

Next Generation Sequencing

Theory and Platforms

Dr. Walter Verweij

walter.verweij@earlham.ac.uk



Earlham
Institute



At the Norwich
Research Park
(NRP)



Sequencing Platforms



Illumina HiSeq 2500 x 2
Illumina HiSeq 2000 x 1
Illumina HiSeq 4000 x 2



Illumina MiSeq x 3

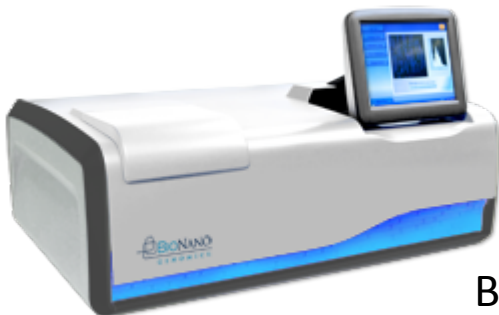


PacBio RSII x 1
PacBio Sequel x 1

Sequencing Platforms



OpGen Argus x 1

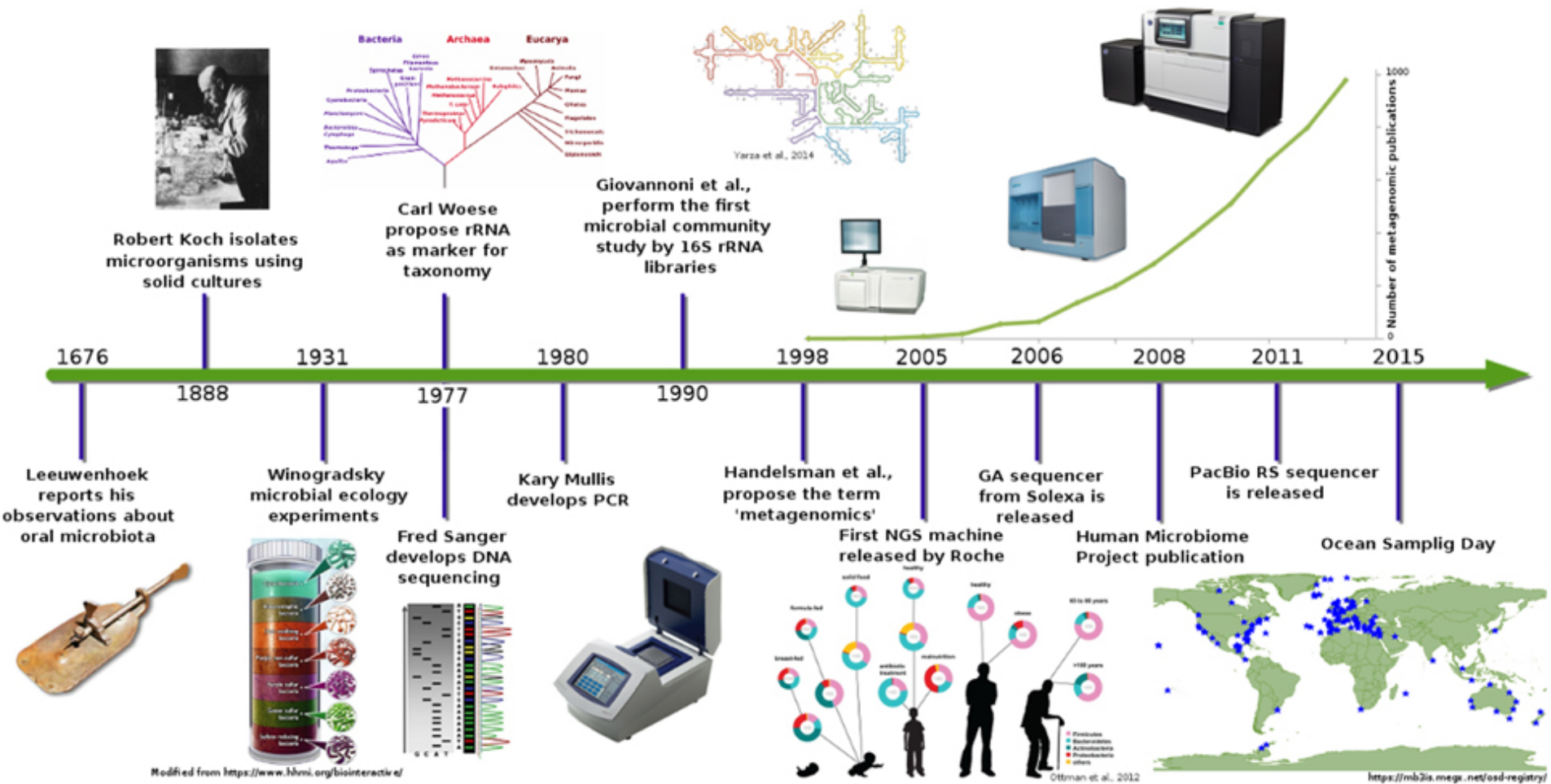


Bionano Optical mapping x 1



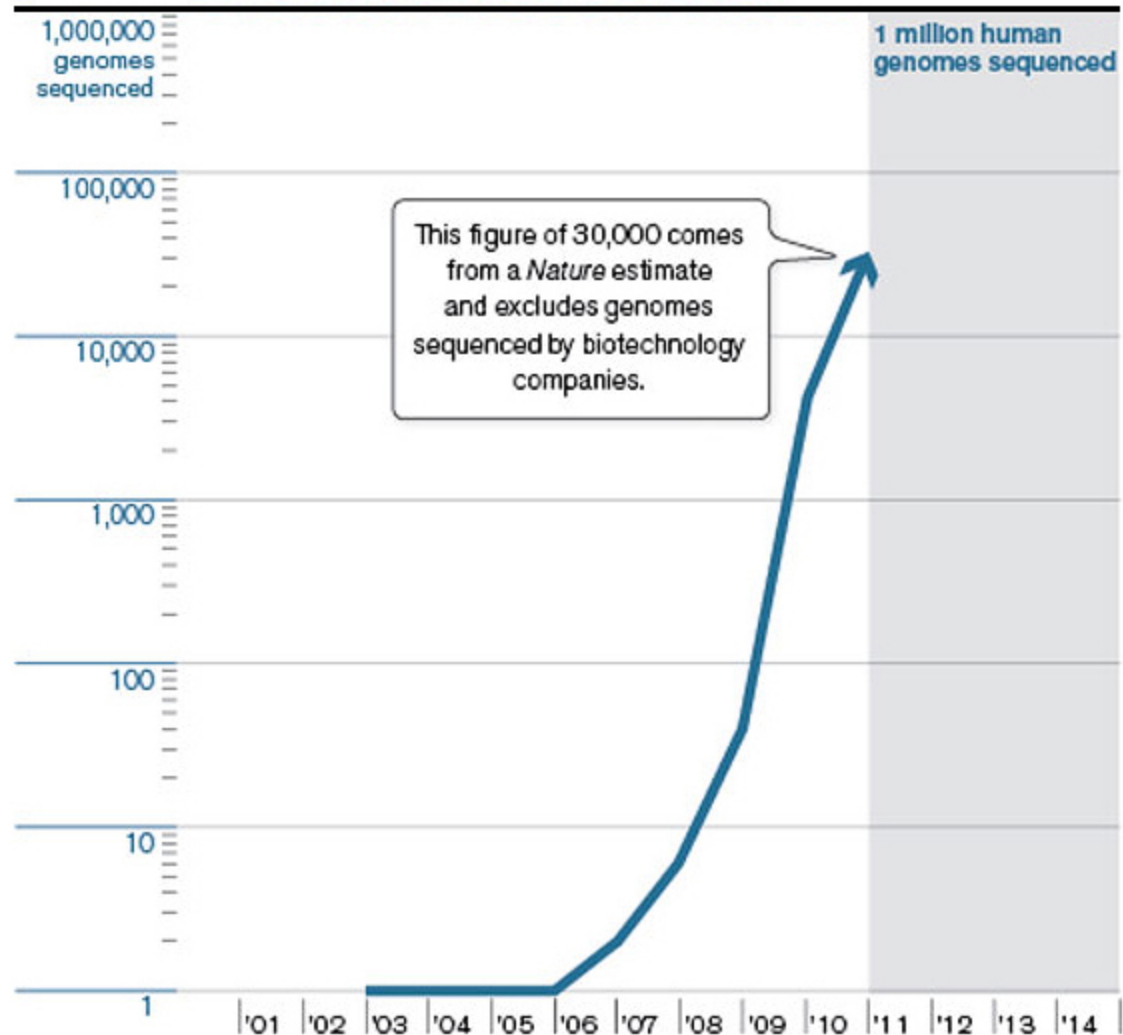
Oxford Nanopore Minlon x ...

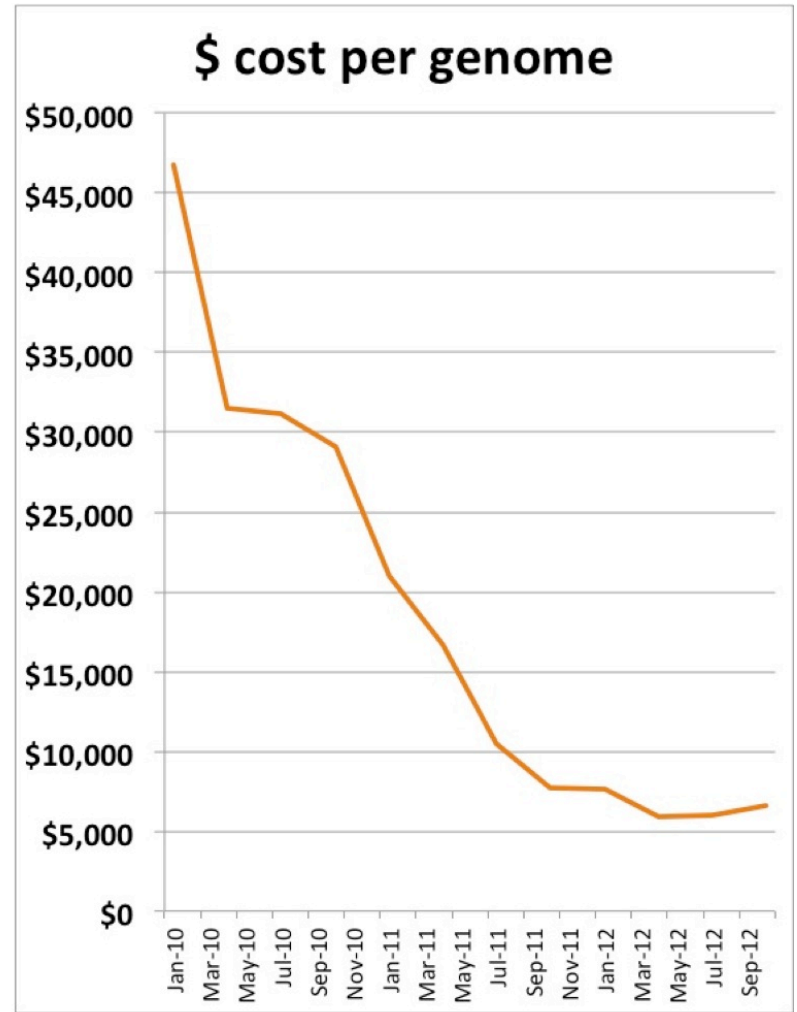
Sequence history



Output Skyrocketing

Number sequenced

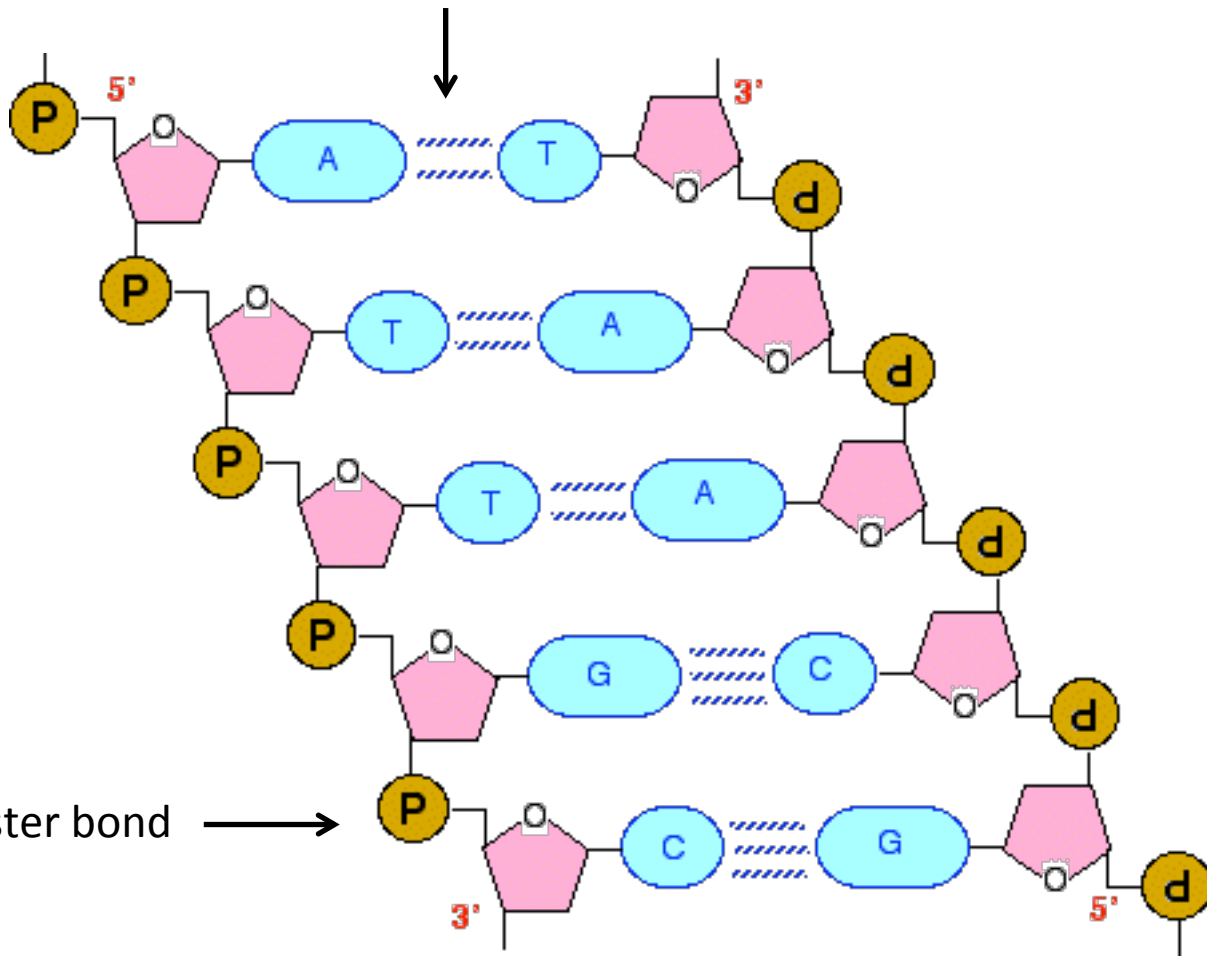




DNA sequencing

- Capillary Sequencing (Sanger sequencing)
- Next Gen Sequencing Platforms
 - Illumina
 - Pacific Biosciences

Hydrogen bond



phosphodiester bond

Why sequence?

1. *De novo*

A species genomic sequence allows genetic analysis.

2. Re-sequencing

Understand the nature of genetic variation – population genetics

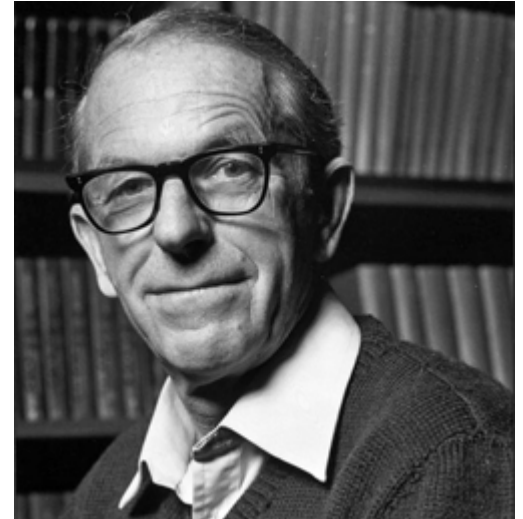
3. Counting

Functional genomics e.g. gene expression etc.

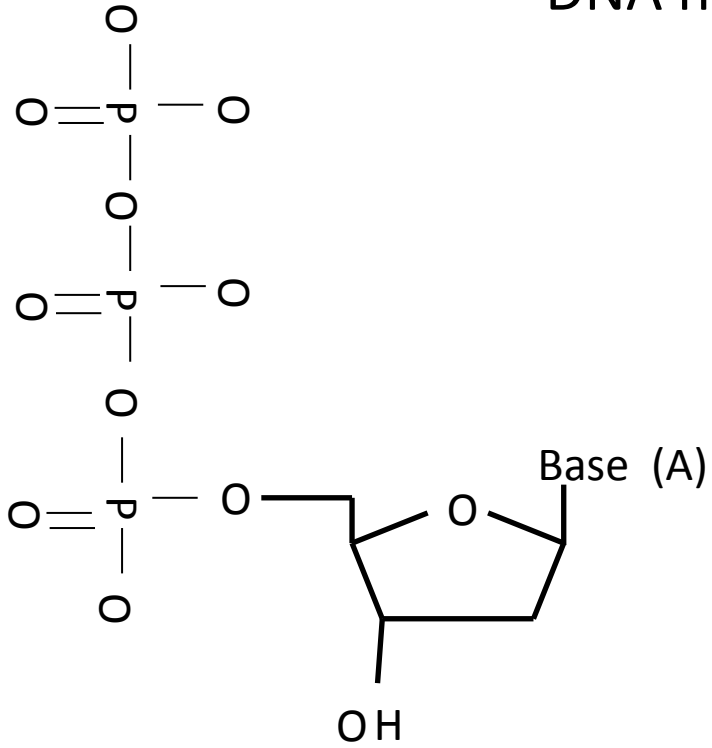
Capillary Sequencing

The DNA sequencing technology developed by Fred Sanger in 1977 is known as the dideoxy method, the chain termination method, or the Sanger sequencing method.

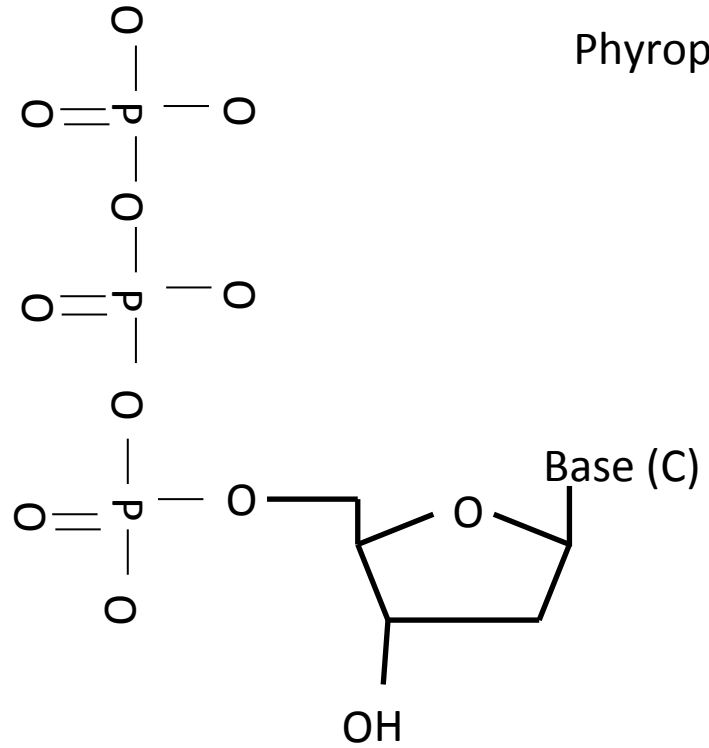
Based on the technique of separation by electrophoresis.



