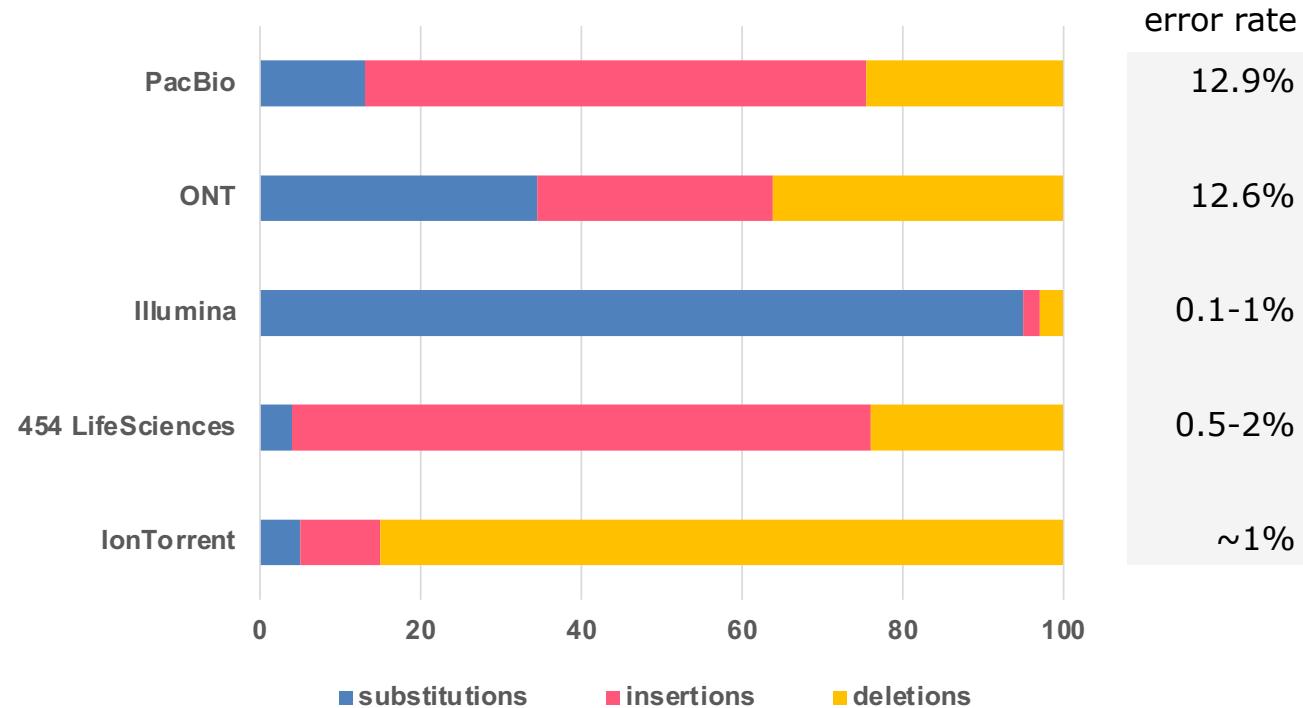
B  
Error rate distribution  
measured in middle of reads



| Error rate (%) |                   |        |                    |
|----------------|-------------------|--------|--------------------|
| Platform       | Number of samples | Median | Standard deviation |
| MiSeq          | 212               | 0.473  | 0.938              |
| MiniSeq        | 40                | 0.613  | 0.459              |
| NextSeq 500    | 160               | 0.429  | 0.827              |
| NextSeq 550    | 171               | 0.593  | 0.435              |
| HiSeq 2500     | 141               | 0.112  | 0.544              |
| NovaSeq 6000   | 239               | 0.109  | 0.350              |
| HiSeq X Ten    | 163               | 0.087  | 0.126              |

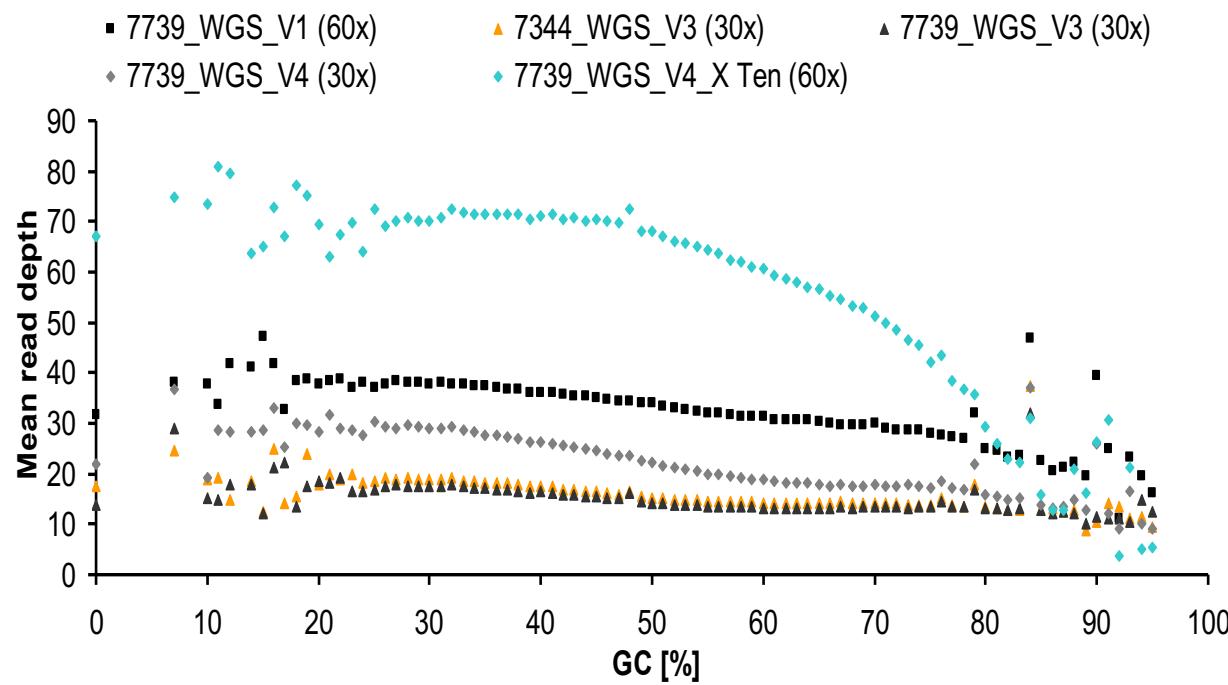
NAR Genom Bioinform, Volume 3, Issue 1, March 2021, lqab019,  
<https://doi.org/10.1093/narqab/lqab019>

