

Next Generation Sequencing

Wouter Bossuyt

Interuniversity Course in Human Genetics

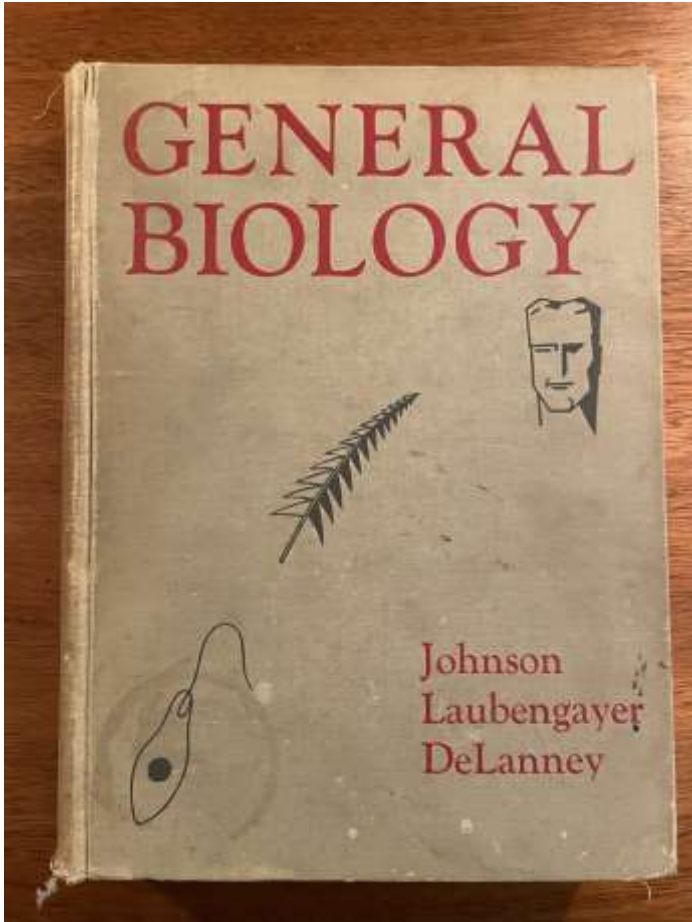
13/10/2023

Sequencing: a history

Landmarks in DNA sequencing

- 1911:
 - Thomas Hunt Morgan disproves himself and find chromosomes as basis of hereditary
- 1944-1952
 - Avery–MacLeod–McCarty experiment (DNA in bacteria)
 - Hershey–Chase experiment (DNA in phages)
- 1953
 - Rosalind Franklin, Watson, Crick: Discovery of DNA double helix structure

1956



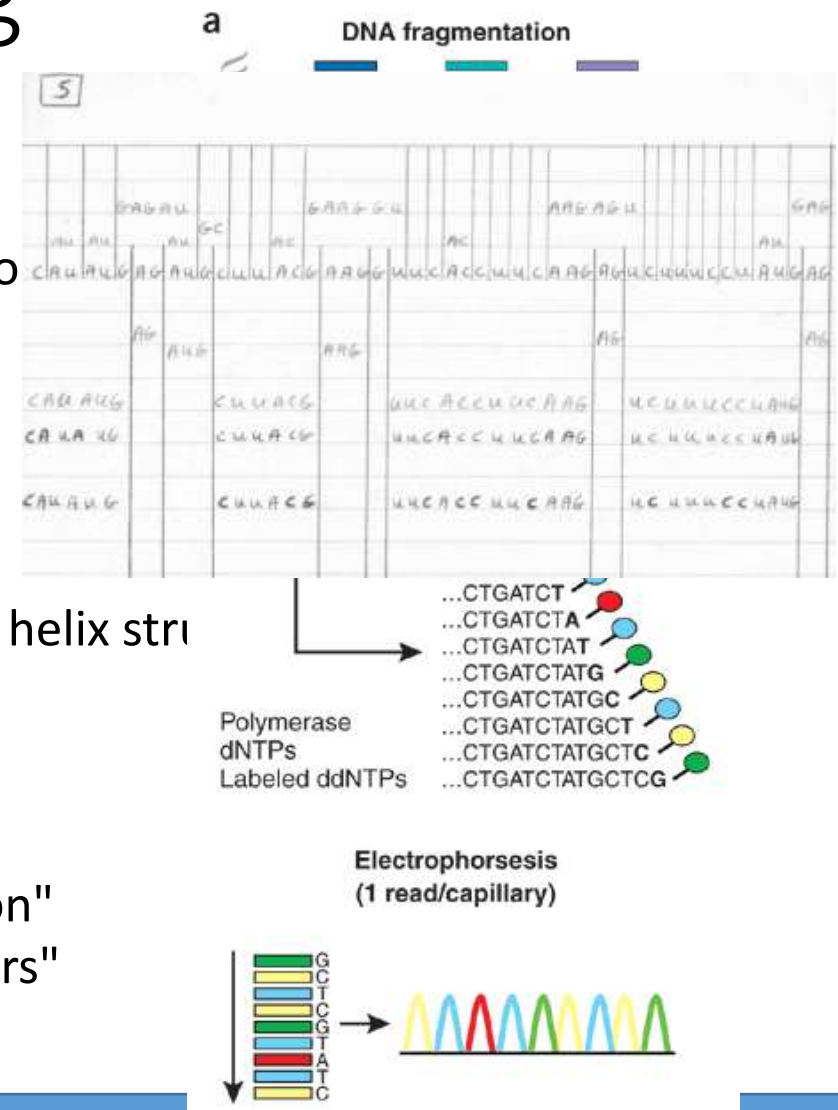
sible. However, a few of the most salient facts that have been discovered and a few of the more plausible hypotheses should be mentioned.

There is evidence from a number of sources that genes are nucleoproteins. Chemical analyses of chromosomes show them to be largely nucleoprotein in nature. A second bit of evidence is that viruses, which are much like genes in certain respects, are nucleoprotein in nature.

range of s
large-sized
probably i
divided th
by the est
some nong
The mo
their abilit
been empl
normal ce
lished no

Landmarks in DNA sequencing

- 1911:
 - Thomas Hunt Morgan disproves himself and find chromoso
- 1944-1952
 - Avery–MacLeod–McCarty experiment (DNA in bacteria)
 - Hershey–Chase experiment (DNA in phages)
- 1953
 - Rosalind Franklin, Watson, Crick: Discovery of DNA double helix stru
- 1953-1977: the 'desperate' era
 - Walter Fiers use RNase digest in competition with Sanger
- 1977
 - A Maxam and W Gilbert "DNA seq by chemical degradation"
 - F Sanger "DNA sequencing with chain-terminating inhibitors"



Landmarks in DNA sequencing

- 1984
 - DNA sequence of the Epstein-Barr virus, 172 kb
- 1987
 - Applied Biosystems - first automated sequencer
- 1991
 - Sequencing of human genome in Venter's lab
- 1996
 - P. Nyrén and M Ronaghi - pyrosequencing
- 2001
 - A draft sequence of the human genome
- 2003
 - human genome completed
- 2004
 - 454 Life Sciences markets first NGS machine



ARTICLES

Genome sequencing in microfabricated high-density picolitre reactors

Marcel Margulies^{1*}, Michael Egholm^{1*}, William E. Altman¹, Said Attiya¹, Joel S. Bader¹, Lisa A. Bembgen¹, Jan Berka¹, Michael S. Braverman¹, Yi-Ju Chen¹, Zhoutao Chen¹, Scott B. Dewell¹, Lei Du¹, Joseph M. Fierro¹, Xavier V. Gomes¹, Brian C. Godwin¹, Wen He¹, Scott Helgesen¹, Chun He Ho¹, Gerard P. Irzyk¹, Szilveszter C. Jando¹, Maria L. I. Alenquer¹, Thomas P. Jarvie¹, Kshama B. Jirage¹, Jong-Bum Kim¹, James R. Knight¹, Janna R. Lanza¹, John H. Leamon¹, Steven M. Lefkowitz¹, Ming Lei¹, Jing Li¹, Kenton L. Lohman¹, Hong Lu¹, Vinod B. Makhijani¹, Keith E. McDade¹, Michael P. McKenna¹, Eugene W. Myers², Elizabeth Nickerson¹, John R. Nobile¹, Ramona Plant¹, Bernard P. Puc¹, Michael T. Ronan¹, George T. Roth¹, Gary J. Sarkis¹, Jan Fredrik Simons¹, John W. Simpson¹, Maithreyan Srinivasan¹, Karrie R. Tartaro¹, Alexander Tomasz³, Kari A. Vogt¹, Greg A. Volkmer¹, Shally H. Wang¹, Yong Wang¹, Michael P. Weiner⁴, Pengguang Yu¹, Richard F. Begley¹ & Jonathan M. Rothberg¹

The proliferation of large-scale DNA-sequencing projects in recent years has driven a search for alternative methods to reduce time and cost. Here we describe a scalable, highly parallel sequencing system with raw throughput significantly greater than that of state-of-the-art capillary electrophoresis instruments. The apparatus uses a novel fibre-optic slide of individual wells and is able to sequence 25 million bases, at 99% or better accuracy, in one four-hour run. To achieve an approximately 100-fold increase in throughput over current Sanger sequencing technology, we have developed an emulsion method for DNA amplification and an instrument for sequencing by synthesis using a pyrosequencing protocol optimized for solid support and picolitre-scale volumes. Here we show the utility, throughput, accuracy and robustness of this system by shotgun sequencing and *de novo* assembly of the *Mycoplasma genitalium* genome with 96% coverage at 99.96% accuracy in one run of the machine.

DNA Sequencing – the next generation

- NGS refers to non-Sanger-based high-throughput DNA sequencing technologies.
- NGS technologies constitute various strategies that rely on a combination of
 - Library/template preparation
 - Parallel sequencing

Different technologies

- Illumina
- Nanopore
- Pacbio
- MGI
- AVITI Elements
- Ultima
- ...

Illumina



Nanopore



AVITI Element



MGI



Pacbio



Ultima Genomics

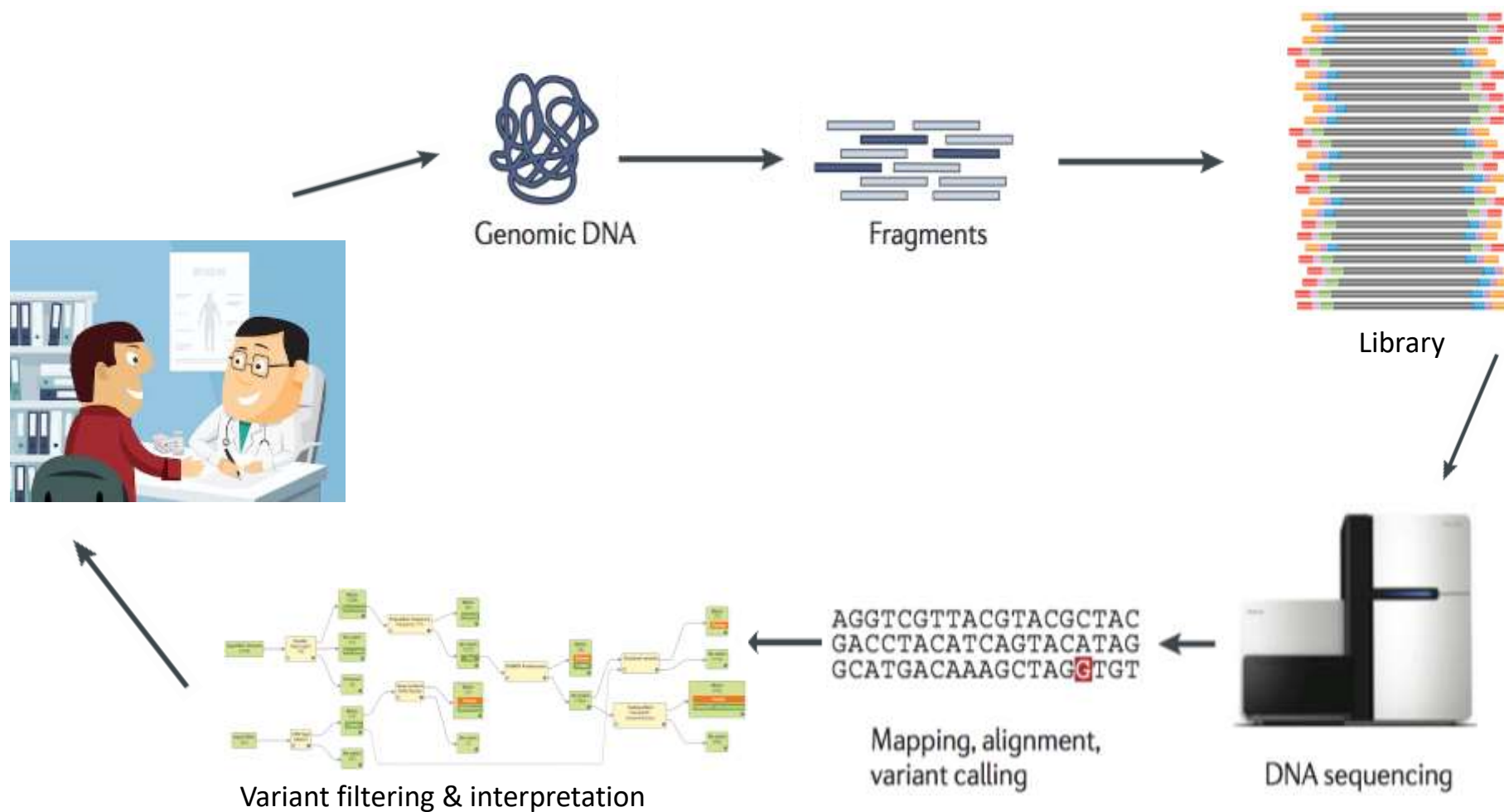


List of NGS technologies and their specs:

<https://docs.google.com/spreadsheets/d/1GMMfhyLK0-q8Xklo3YxlWaZA5vVMuhU1kg41g4xLkXc/edit?usp=sharing>

Sequencing workflows

NGS workflow



What do I chose? Long read, short read,...

	illumina	Nanopore	Pacbio
Read length	35 bp to 600 bp	Anything goes	250bp to 25kb (or 100kb)
Accuracy	High	Medium and improving	high
Capacity	Small to very large	Tiny to large	medium
Biases	Fragment size and GC		DNA modifications
Applications			
WGS	+++	++	+++
RNAseq	+++	+	++
Targeted resequencing	+++	Only large fragments	++
Single cell sequencing	+++		

Illumina (solexa) sequencing

- Illumina MiSeq, NextSeq 500, Nextseq2000, HiSeq4000 & NovaSeq 6000



Illumina flowcells



MiSeq



NextSeq500



NextSeq2000

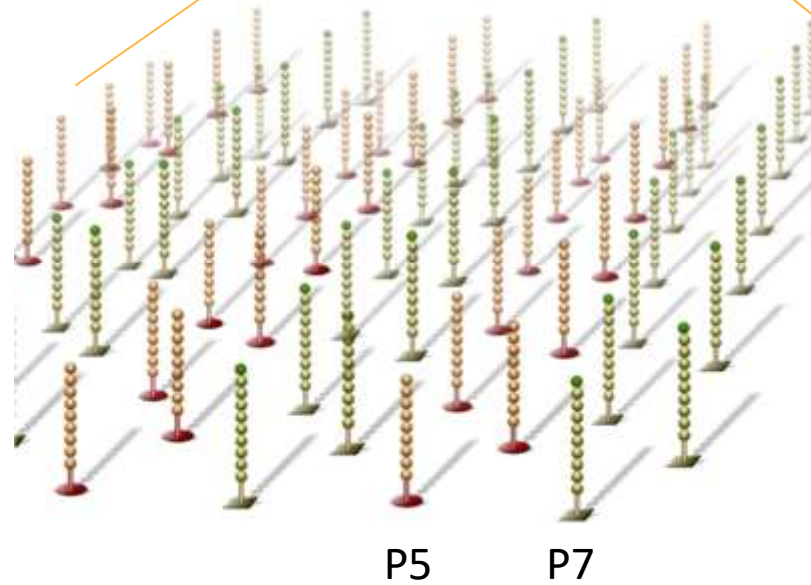
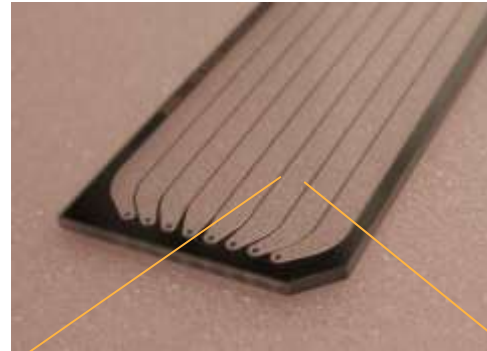


NovaSeq



HiSeq4000

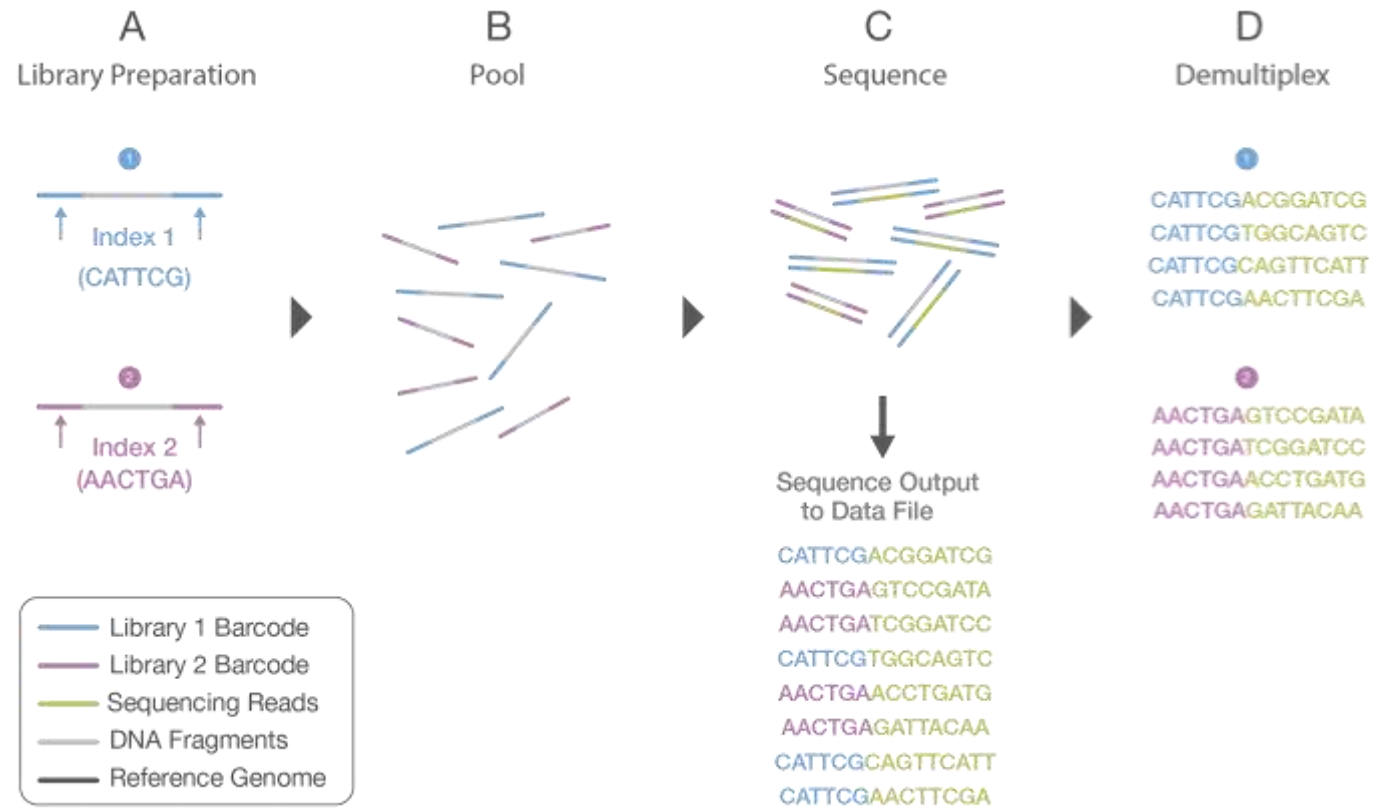
NGS Library Prep



NGS Library Prep

- Indexing

- Sample barcodes
- High diversity necessary
- Unique dual indexing is top



NGS Library Prep

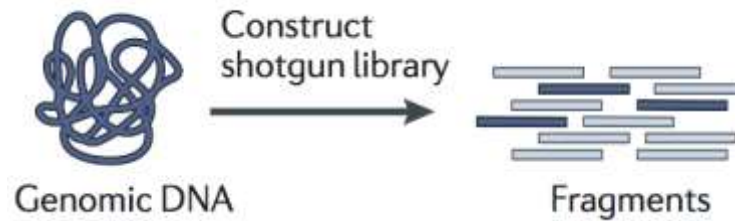
- General DNA library prep
- General RNA library prep
- Targeted library prep approaches

NGS Library Prep

- General DNA library prep
- General RNA library prep
- Targeted library prep approaches