DNA nucleotide binding

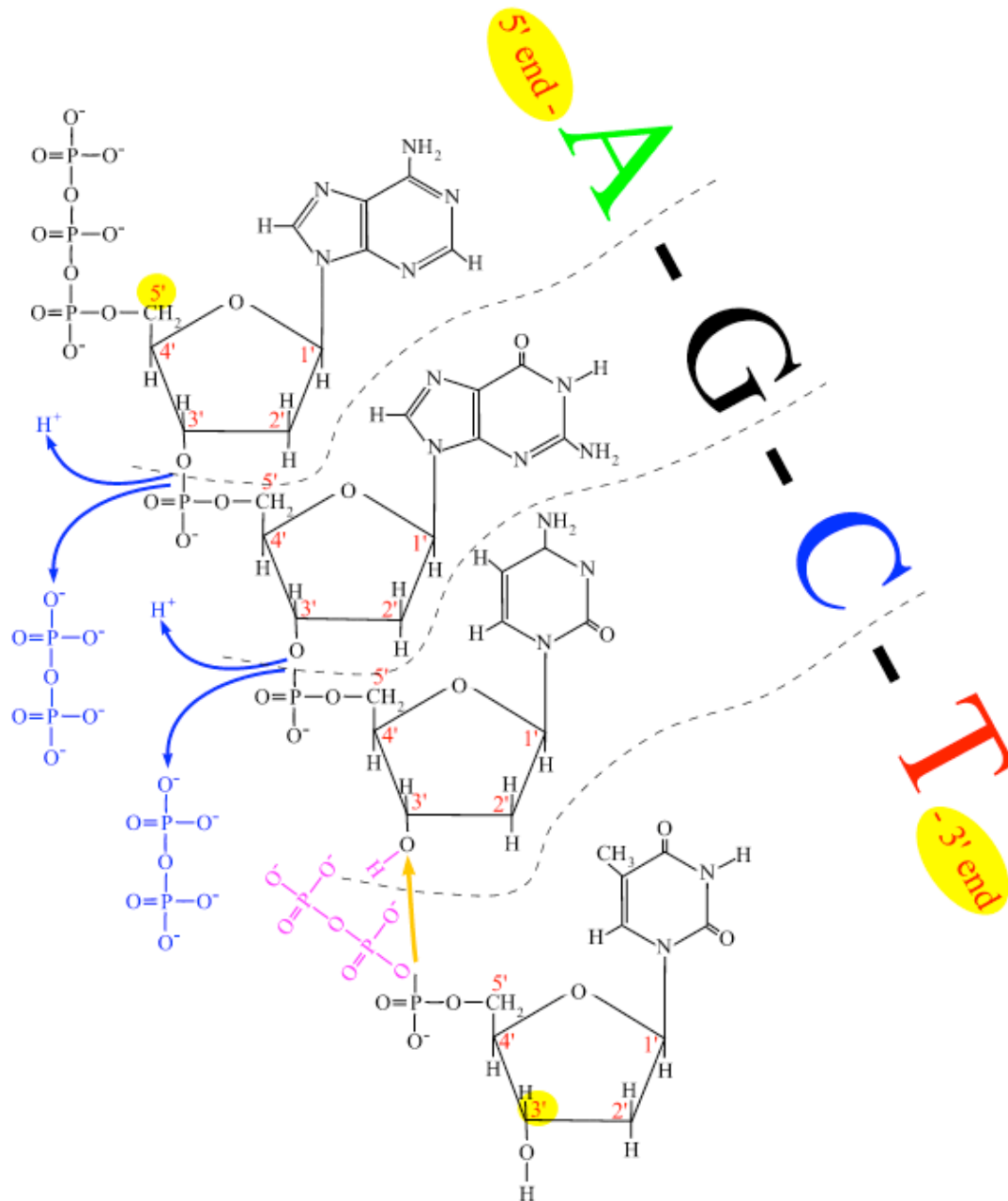


hydrogen



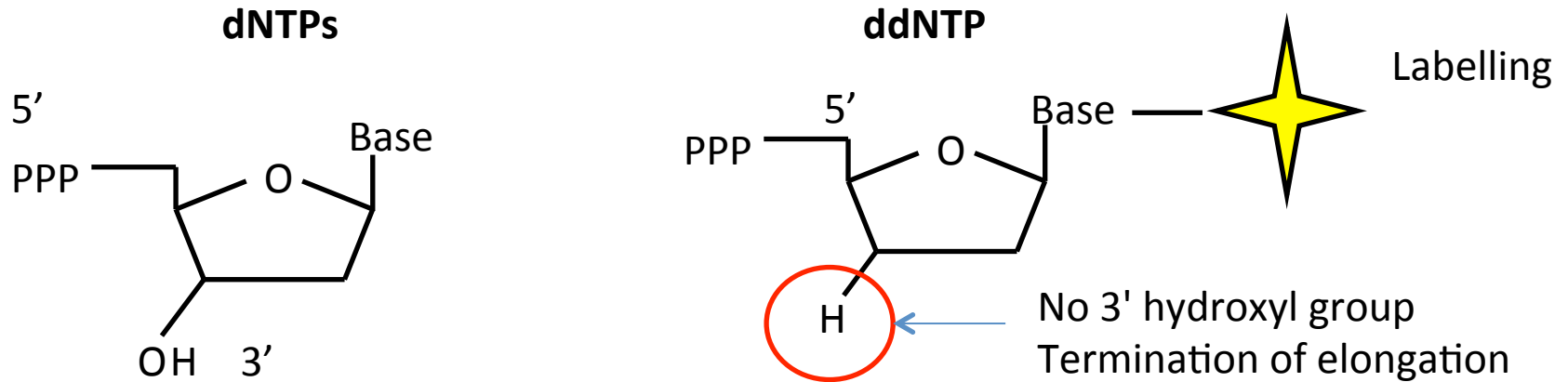
Pyrophosphate

From nucleotide to DNA



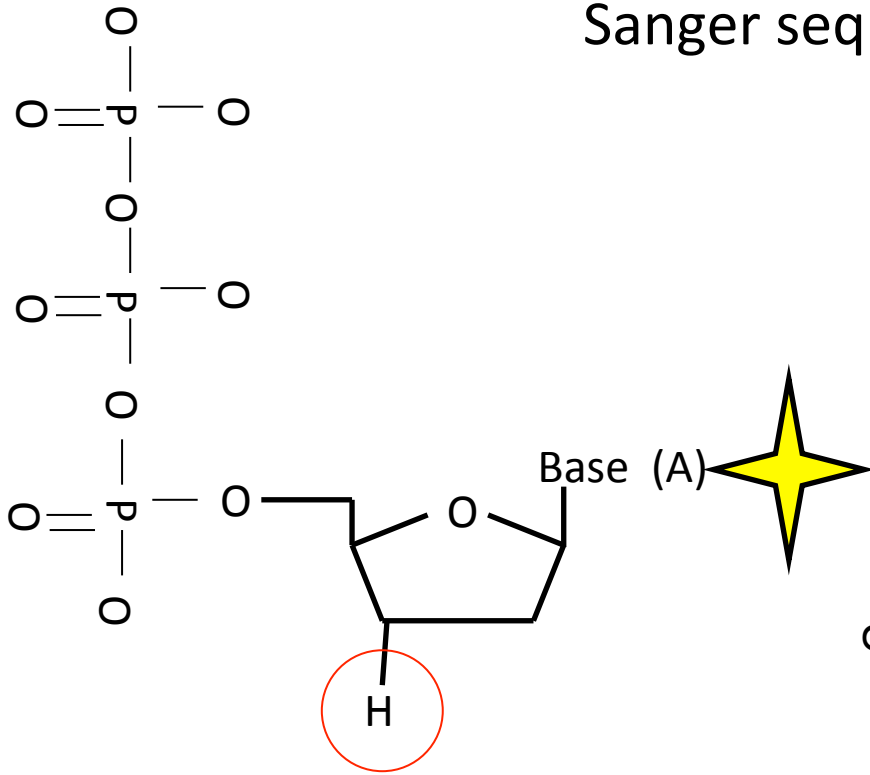
Capillary Sequencing – Chain Termination

Key principle of Sanger sequencing is the use of (radioactive or fluorescence) labelled dideoxy NTPs (ddNTPs) which terminate elongation of a DNA fragment

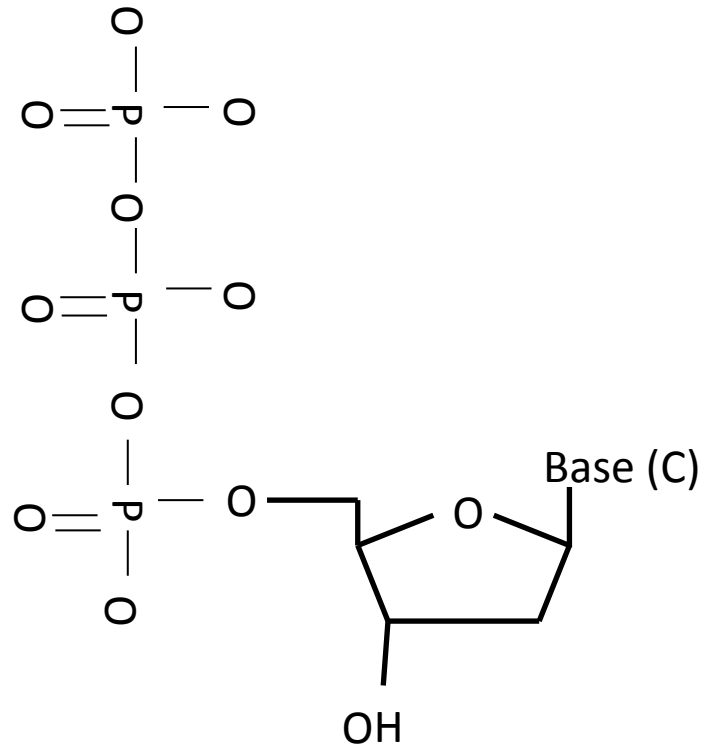


Normal and labelled ddNTPs in reaction

Sanger sequencing principle



--- No phosphodiester bond ---



Amplification is stopped for THIS molecule

||||| GTACGTTG*
AGCTGTCATGCAACGTCGTATGAC



———— GTACGTTG* (14nt)

||||| GTT*
AGCTGTCATGCAACGTCGTATGAC



———— GTT* (9nt)

||||| GTACGTTGCA*
AGCTGTCATGCAACGTCGTATGAC



———— GTACGTTGCA* (16nt)

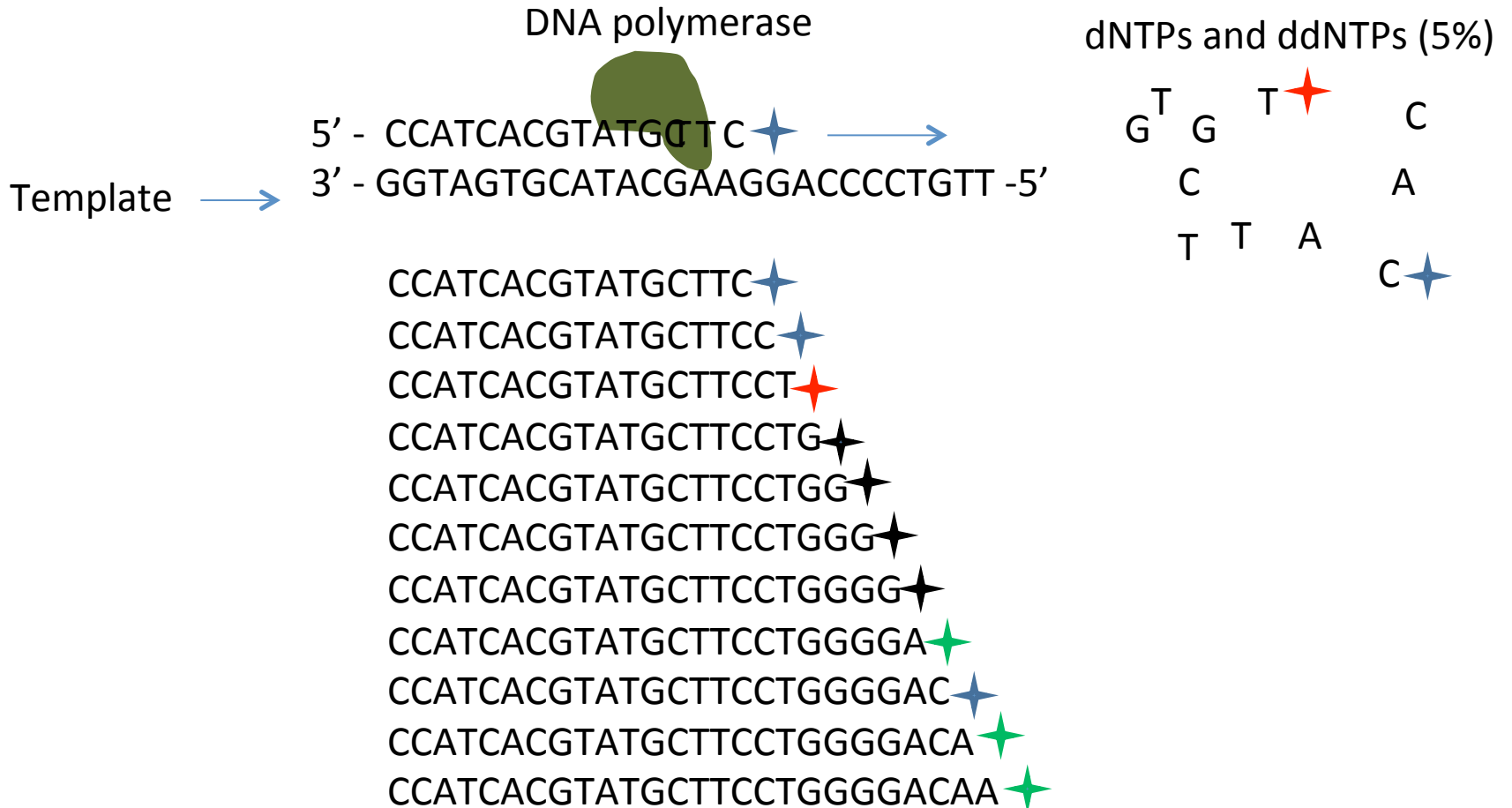
———— *
——— *
———— *
——— *



Separate and detect fluorescent labeled molecules

- Add:
- DNA polymerase
 - dNTPs
 - primer
 - ddNTPs (5%)
 - Buffer

Chain Termination

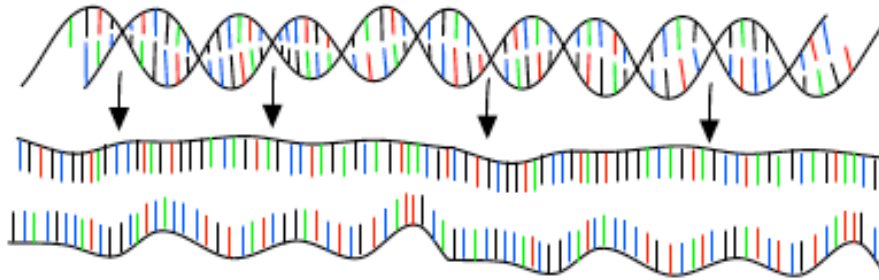


Sequencing

30 cycles of 3 steps :

Step 1 : denaturation

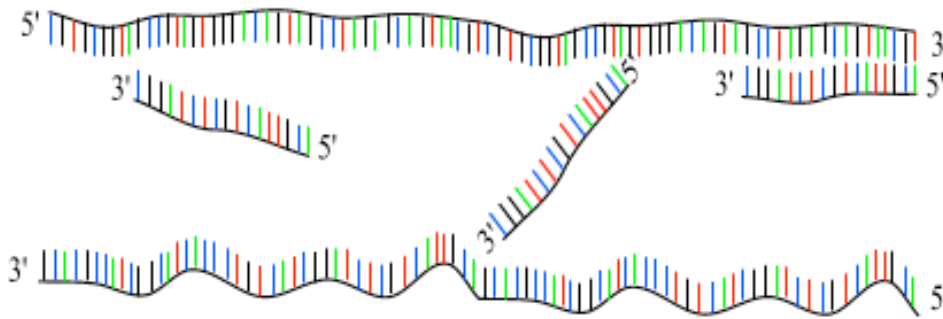
1 minut 94 °C



Step 2 : annealing

15 seconds 50 °C

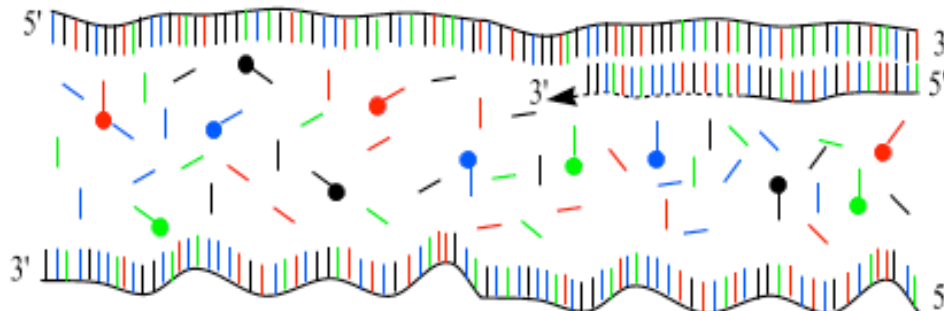
1 primer !!!!



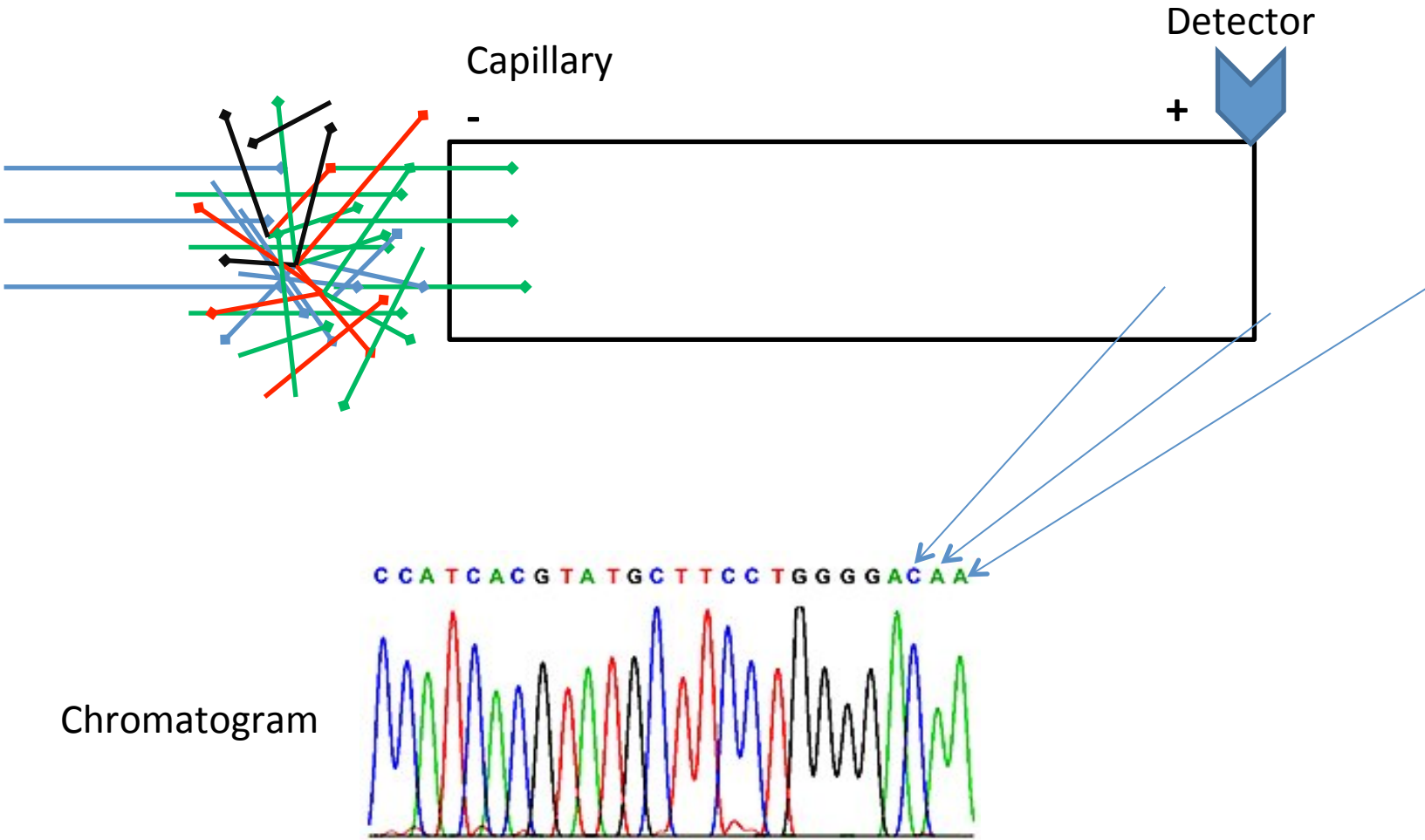
Step 3 : extension

4 minutes 60 °C

**mixture of dNTP's |
and ddNTP's ↓**



Capillary Sequencing

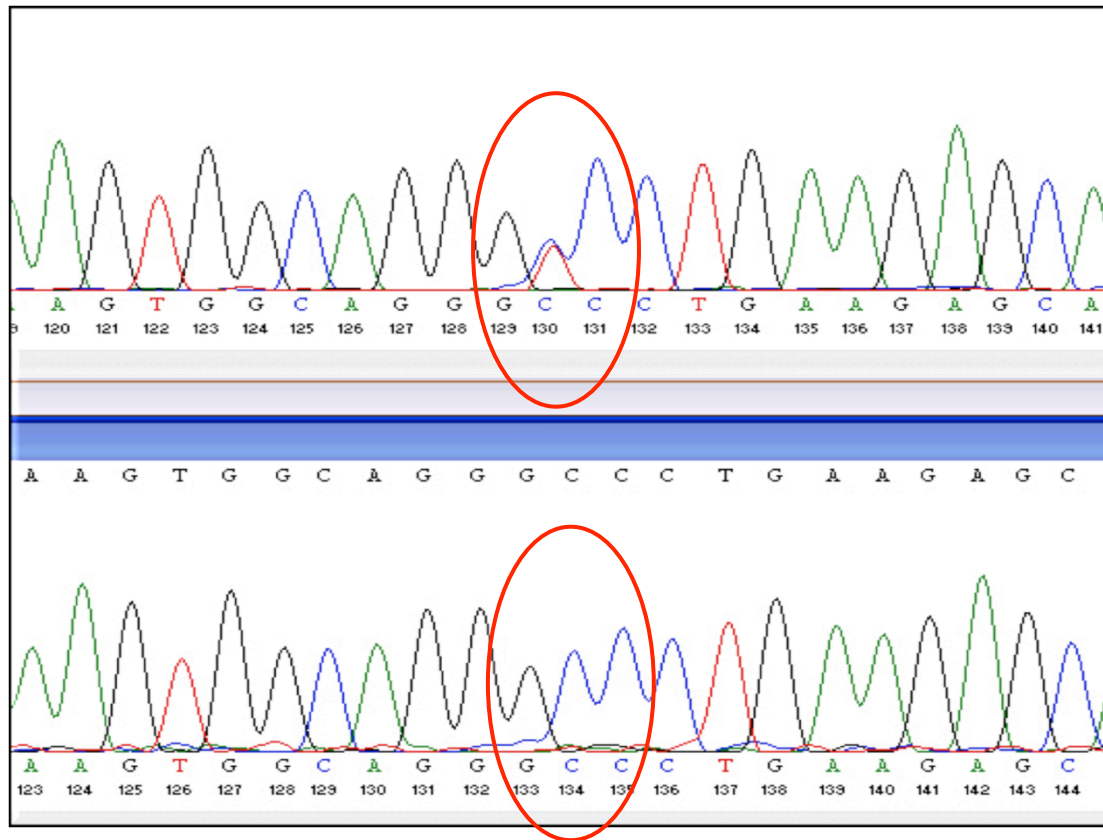


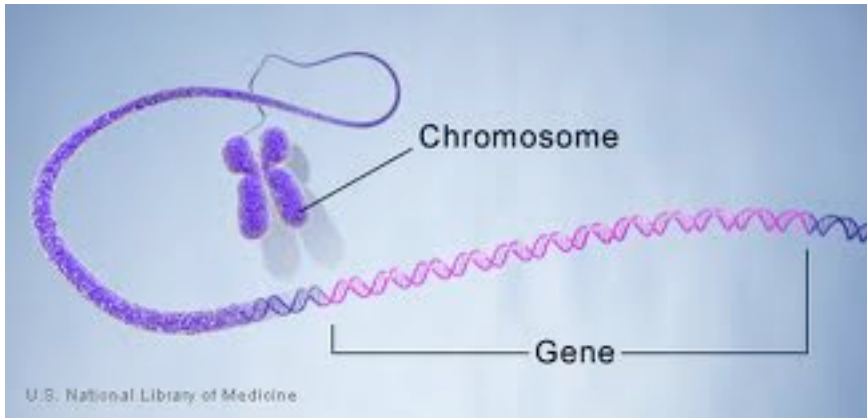
ABI 3730

- Individual reaction
- 0.0001 Gb/Run
- 1000 bp reads
- 1-2 hour run time
- Fast and easy for individual samples
- Accurate
- Easy adoption
- Robust technology
- Out dated!!