# Platform Specific Error Types

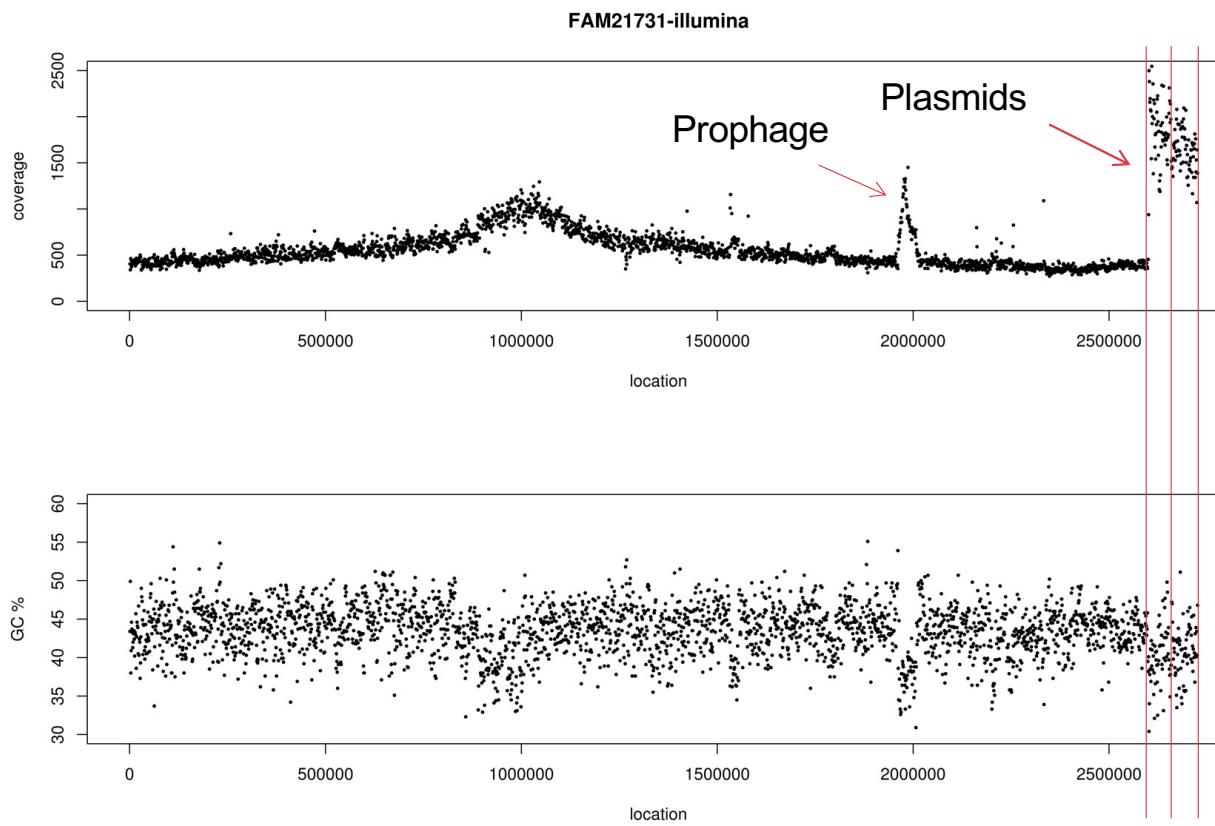


NAR Genom Bioinform, Volume 2, Issue 2, June 2020, lqaa037, <https://doi.org/10.1093/nargab/lqaa037>