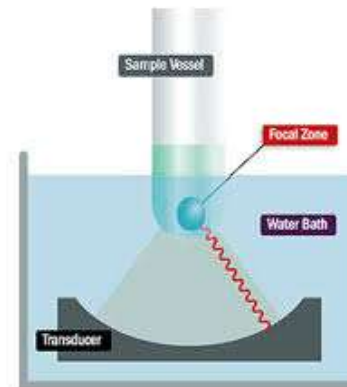
NGS Library prep

- Fragmentation



- acoustic

covaris

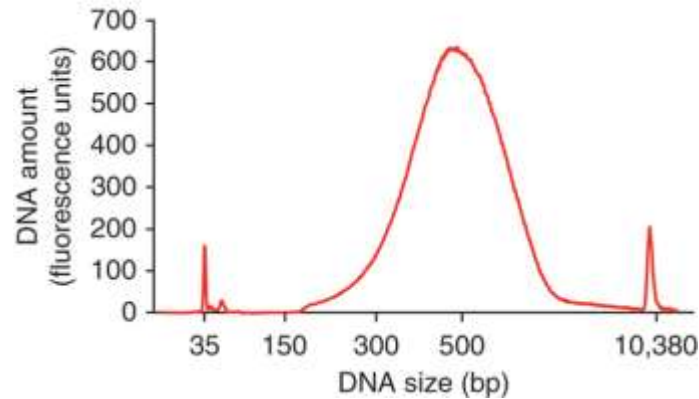


- enzymatic

NGS Library prep

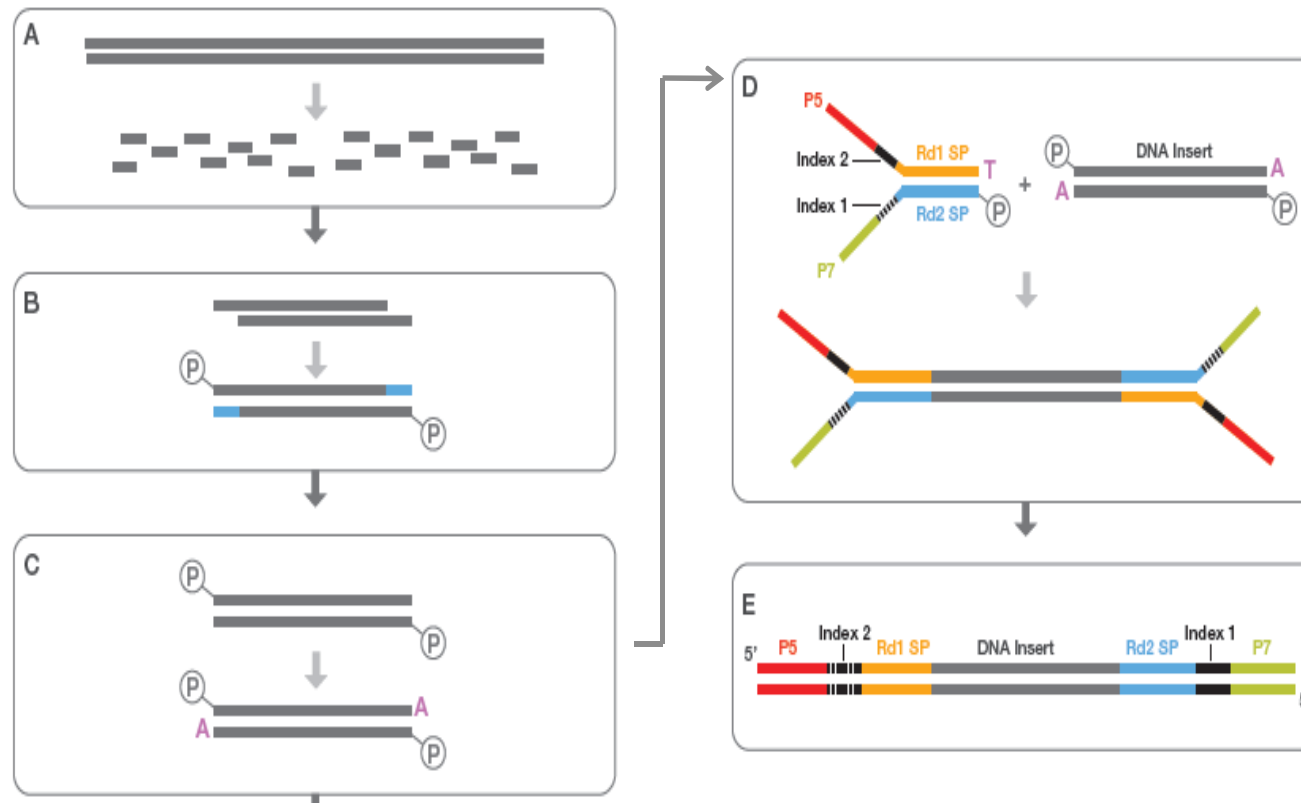
- Fragmentation

- quality check: BioAnalyser / FragmentAnalyser



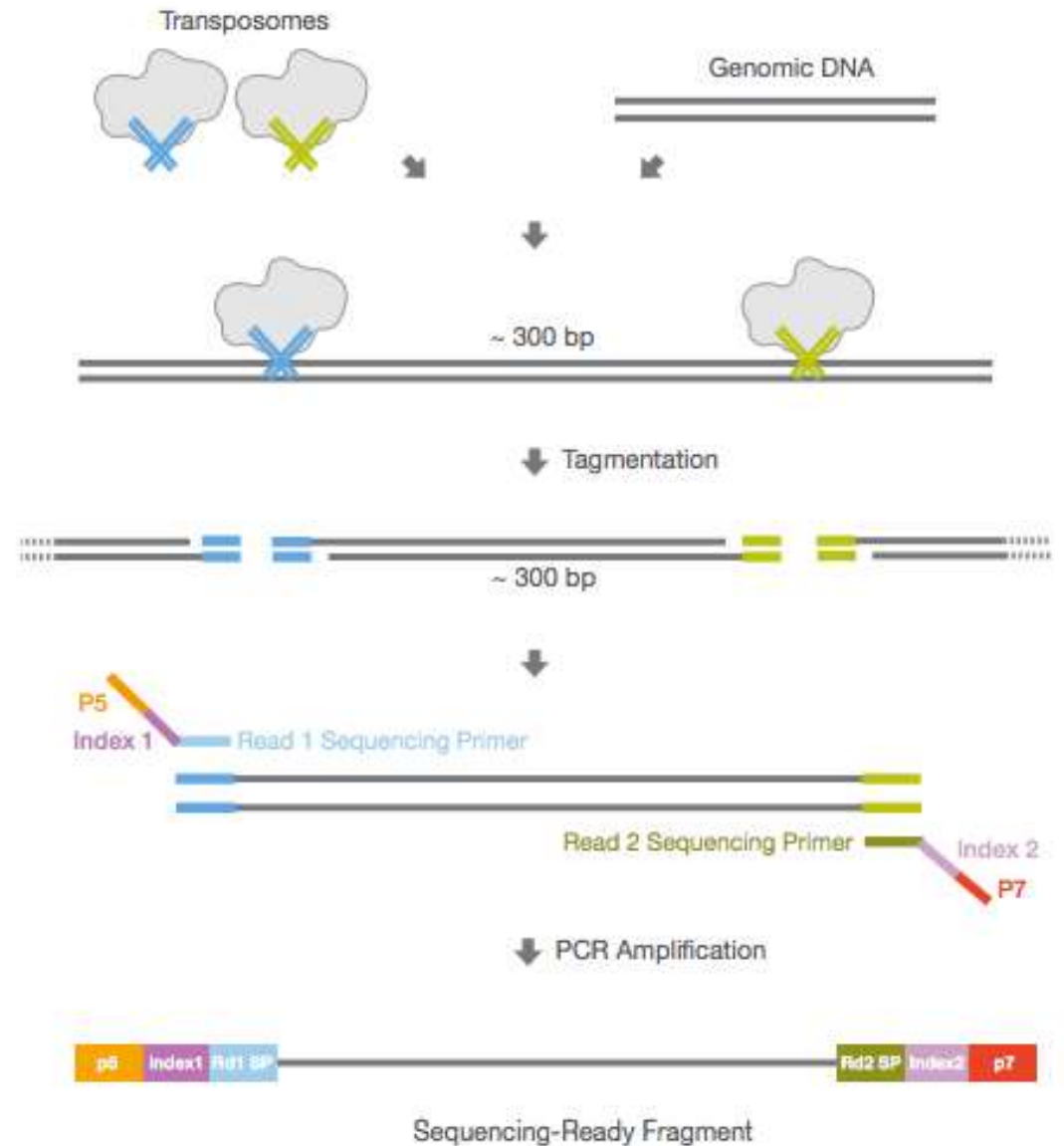
NGS Library prep

- Standard A-tailing & adaptor-ligation (DNA/RNA)



NGS Library Prep

- Illumina Nextera tagmentation
- Transposon-based adapter insertion
- PCR-based indexing



NGS Library Prep

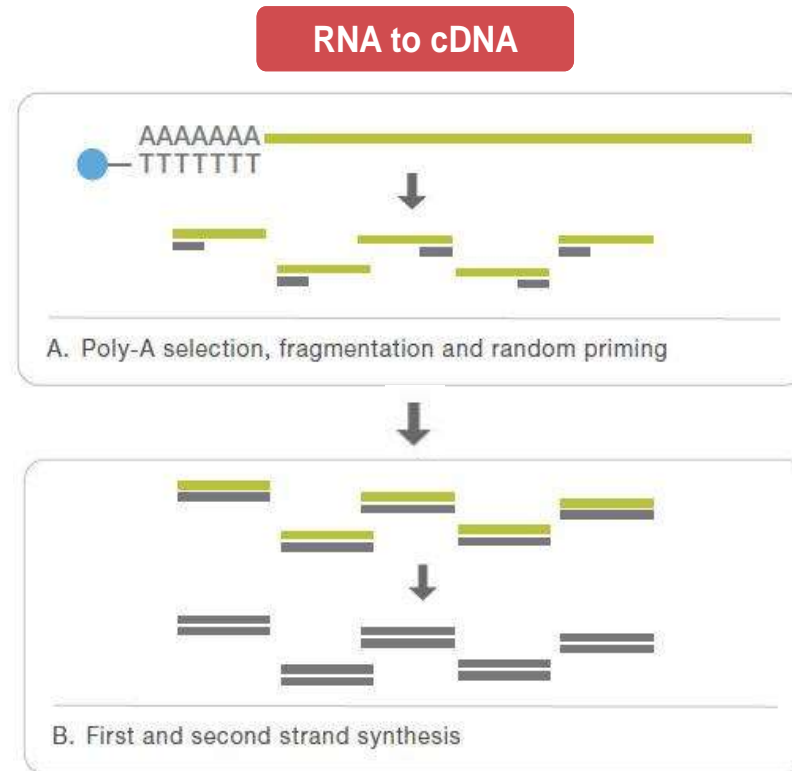
- General DNA library prep
- General RNA library prep
- Targeted library prep approaches

NGS Library Prep

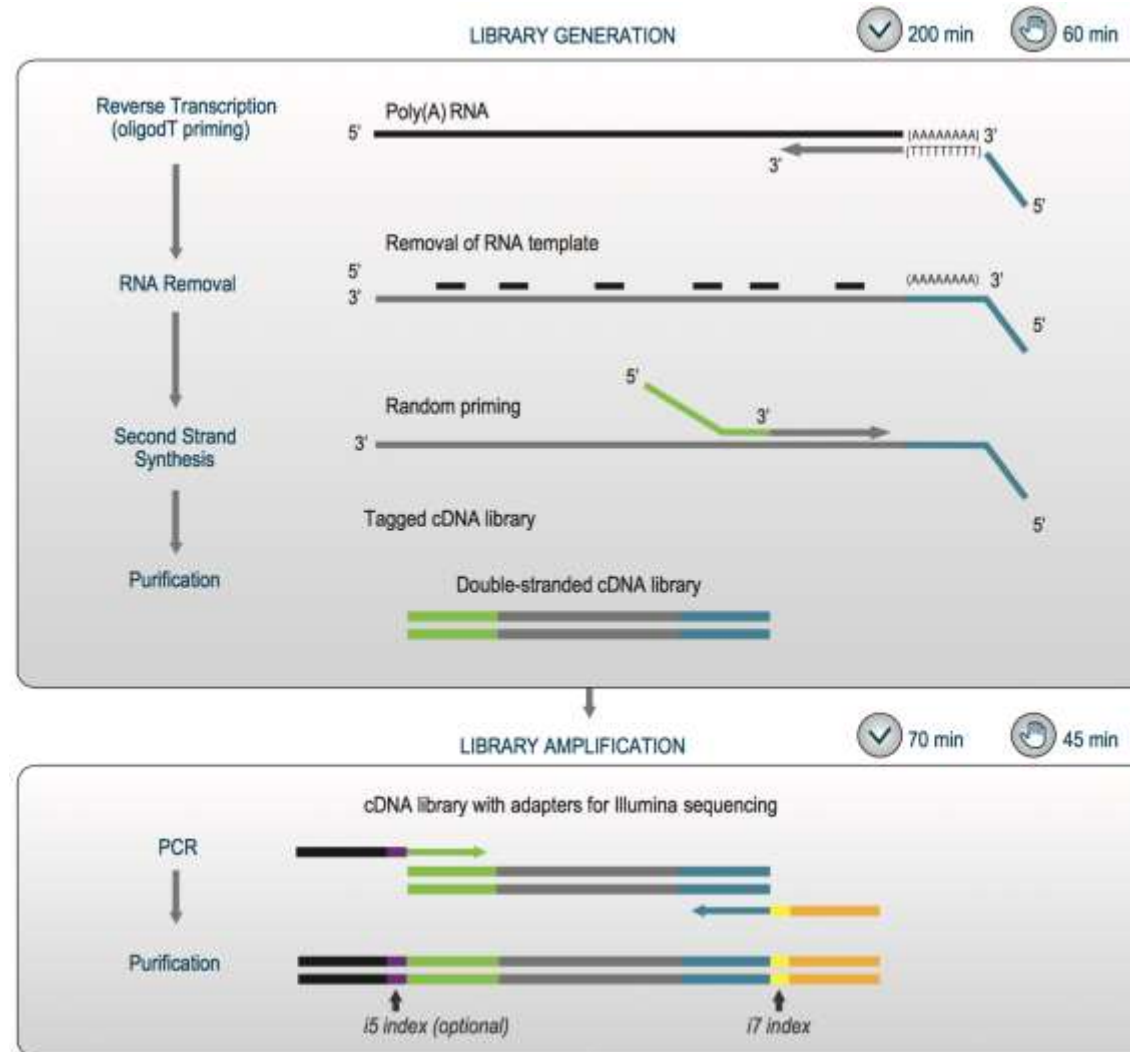
- General DNA library prep
- General RNA library prep
- Targeted library prep approaches

NGS Library Prep

- Illumina TruSeq Stranded mRNA



LexoGen QuantSeq (RNA)



NGS Library Prep

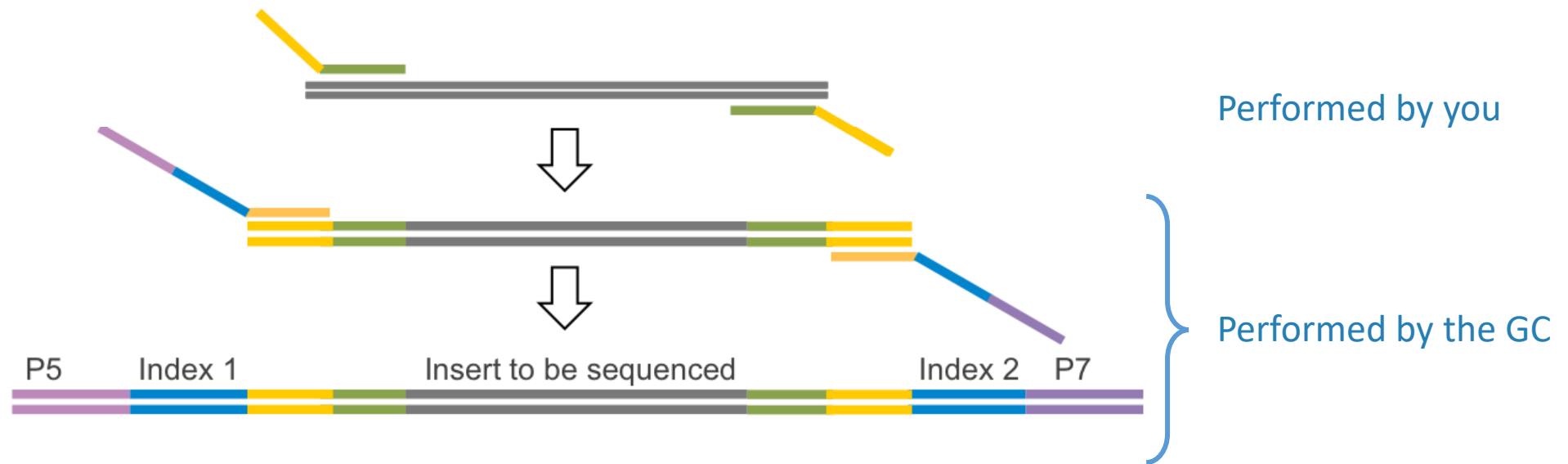
- General DNA library prep
- General RNA library prep
- Targeted library prep approaches

NGS Library Prep

- General DNA library prep
- General RNA library prep
- Targeted library prep approaches
 - amplicon based
 - ligation based
 - enrichment based