Question:

1) How can you explain the following sanger sequence chromatogram?

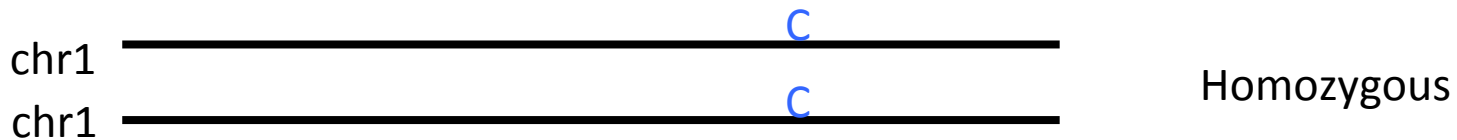
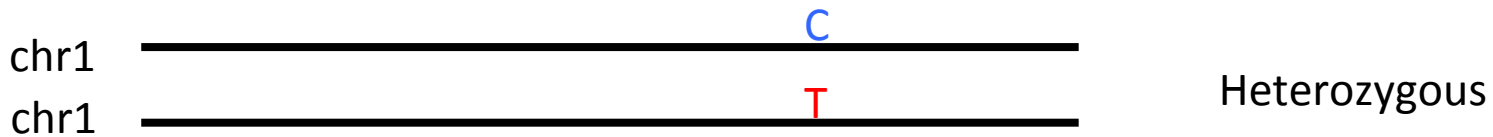




2 copies of the chromosome in one cell
Amplify the gene by PCR

→ homozygous: these two copies are IDENTICAL

→ Heterozygous: these two copies are NOT identical (SNPs)



Next Gen Sequencing (NGS) Platforms

Illumina

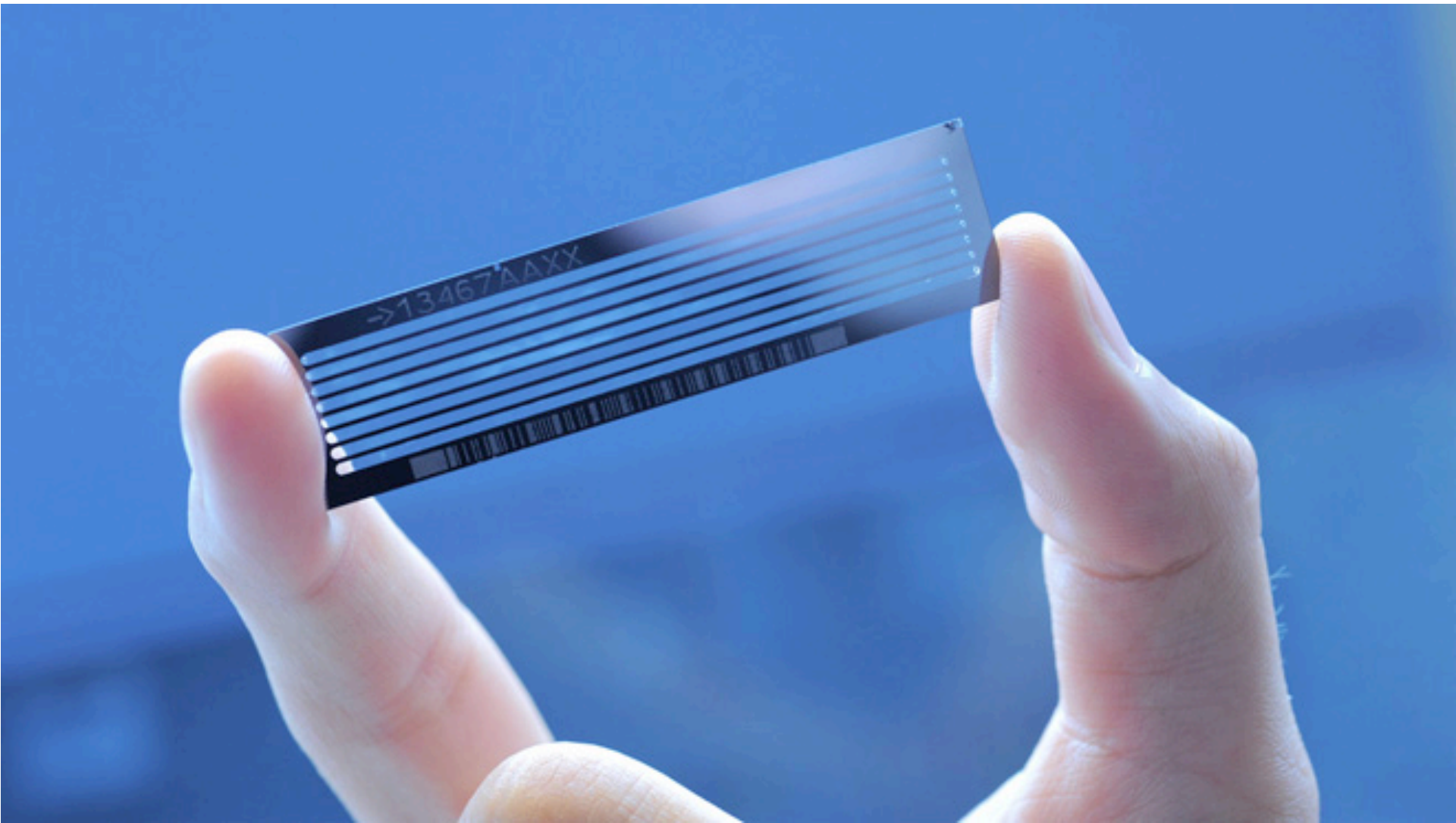
Pacific Biosciences (PacBio)

Illumina Sequencing

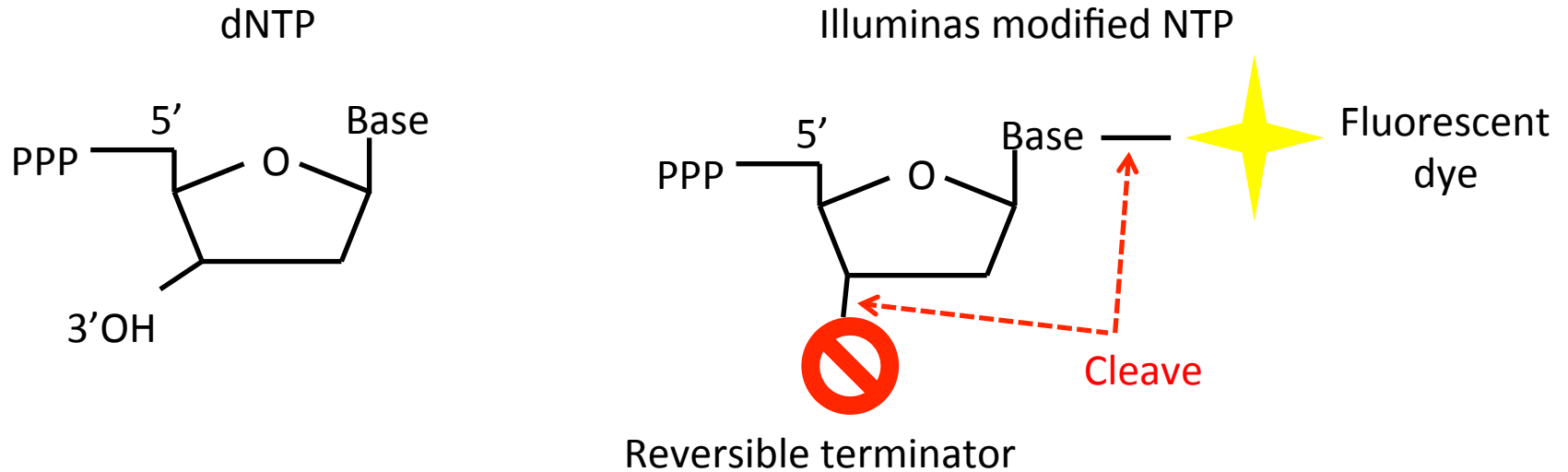
Has become the most popular next generation sequencing platform.

Synthetic sequencing using fluorescent reversible terminator technology.



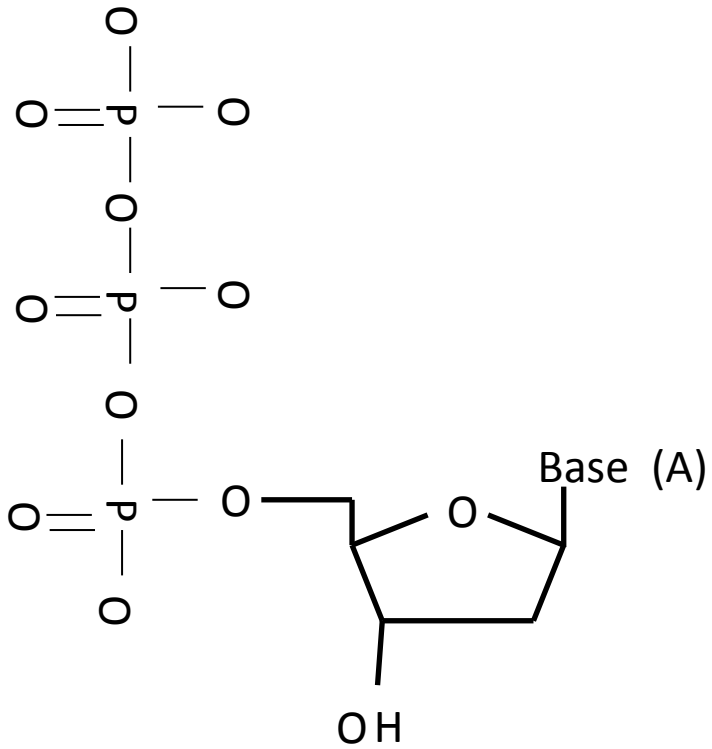


Illumina Sequencing

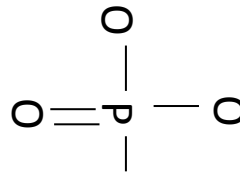


Every single base labelled (Each of the 4 bases have a different fluorescent dye)

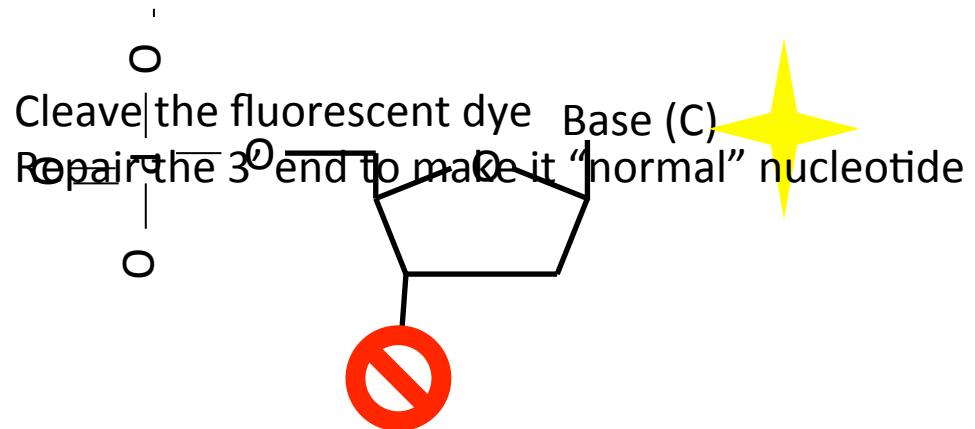
Every base has a terminator group -**cleaved off** after each round of sequencing



hydrogen

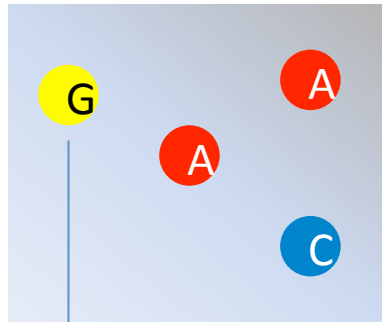


Pyrophosphate

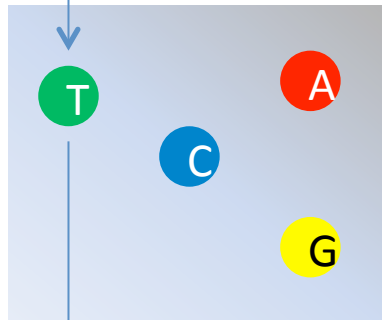


OH

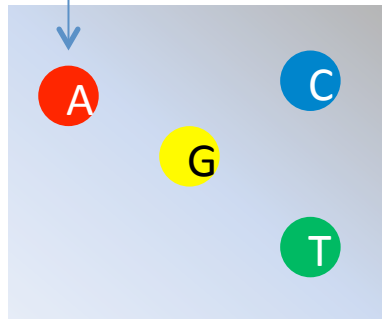
Illumina Base Calling



Cycle 1



Cycle 2



Cycle 3

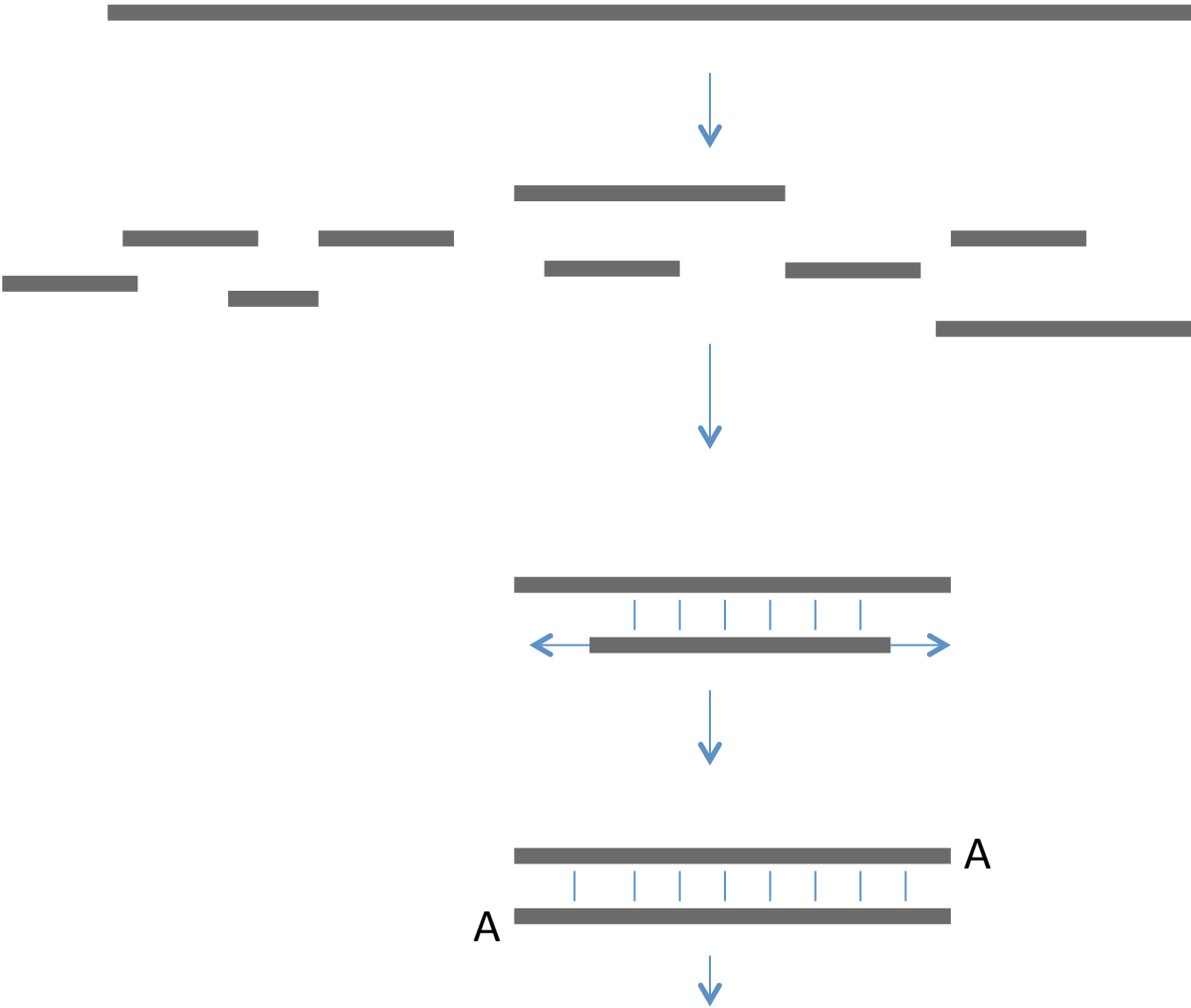
G T A....

```

CAAGGCCGCCAACAATGGTGGTGATAAGCGGGGGTG
AAGGCCGCCAACAATGGTGGTGATAAGCGGGGTGG
AGGCCGCCAACAATGGTGGTGATAAGCGGGGTGGC
AGGCCGCCAACAATGGTGGTGATAAGCGGGGTGGC
AGGCCGCCAACAATGGTGGTGATAAGCGGGGTGG
GGCCGCCAACAATGGTGGTGATAAGCGGGGTGGCG
GGCCGCCAACAATGGTGGTGATAAGCGGGGTGGCGT
GCCGCCAACAATGGTGGTGATAAGCGGGGTGGCGT
CGCCAACAATGGTGGTGATAAGCGGGGTGG
GCCAACAATGGTGGTGATAAGCGGGGTGGCGTGAT
CCAACAATGGTGGTGATAAGCGGGGTGG
CCAACAATGGTGGTGATAAGCGGGGTGGCGTGATG
CCAACAATGGTGGTGATAAGCGGGGTGGCGTGATG
CAACAATGGTGGTGATAAGCGGGGTGGCGTGATG
CAACAATGGTGGTGATAAGCGGGGTGGCGTGATG
CAACAATGGTGGTGATAAGCGGGGTGGCGTGATG
ACAATGGTGGTGATAAGCGGGGTGG
ACAATGGTGGTGATAAGCGGGGTGGCGTGATGCAT
CAATGGTGGTGATAAGCGGGGTGG
AATGGTGGTGATAAGCGGGGTGGCGTGATGCATTG
AATGGTGGTGATAAGCGGGGTGGCGTGAT
ATGGTGGTGATAAGCGGGGTGGCGTGATGCATTCC
AAAGGCGAGCACAAAGGCCGCCAACAATGGTGGTGATAAGCGGGGTGGCGTGATGCATTCCGTCCTCTTCTCTGGTGGI
    
```