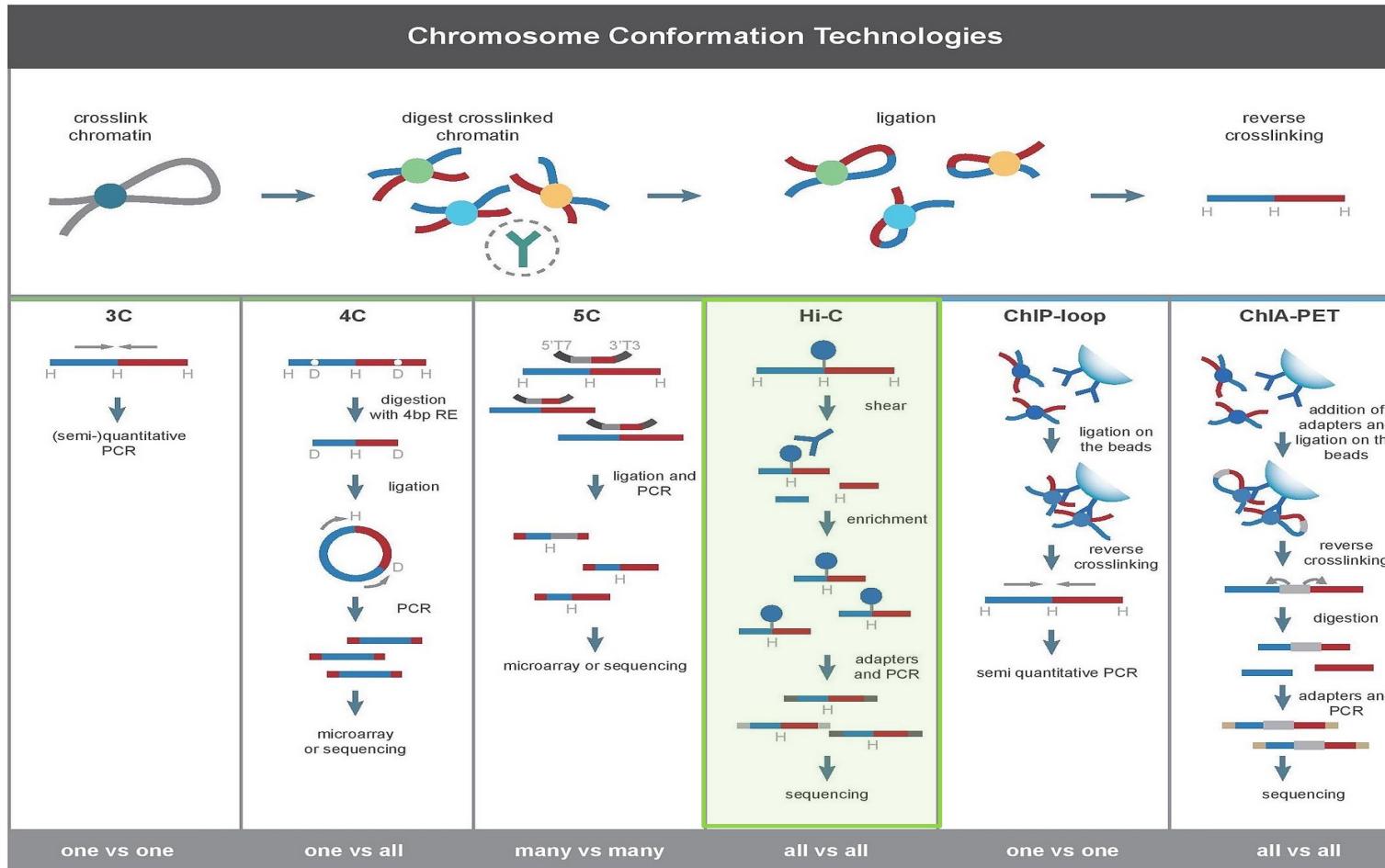
- Bias PCR amplification due to differences in the GC content



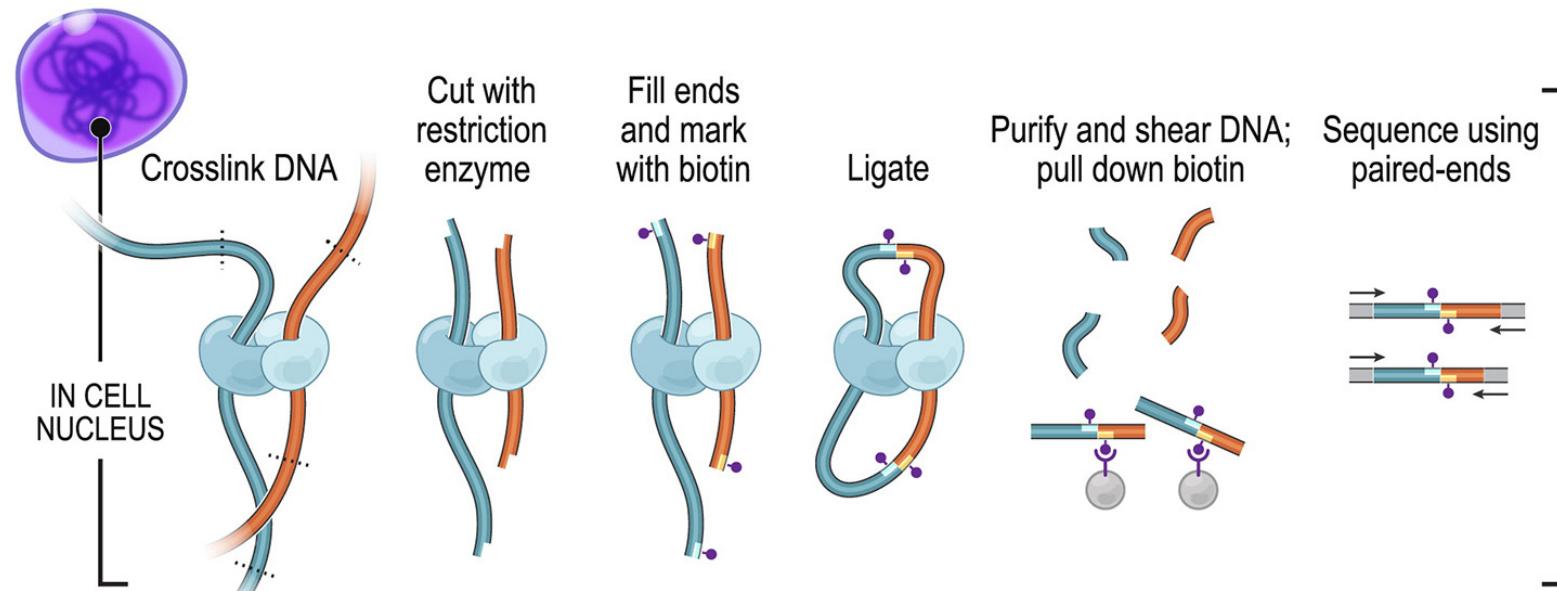
Meienberg J, Zerjavic1 K, Keller I, Okoniewski M, Patrignani A, Ludin K, Xu Z, Steinmann B, Carrel T, Röthlisberger B, Schlapbach R, Bruggmann R and Matyas G; Nucleic Acids Res . 2015 Jun 23;43(11):e76. doi: 10.1093/nar/gkv216. Epub 2015 Mar 27.