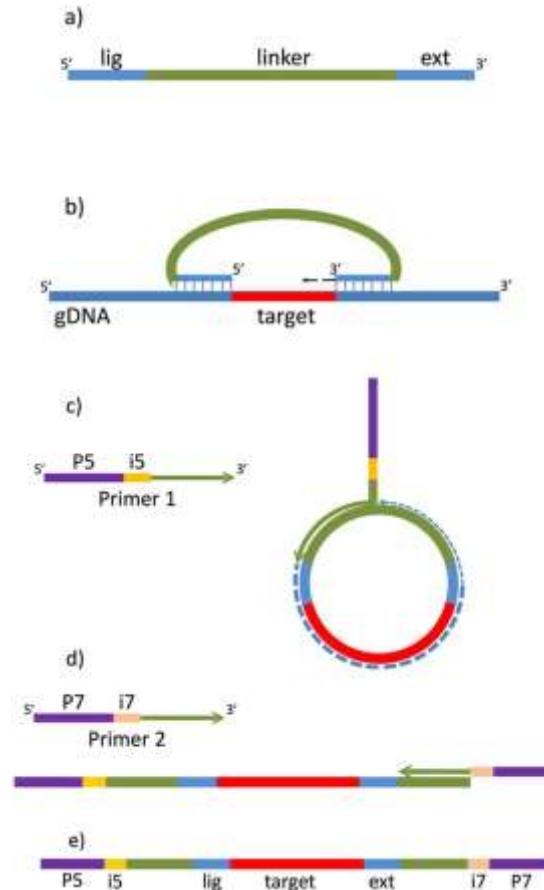
NGS Library Prep

- Custom two step PCR



NGS Library Prep

- MIPs



Multiplex probes

Hybridization and fill-in reaction

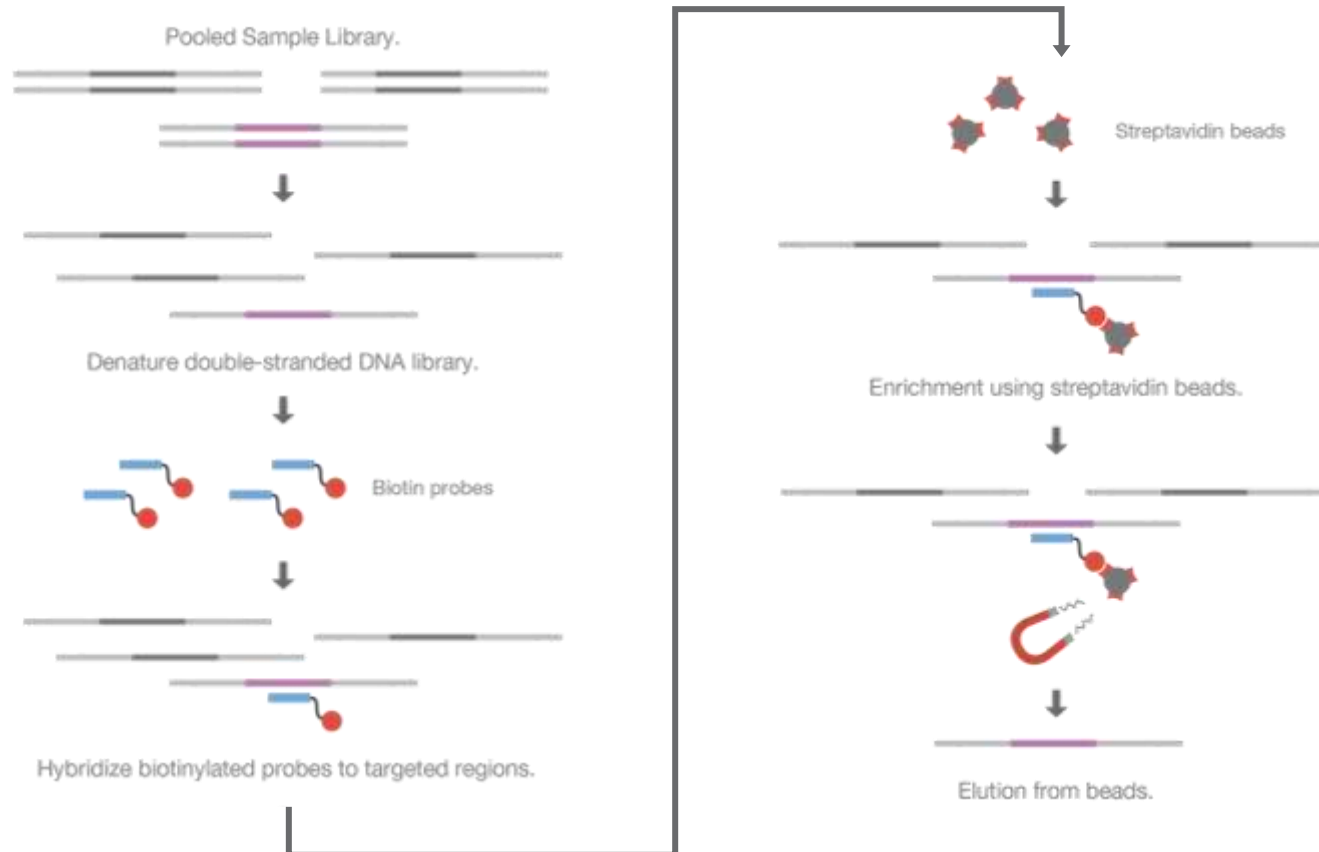
PCR linearization

PCR amplification

Niedzicka *et al.* 2016

NGS Enrichment

- Sequence capture



Illumina sequencing

Sample prep

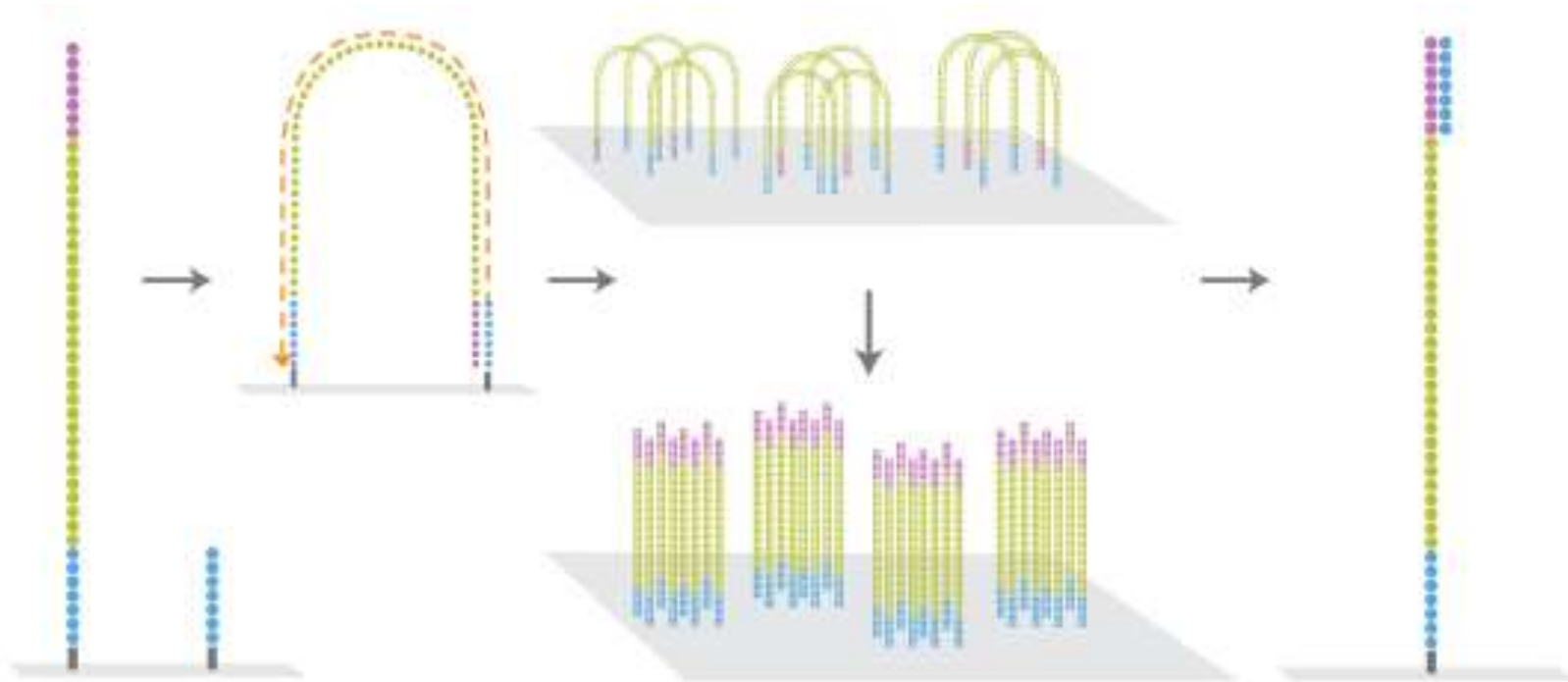


Clonal Amplification



Parallel sequencing

NGS Illumina Clustering



Sample prep



Clonal Amplification



Parallel sequencing

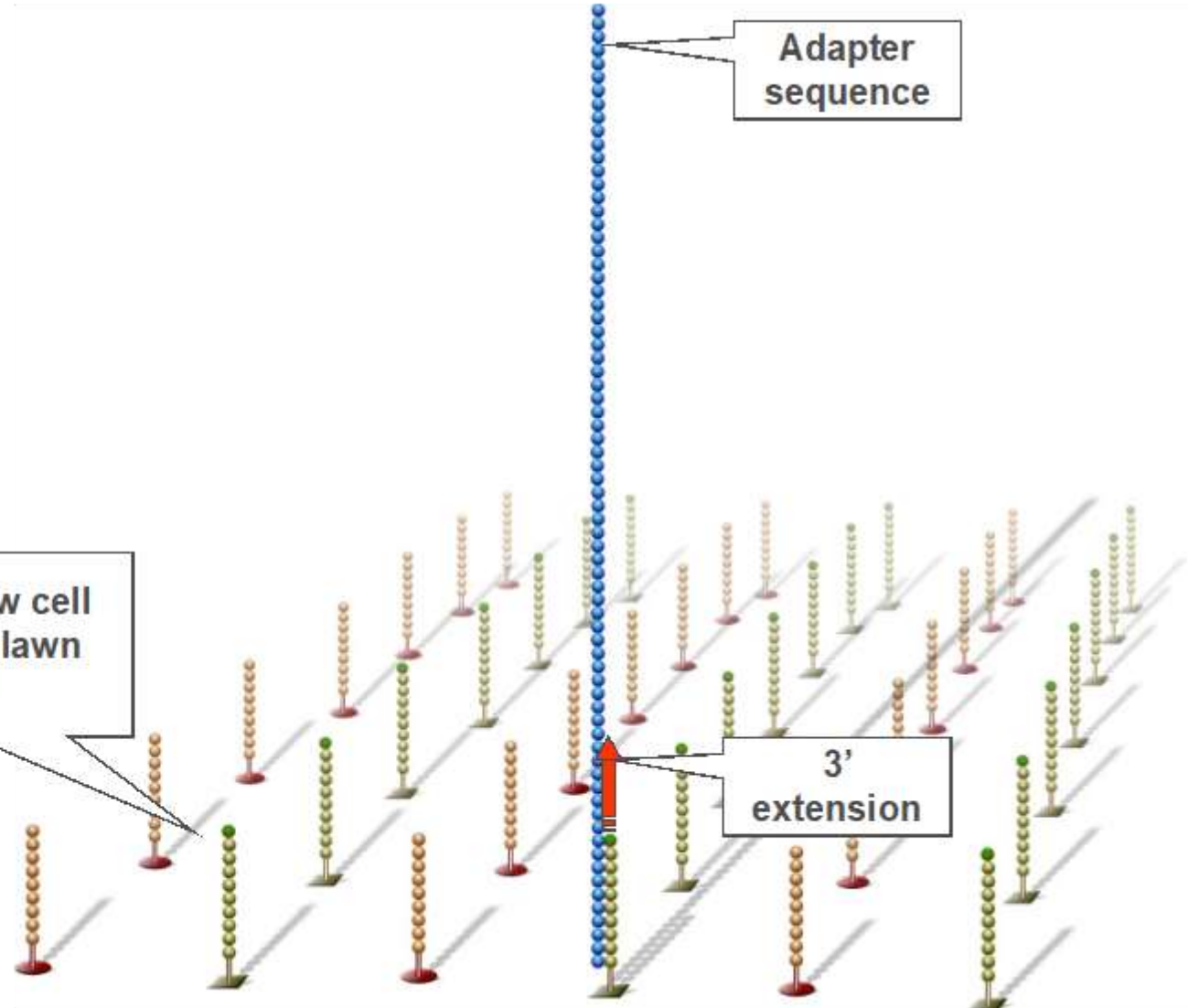
Single DNA libraries are hybridized to primer lawn

Bound libraries then extended by polymerases

Surface of flow cell coated with a lawn of oligo pairs

Adapter sequence

3' extension



Sample prep

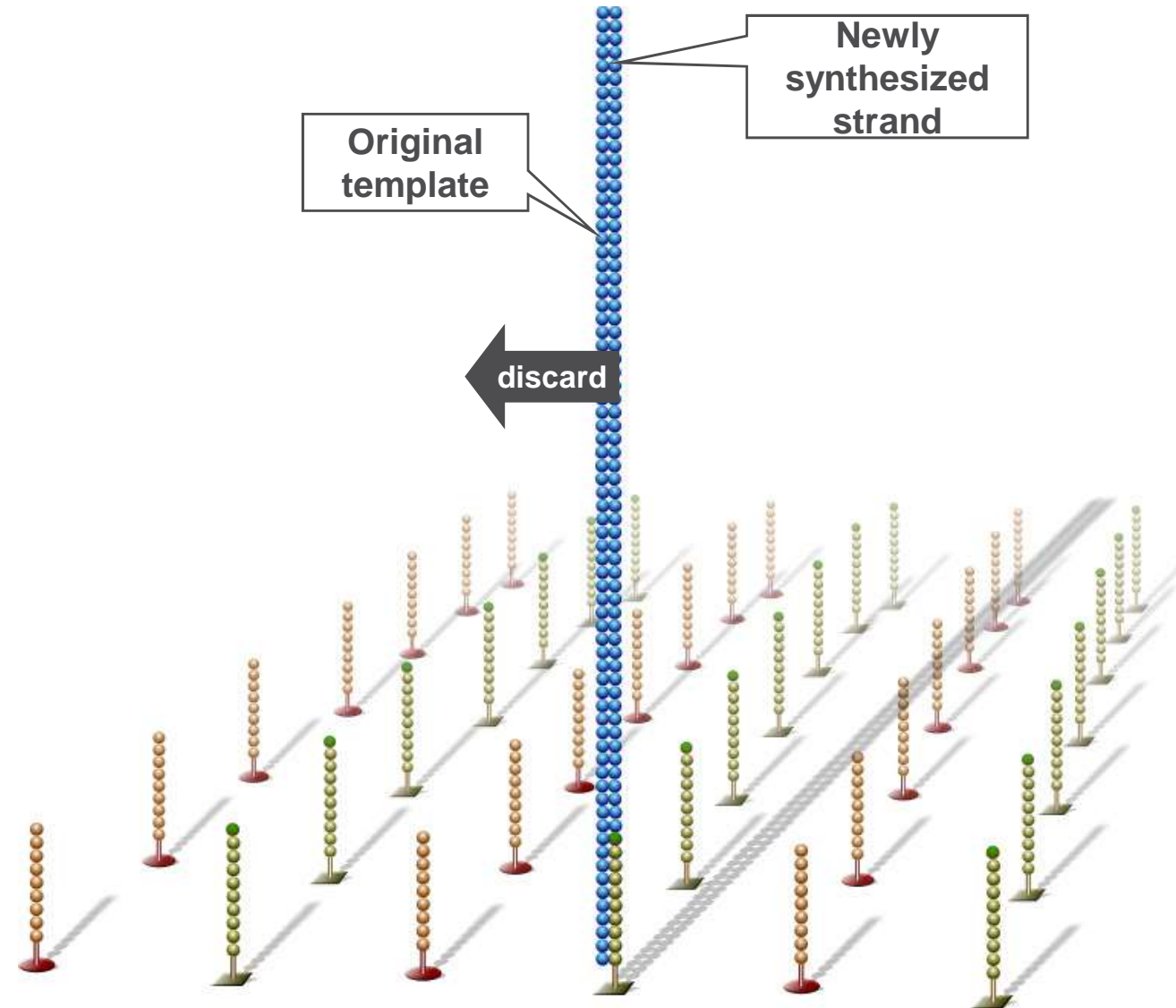
Clonal Amplification

Parallel sequencing

Double-stranded molecule is denatured

Original template washed away

Newly synthesized strand is covalently attached to flow cell surface



Sample prep



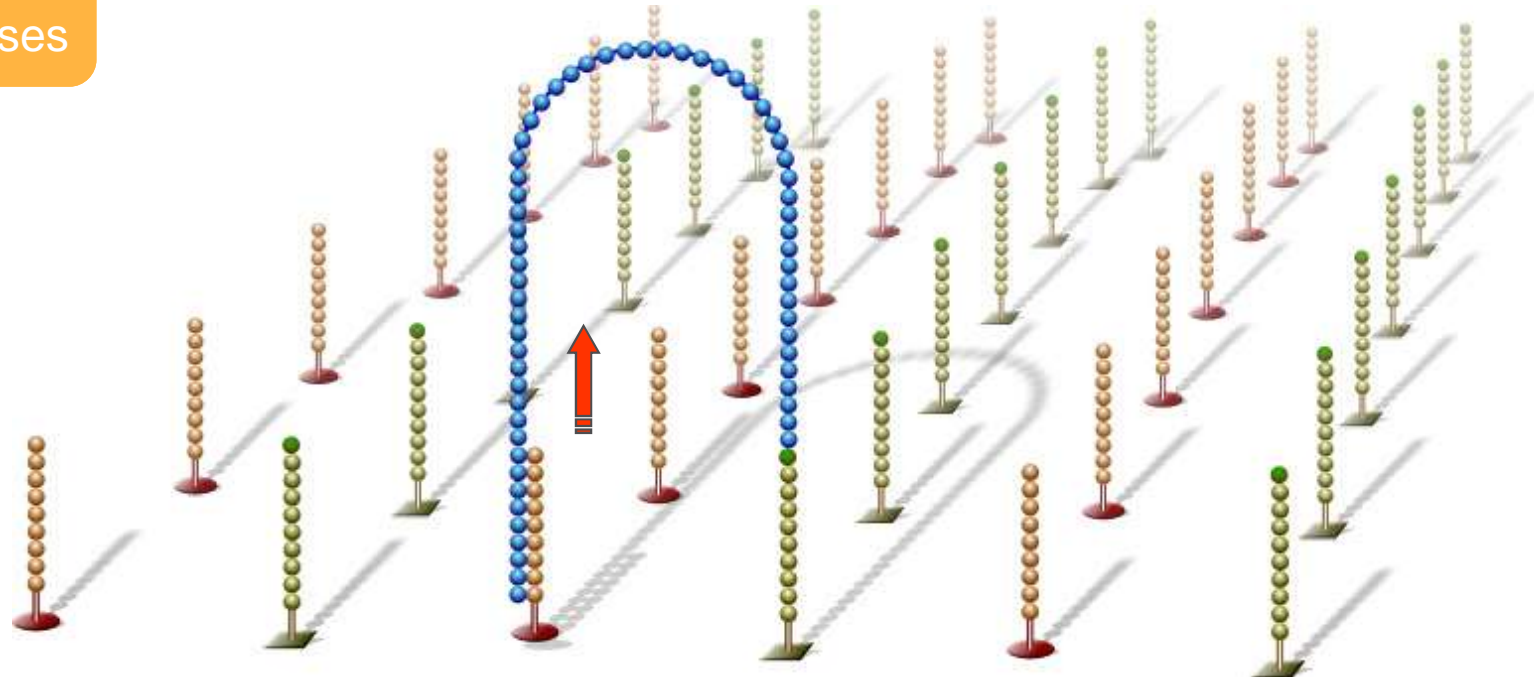
Clonal Amplification



Parallel sequencing

Single-stranded molecule flips over and forms a bridge by hybridizing to adjacent, complementary primer