Aligned to consensus

Illumina library



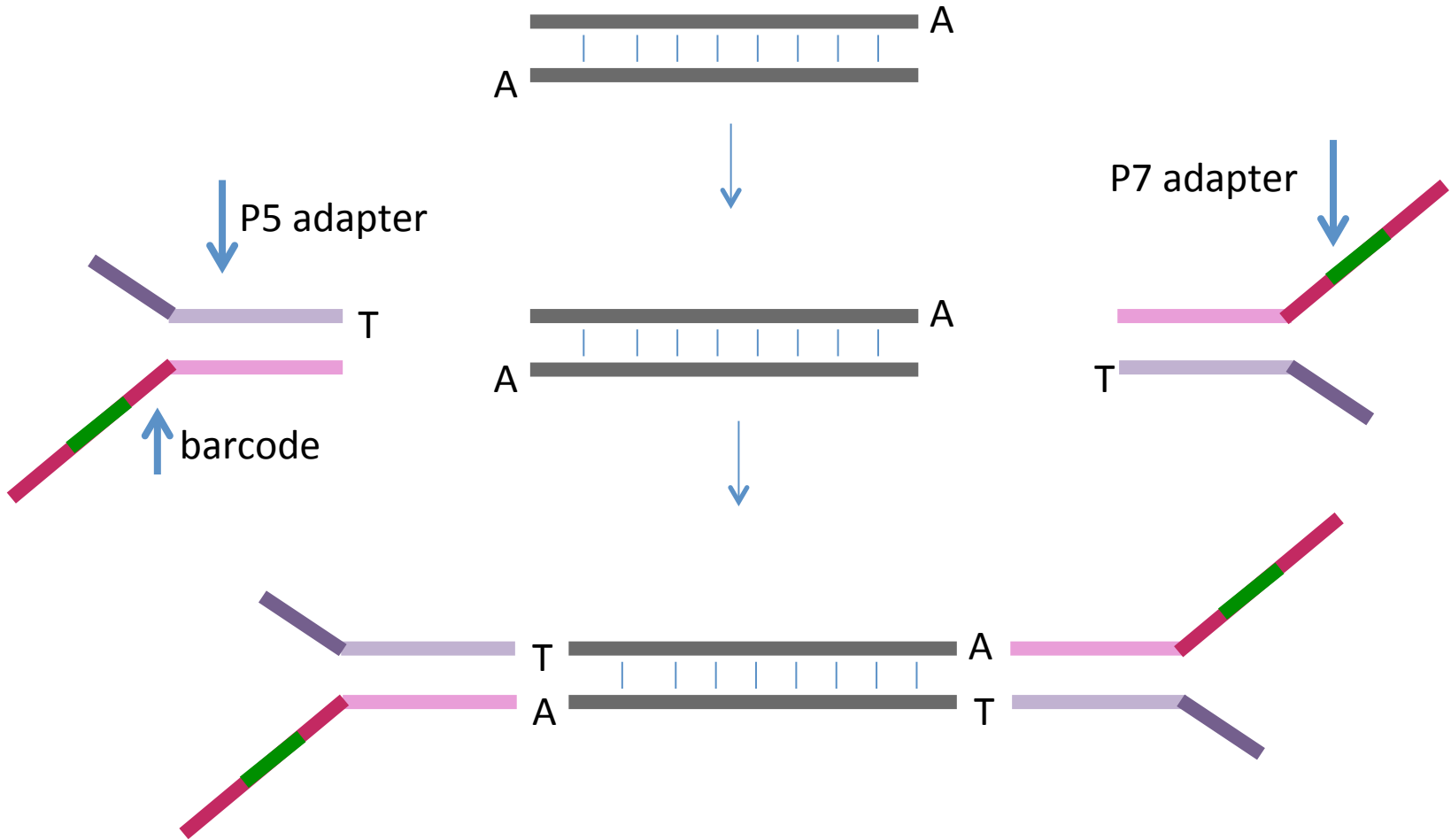
Genome

Fragment Genome
50 to 500bp

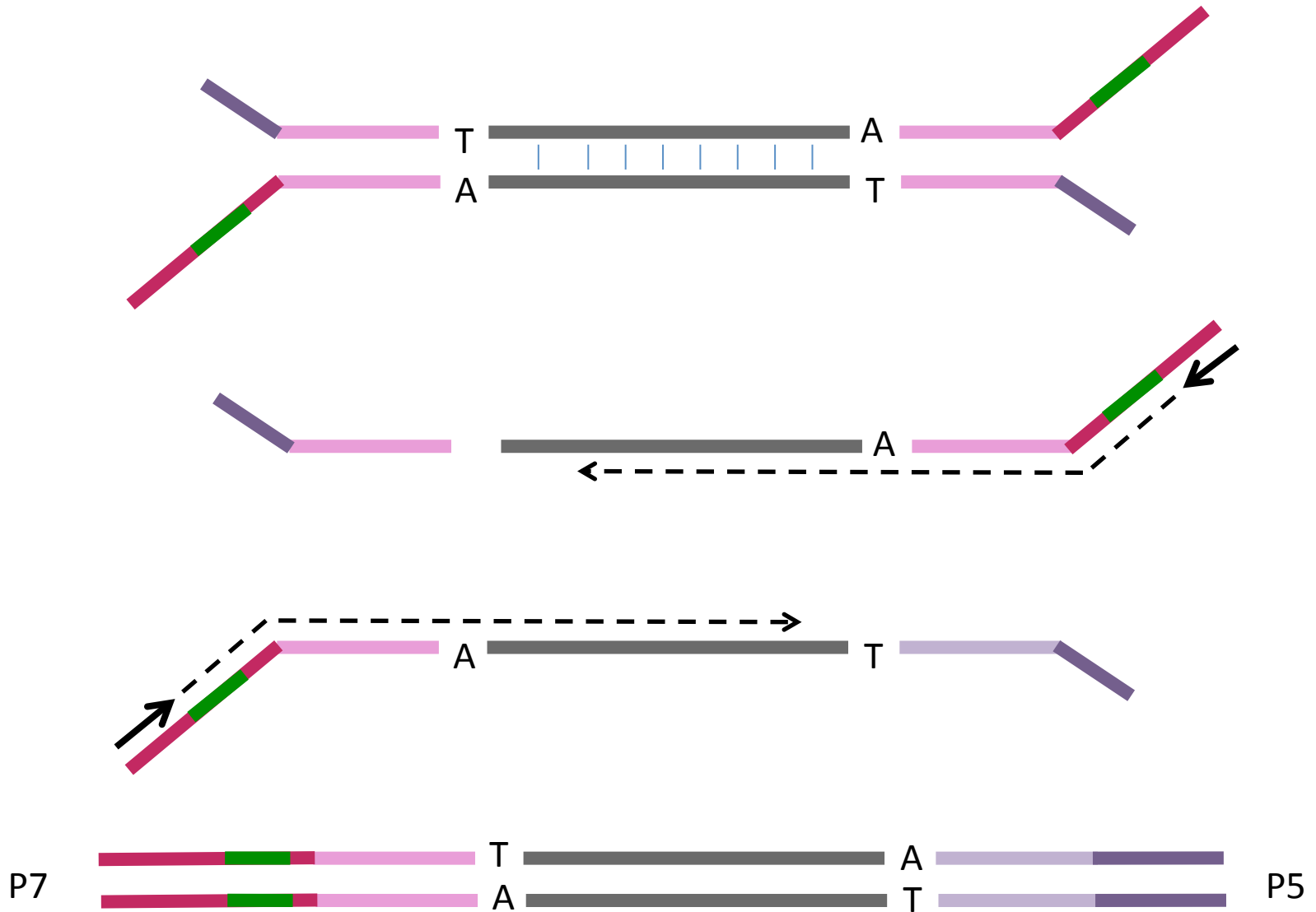
End Repair :DNA ends
are blunted and
phosphorylated

A tailing to the 3'-
end of each DNA
fragment

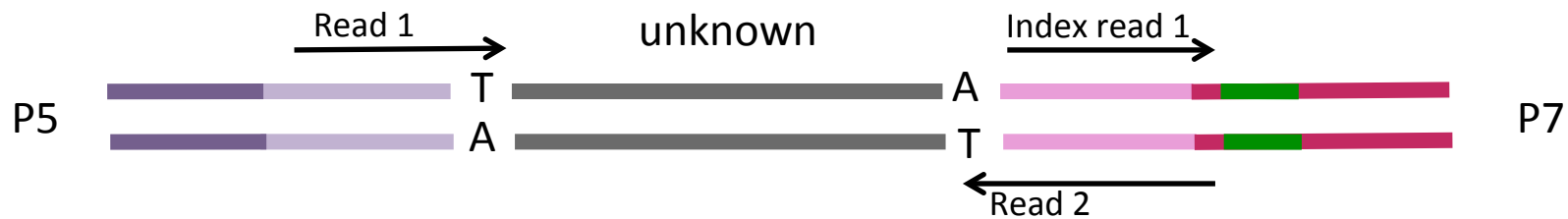
Illumina library



Illumina library



A billion times, coverage the entire genome



5'- AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATC* T 3'

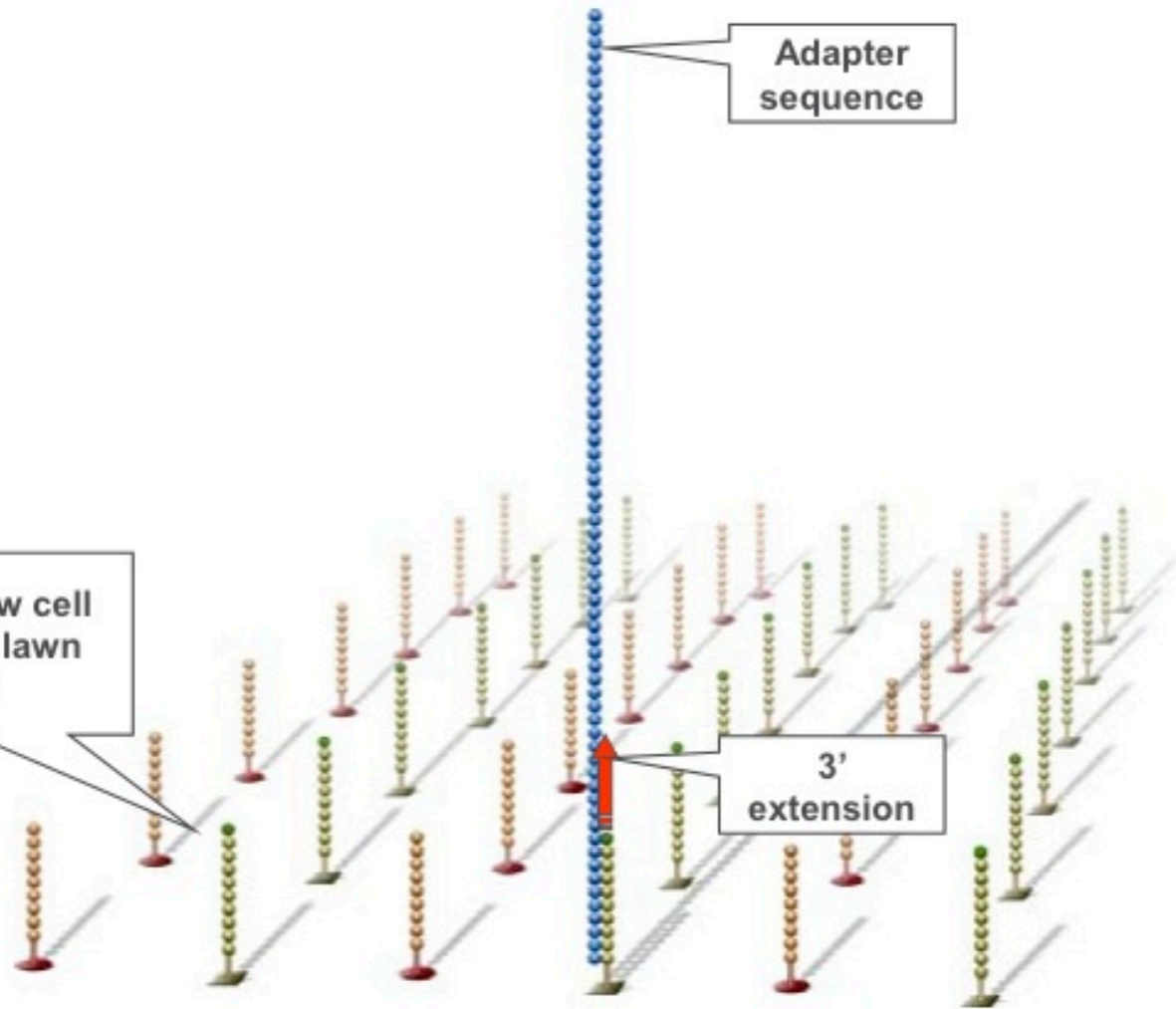
5'- /5Phos/GATCGGAAGAGCACACGTCTGAACTCCAGTCAC BC ATCTCGTATGCCGTCTTCTGCTT*G 3'

Hybridize Fragment & Extend

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

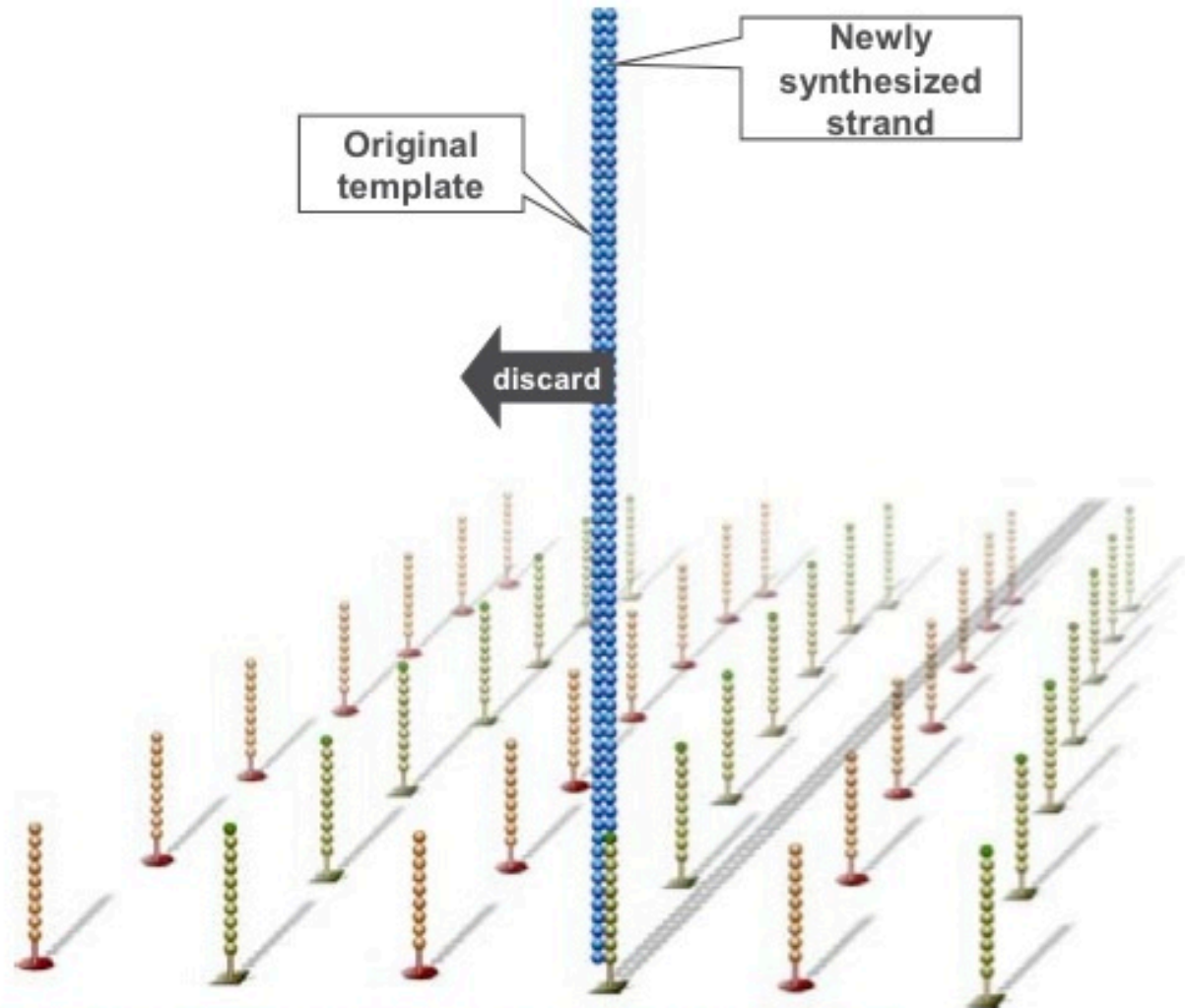


Denature Double-Stranded DNA

Double-stranded molecule is denatured

Original template washed away

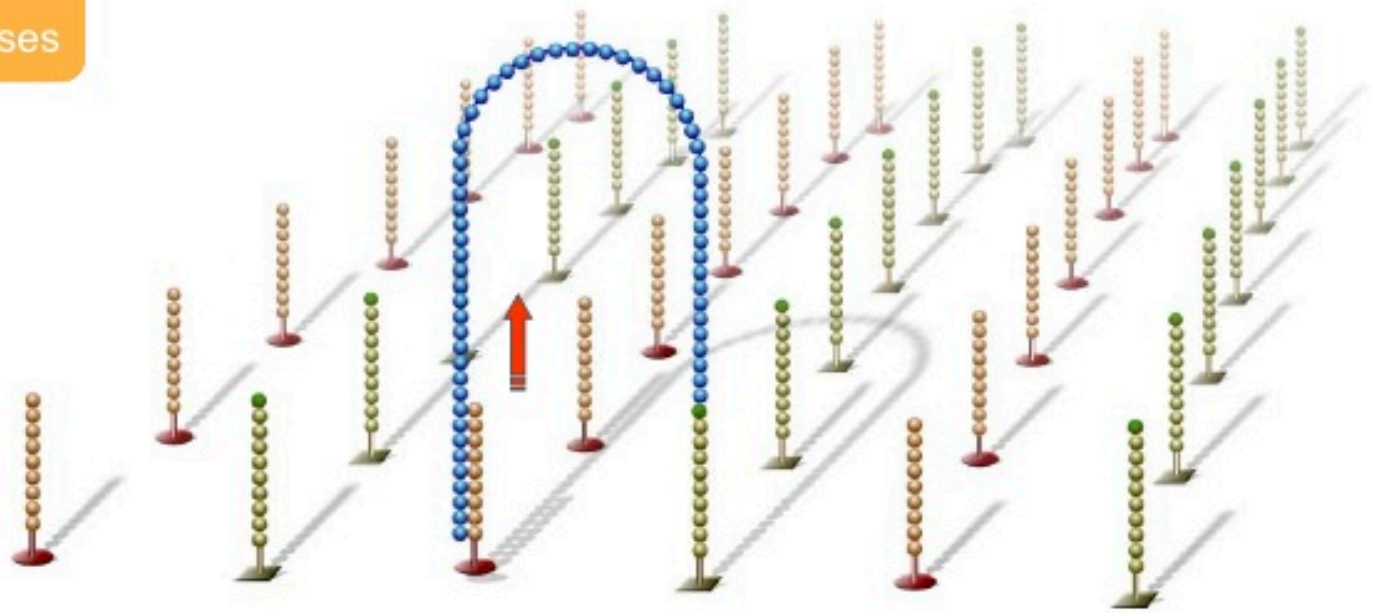
Newly synthesized strand is covalently attached to flow cell surface



Bridge Amplification

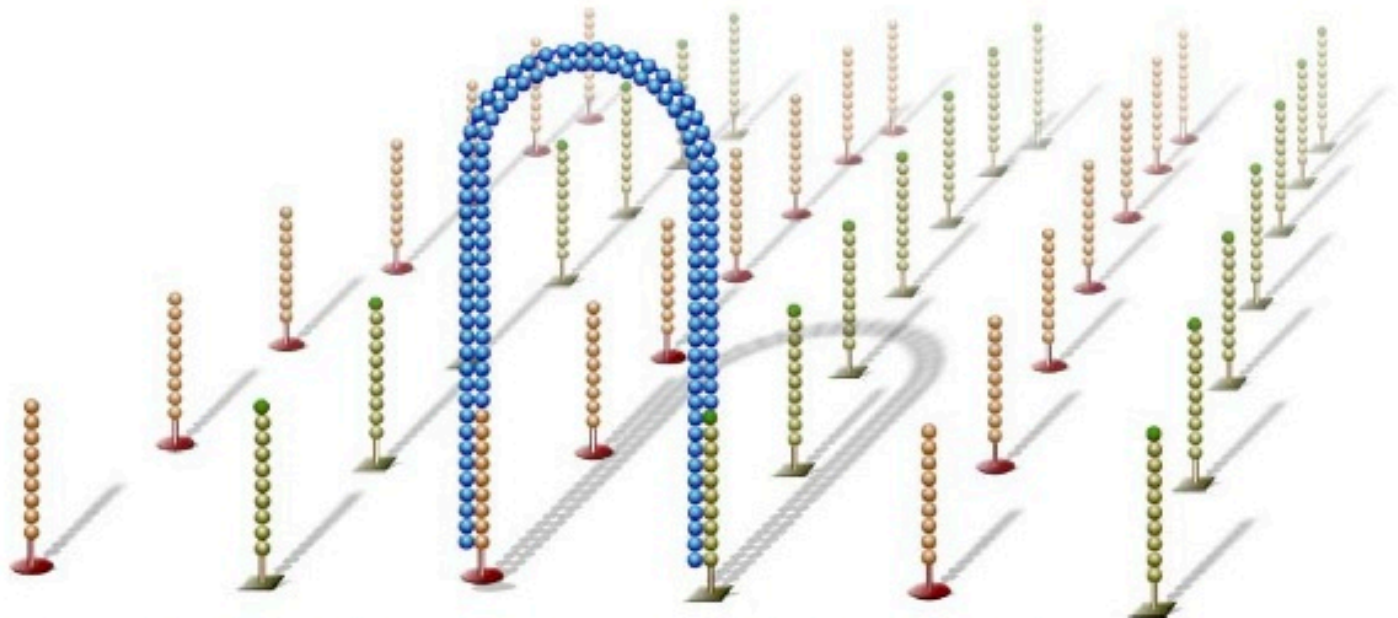
Single-stranded molecule flips over and forms a bridge by hybridizing to adjacent, complementary primer