# Hi-C | Chromosome Conformation Capturing



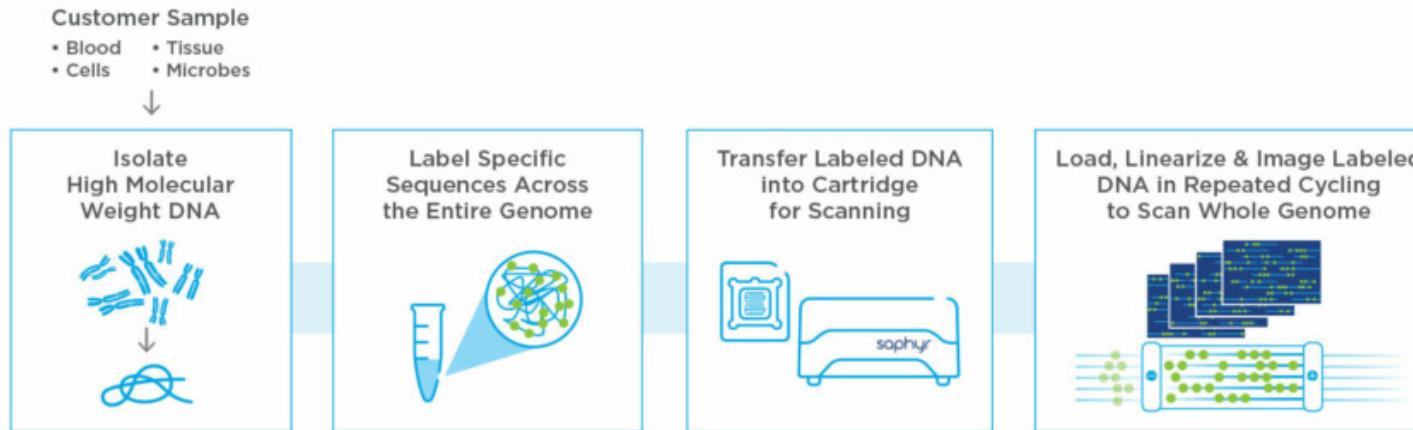
# Hi-C | Chromosome Conformation Capturing



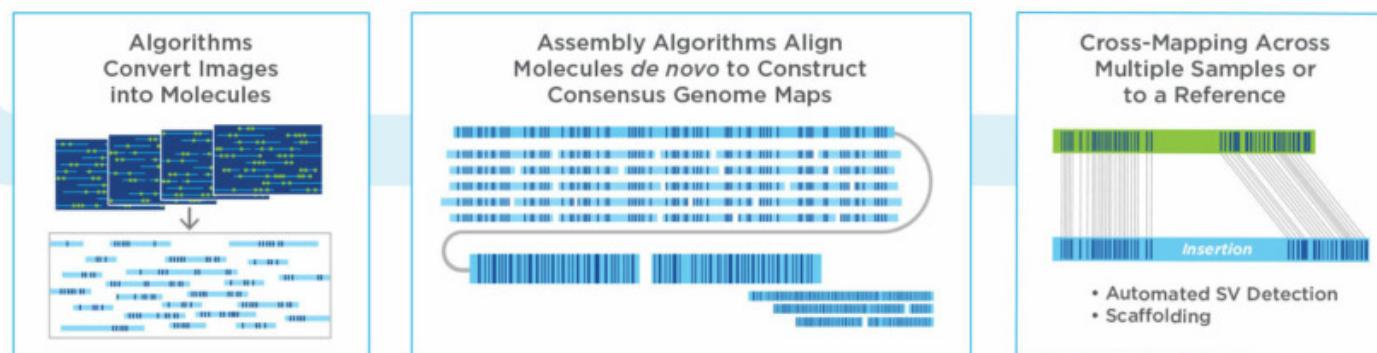
Distance between reads can be up to 500 Mb

Rao et al., 2014. Cell 159: 1665-1680.