Hybridized primer is extended by polymerases



Sample prep

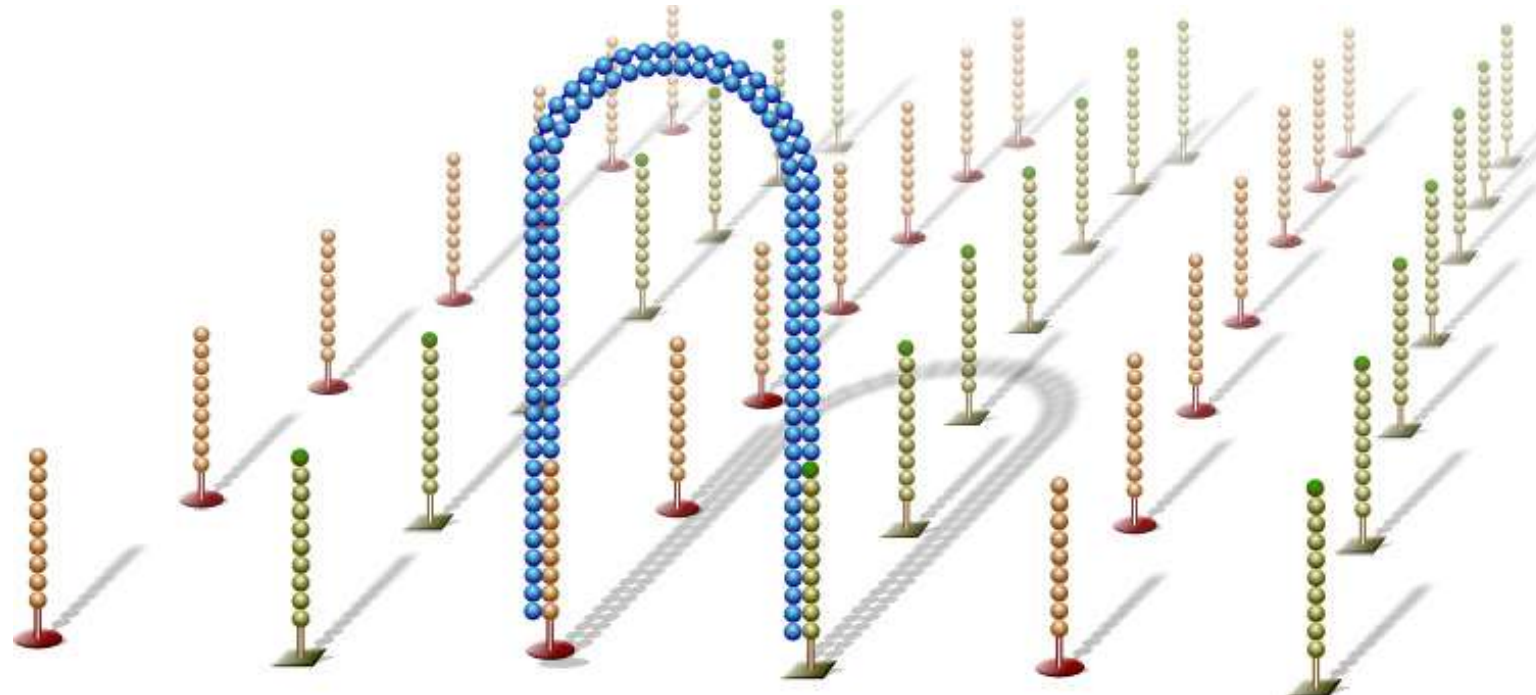


Clonal Amplification



Parallel sequencing

Double-stranded bridge is formed



Sample prep



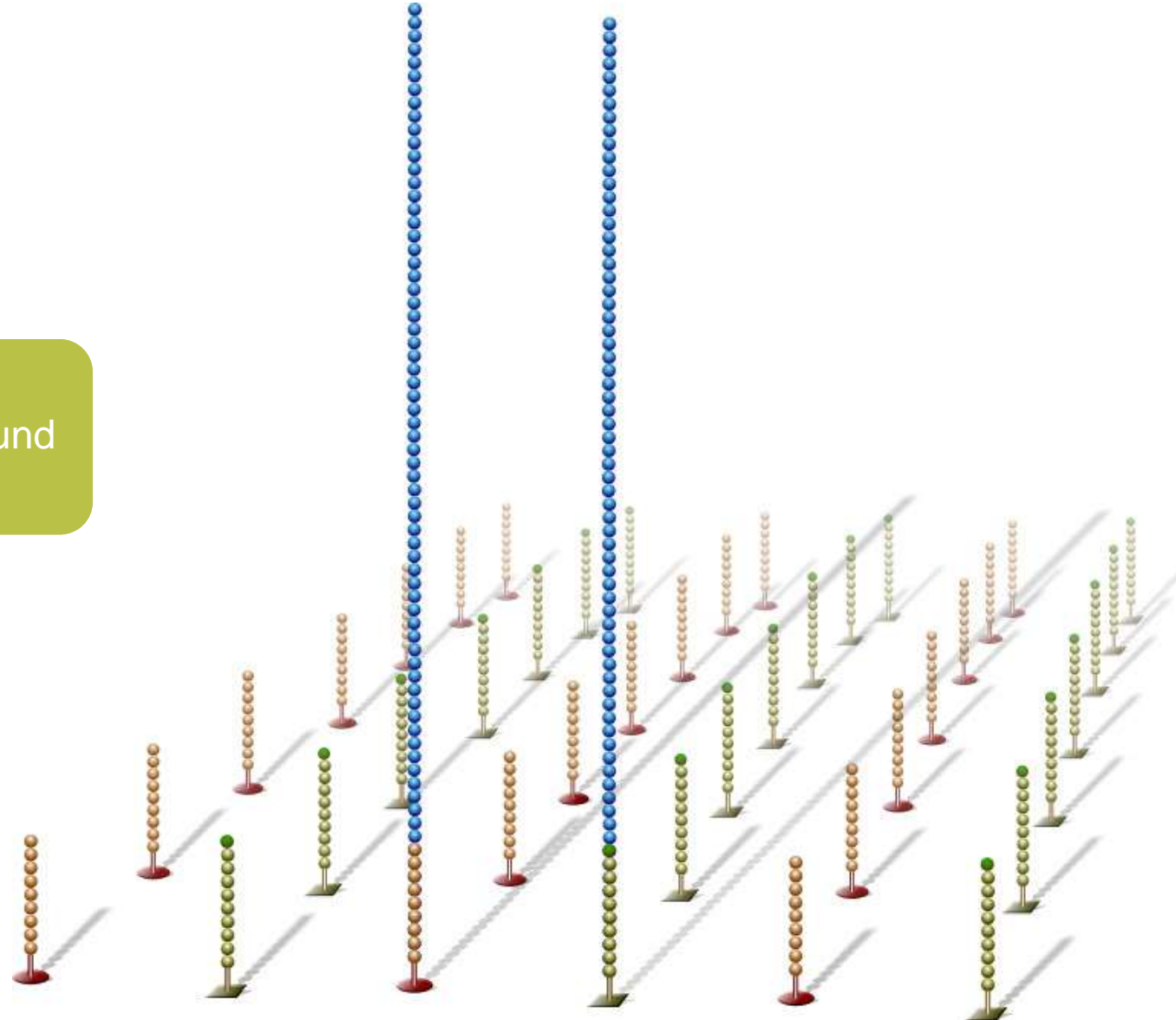
Clonal Amplification



Parallel sequencing

Double-stranded bridge is denatured

Result:
Two copies of covalently bound single-stranded templates



Sample prep



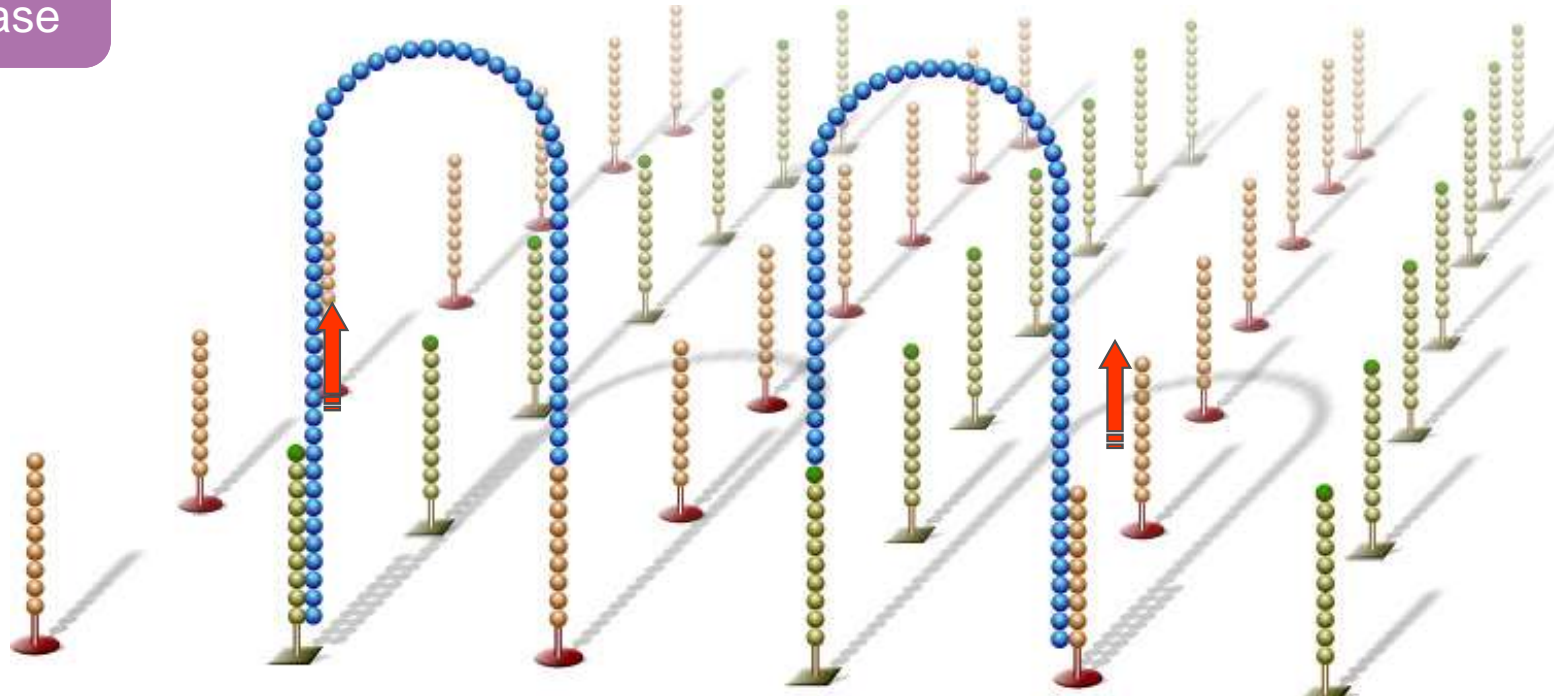
Clonal Amplification



Parallel sequencing

Single-stranded molecules flip over to hybridize to adjacent primers

Hybridized primer is extended by polymerase



Sample prep

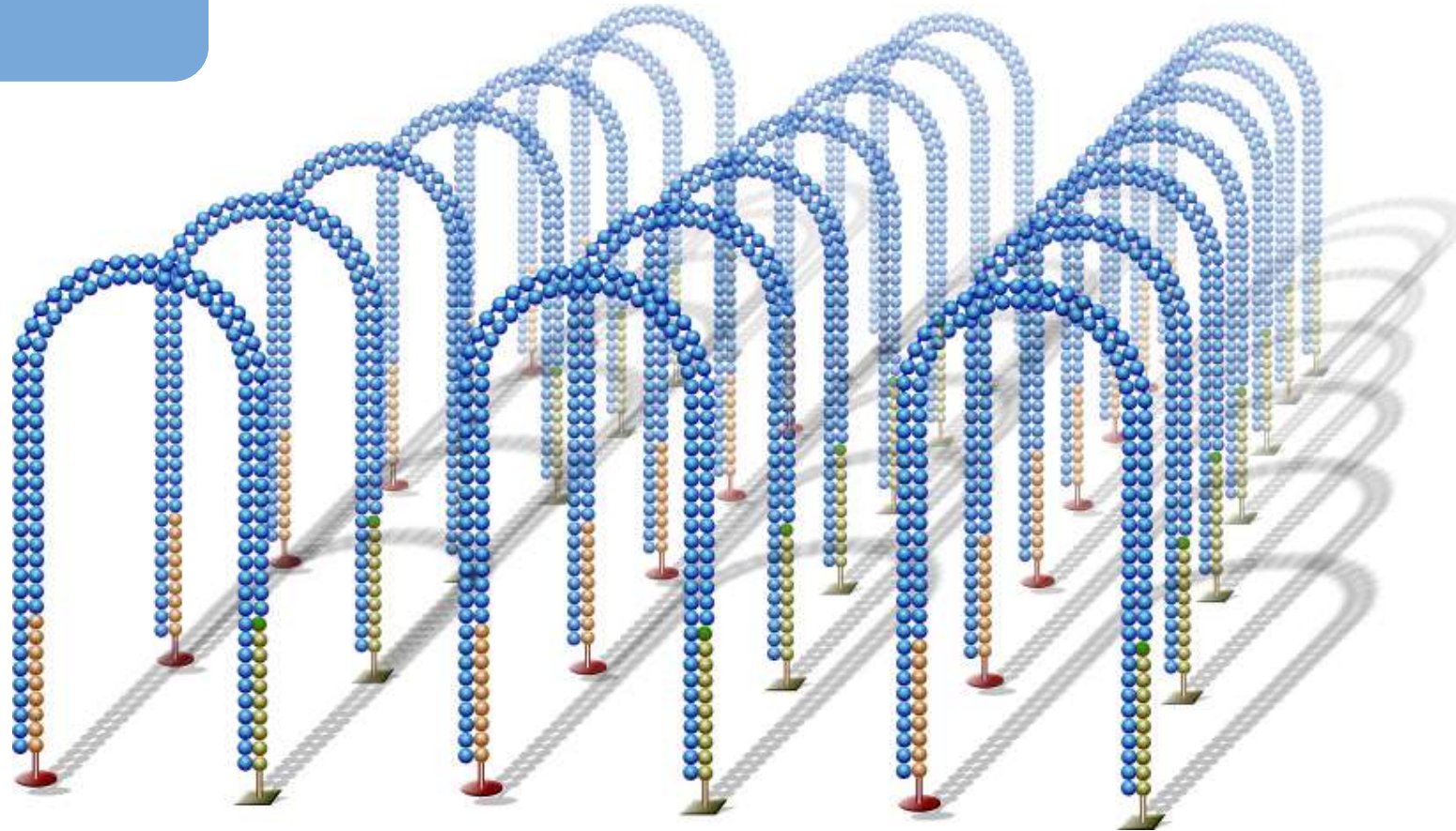


Clonal Amplification



Parallel sequencing

Bridge amplification cycle repeated until multiple bridges are formed



Sample prep

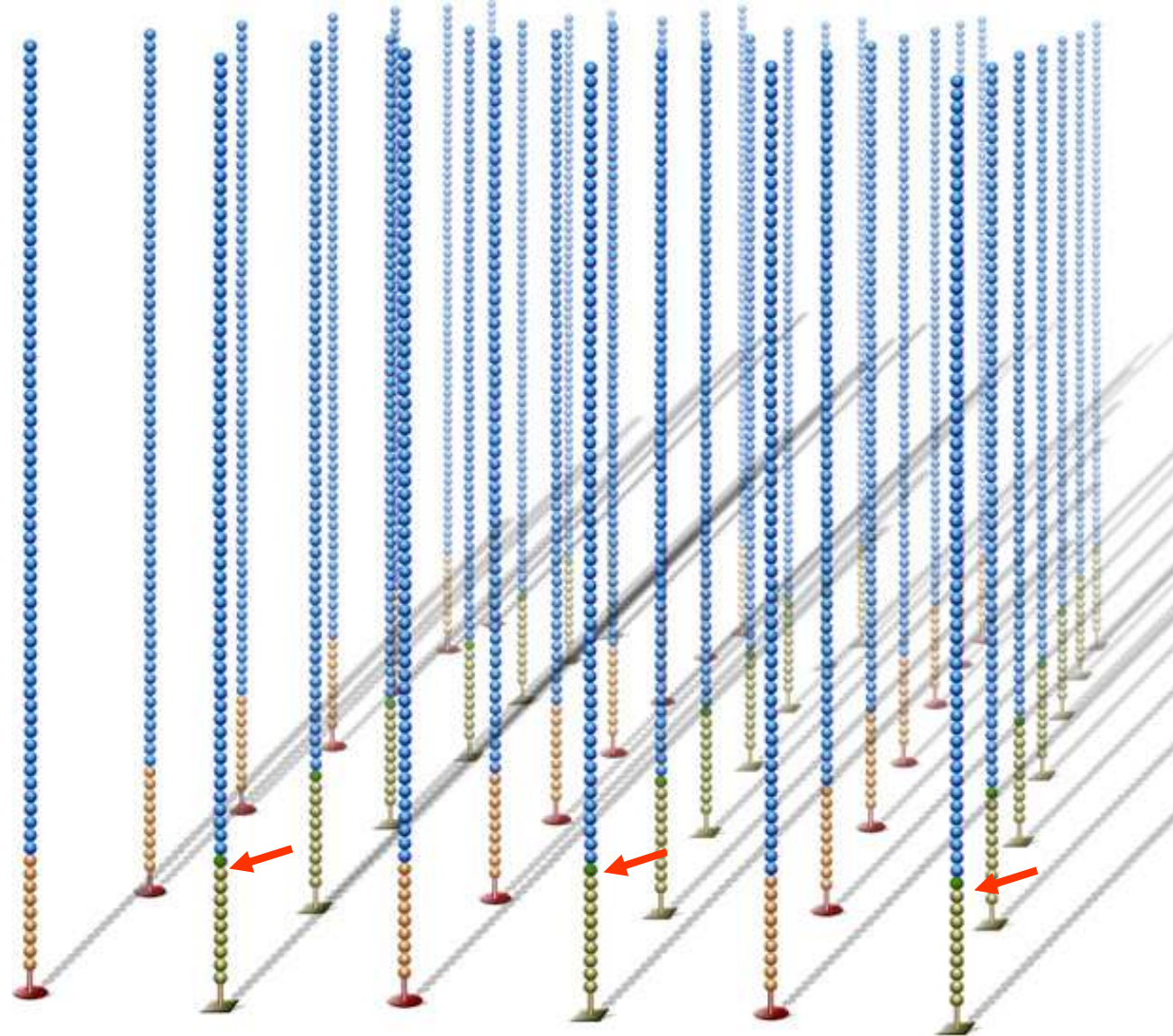


Clonal Amplification



Parallel sequencing

dsDNA bridges are denatured



Sample prep

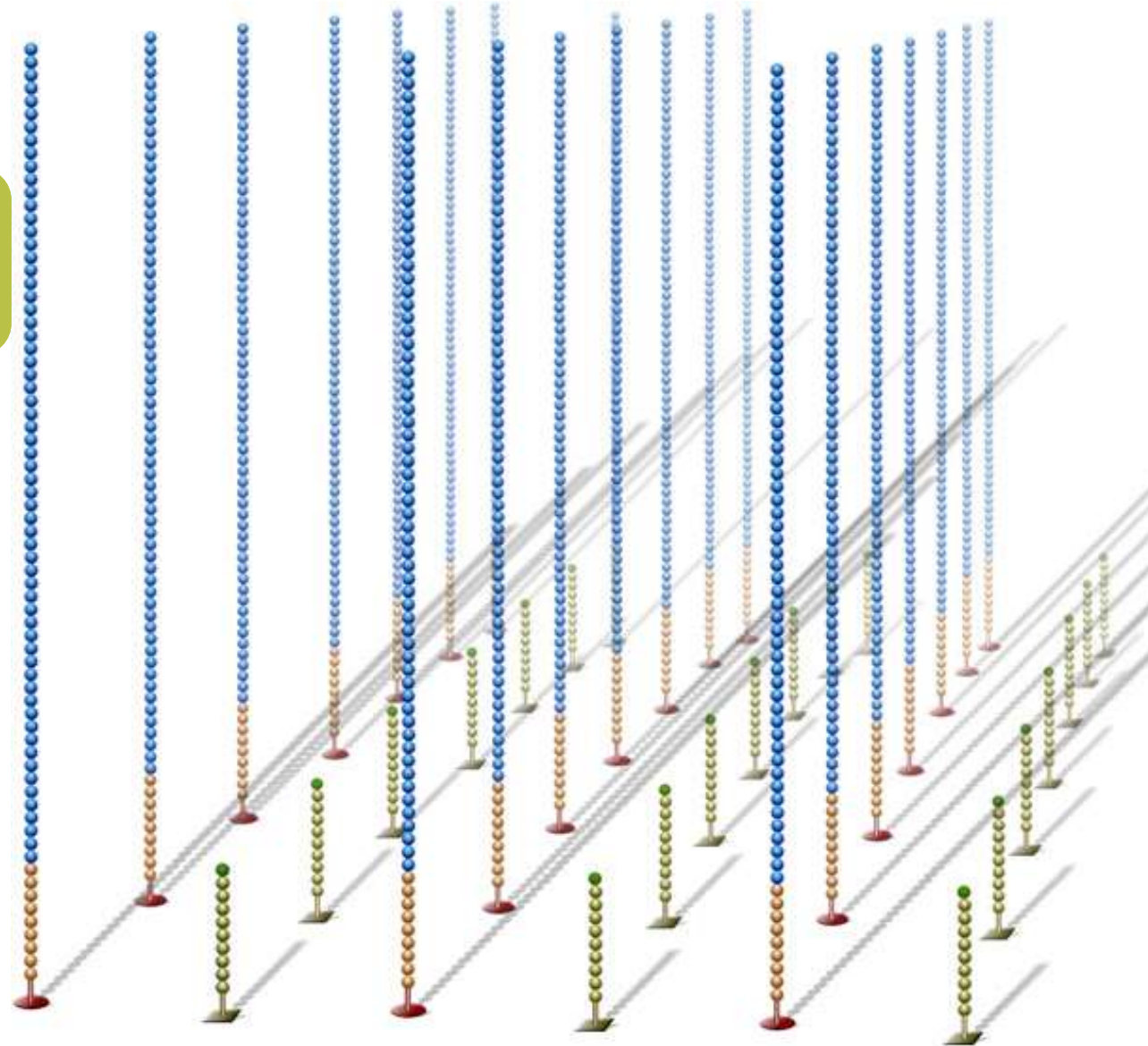


Clonal Amplification



Parallel sequencing

Reverse strands cleaved and washed away, leaving a cluster with forward strands only