Hybridized primer is extended by polymerases



Bridge Amplification

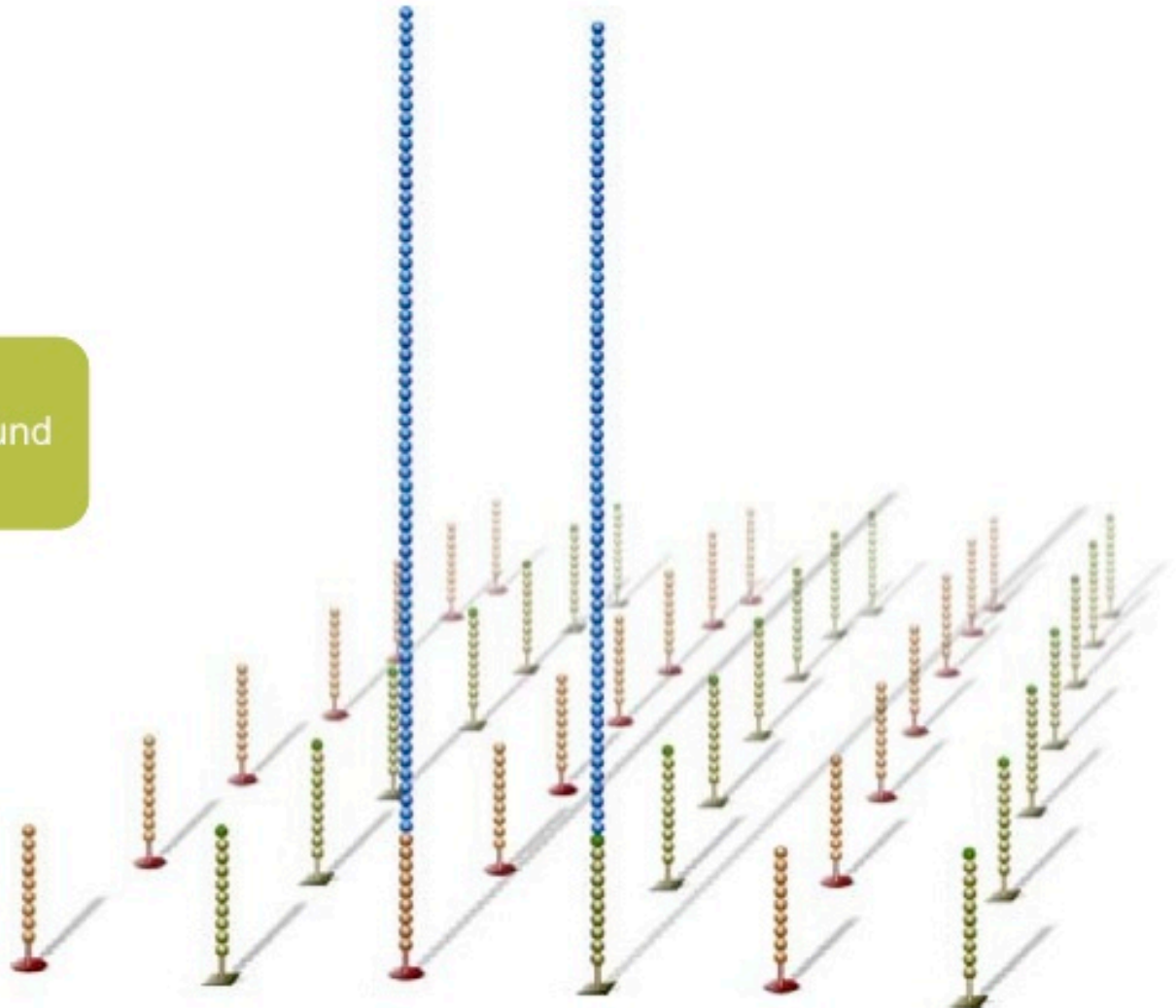
Double-stranded bridge is formed



Denature Double-Stranded Bridge

Double-stranded bridge is denatured

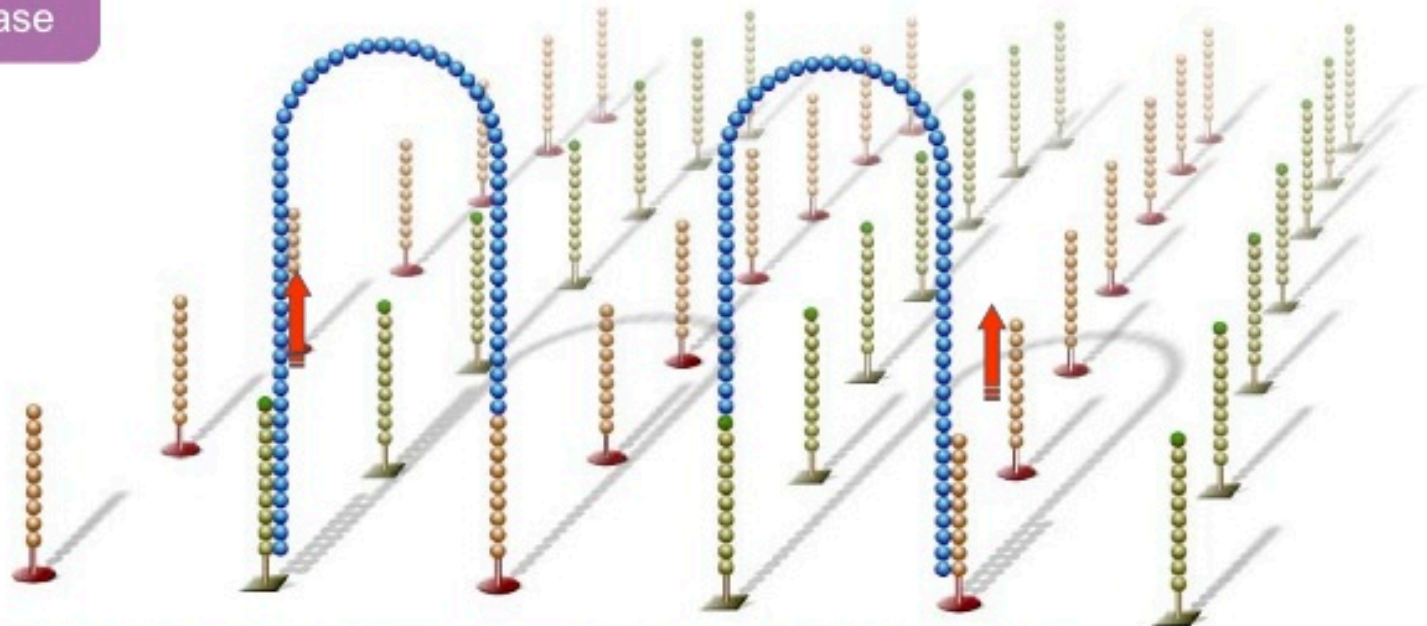
Result:
Two copies of covalently bound single-stranded templates



Bridge Amplification

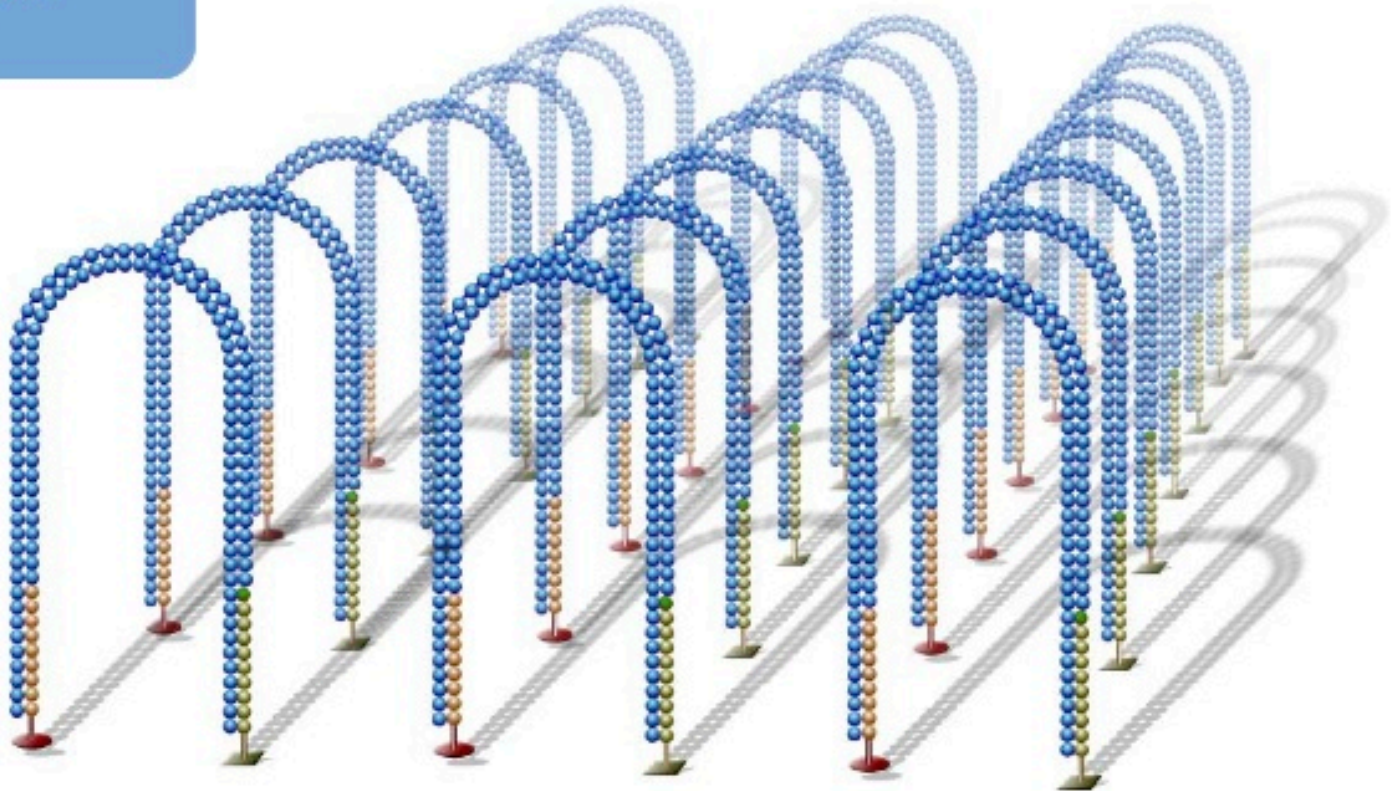
Single-stranded molecules flip over to hybridize to adjacent primers