# Optical Mapping | Bionano

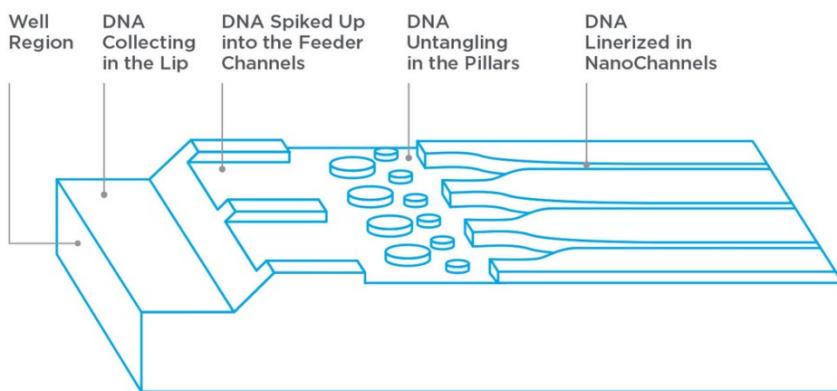


High-throughput, High-resolution Imaging of Megabase Length Molecules

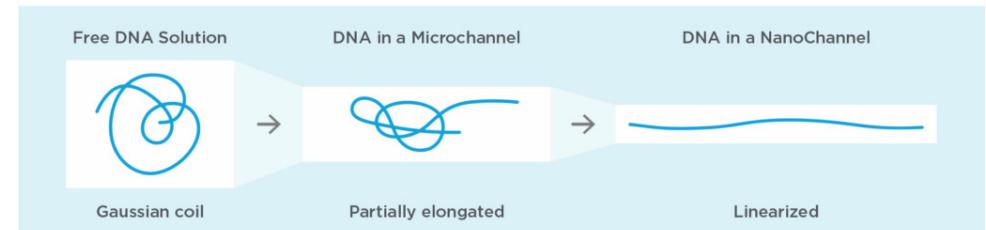


# Optical Mapping | Bionano

SAPHYRCHIP™ SCHEMATIC

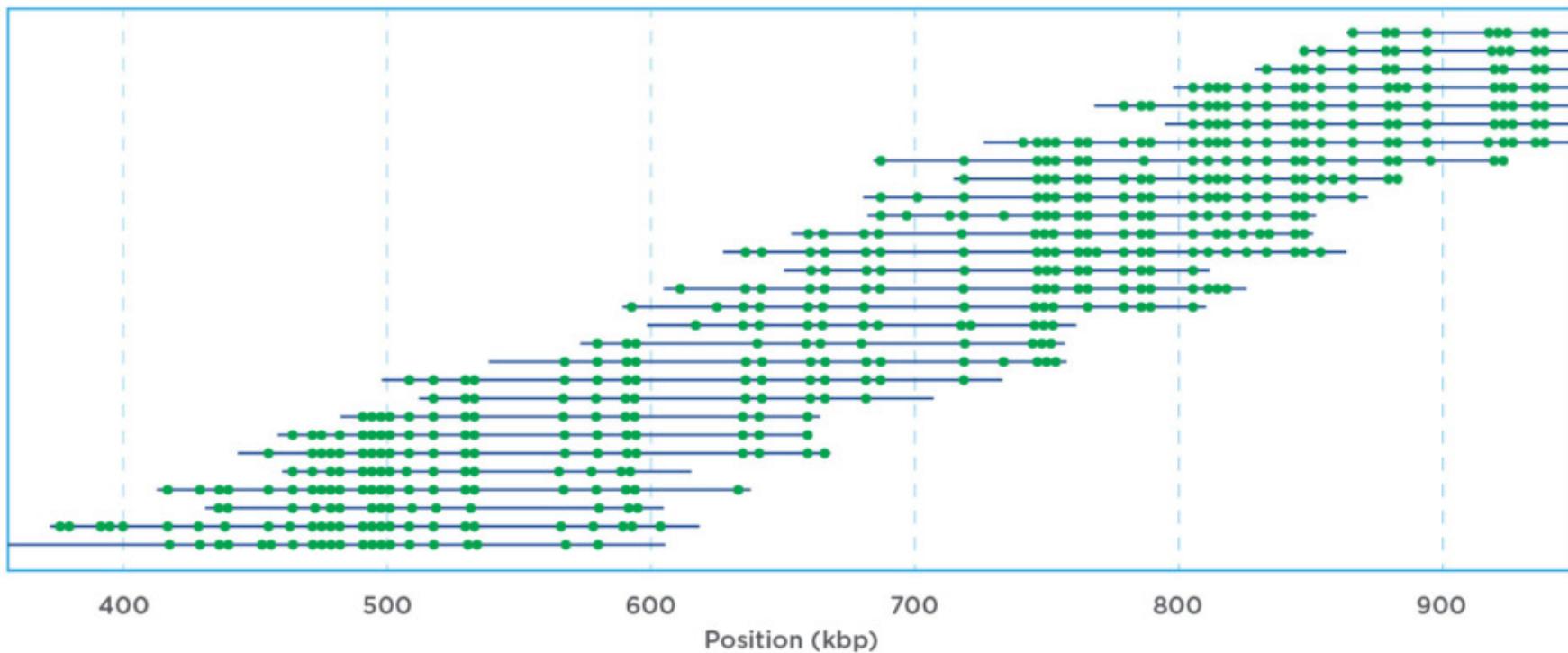


SINGLE DNA MOLECULE LINEARIZATION IN NANOCHANNEL



# Optical Mapping | Bionano

## DIGITAL REPRESENTATION OF LABELED LONG DNA



# Towards high Quality Assemblies

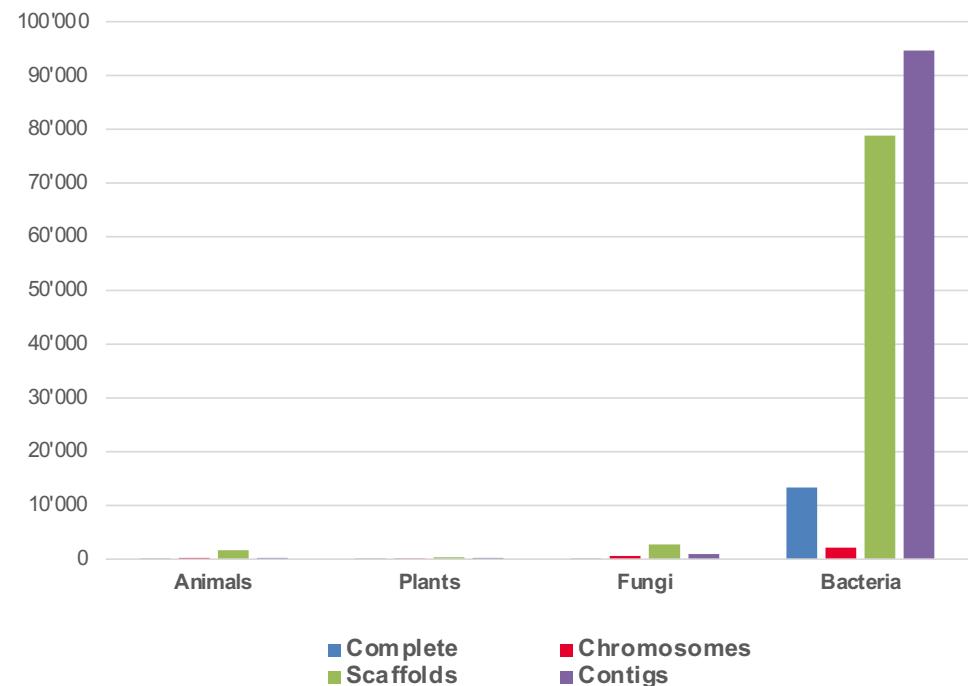


# Genomes in Public Databases (NCBI)

| Assembly Level | Total   | Complete    |           | Chromosomes |           | Scaffolds   |           | Contigs     |           |
|----------------|---------|-------------|-----------|-------------|-----------|-------------|-----------|-------------|-----------|
|                |         | Chromosomes | Scaffolds | Chromosomes | Scaffolds | Chromosomes | Scaffolds | Chromosomes | Scaffolds |
| Animals        | 3'853   | 3           | (0.1%)    | 695         | (18.0%)   | 2'773       | (72.0%)   | 382         | (9.9%)    |
| Plants         | 1'488   | 3           | (0.2%)    | 416         | (28.0%)   | 665         | (44.7%)   | 404         | (27.2%)   |
| Fungi          | 6'886   | 86          | (1.2%)    | 826         | (12.0%)   | 4'390       | (63.8%)   | 1'584       | (23.0%)   |
| Protists       | 943     | 18          | (1.9%)    | 119         | (12.6%)   | 593         | (62.9%)   | 213         | (22.6%)   |
| Eukaryotes     | 13'170  | 110         | (0.8%)    | 2'056       | (15.6%)   | 8'421       | (63.9%)   | 2'583       | (19.6%)   |
| Bacteria       | 276'806 | 19'884      | (7.2%)    | 3'270       | (1.2%)    | 102'252     | (36.9%)   | 151'400     | (54.7%)   |
| Archaea        | 5'615   | 401         | (7.1%)    | 25          | (0.4%)    | 2'163       | (38.5%)   | 3'026       | (54.7%)   |
| Σ Prokaryotes  | 282'421 | 20'285      | (7.2%)    | 3'295       | (1.2%)    | 104'415     | (37.0%)   | 154'426     | (54.7%)   |