Sample prep

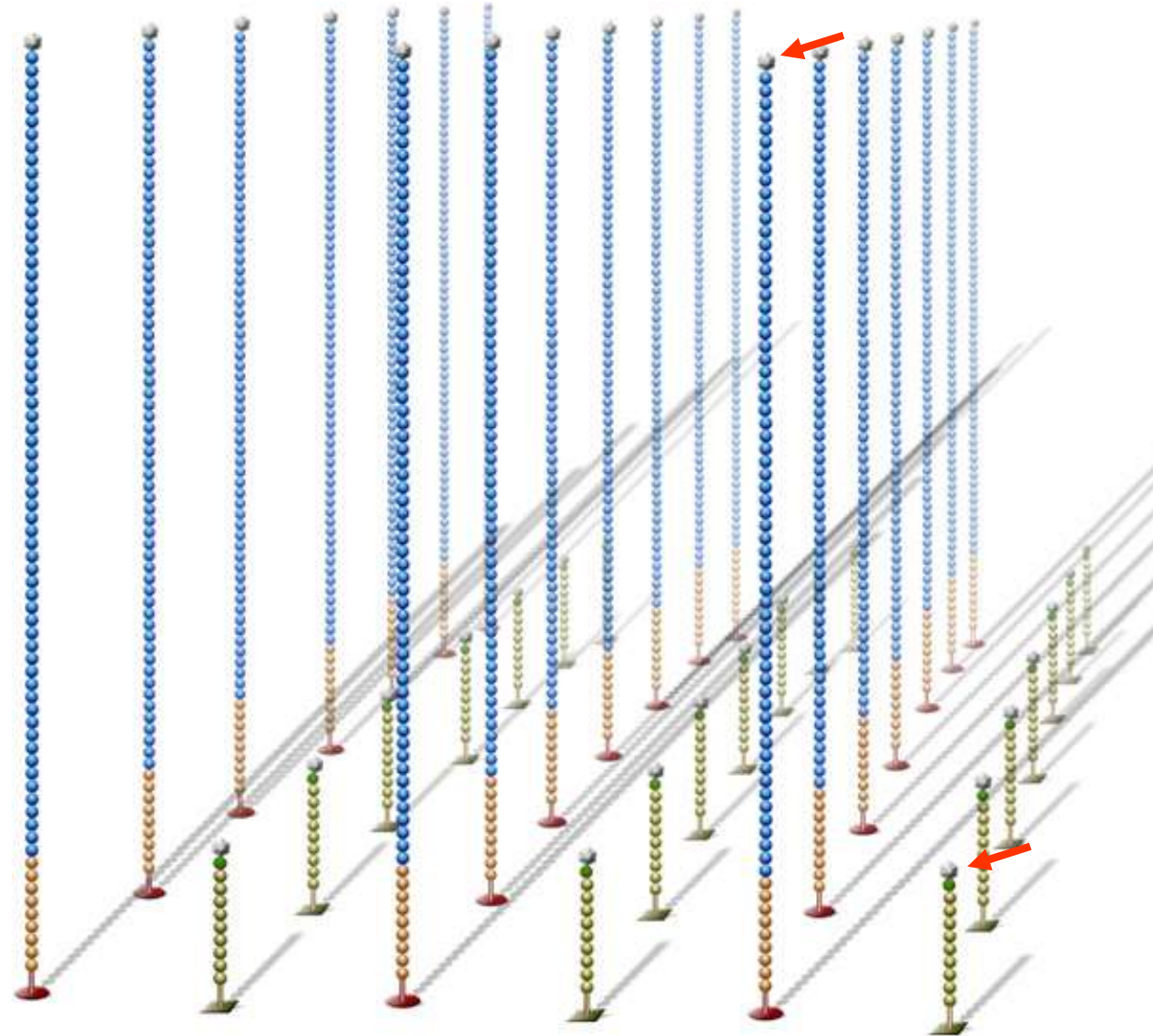


Clonal Amplification



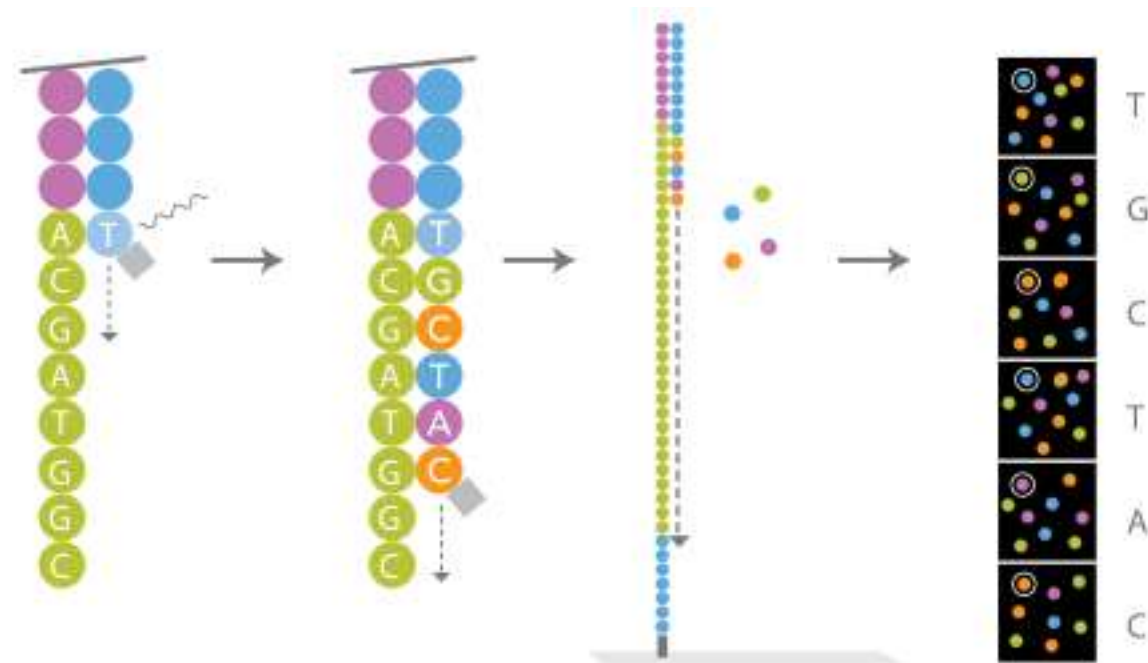
Parallel sequencing

Free 3' ends are blocked to prevent unwanted DNA priming



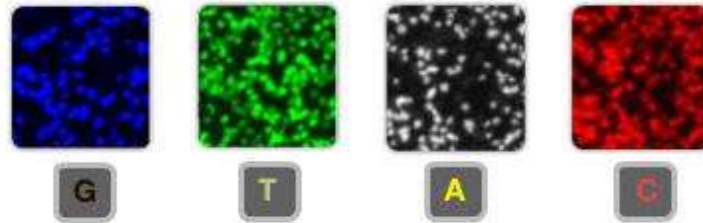
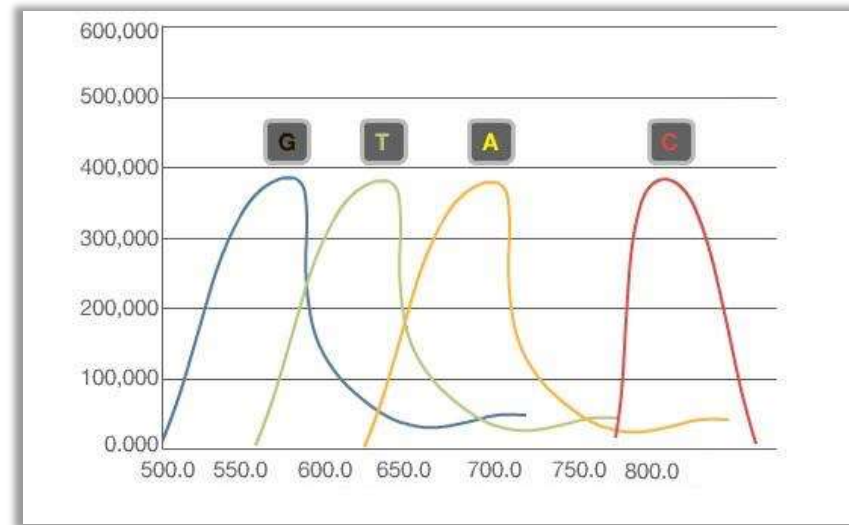
NGS Illumina Sequencing

- Sequencing By Synthesis (SBS)



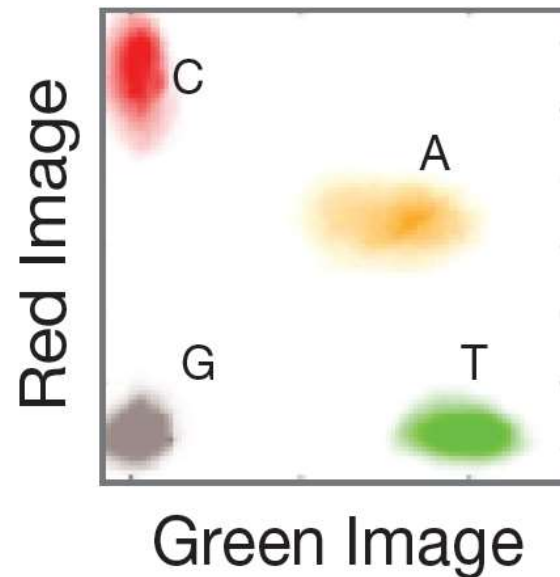
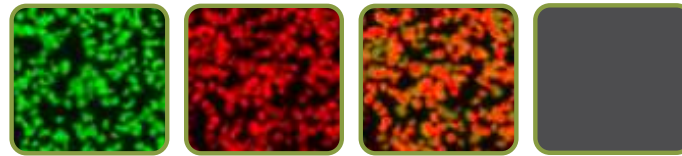
NGS Illumina Sequencing

- Imaging MiSeq, HiSeq4000



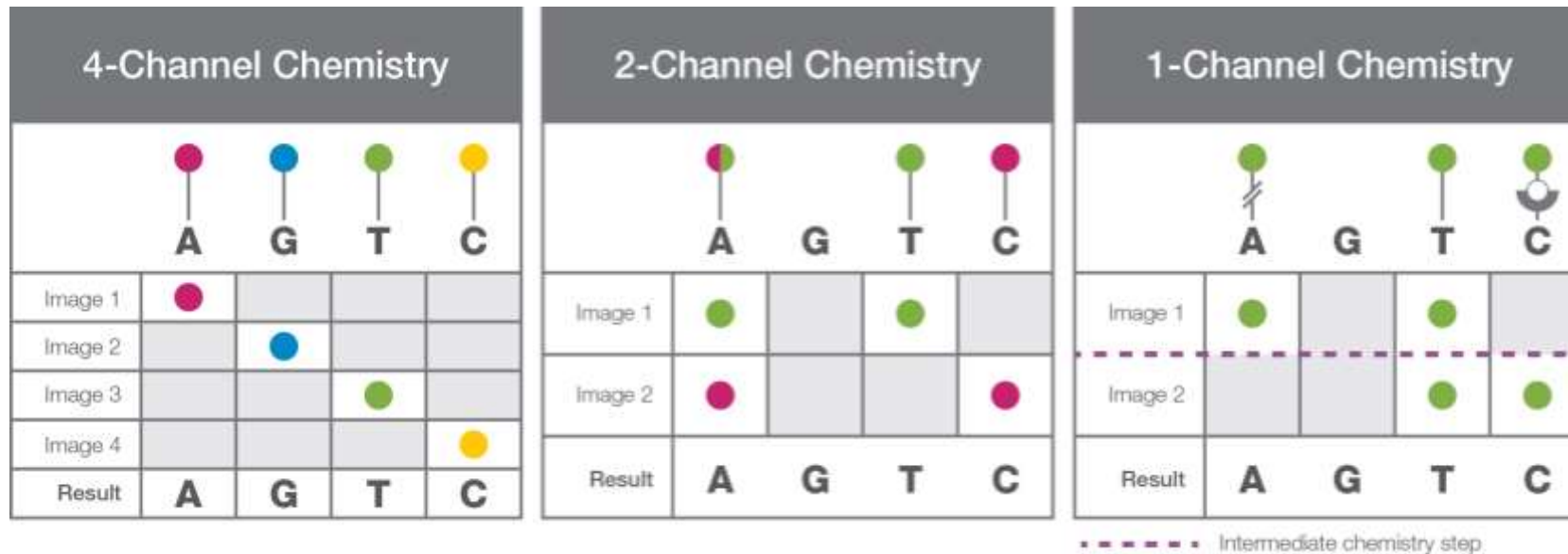
NGS Illumina Sequencing

- Imaging NextSeq500, Nextseq2000* NovaSeq



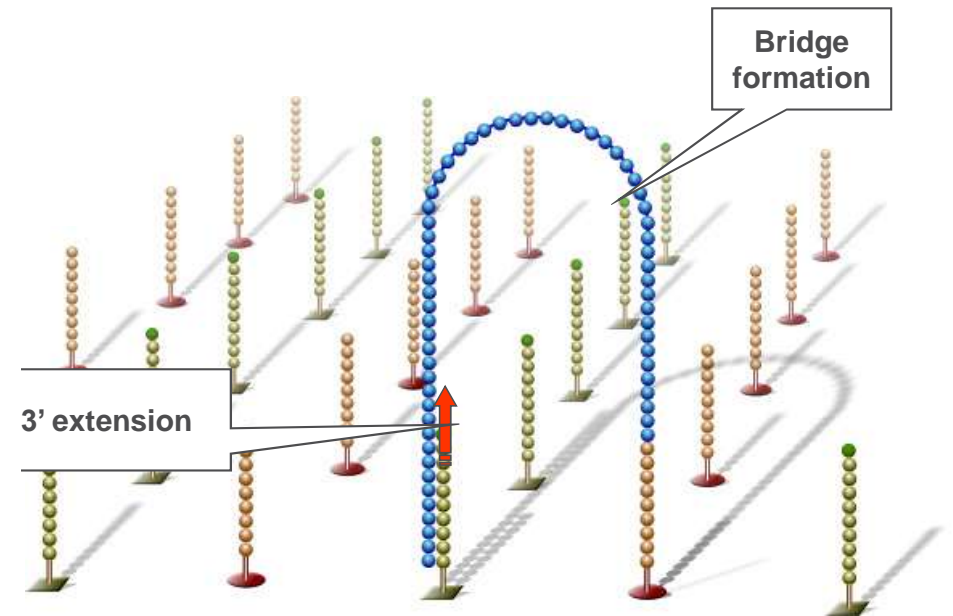
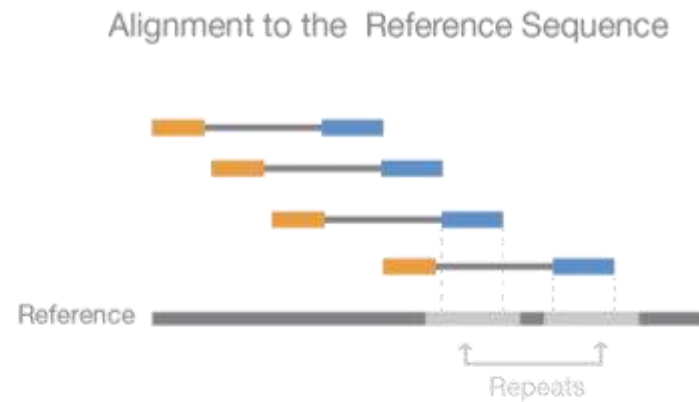
NGS Illumina Sequencing

- Different SBS dyes



NGS Illumina Sequencing

- Paired-end sequencing



Illumina sequencing QC

Quality Control

Run Summary

Level	Yield Total (G)	Projected Total Yield (G)	Aligned (%)	Error Rate (%)	Intensity Cycle 1	% \geq Q30
Read 1	0,3	0,3	33,59	1,76	76	95,3
Read 2 (I)	0,0	0,0	0,00	0,00	77	64,9
Read 3 (I)	0,0	0,0	0,00	0,00	379	97,2
Read 4	0,3	0,3	33,07	1,79	100	
Total	0,6	0,6	33,33	1,78	158	

Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

Q30 by cycle

Run Folder: Z:\HiSeq\210315\NextSeq\FCA\210315_NB501171_0668_AHCKWYBGXH

Browse

Refresh

Analysis Imaging Summary Indexing

Data By Cycle

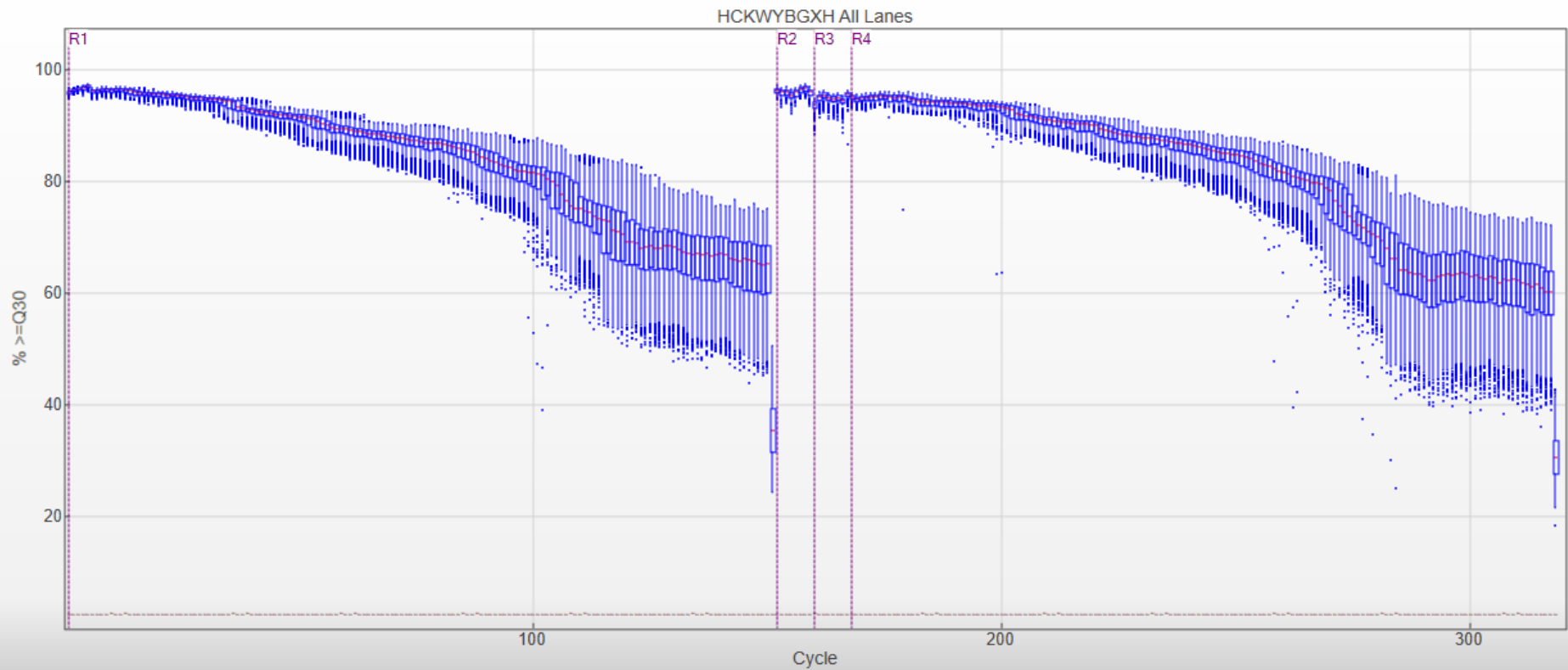
% >=Q30

Lane: All

Surface: All

Fix Scale

Accum



NGS Illumina Sequencing

- Quality Control

Run Summary

Level	Yield Total (G)	Projected Total Yield (G)	Aligned (%)	Error Rate (%)	Intensity Cycle 1	% >= Q30
Read 1	0,3	0,3	33,59	1,76	76	95,3
Read 2 (I)	0,0	0,0	0,00	0,00	77	64,9
Read 3 (I)	0,0	0,0	0,00	0,00	379	97,2
Read 4	0,3	0,3	33,07	1,79	100	88,3
Total	0,6	0,6	33,33	1,78	158	91,5

Read 1

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)	Error Rate 100 cycle (%)	Intensity Cycle 1
1	2	877 +/- 6	94,89 +/- 1,62	0,108 / 0,078	1,23	1,16	95,3	0,3	250	33,59 +/- 0,22	1,76 +/- 0,03	0,16 +/- 0,08	0,20 +/- 0,08	0,26 +/- 0,06	76 +/- 8

Read 2 (I)

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)	Error Rate 100 cycle (%)	Intensity Cycle 1
1	2	877 +/- 6	94,89 +/- 1,62	0,000 / 0,000	1,23	1,16	64,9	0,0	0	0,00 +/- 0,00	0,00 +/- 0,00	0,00 +/- 0,00	0,00 +/- 0,00	0,00 +/- 0,00	77 +/- 0

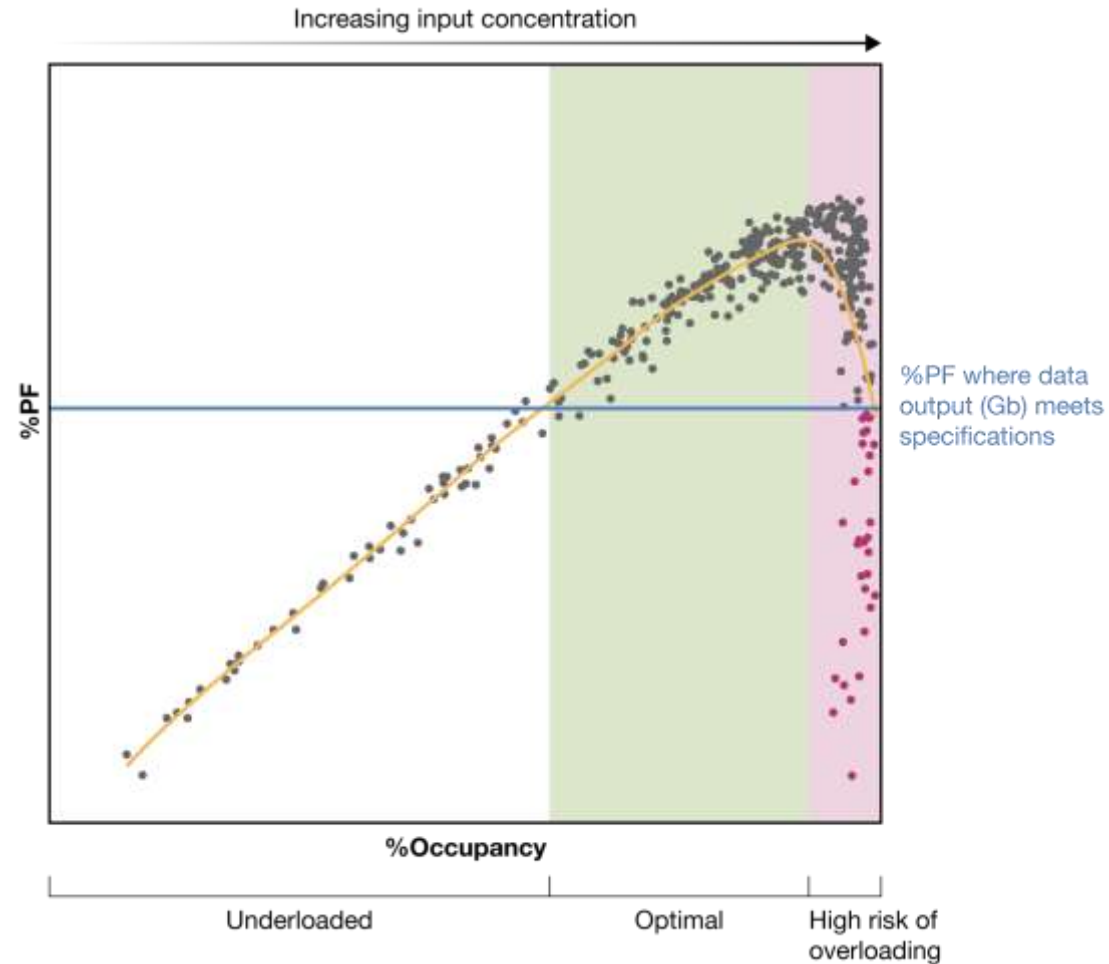
Read 3 (I)

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)	Error Rate 100 cycle (%)	Intensity Cycle 1
1	2	877 +/- 6	94,89 +/- 1,62	0,000 / 0,000	1,23	1,16	97,2	0,0	0	0,00 +/- 0,00	0,00 +/- 0,00	0,00 +/- 0,00	0,00 +/- 0,00	0,00 +/- 0,00	379 +/- 6