Hybridized primer is extended by polymerase



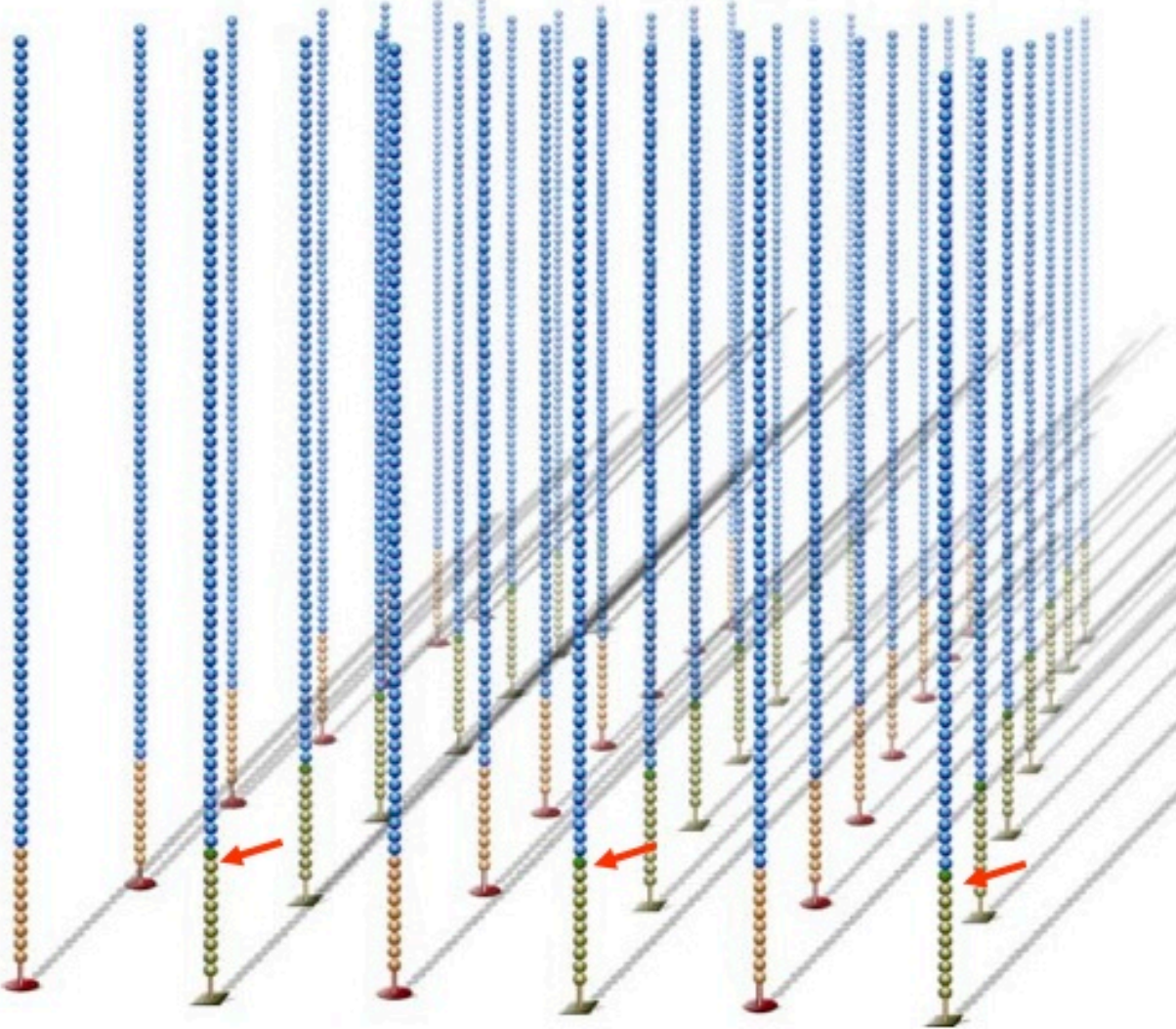
Bridge Amplification

Bridge amplification cycle repeated until multiple bridges are formed



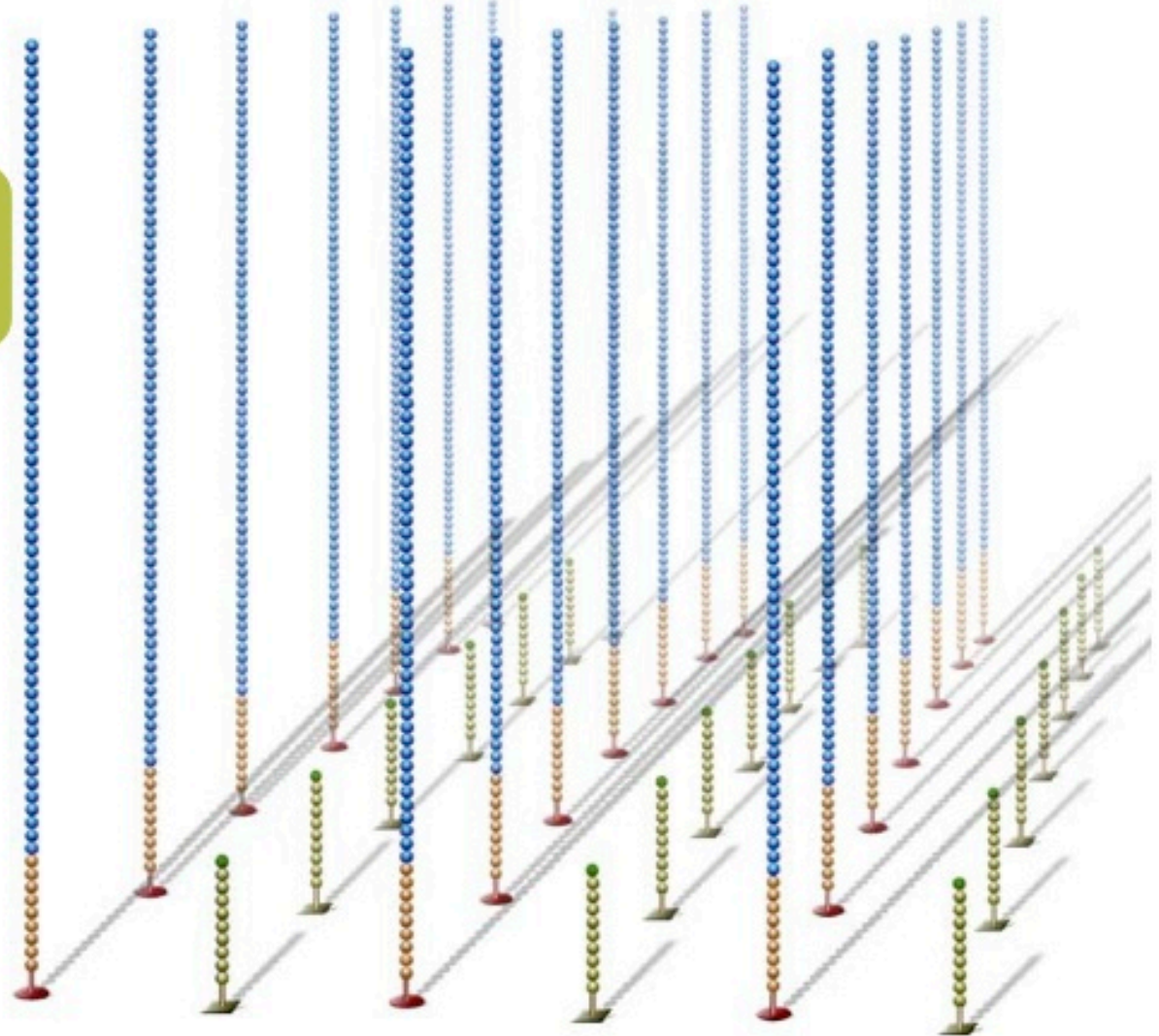
Linearization

dsDNA bridges are denatured



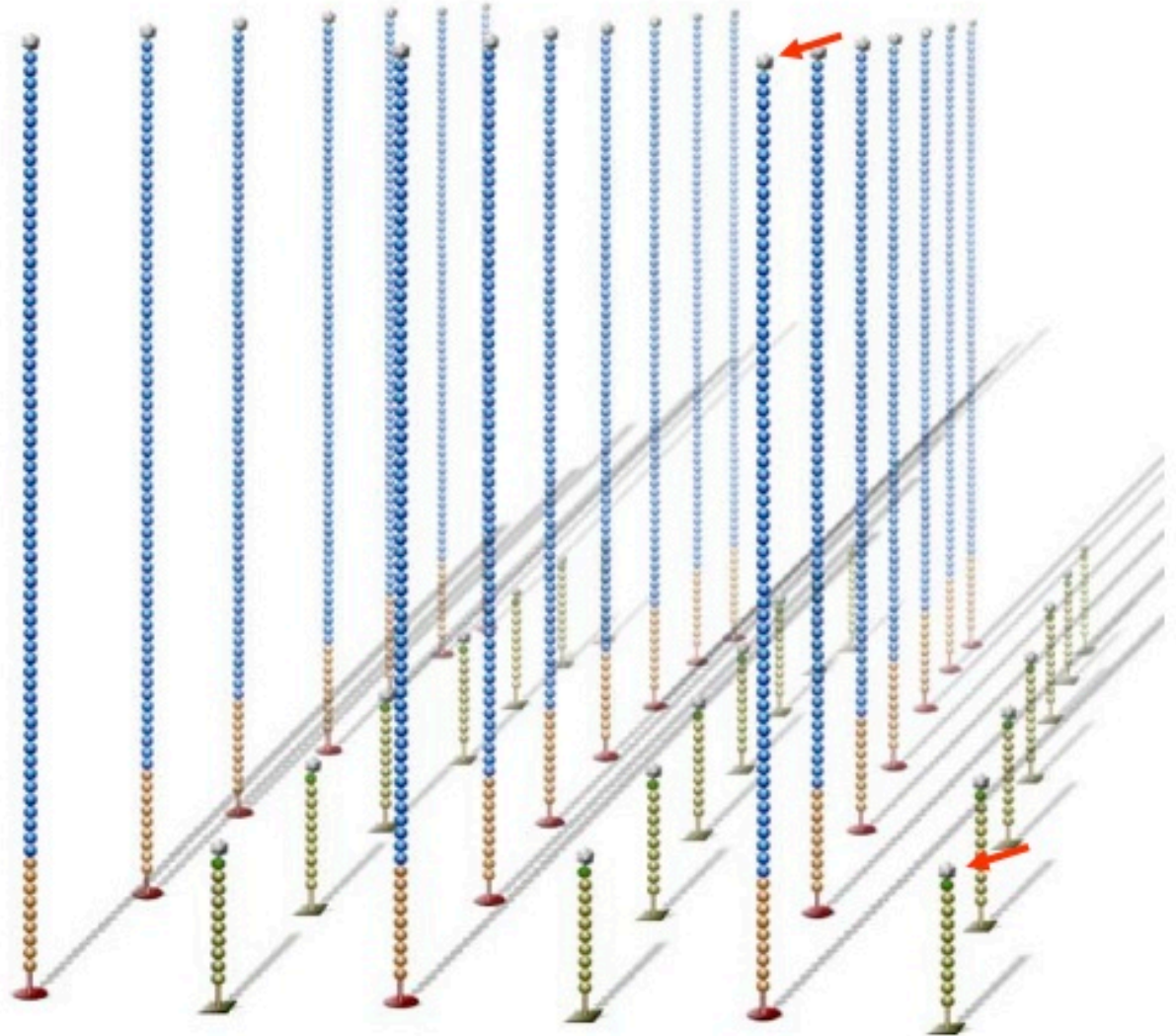
Reverse Strand Cleavage

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



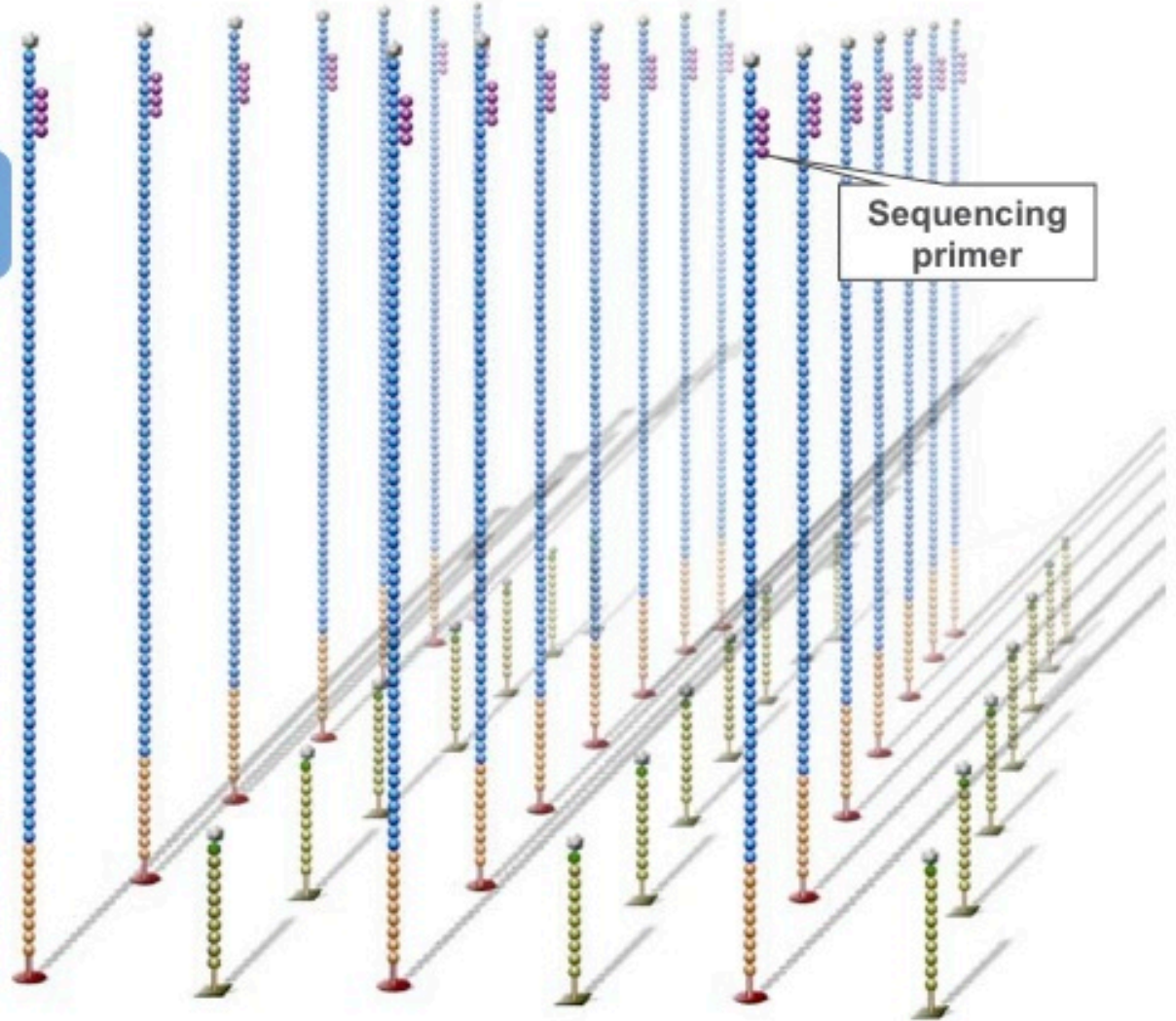
Blocking

Free 3' ends are blocked to prevent unwanted DNA priming

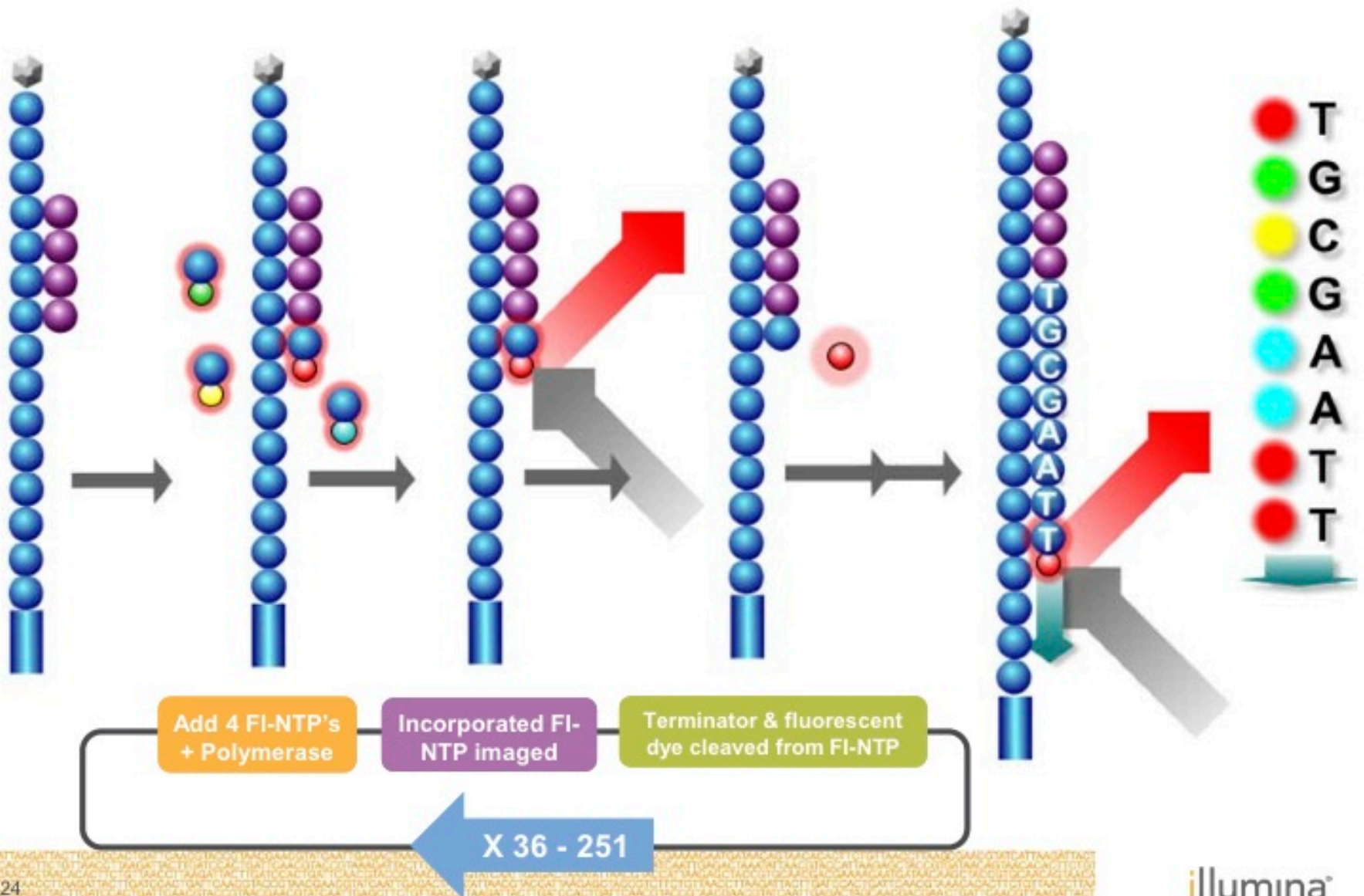


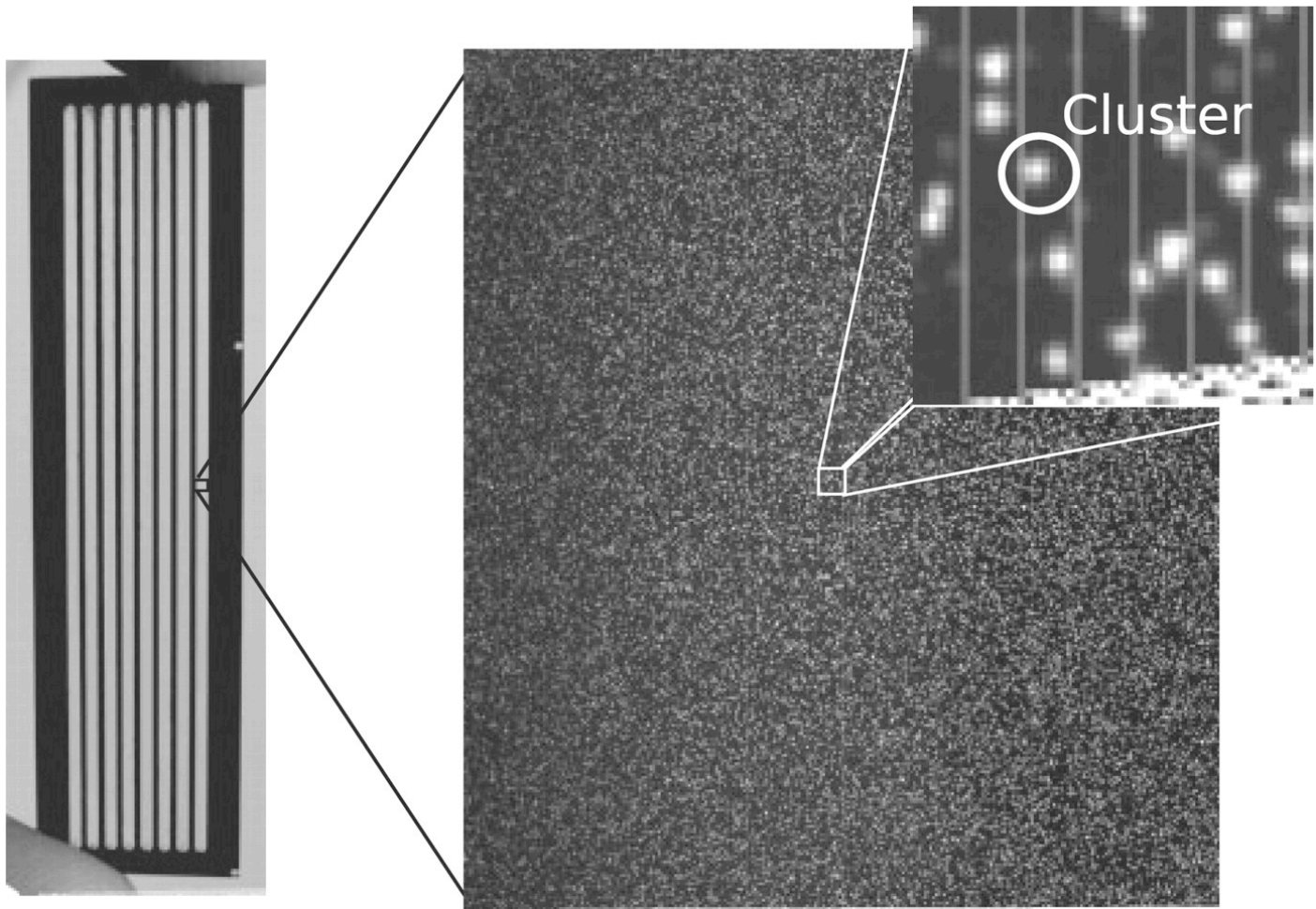
Read 1 Primer Hybridization

Sequencing primer is hybridized to adapter sequence



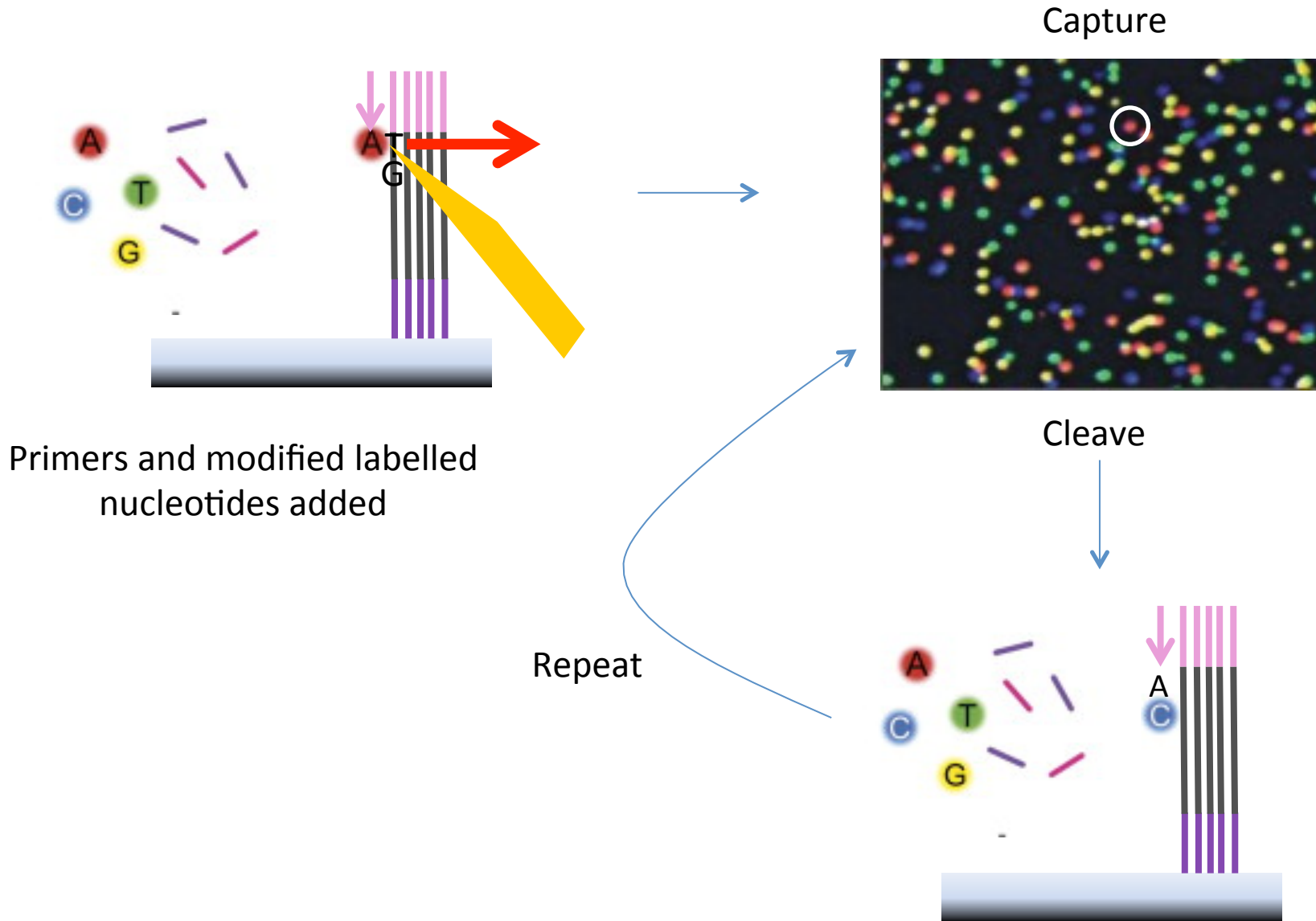
Sequencing by Synthesis





Tens of millions clusters per lane

Illumina Sequencing



Illumina - HiSeq

- 2 Tb/Run (8 lanes, 2 flow cells)
- up to 150 bp reads
- 12 day run time / 24hrs on Rapid Mode
- High throughput
- Low error rate (0.1%)
- Simple data format

Illumina - Miseq

The background of the slide features a grayscale image of an Illumina MiSeq sequencer. On the left, the physical instrument is shown, with the 'MiSeq' logo visible on its side. On the right, a computer monitor displays the 'Illumina Control Software' interface, which includes a large central button and several smaller icons below it.

- 15 Gb/Run
- up to 300 bp reads (longer reads)
- 1 day run time (shorter)
- Low error rate (0.1%)
- Simple data format

???

My supervisor gave me the task to find out how many samples (*Arabidopsis*, ~135Mb genome) he can run on one HiSeq lane in order to get 40x coverage for each genotype. He wants to use it for de novo assembly.

- 1 Hiseq lane produces 50Gbase / lane
- How many *Arabidopsis* samples can he pool??

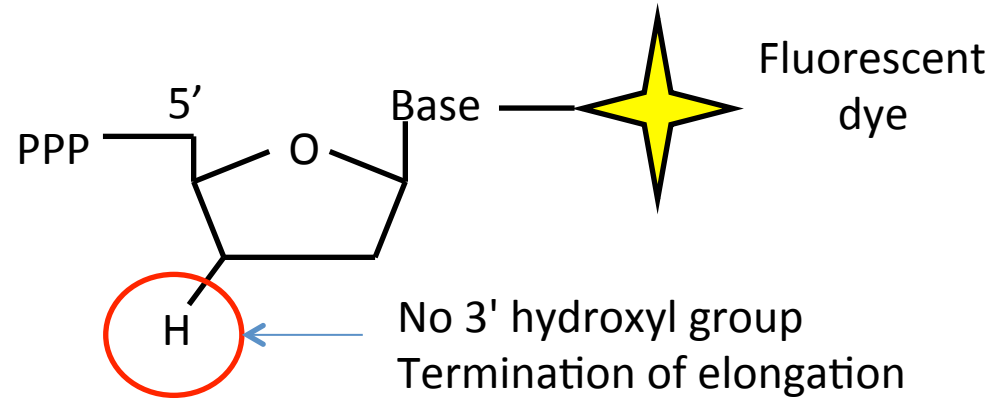
$$135\text{Mb} \times 40 = 5.4 \text{ Gb} \quad 50/5.4 = \underline{9.25 \text{ samples/lane}}$$

Pacific Biosciences (PacBio)

- * Pyrosequencing in a zero mode waveguide (ZMW)
- * Single molecule real time sequencing (SMRT)
- * Accuracy ~85%
- * Sequencing of long DNA molecules (20kb)
- * Very suitable for sequencing repeats, cap filling and GC rich genomes

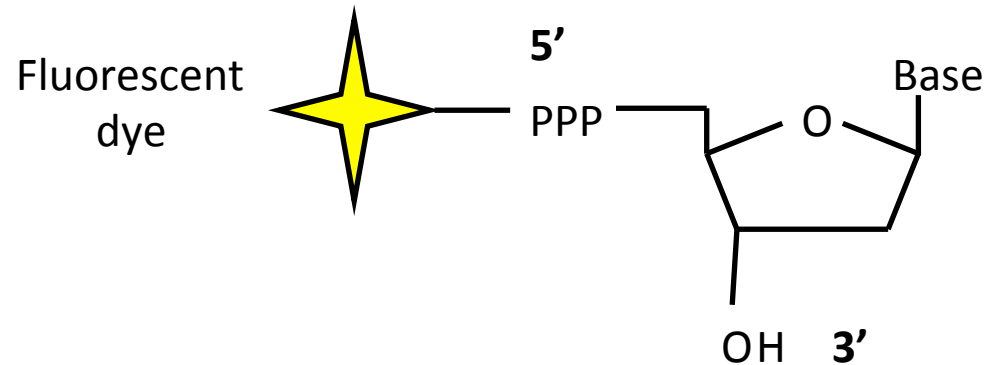


PacBio vs Sanger

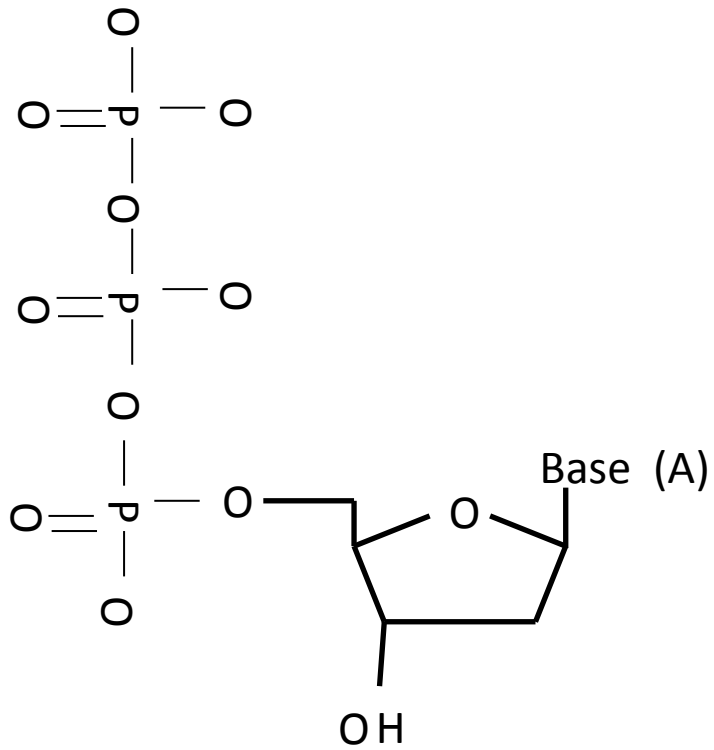


Sanger ddNTPs

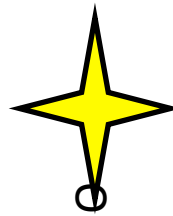
PacBio Phospho-Linked Nucleotides



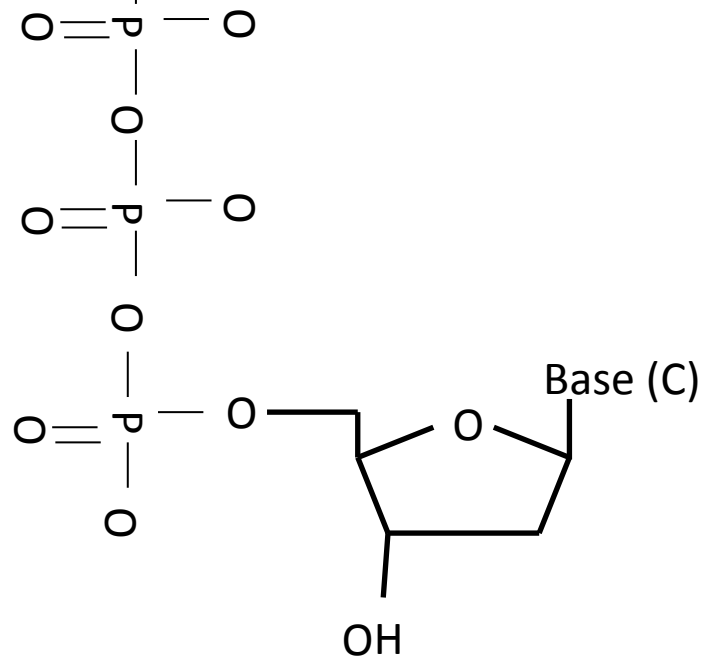
During polymerase cleaves away fluorescent label, leaving natural DNA