| <hr/>          |         |             |           |             |           |             |           |             |           |
| Overall        | 295'591 | 20'395      | (6.9%)    | 5'351       | (1.8%)    | 112'836     | (38.2%)   | 157'009     | (53.1%)   |

\*data accessed October 2020 at NCBI genomes  
<https://www.ncbi.nlm.nih.gov/genome/browse#!/overview/>

Genome Assemblies at NCBI



# The Human Genome

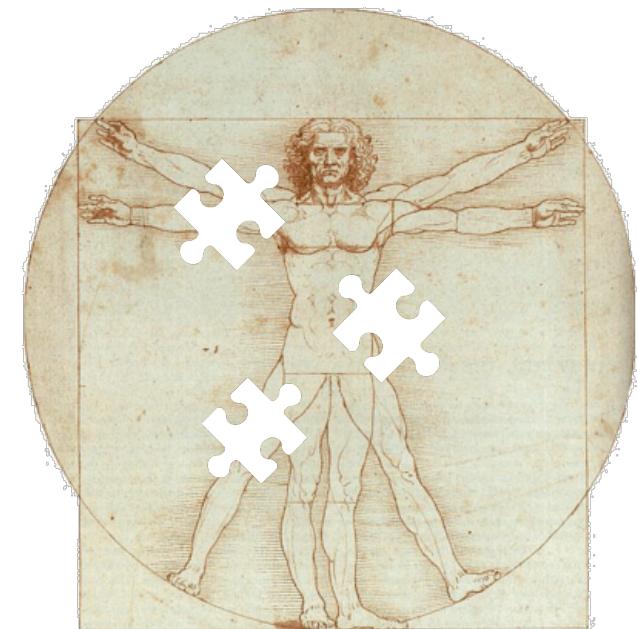
**Psst, the human genome was never completely sequenced. Some scientists say it should be**

By SHARON BEGLEY @sxbegley / JUNE 20, 2017

“As a matter of truth in advertising, **the ‘finished’ sequence isn’t finished**,” said Eric Lander, who led the lab at the Whitehead Institute that deciphered more of the genome for the government-funded Human Genome Project than any other. “I always say ‘finished’ is a term of art.”

**“It’s very fair to say the human genome was never fully sequenced,”** Craig Venter, another genomics luminary, told STAT.

**“The human genome has not been completely sequenced and neither has any other mammalian genome** as far as I’m aware,” said Harvard Medical School bioengineer George Church, who made key early advances in sequencing technology.



***still missing sequences of the human genome  
4-9%***  
*(Georg Church, Karen Miga)*

# Telomere-to-telomere Initiative

## Article

### Telomere-to-telomere assembly of a complete human X chromosome

<https://doi.org/10.1038/s41586-020-2547-7>

Received: 30 July 2019

Accepted: 29 Ma

Published online:

Open access

Check for update

Karen H. Miga<sup>1,24</sup>, Sergey Koren<sup>2,24</sup>, Arang Rhie<sup>2</sup>, Mitchell R. Vollger<sup>3</sup>, Ariel Gershman<sup>4</sup>, Andrey Bzikadze<sup>5</sup>, Shelise Brooks<sup>6</sup>, Edmund Howe<sup>7</sup>, David Porubsky<sup>8</sup>, Glennis A. Logsdon<sup>3</sup>, Valerie A. Schneider<sup>8</sup>, Tamara Potapova<sup>7</sup>, Jonathan Wood<sup>9</sup>, William Chow<sup>9</sup>, Joel Armstrong<sup>1</sup>.