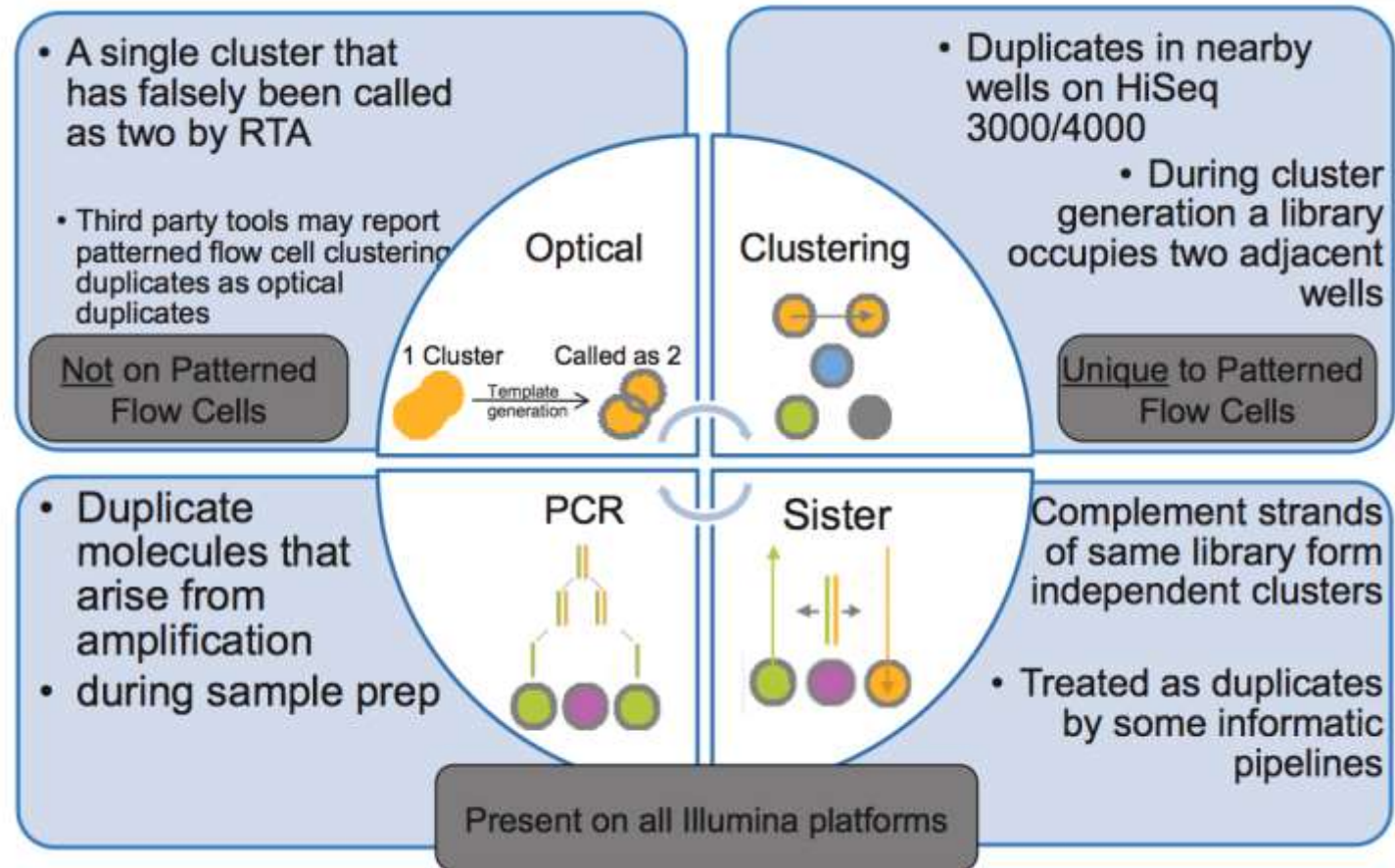
Read 4

Lane	Tiles	Density (K/mm2)	Cluster PF (%)	Phas/Prephas (%)	Reads (M)	Reads PF (M)	% >= Q30	Yield (G)	Cycles Err Rated	Aligned (%)	Error Rate (%)	Error Rate 35 cycle (%)	Error Rate 75 cycle (%)	Error Rate 100 cycle (%)	Intensity Cycle 1
1	2	877 +/- 6	94,89 +/- 1,62	0,147 / 0,068	1,23	1,16	88,3	0,3	250	33,07 +/- 0,61	1,79 +/- 0,01	0,20 +/- 0,10	0,25 +/- 0,08	0,30 +/- 0,06	100 +/- 5

Optimal loading



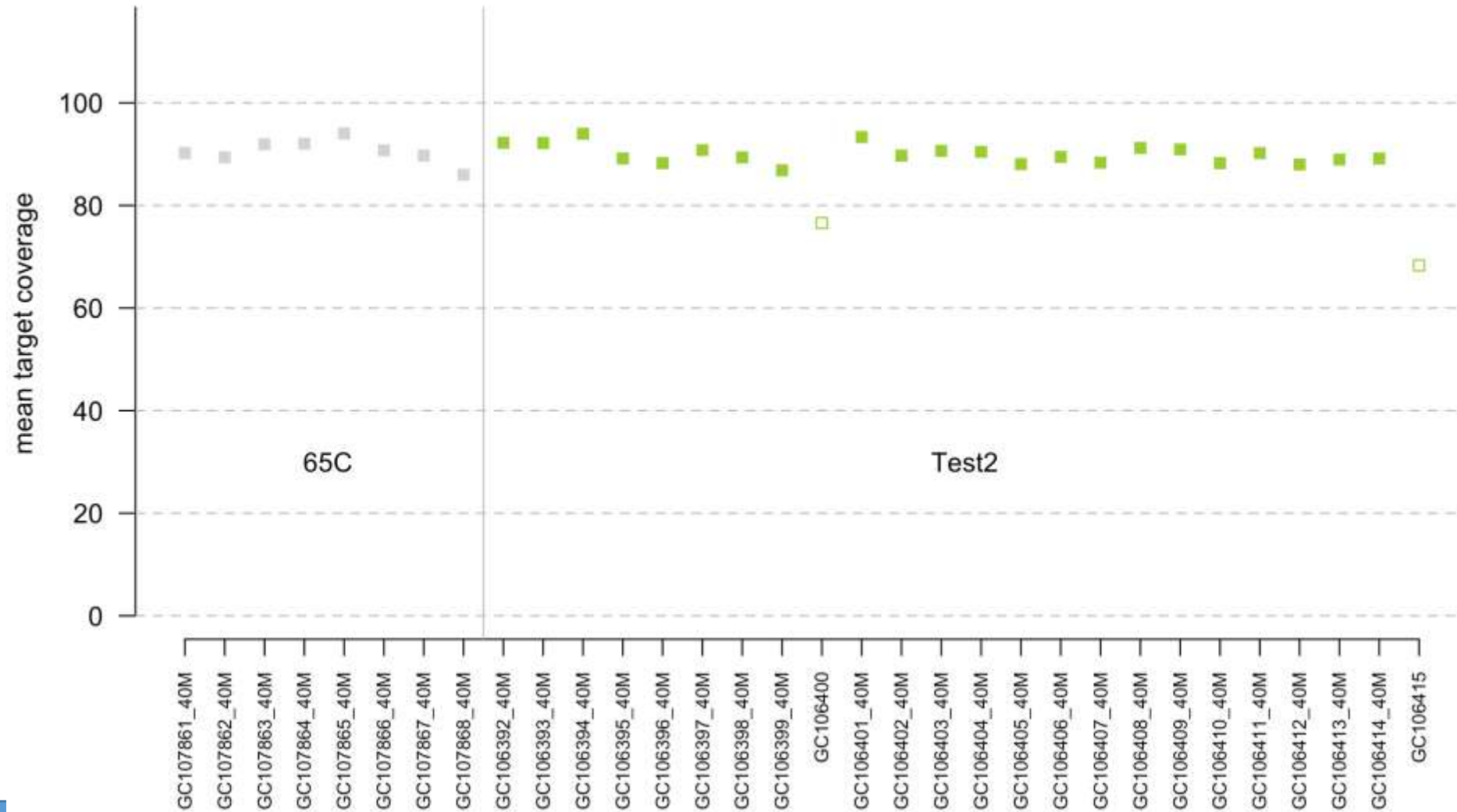
Duplicates



Common loss in Illumina

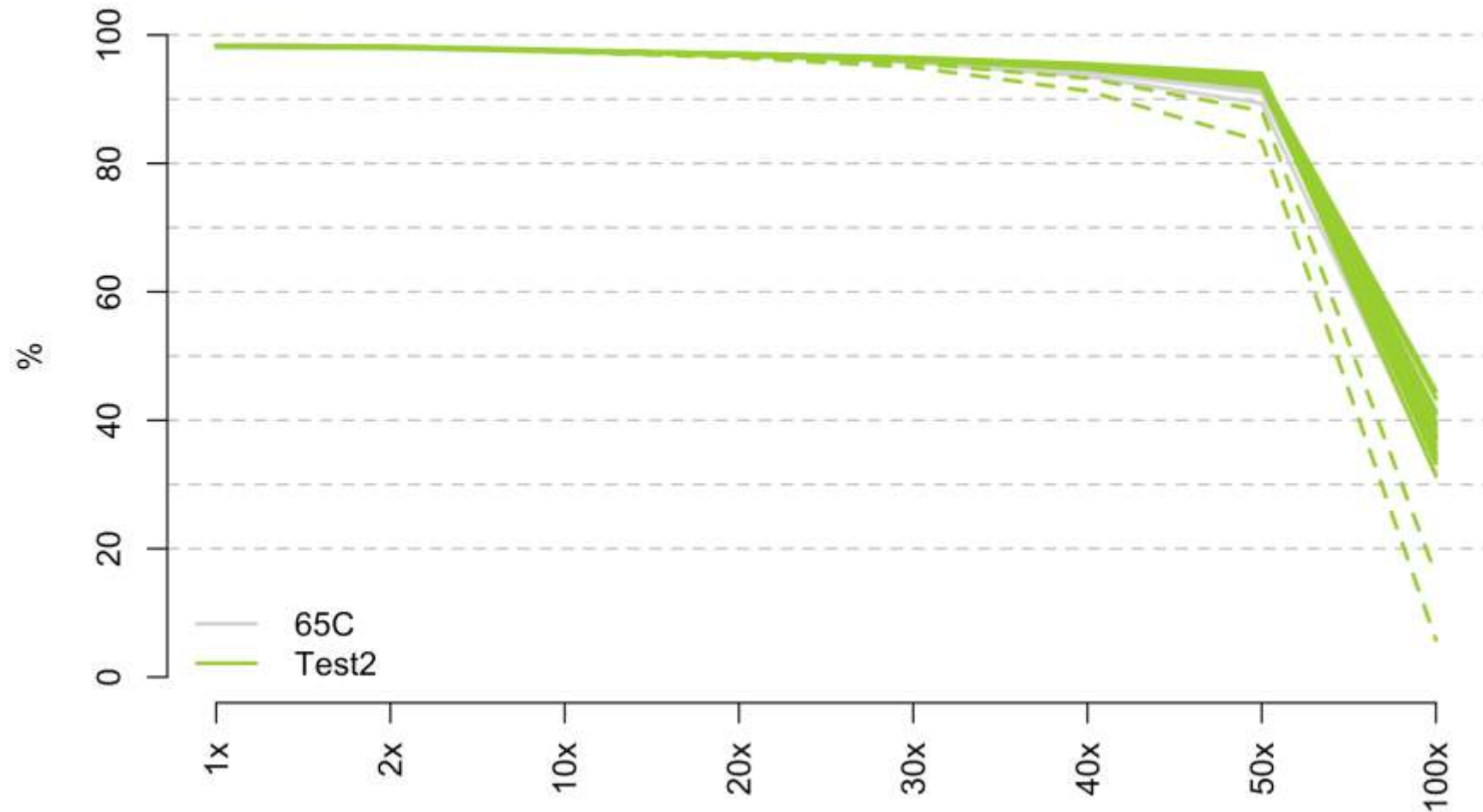
- Insufficient coverage: suboptimal loading
- Low complexity library
- Overloading
- Repeats
- Duplicates
- Too short fragment size
 - Overlap
 - High % adapter
- Low Q30

Metrics Picard – coverage technical target

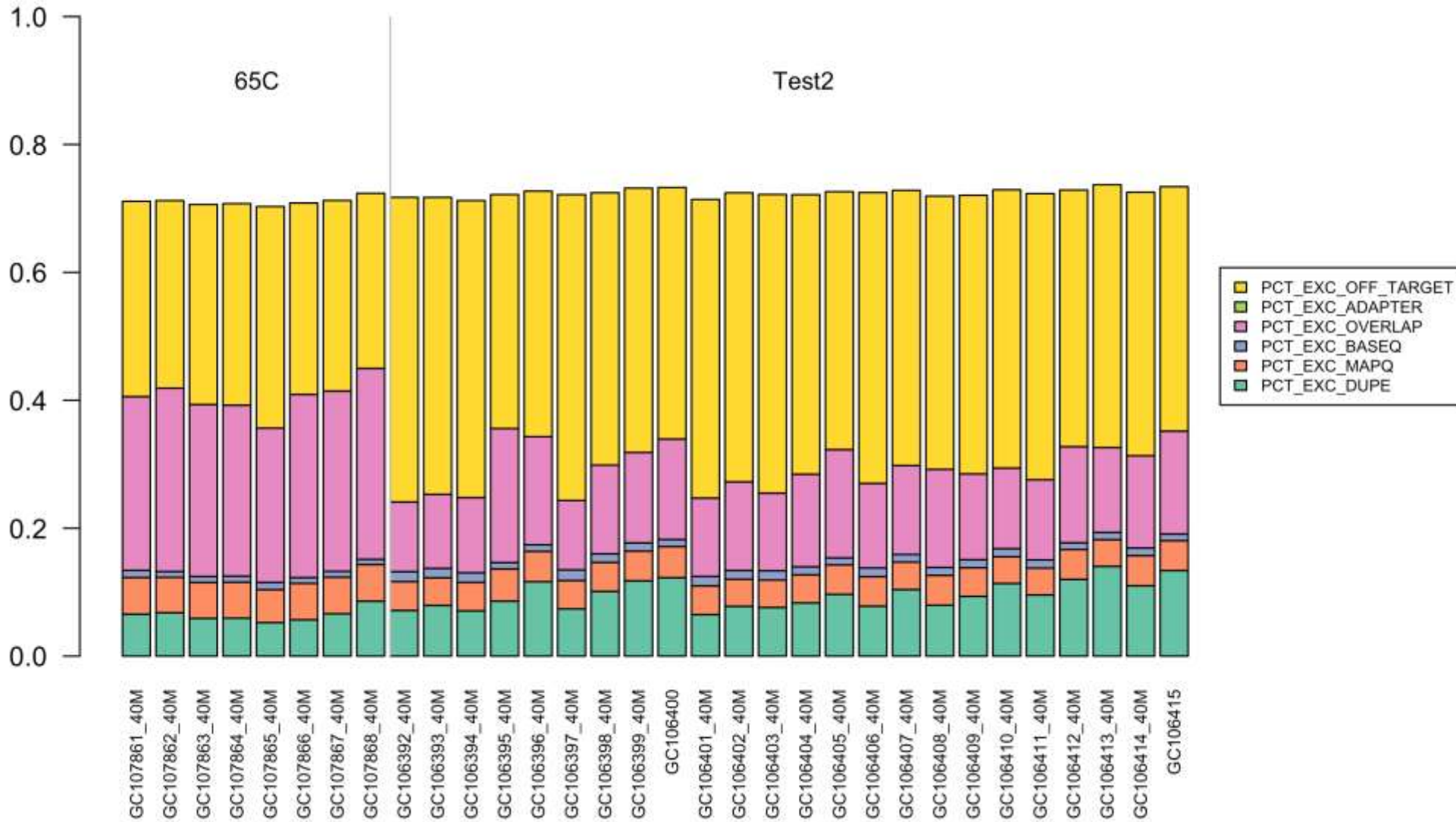


Metrics Picard – uniformity

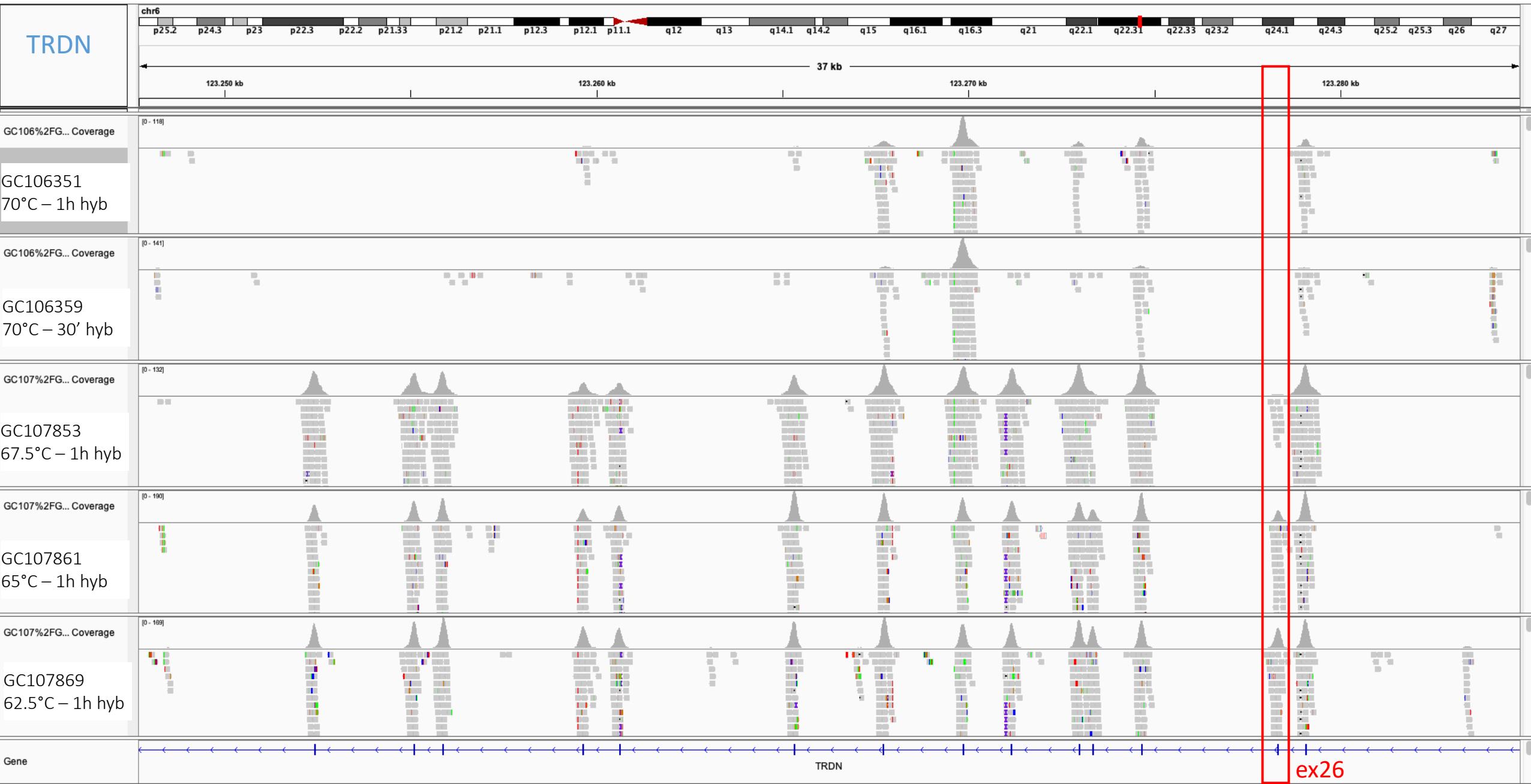
% TARGET BASES at ..x



Metrics Picard – Excluded bases



IGV screenshot of part of TRDN (chr6:123,247,689-123,284,977)



Nanopore sequencing

- MinION
- GridION
- PromethION
- SmidgION



Oxford Nanopore sequencing

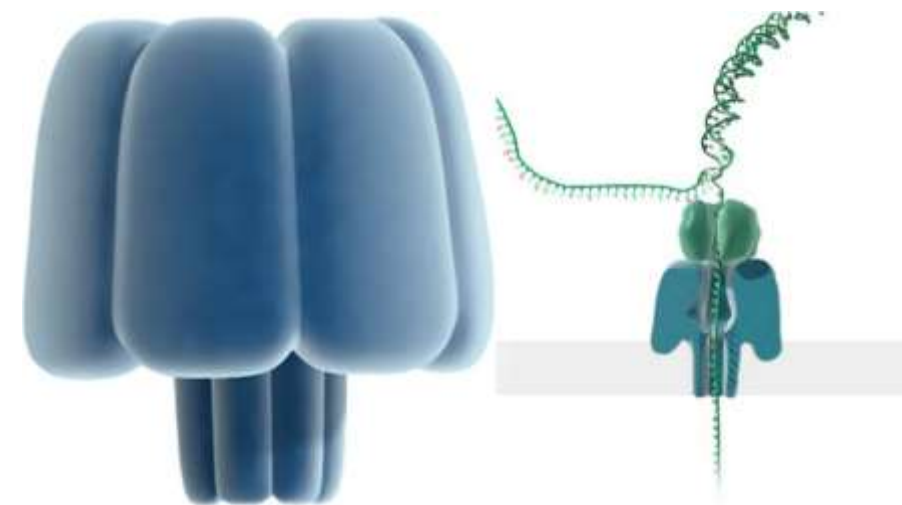


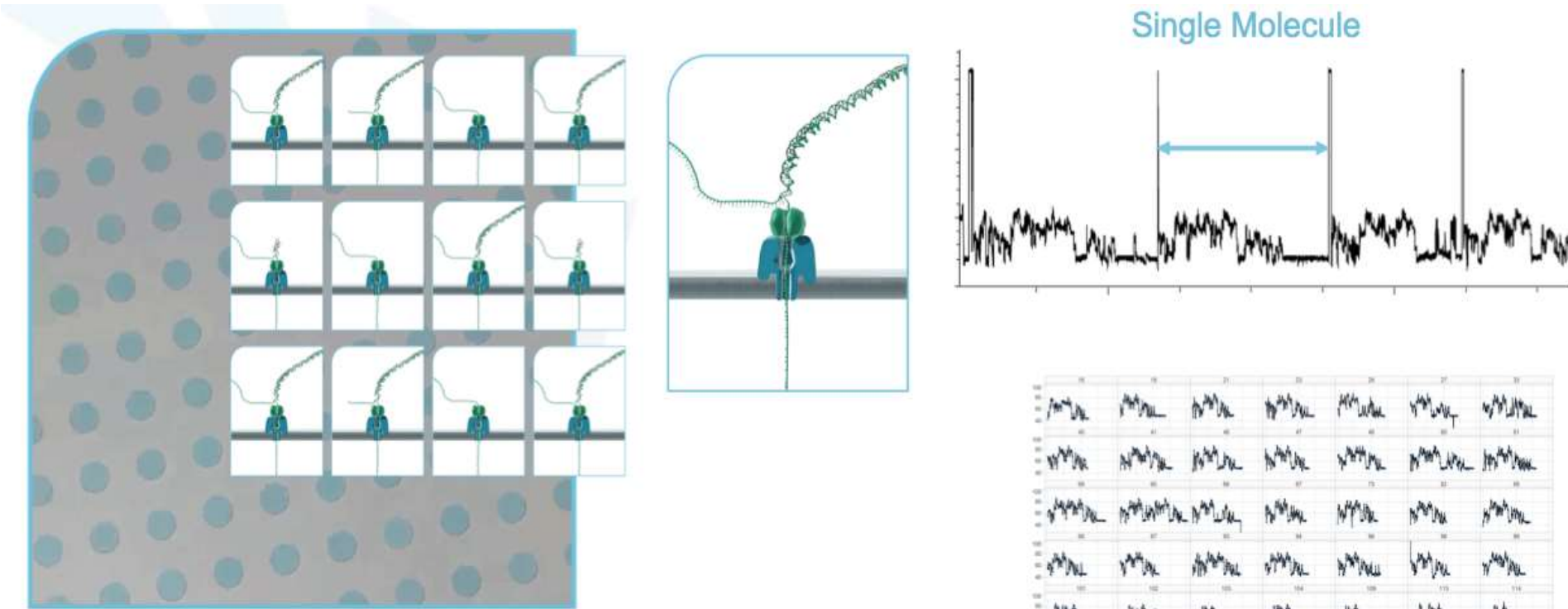
Tether keeps DNA fragment on the membrane leading to a ~20K fold higher DNA concentration close to the pore.