hydrogen

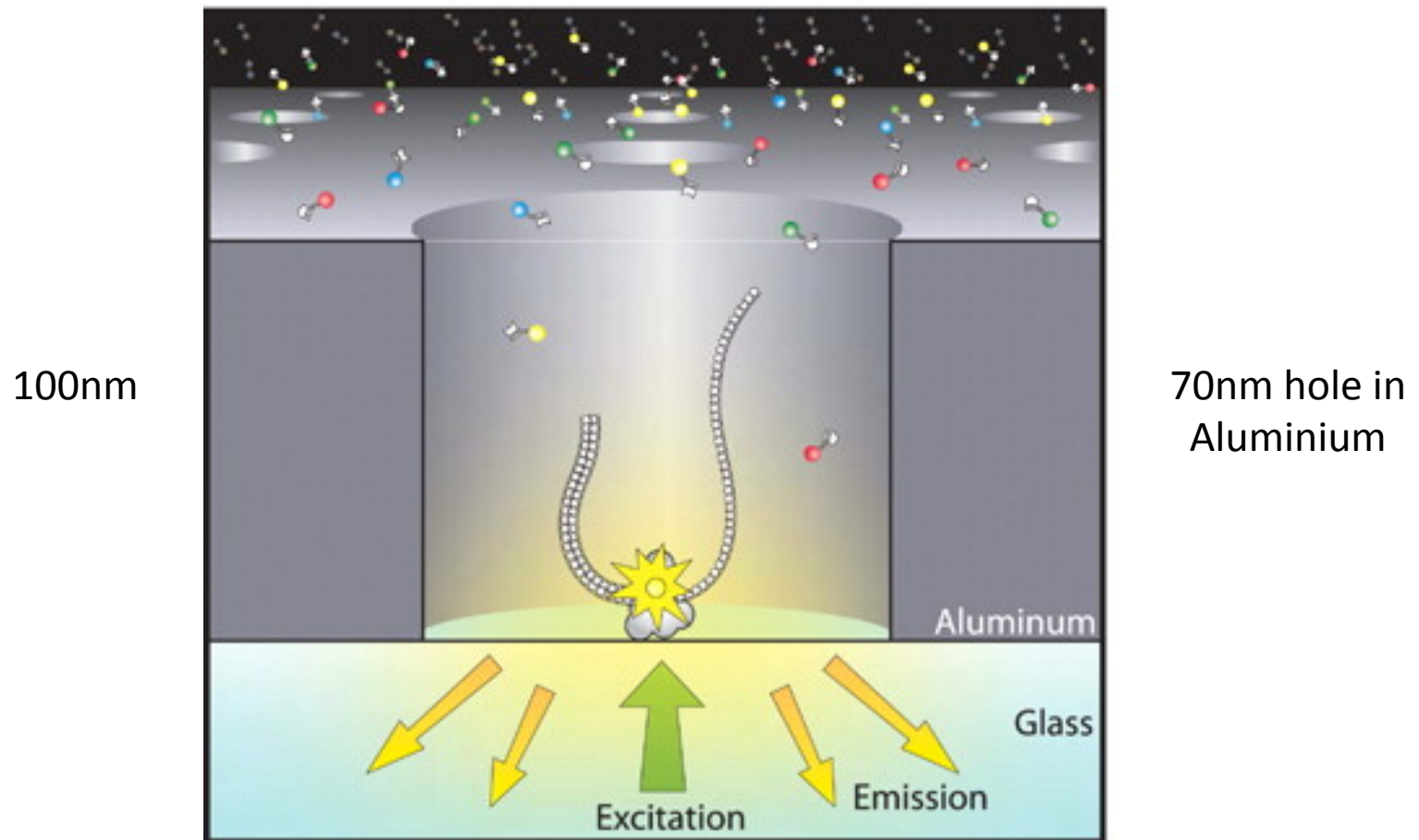


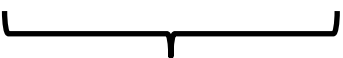
Phyrophosphate



Pacific Biosciences

Uses a zero mode waveguide (ZMW) to guide light energy to a single DNA polymerase with a single DNA template

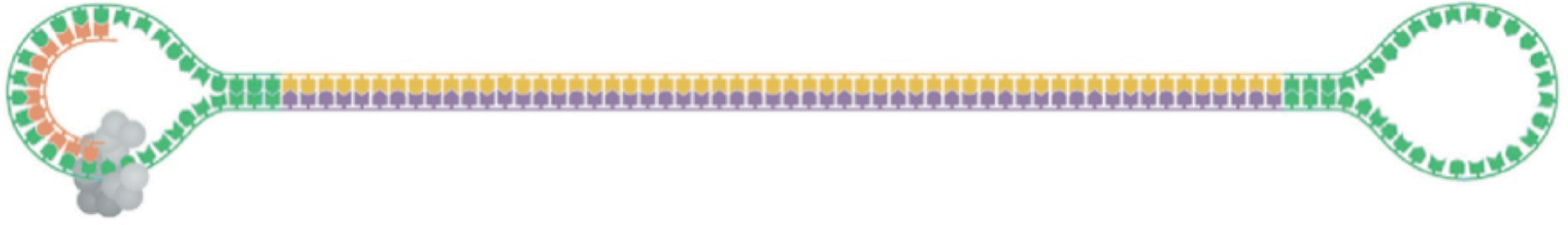


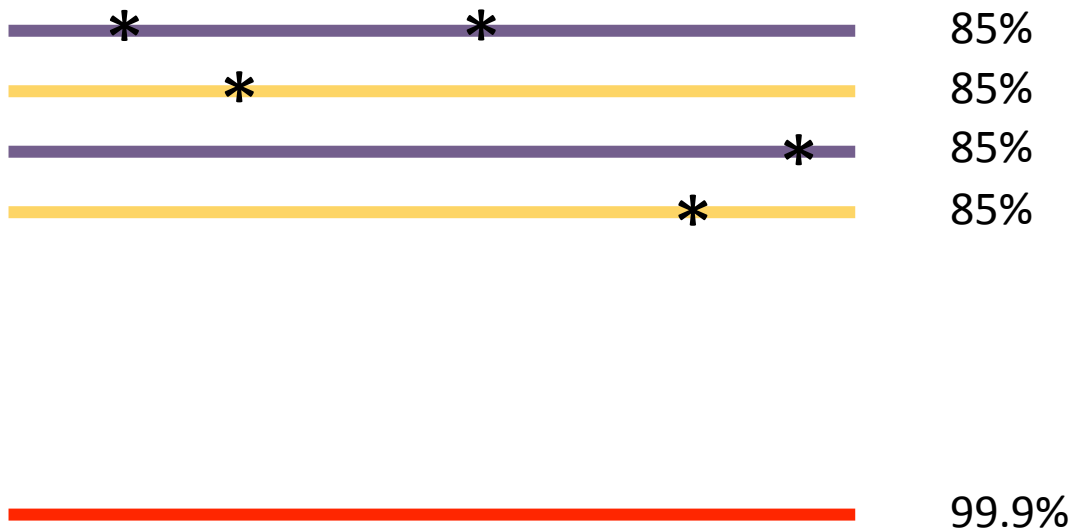
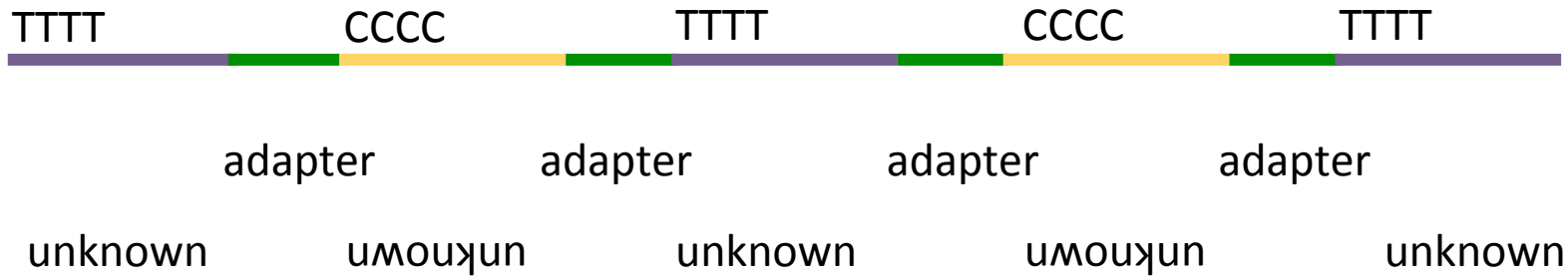


SMRT adapter



SMRT adapter





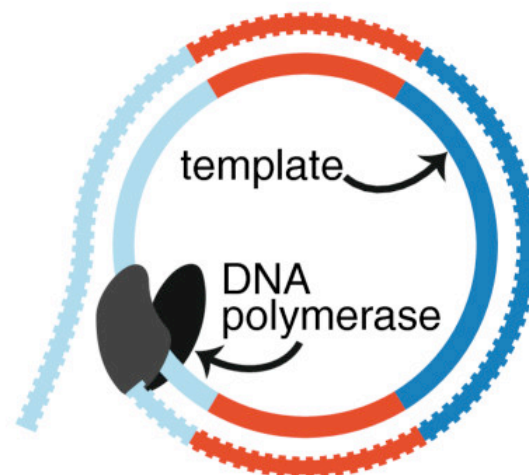
1. generate amplicon

5' forward strand 3'
3' reverse strand 5'

2. ligate adaptors



3. sequence

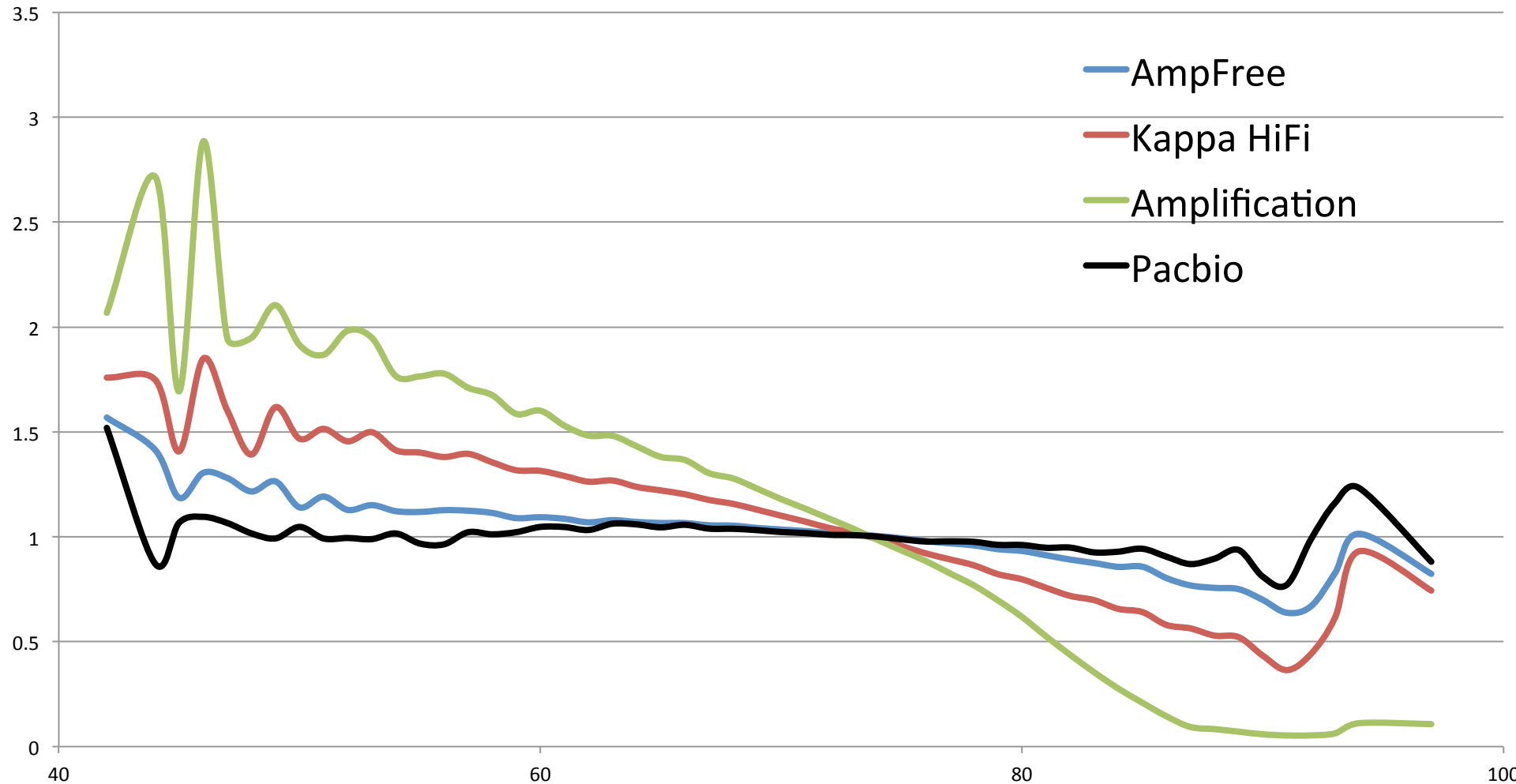


4. data analysis



1° analysis

GC bias – relative coverage



Pacific Biosciences

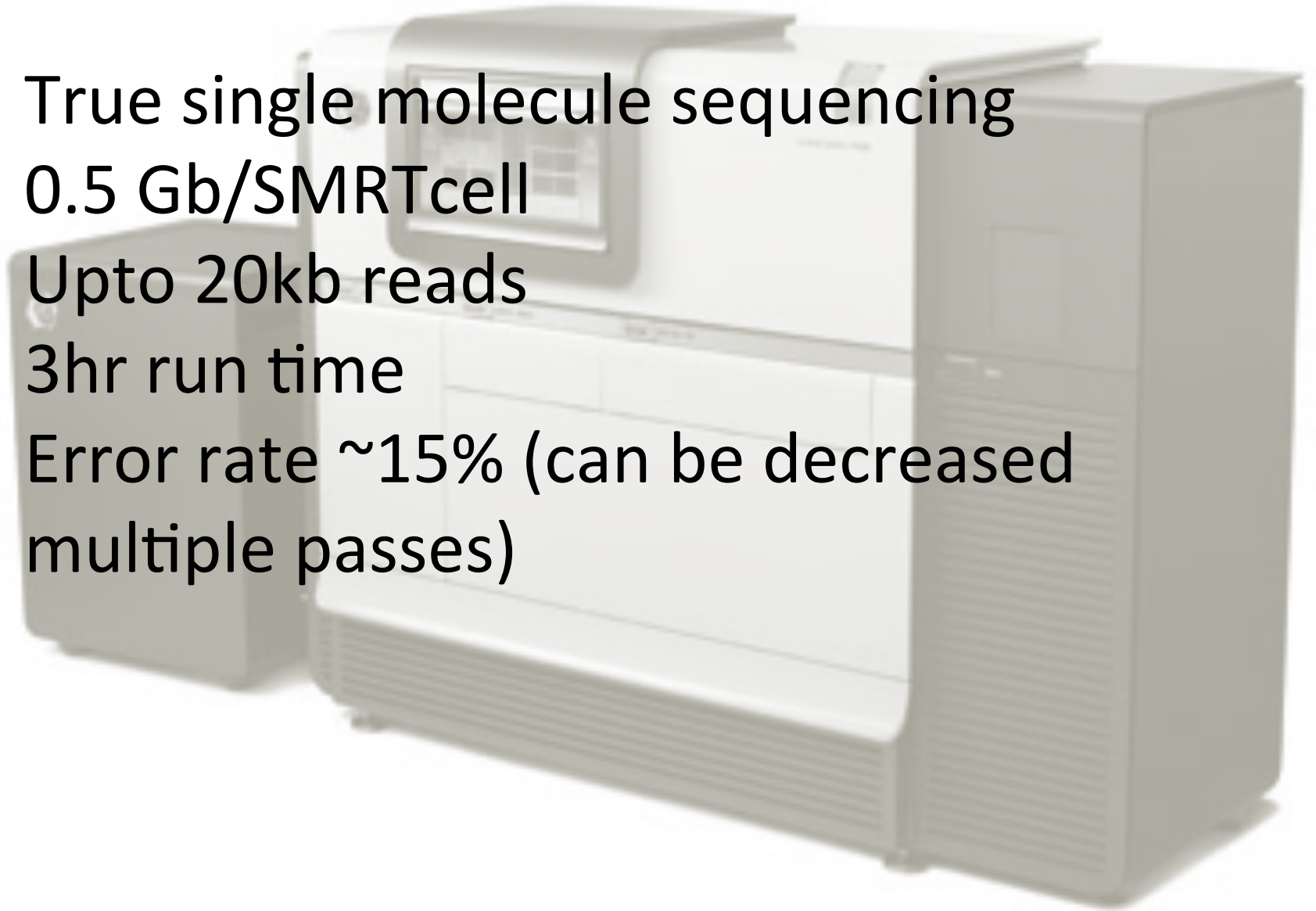
True single molecule sequencing

0.5 Gb/SMRTcell

Upto 20kb reads

3hr run time

Error rate ~15% (can be decreased
multiple passes)



[http://www.youtube.com/watch?v= B_cUZ8hSYU](http://www.youtube.com/watch?v=B_cUZ8hSYU)

<http://www.youtube.com/watch?v=v8p4ph2MAvI>

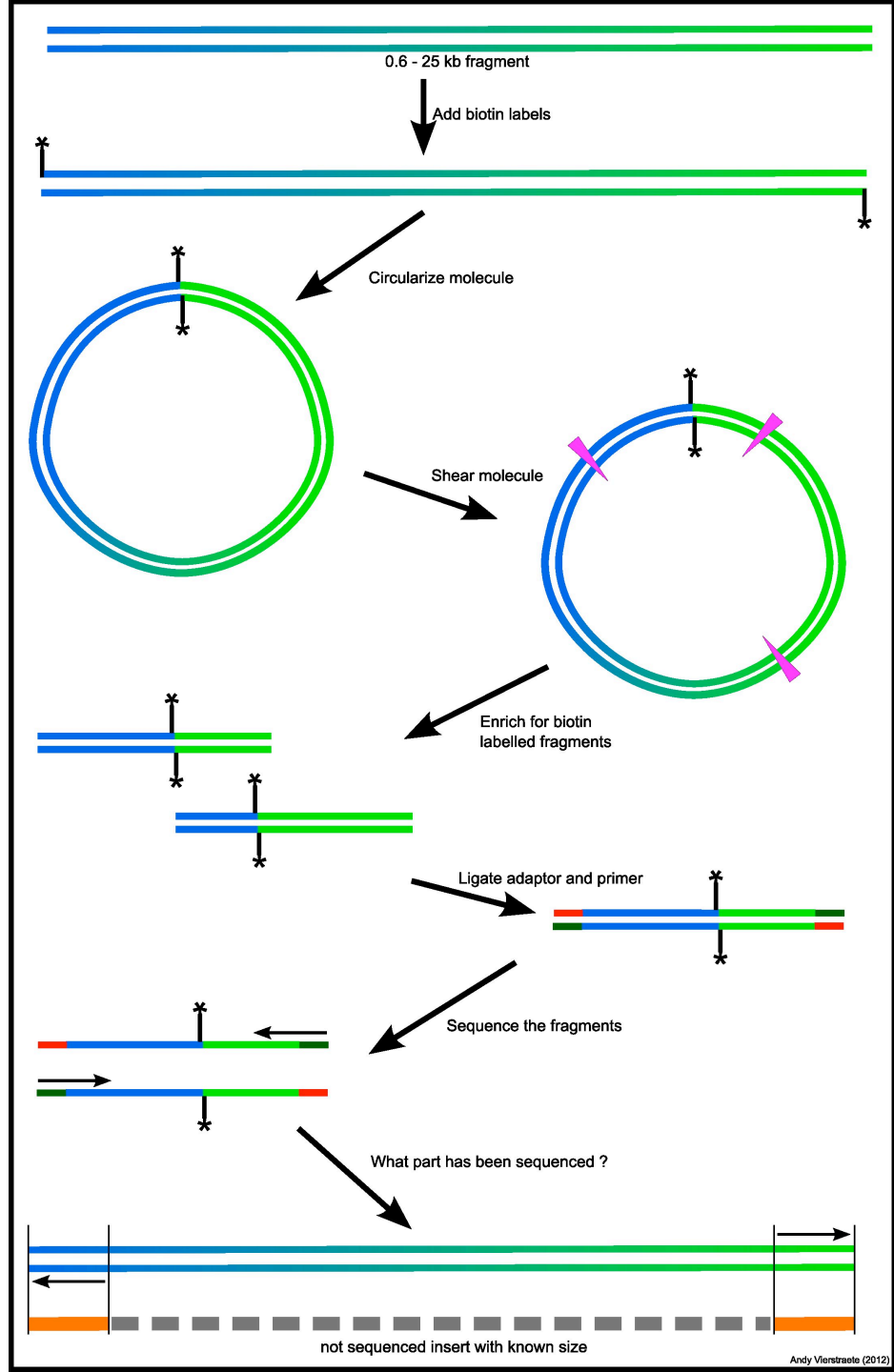
Summary

- Capillary sequencing
- NGS platforms
 - Illumina
 - Pacific Biosciences

Different types of libraries

Mate pair library

Mate pair library is very helpful in de novo sequencing; assembly genomes to sequence repeats.