## Article

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

<https://doi.org/10.1038/s41586-021-03420-7>

Received: 4 September 2020

Accepted: 4 March 2021

Published online: 07 April 2021

Open access

Check for updates

Glennis A. Logsdon<sup>1</sup>, Mitchell R. Vollger<sup>1</sup>, PingHsun Hsieh<sup>1</sup>, Yafei Mao<sup>1</sup>, Mikhail A. Liskovskykh<sup>2</sup>, Sergey Koren<sup>3</sup>, Sergey Nurk<sup>3</sup>, Ludovica Mercuri<sup>4</sup>, Philip C. Dishuck<sup>1</sup>, Arang Rhie<sup>3</sup>, Leonardo G. de Lima<sup>5</sup>, Tatiana Dvorkina<sup>6</sup>, David Porubsky<sup>1</sup>, William T. Harvey<sup>1</sup>, Alla Mikheenko<sup>6</sup>, Andrey V. Bzikadze<sup>7</sup>, Milinn Kremitzki<sup>8</sup>, Tina A. Graves-Lindsay<sup>8</sup>, Chirag Jain<sup>9</sup>, Kendra Hoekzema<sup>1</sup>, Shwetha C. Murali<sup>10</sup>, Katherine M. Munson<sup>1</sup>, Carl Baker<sup>1</sup>, Melanie Sorensen<sup>1</sup>, Alexandra M. Lewis<sup>1</sup>, Urvashi Surti<sup>10</sup>, Jennifer L. Gerton<sup>5</sup>, Vladimir Larionov<sup>2</sup>, Mario Ventura<sup>4</sup>, Karen H. Miga<sup>11</sup>, Adam M. Phillippy<sup>12</sup> & Evan E. Eichler<sup>1,24</sup>, Karen H. Miga<sup>1,11</sup>, Adam M. Phillippy<sup>1,12</sup>

The complete assembly of each human chromosome is essential for understanding human biology and evolution<sup>1,2</sup>. Here we use complementary long-read sequencing technologies to complete the linear assembly of human chromosome 8. Our assembly resolves the sequence of five previously long-standing gaps, including a 2.08-Mb centromeric  $\alpha$ -satellite array, a 644-kb copy number polymorphism in the  $\beta$ -defensin gene cluster that is important for disease risk, and an 863-kb variable number tandem repeat at chromosome 8q11.2 that can function as a centromere. We show that

## The complete sequence of a human genome

Sergey Nurk<sup>1,1</sup>, Sergey Koren<sup>1,1</sup>, Arang Rhie<sup>1,1</sup>, Mikko Rautiainen<sup>1,1</sup>, Andrey V. Bzikadze<sup>2</sup>, Alla Mikheenko<sup>3</sup>, Mitchell R. Vollger<sup>4</sup>, Nicolas Altemose<sup>5</sup>, Lev Uralsky<sup>4,7</sup>, Ariel Gershman<sup>8</sup>, Sergey Aganezov<sup>9</sup>, Savannah J. Hoyt<sup>10</sup>, Mark Diekhans<sup>11</sup>, Glennis A. Logsdon<sup>1</sup>, Michael Alonso<sup>12</sup>, Stylianos E. Antonarakis<sup>12</sup>, Matthew Borchers<sup>13</sup>, Gerard G. Bouffard<sup>14</sup>, Shelise Y. Brooks<sup>14</sup>, Gina V. Caldas<sup>15</sup>, Haoyu Cheng<sup>16,17</sup>, Chen-Shan Chin<sup>18</sup>, William Chow<sup>19</sup>, Leonardo G. de Lima<sup>13</sup>, Philip C. Dishuck<sup>1</sup>, Richard Durbin<sup>21</sup>, Tatiana Dvorkina<sup>2</sup>, Ian T. Fiddes<sup>22</sup>, Giulio Formenti<sup>23,24</sup>, Robert S. Fulton<sup>25</sup>, Arkarachai Fungtammasan<sup>18</sup>, Erik Garrison<sup>11,26</sup>, Patrick G.S. Grady<sup>10</sup>, Tina A. Graves-Lindsay<sup>27</sup>, Ira M. Hall<sup>19</sup>, Nancy F. Hansen<sup>29</sup>, Gabrielle A. Hartley<sup>10</sup>, Marina Haukness<sup>1</sup>, Kerstin Howe<sup>30</sup>, Michael W. Hunkapiller<sup>30</sup>, Chirag Jain<sup>1,31</sup>, Miten Jain<sup>11</sup>, Erich D. Jarvis<sup>23,24</sup>, Peter Kerepeldiev<sup>32</sup>, Melanie Kirsch<sup>6</sup>, Mikhail Kolmogorov<sup>33</sup>, Jonas Korlach<sup>30</sup>, Milinn Kremitzki<sup>27</sup>, Heng Li<sup>16,17</sup>, Valerie V. Maduro<sup>34</sup>, Tobias Marschall<sup>35</sup>, Ann M. McCartney<sup>1</sup>, Jennifer McDaniel<sup>36</sup>, Danny E. Miller<sup>4,37</sup>, James C. Mullikin<sup>14,21</sup>, Eugene W. Myers<sup>38</sup>, Nathan D. Olson<sup>39</sup>, Benedict Paten<sup>11</sup>, Paul Peluso<sup>30</sup>, Pavel A. Pevzner<sup>33</sup>, David Porubsky<sup>1</sup>, Tamara Potapova<sup>13</sup>, Evgeny I. Rogae<sup>4,7,39,40</sup>, Jeffrey A. Rosenfeld<sup>41</sup>, Steven L. Salzberg<sup>4,42</sup>, Valerie A. Schneider<sup>43</sup>, Fritz J. Sedlacek<sup>1</sup>, Kishwar Shafin<sup>10</sup>, Colin J. Shew<sup>20</sup>, Alaina Shumate<sup>42</sup>, Yumi Sims<sup>19</sup>, Arian F. A. Smit<sup>45</sup>, Daniela C. Soto<sup>20</sup>, Ivan Sovic<sup>10,46</sup>, Jessica M. Storer<sup>45</sup>, Aaron Streets<sup>5,47</sup>, Beth A. Sullivan<sup>48</sup>, Francois Thibaud-Nissen<sup>43</sup>, James Torrance<sup>10</sup>, Justin Wagner<sup>39</sup>, Brian P. Walenz<sup>1</sup>, Aaron Wenger<sup>30</sup>, Jonathan M. D. Wood<sup>19</sup>, Chunlin Xiao<sup>43</sup>, Stephanie M. Yan<sup>49</sup>, Alice C. Young<sup>14</sup>, Samantha Zarate<sup>1</sup>, Urvashi Surti<sup>10</sup>, Rajiv C. McCoy<sup>49</sup>, Megan Y. Dennis<sup>40</sup>, Ivan A. Alexandrov<sup>3,7,51</sup>, Jennifer L. Gerton<sup>13</sup>, Rachel J. O'Neill<sup>10</sup>, Winston Tim<sup>4,42</sup>, Justin M. Zook<sup>36</sup>, Michael C. Schatz<sup>9,49</sup>, Evan E. Eichler<sup>1,24</sup>, Karen H. Miga<sup>1,11</sup>, Adam M. Phillippy<sup>1,12</sup>