Motor protein unwinds DNA and ratchets it through the pore.

Abasic nucleotides in the hairpin are a recognition point.

Brake protein prevents the motor protein from zipping through the complement strand.





- ⦿ Data acquired as full length reads – real time
- ⦿ Data throughput = No. pores x average speed/pore

From squiggles to sequencing

- New basecaller: extract more correct information from squiggles
- Training of base caller for methylation data

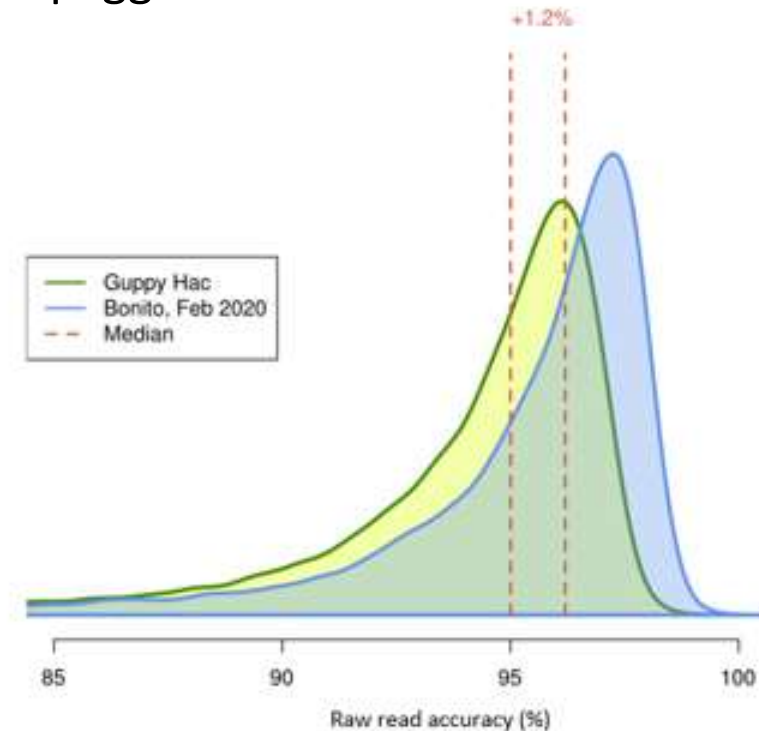
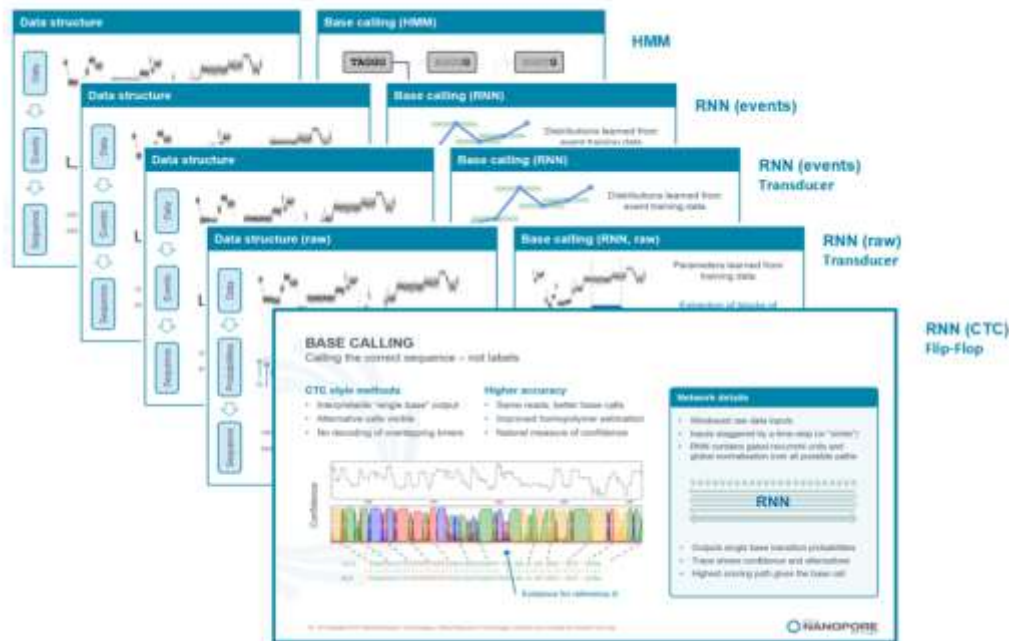
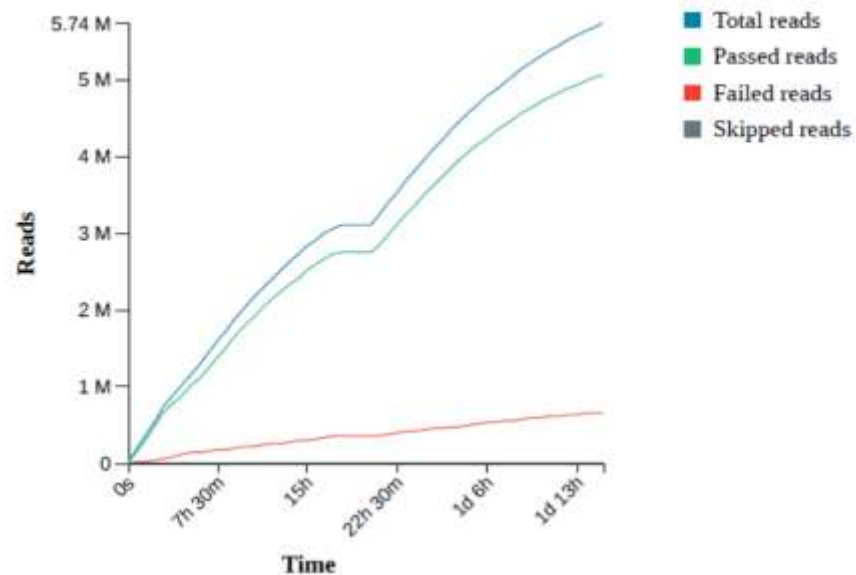


Figure 1: Raw read accuracy of Bonito basecaller on the human reference genome NA12878 against high-accuracy Guppy, currently integrated into MINKNOW onboard nanopore devices.

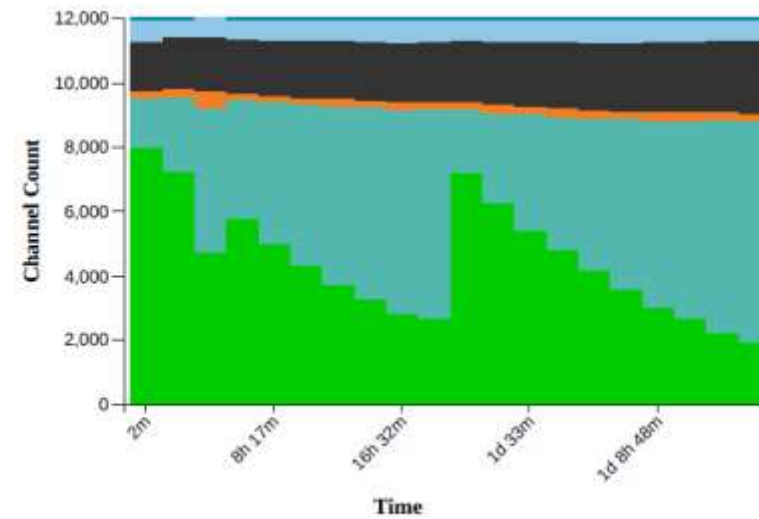
Run QC on Nanopore

- Every flowcell is different
- Nuclease wash (and refueling) increases output

Cumulative Output Reads

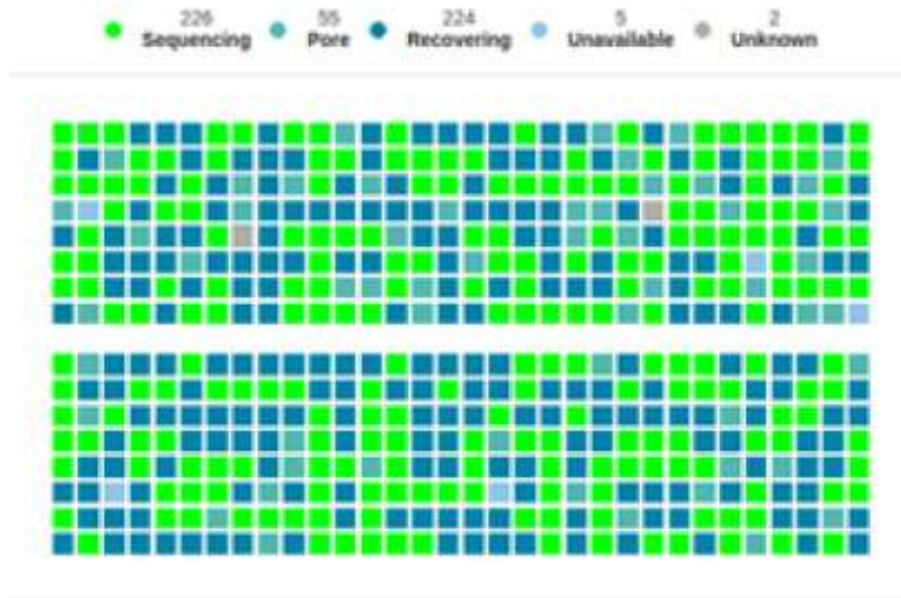


Mux Scan Categorized

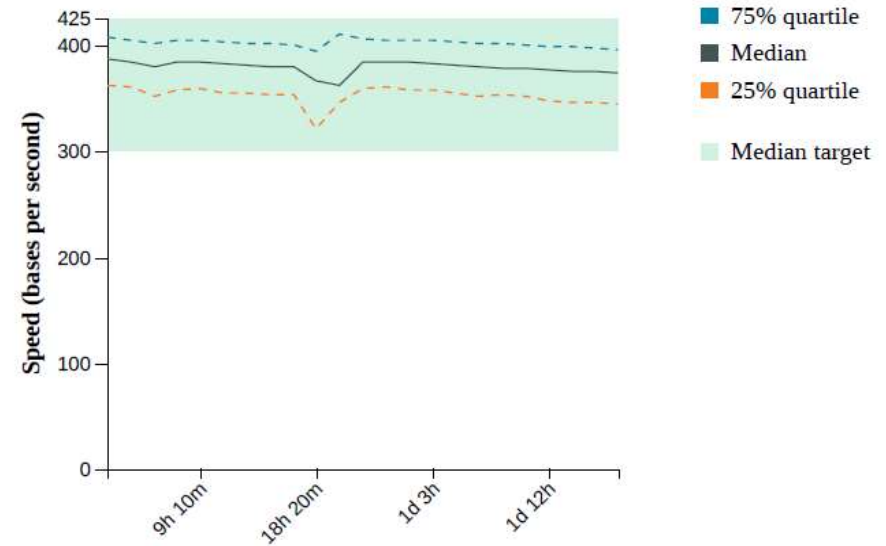


● Single Pore ● Reserved Pore ● Unavailable ● Multiple ● Saturated ● Zero ● Other

Translocation speed and pore status



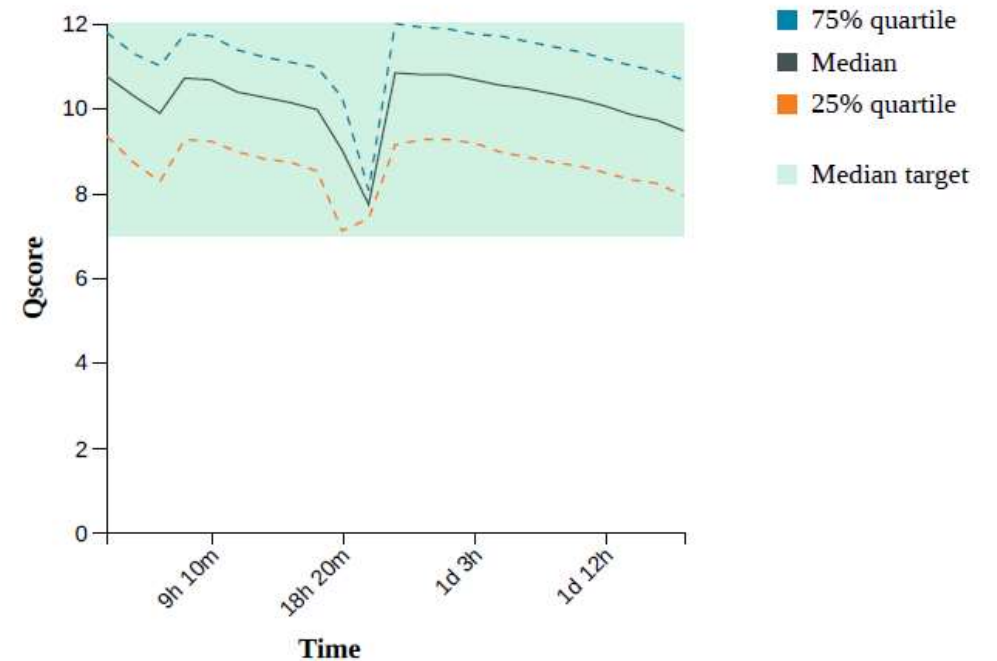
Translocation Speed



Run QC on Nanopore

- Q-score need to be stable
- New flowcell chemistry improves Q-score
- Barcode selection → select high quality door

QScore

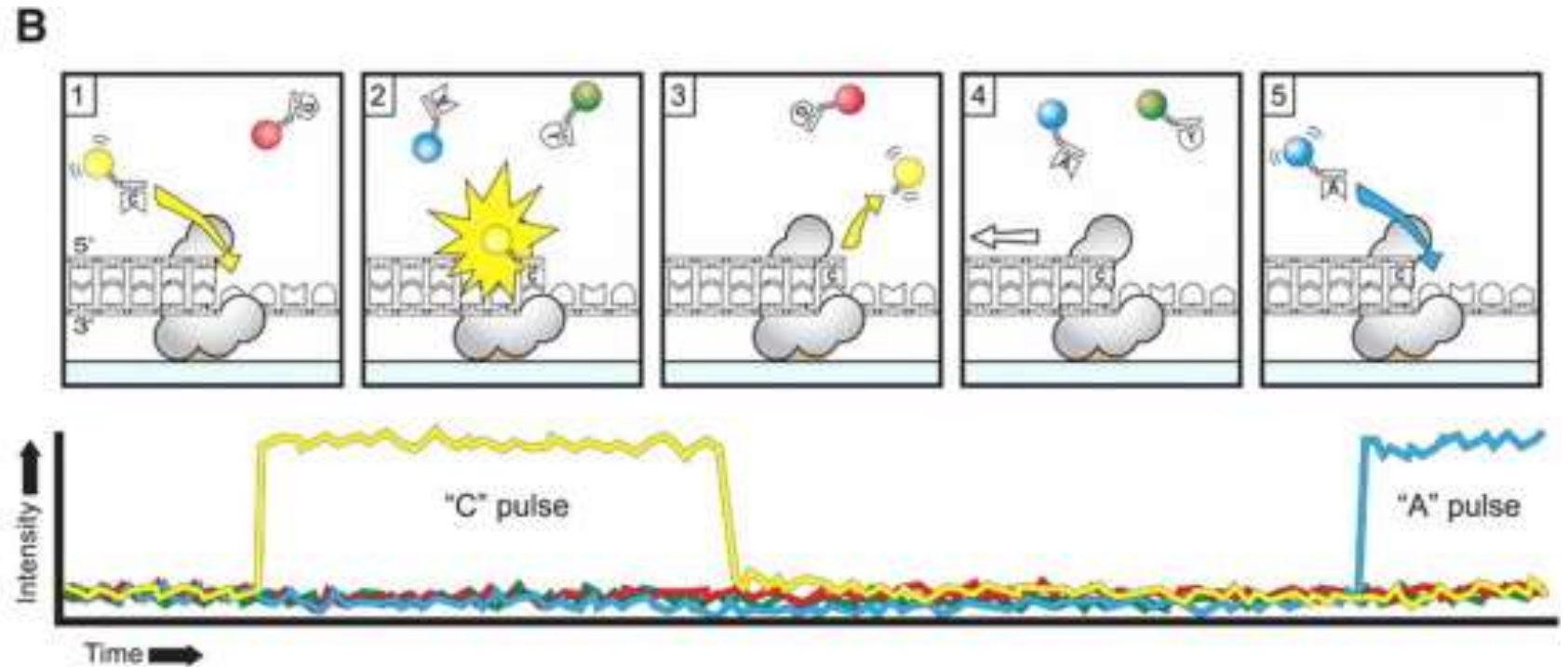
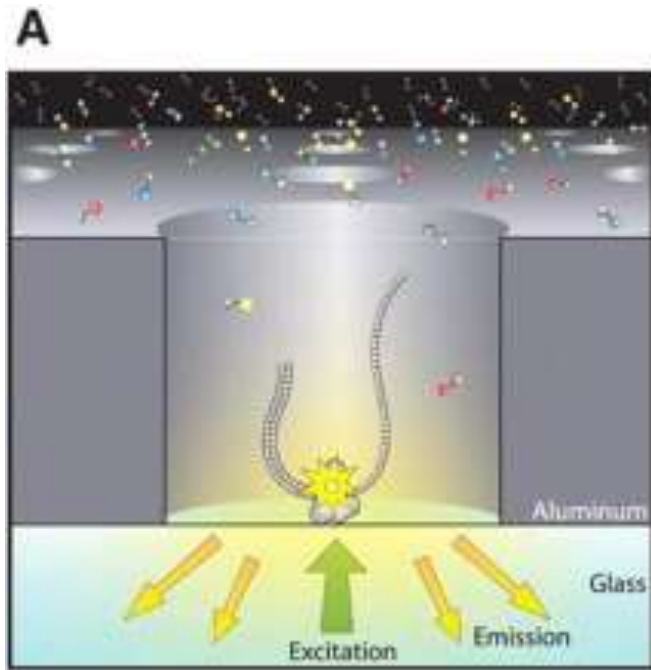


SMRT Sequencing

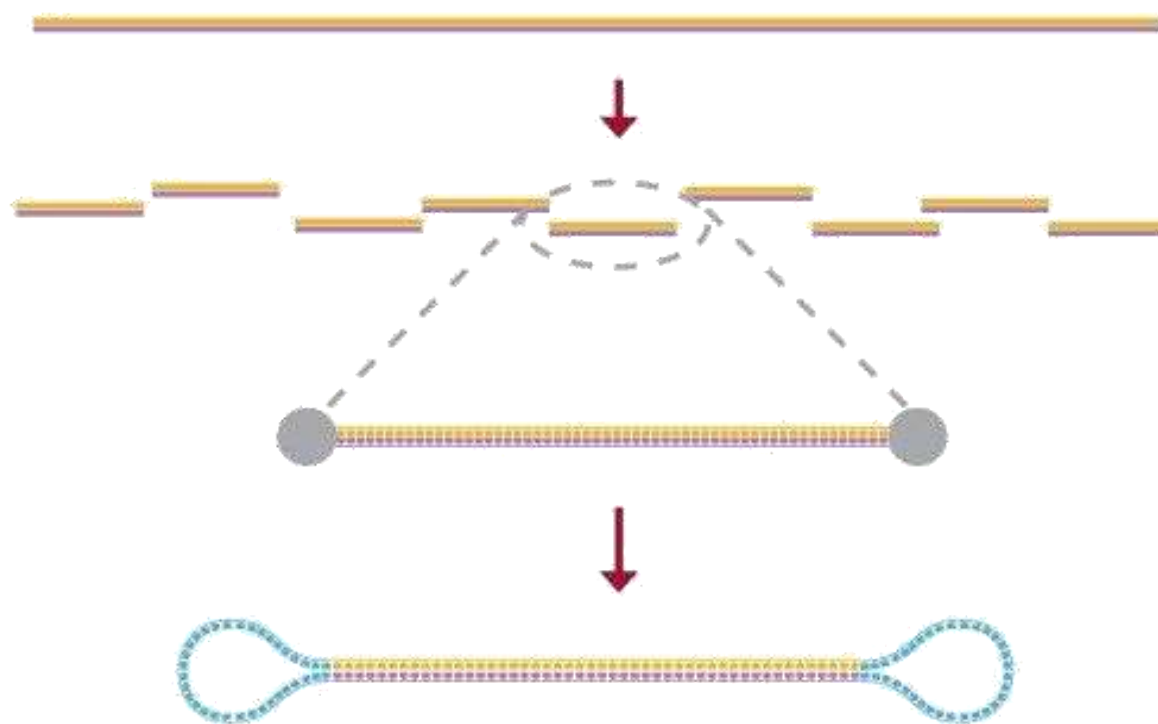
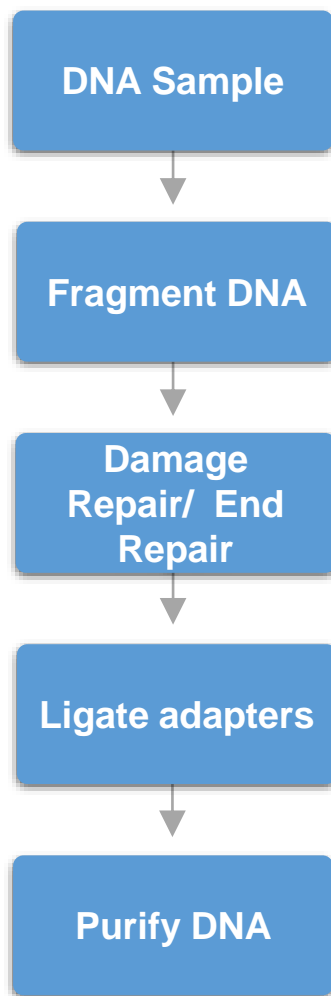
- Pacific Biosciences
PacBio Sequel IIe