a



b

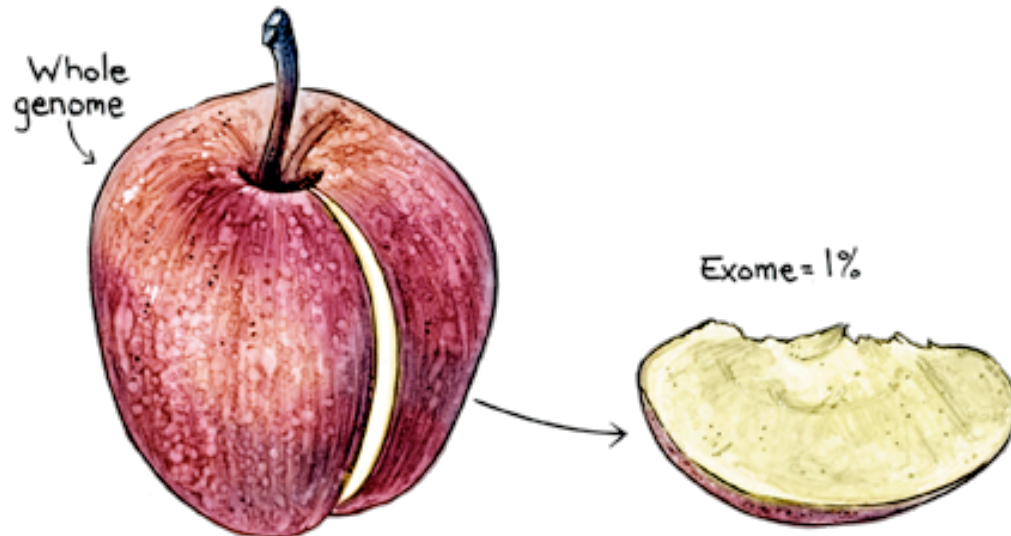
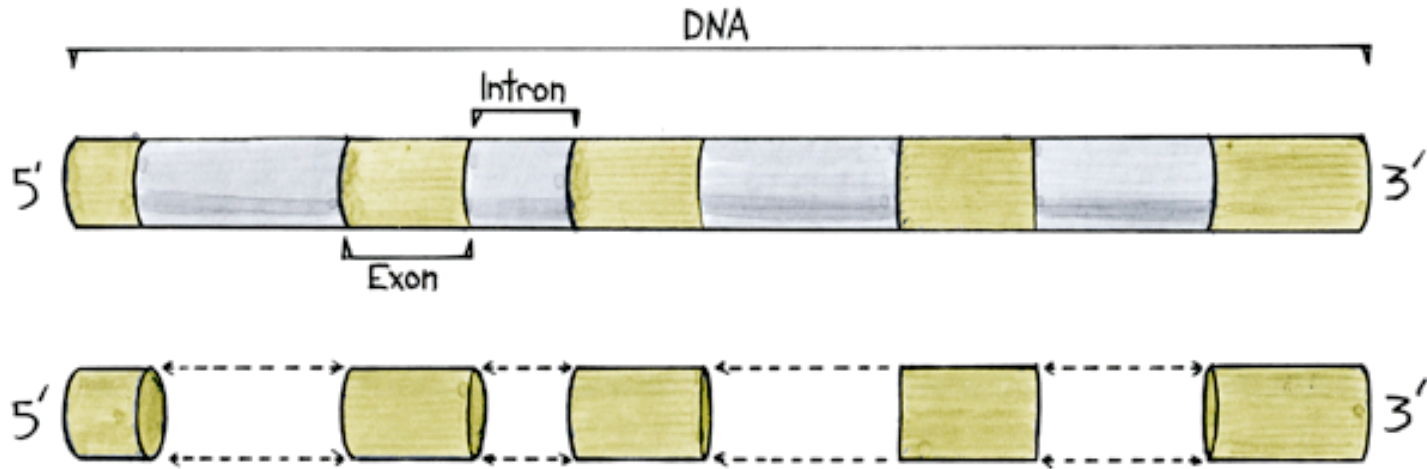
Scaffold

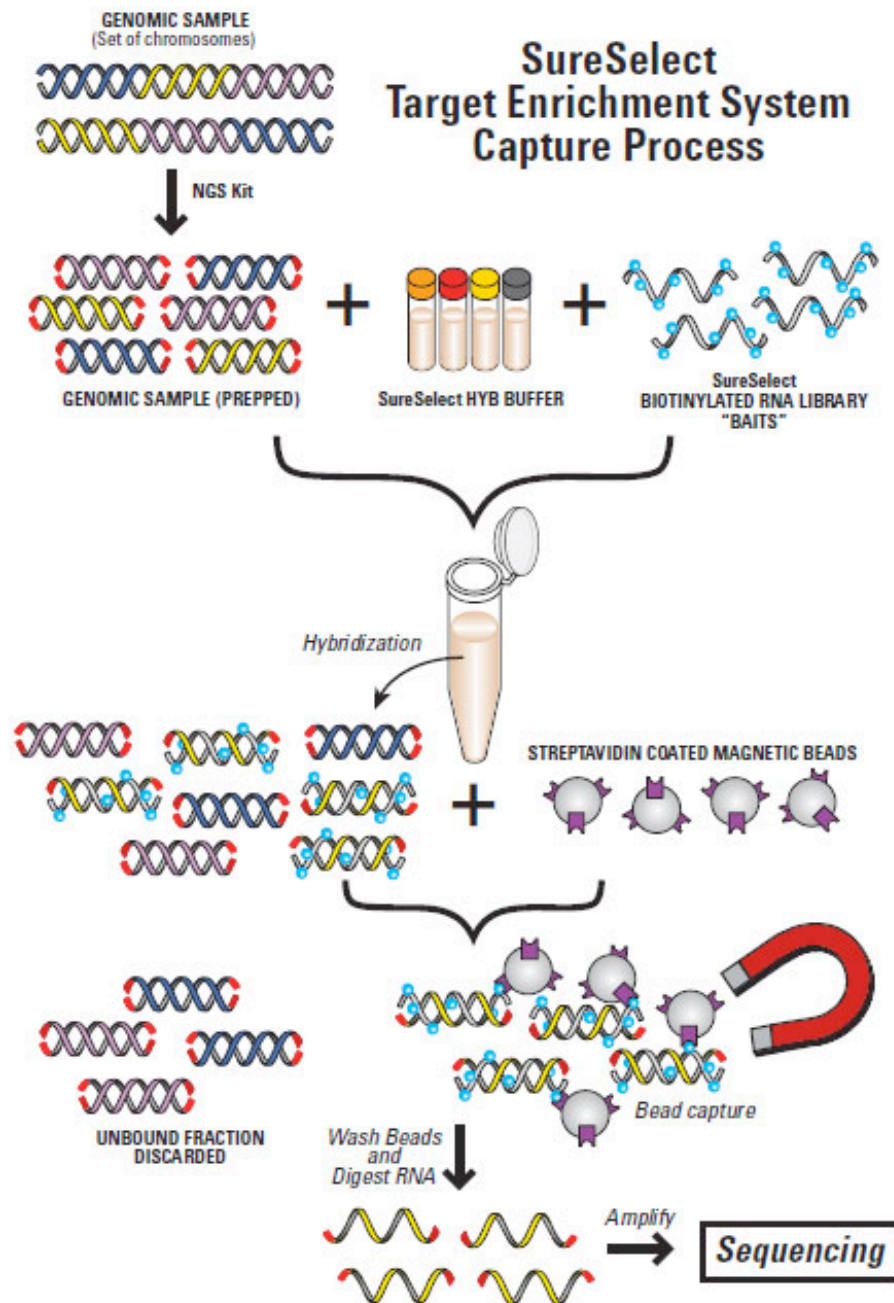


Table 2 Scaffolding value of different paired-read library combinations

No. of libraries*	Most efficient	Scaffold N50	Least efficient	Scaffold N50
1	15 kb	163,475	PE	37,694
2	5 kb + 25 kb	522,027	PE + 3 kb	46,699
2	5 kb + 20 kb	474,308	PE + 5 kb	141,403
2	8 kb + 25 kb	470,890	PE + 25 kb	142,007
3	5 kb + 20 kb + 25 kb	834,964	PE + 3 kb + 5 kb	158,525
3	5 kb + 15 kb + 25 kb	789,954	PE + 3 kb + 8 kb	171,253
3	8 kb + 20 kb + 25 kb	726,289	PE + 3 kb + 25 kb	198,696
7	ALL	1,287,609	N/A	N/A

Exome sequencing





Only sequence data from regions of interest:

- Certain chromosomes
- Gene families
- Genes involved in cancer
- Etc...

RenSeq (Resistance Enrichment Sequencing)

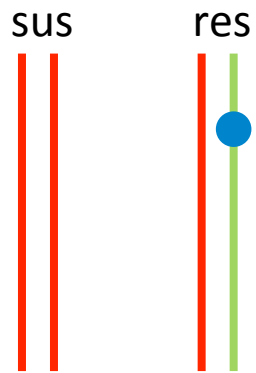
- In potato ~ 475 resistant gene present
- One or few R genes give resistance to Pi
- RenSeq = a method to quickly identify the R genes that gives resistance to *phytophthora infestans* (or any other pathogen)

resistant



susceptible

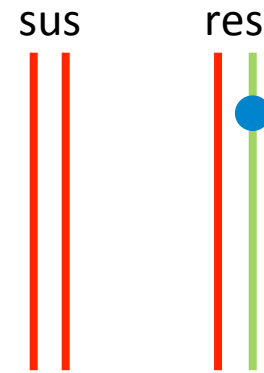




parents

F1 cross

x	r	r
R	Rr	Rr
r	rr	rr



BR

ccacgtgtagtcgtagc**ct**acgt
 tagtcgtagc**ct**acgtgatgataa
 tcgtagc**ct**acgtgatgataaata
 agtctgaccacgtgtagtcgtagc**gt**acgtgatgataaata
 cgtgtagtcgtagc**gt**acg
 acgtgtagtcgtagc**gt**ac

50%

50%

BS

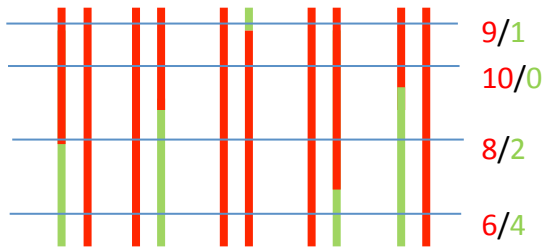
ccacgtgtagtcgtagc**gt**acgt
 tagtcgtagc**gt**acgtgatgataa
 tcgtagc**gt**acgtgatgataaata
 agtctgaccacgtgtagtcg agc**gt**acgtgatgataaata
 cgtgtagtcgtagc**gt**acg
 acgtgtagtcgtagc**gt**ac

100%



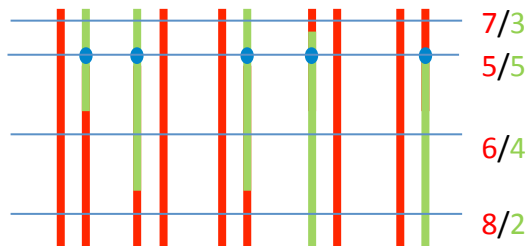
recombination

BS

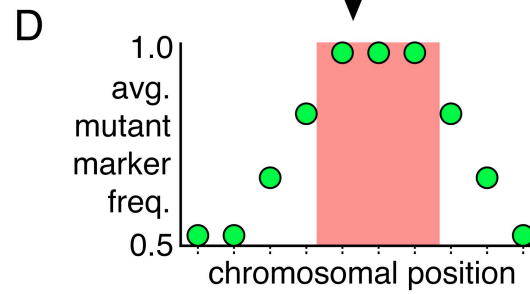
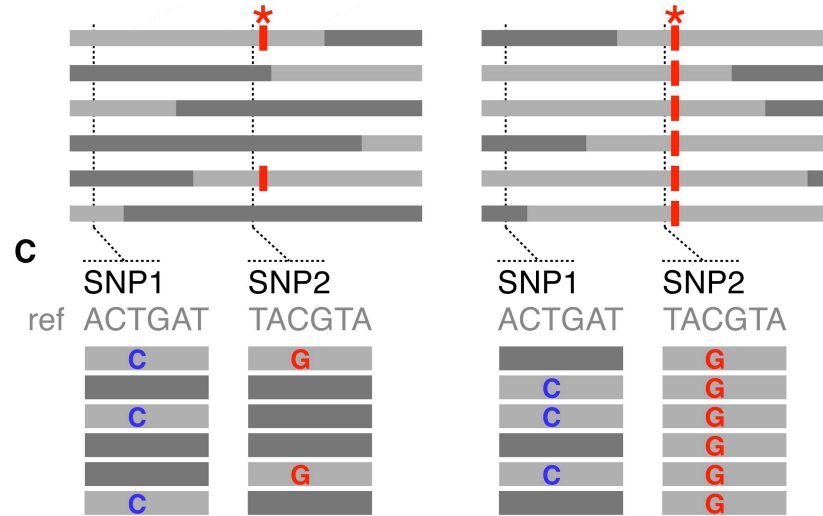


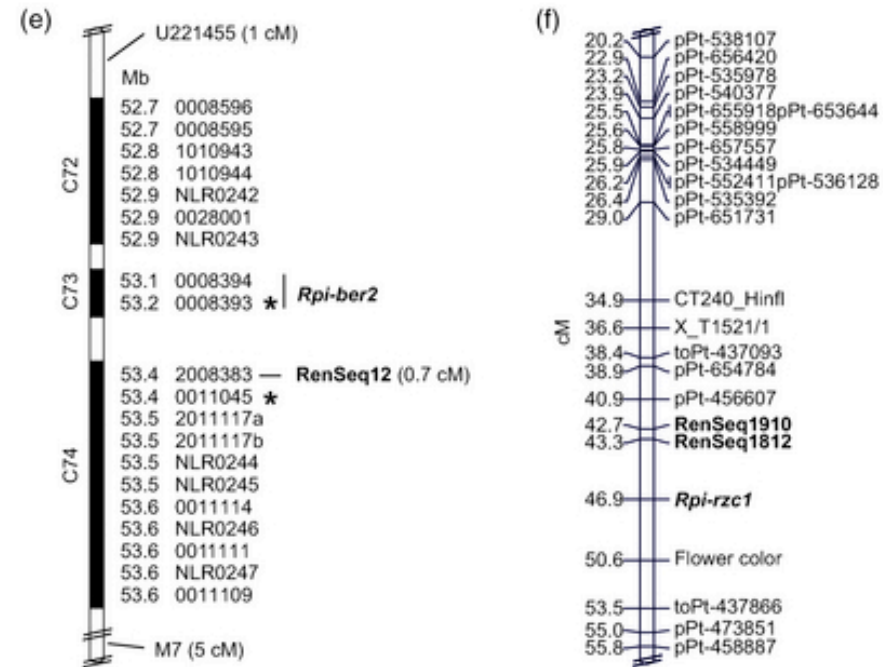
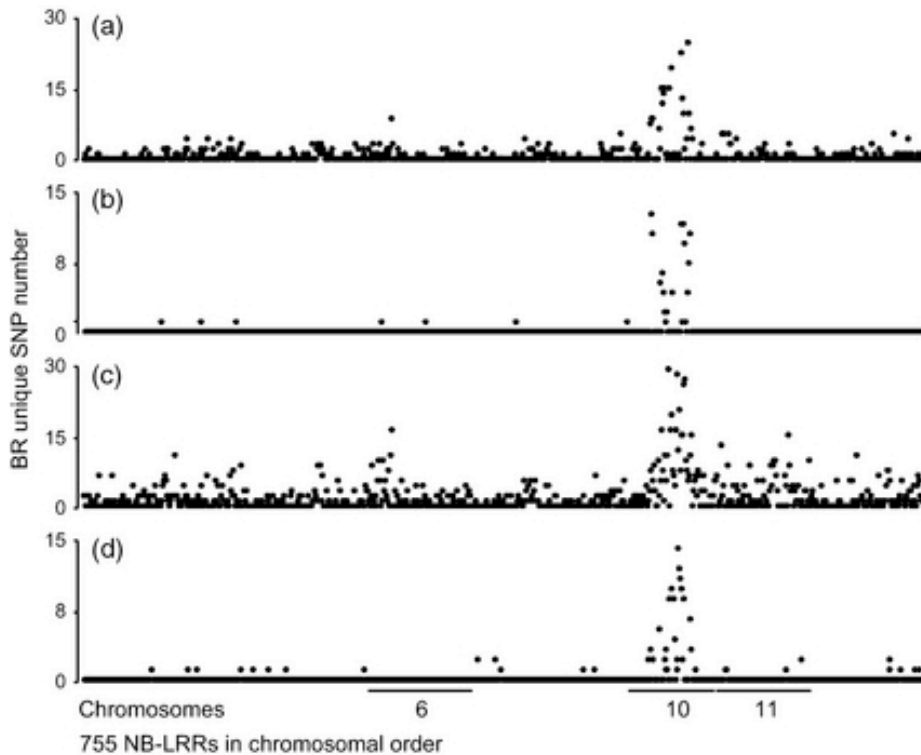
sus : res ratio = 100 : 0

BR



sus : res ratio = 50 : 50

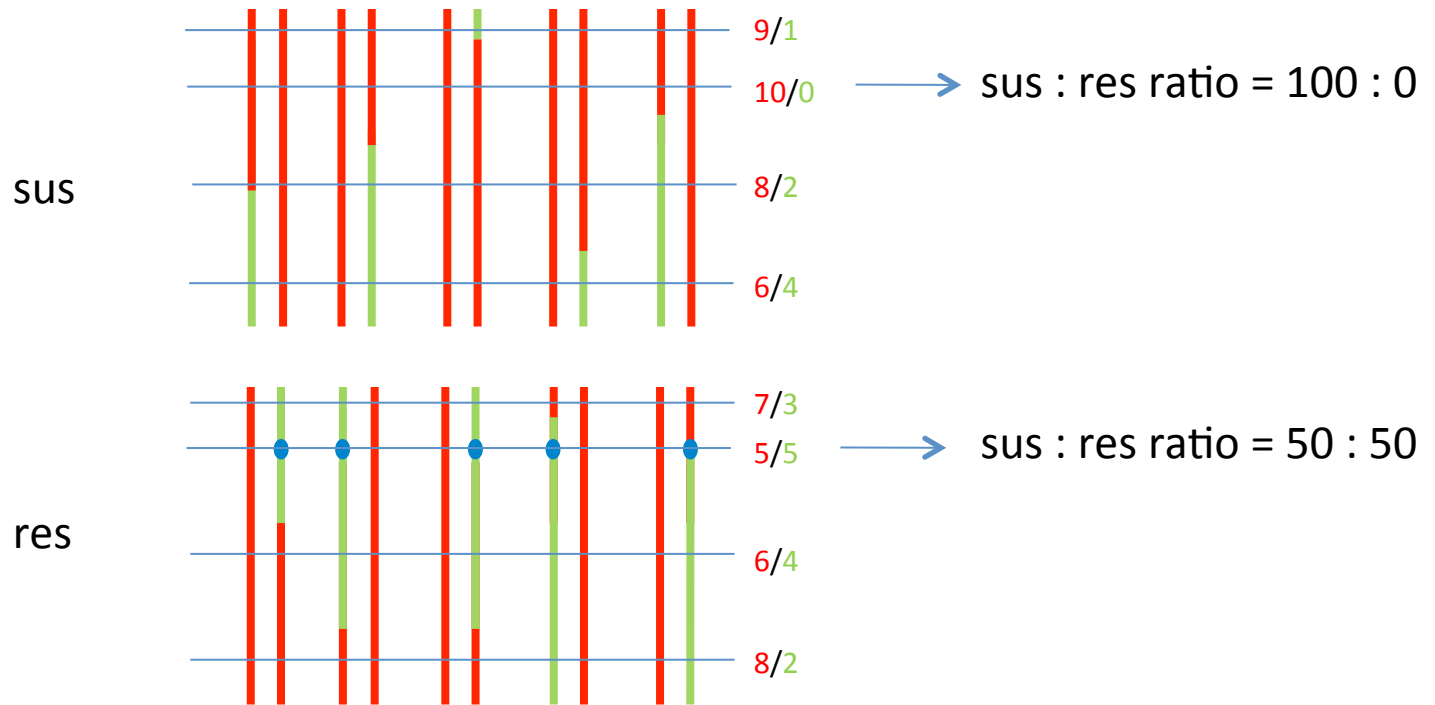




(a + c) = quick mapping; map reads to DM R genes

(b + d) = gene specific mapping; reads from BS were used for de novo assembly

recombination



Questions?