<sup>1-51</sup> Affiliations are listed at the end

<sup>1</sup> Equal contribution

<sup>2</sup> Corresponding authors: Evan E. Eichler (ehee@gs.washington.edu); Karen H. Miga (khmiga@ucsc.edu); Adam M. Phillippy (adam.phillippy@nih.gov)

## Abstract

In 2001, Celera Genomics and the International Human Genome Sequencing Consortium published their initial drafts of the human genome, which revolutionized the field of genomics. While these drafts and the updates that followed effectively covered the euchromatic fraction of the genome, the heterochromatin and many other complex regions were left unfinished or erroneous. Addressing this remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium has finished the first truly complete 3.055 billion base pair (bp) sequence of a human genome, representing the largest improvement to the human reference genome since its initial release. The new T2T-CHM13 reference includes gapless assemblies for all 22 autosomes plus Chromosome X, corrects numerous errors, and introduces nearly 200 million bp of novel sequence containing 2,226 paralogous gene copies, 115 of which are predicted to be protein coding. The newly completed regions include all centromeric satellite arrays and the short arms of all five acrocentric chromosomes, unlocking these complex regions of the genome to variational and functional studies for the first time.

Nurk S, Koren S, Rieh A, Rautiainen M, et al. The complete sequence of a human genome. bioRxiv, 2021.

# Is Perfect Assembly Possible?

Theorem: Perfect assembly possible if *(Gene Myers)*

- i. errors random
- ii. sampling is Poisson
- iii. reads long enough to solve repeats

(Note: e-rate not needed)

- computationally expensive
- depth of coverage ~100x required

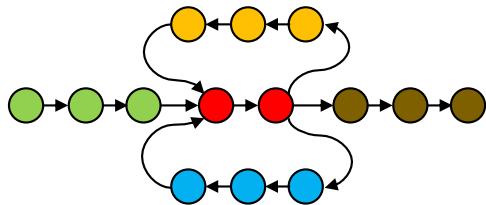
- error correction
  - self-alignment of long reads to correct errors
  - error correction using high accurate illumina reads  
(method based on variable length k-mers)

# Repeats in Genome Assembly

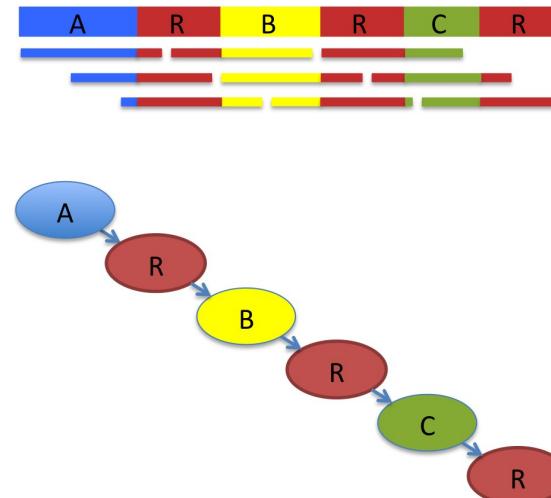
a)



b)



- (a) Genomic region with a triple repeat and 4 unique regions A, B, C, D. Mapped short reads are shown below this region.
- (b) Corresponding de Bruijn graph representation of the genomic region if the triple repeat is longer than the k-mer size k. Multiple paths are possible through the graph.



adapted from Michael Schatz

# Complete Genomes

- for universal applicability technologies should deliver
- extremely long reads → ideally chromosome size
- extremely low-cost data
- extremely fast and real-time
- very flexible instruments instead of monolithic ones
- extremely sensitive / low input requirements
- error free

# European Reference Genome Atlas (ERGA)

- Aim of ERGA is to generate *de novo* genome assemblies that are
  - reference-quality
  - complete
  - error-free
- for species across the whole of European biodiversity
  - includes at least 200,000 species of higher organisms



**Richard Durbin**

*"We are entering a new era of genomics, in which reference sequences for all species will change the way we interact with the whole living world. The ERGA will coordinate European activities in this area and ensure that Europe plays a leading role in this fundamental initiative for 21st century science."*



**Gene W. Myers**

*"An European-wide initiative like the ERGA is absolutely necessary given the scale of the proposed atlas. I believe this atlas will revolutionize conservation science, agricultural and animal science, and even contribute significantly to medicine."*

# European Reference Genome Atlas (ERGA)

- Methods used to achieve the goal:
  - PacBio HiFi Reads
  - HiC reads
- Challenges to sequence > 200k genomes to completion
  - sampling and sampling permissions
  - DNA extraction (HMW DNA)
  - vast sequencing capacity
  - vast compute and storage resources required
  - standards to ensure quality and comparability
  - choice of methods
  - gene annotation not even discussed yet
  - how to assess assembly quality

## Now it is your turn – Your Project in this Course

- Perform all steps of a de novo genome assembly by yourself
- Start are subsampled read data which are from a Nature Communications publication
  - Nat Commun. 2020 Feb 20;11(1):989. doi: 10.1038/s41467-020-14779-y.
- Each student gets their own data set (PacBio reads)
- Each student will write an own report.
  - report has two parts: (i) assembly & (ii) annotation  
(annotation will be performed in the next course by Christian Parisod (6 weeks))
- At the end of this course (last day), each student presents the own work (15 minutes + discussion)