SMRT[®] Technology



Template preparation



- from 250 bp to 400 kb
- sequences of both forward and reverse strands in the same trace

Universal SMRTbell™ Template



Large Insert Sizes

- Recommended Insert Size: > 3 kb
- Maximum length over 300 kb



Generates one pass on each molecule sequenced

Circular Consensus Sequencing (CCS)



Small Insert Sizes

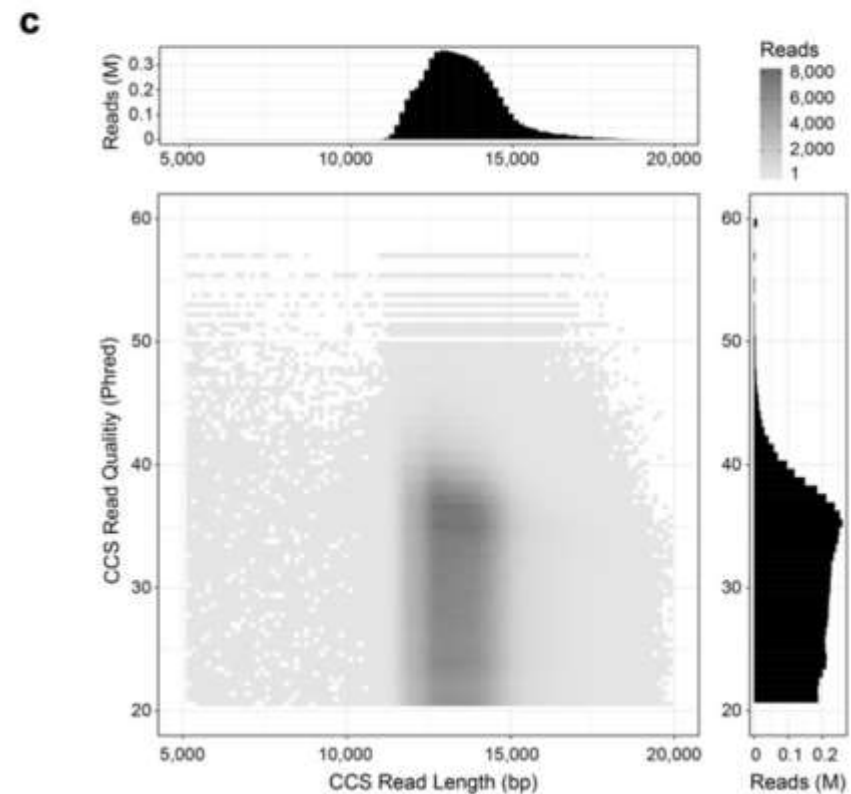
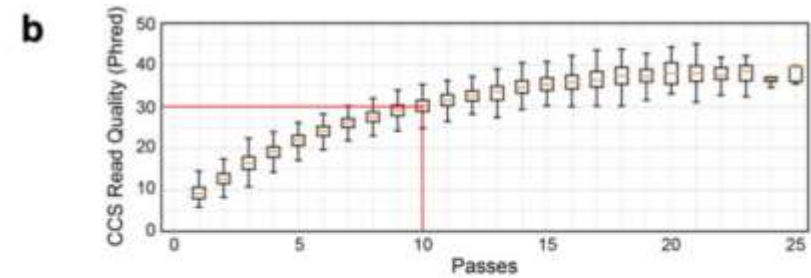
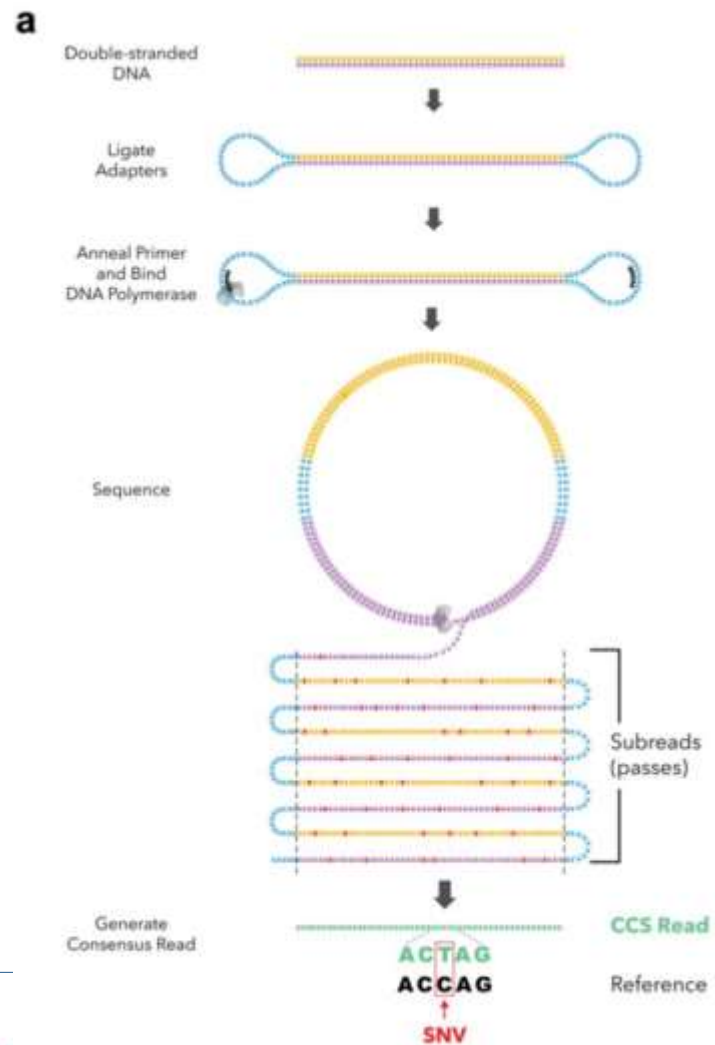
- Recommended Insert Size: 500 bp-20kb



Continued generation of reads per insert size

Generates multiple passes on each molecule sequenced

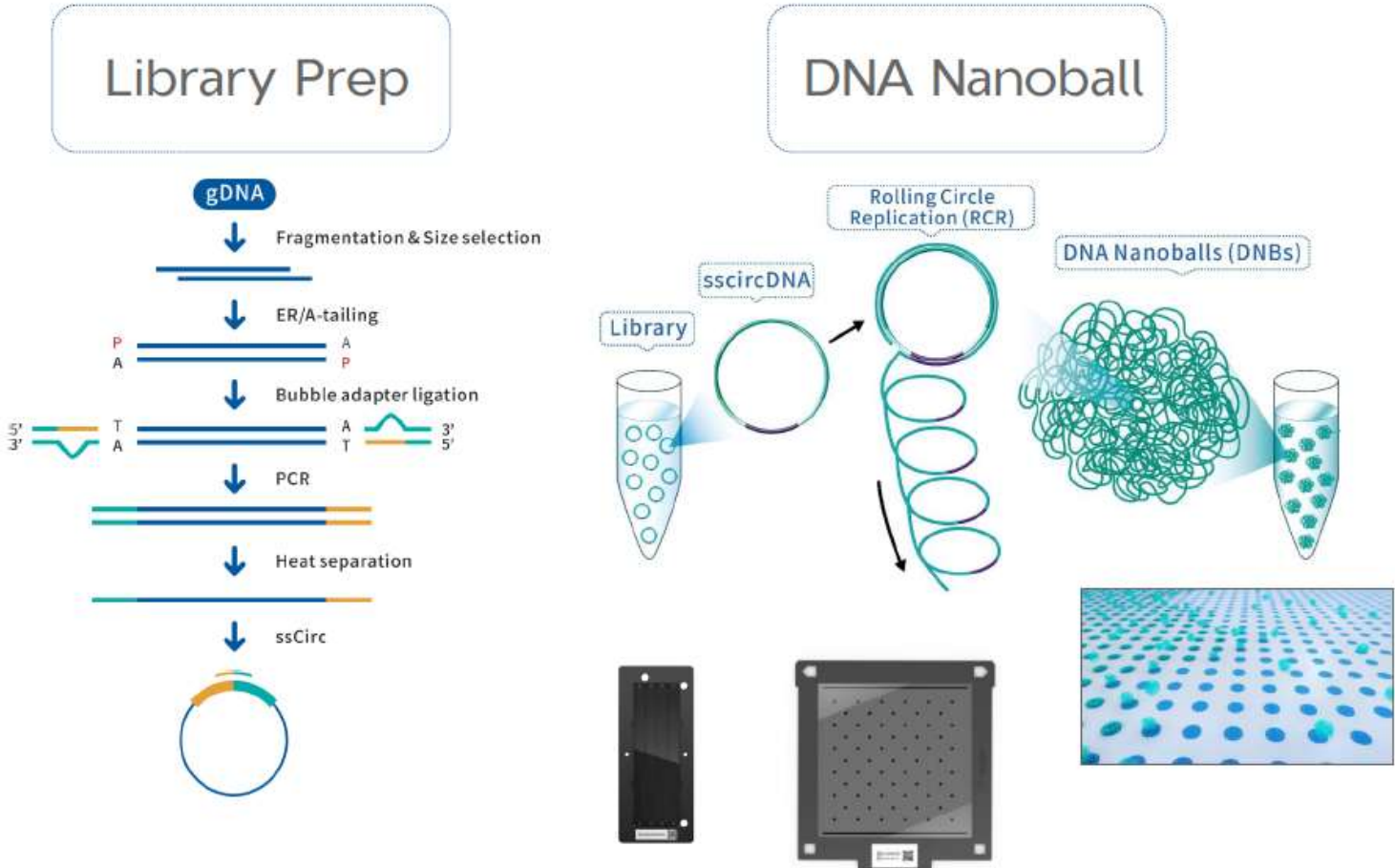
SMRT[®] Sequencing Accuracy



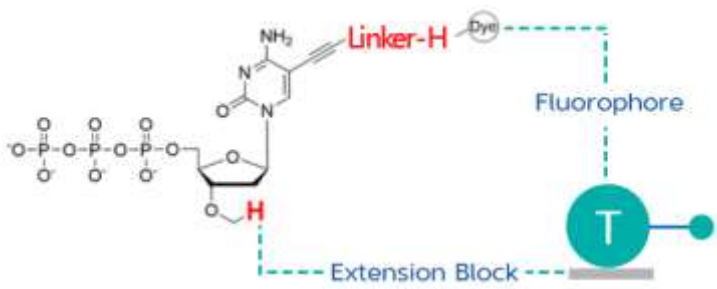
Benefits of SMRT[®] Sequencing

- Produce reads with average lengths of 6000 to 10000, with longest reads over 175,000 base pairs
- Greater than **99.999% (QV 50)** accurate sequencing results
- Sensitivity to detect minor variants at frequency less than 0.1%
- Detect broad spectrum of base modification events in the same sequencing run that reads canonical base sequence
- No amplification bias and least GC bias for improved coverage uniformity

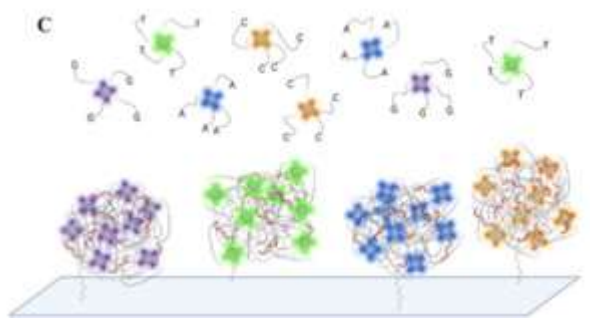
MGI and Aviti Element



MGI: SBS



AVITI: avidite binding



PHASE TWO: INTERPRETATION

SEEDMAN the New Yorker



NGS Applications

NGS as a tool for studying genome variation and regulation

DNA

- Targeted resequencing
 - Amplicon
 - MIPs
 - Capture panels
- *de novo* assembly
- Bacterial WGS
- Vertebrate WGS
- Long read sequencing

RNA

- Truseq stranded mRNA
- Lexogen quantseq
- IsoSeq (Pacbio)

Single cell genomics

- Various single cell library prep methods for DNA and RNA

Whole genome sequencing

- **Copy number variation analysis**
 - Sequencing a genome at 0.1-0.3x
 - Sequencing a genome at 1-3x
- **Structural variation analysis**
 - Sequencing a genome at 5-10x
- **Whole genome re-sequencing**
 - Sequencing a genome at >30x
 - yeast, fruit fly, bacterial genomes, human...



De novo assembly

- Assembling a genome from scratch
- Extremely computationally heavy
- No reference to distinguish variation from artefacts
- Combination of multiple sequencing and optical mapping techniques required

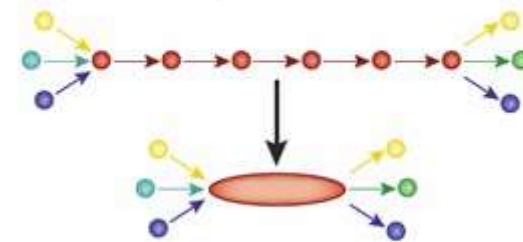
1. Fragment DNA and sequence



2. Find overlaps between reads

...AGCCTAGACCTACA**GGATGCGGACACGT**
GGATGCGGACACGTCGCATATCOGGT...

3. Assemble overlaps into contigs

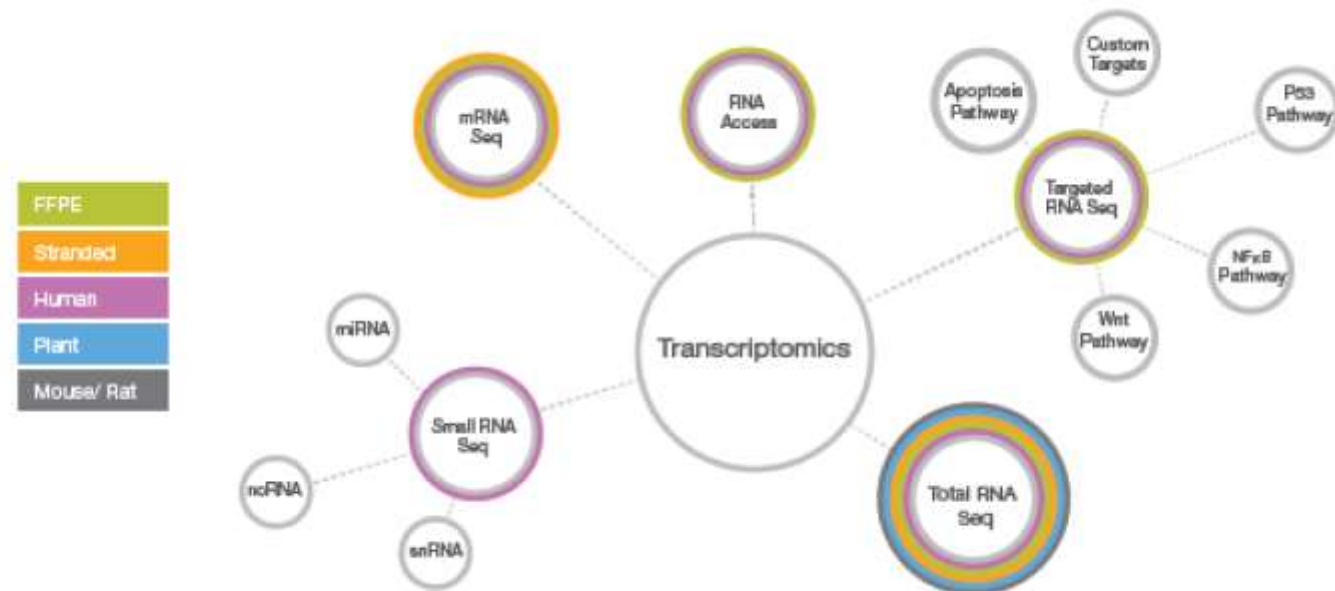


4. Assemble contigs into scaffolds



RNA SEQUENCING

- Rapid expression profiling, transcriptome sequencing and small RNA's

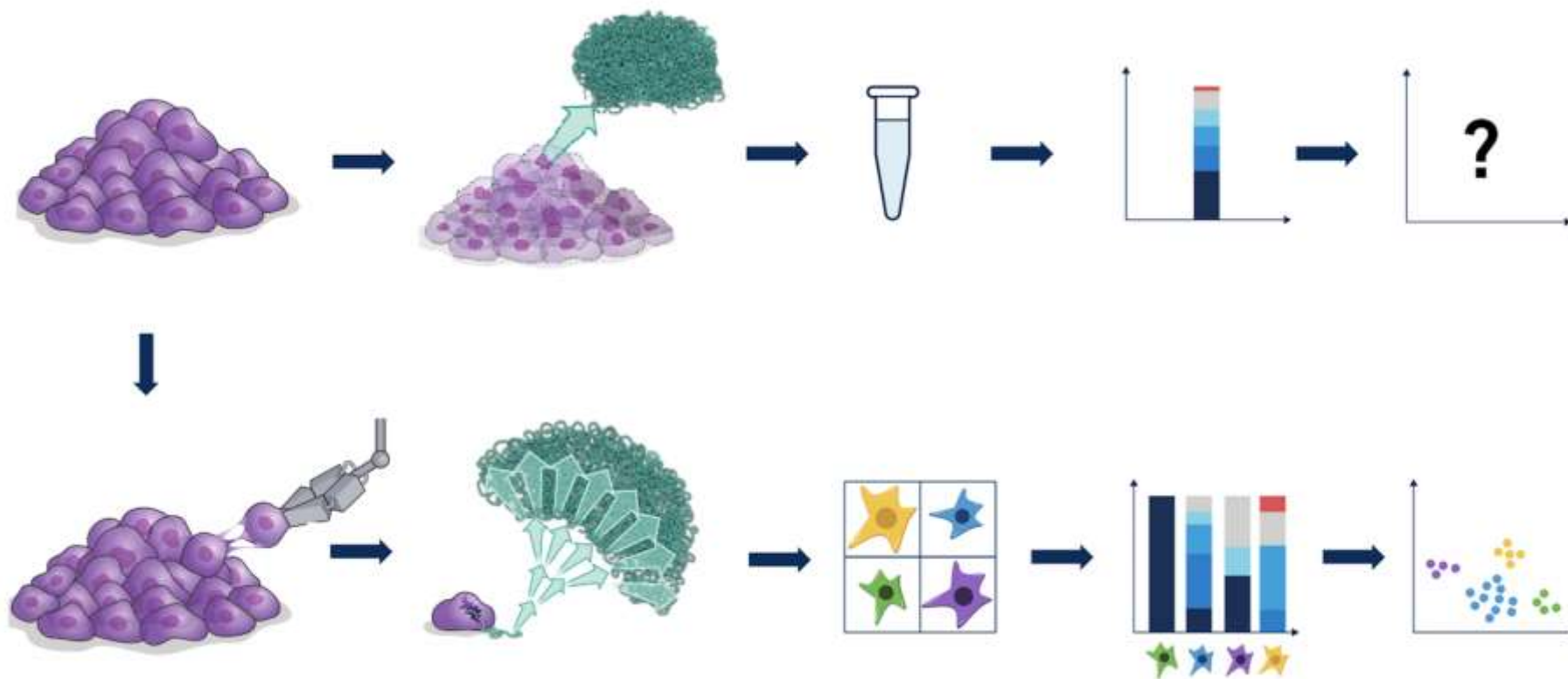


RNAseq library preps

	Differential expression	Whole transcript, fusion, isoforms	Small RNA	Illumina compatible	Low input
Lexogen QuantSeq 3' mRNA	✓			✓	✓
Lexogen Small RNA seq	✓		✓	✓	✓
Illumina TruSeq stranded mRNA	✓	✓		✓	
Illumina TruSeq stranded total RNA	✓	✓	✓	✓	
IsoSeq	(✓)	✓			
Smart-Seq2	✓	✓		✓	✓✓✓

Single cell RNA-seq as a complementary technique to bulk RNA-seq

Side note



Single cell RNA-seq as a complementary technique to bulk RNA-seq

	RNA type	Transcript targeted	sensitivity	throughput	sequencing	Specific property
SMART SEQ2	mRNA	full transcript	sensitive	low	deep	FACS sorting specific populations
10x genomics 3' RNA seq	mRNA	3' end	medium	High (> 10000)	shallow	Can be combined with surface markers

