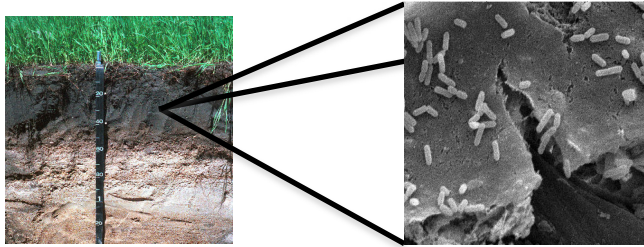
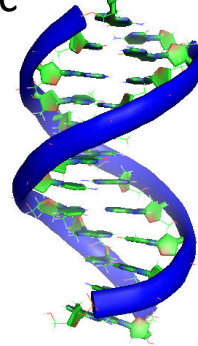




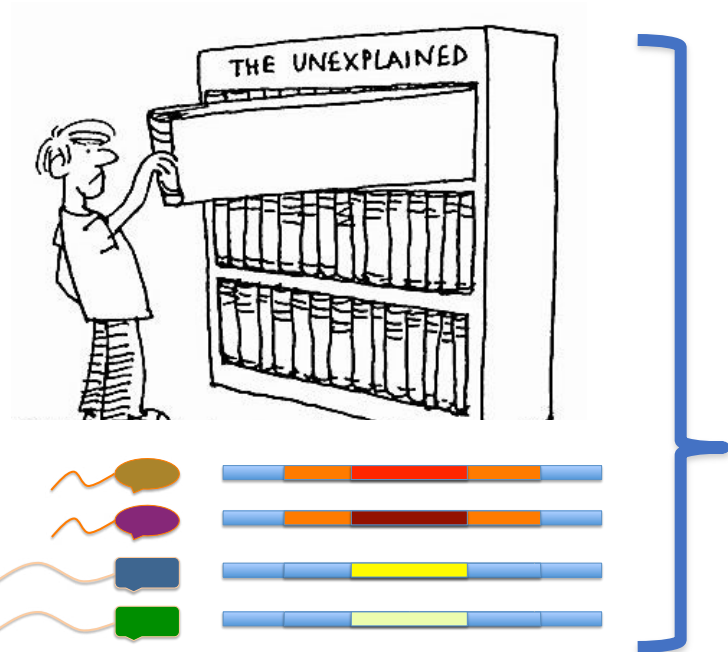
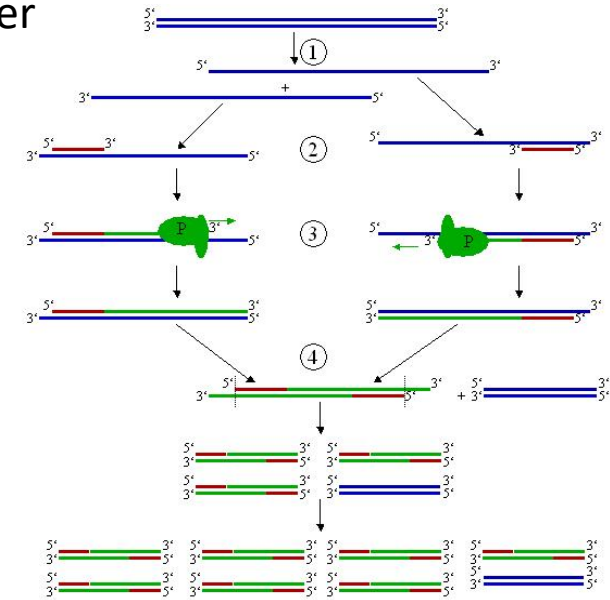
Microbial diversity with NGS-PMGA at a glance



1. Sampling and nucleic acids extraction



2. PCR amplification of phylogenetic marker



4. Data analysis



3. Sequencing

Microbial diversity with NGS-PMGA

1. Sampling and nucleic acids extraction



1. Sampling and nucleic acids extraction

- Sampling
 - Representative
 - Number of replicates (be aware of power tests etc, or consult your statistician!!! E.g. PERMANOVA requires ≥ 4 reps for obtaining P -values lower than 0.05 in pairwise cmps)
 - Controls!!!
 - Selective extraction/sampling (in case of environments with same gene carrying off-target organisms; e.g. sampling of epiphytes)
 - ...
- Nucleic acids extraction... steps:
 - Cell wall/membrane disruption
 - Nuclease inactivation denaturation
 - Removal of proteins and PCR inhibitors

Microbial diversity with NGS-PMGA

2a. Phylogenetic marker choice



2a. Phylogenetic marker choice

- Has to exist in the target organisms:
 - 16S rRNA gene: prokaryotes
 - ITS: fungi
 - 18S rRNA gene: protists, mycorrhizae, eukaryotes
 - amoA: ammonia oxidizers
 - ...
- Suitable conserved sequences (exist in all target organisms) that can facilitate primer designing around hypervariable regions (differ between target organisms like barcodes)
- PCR products compatible with sequencing technologies

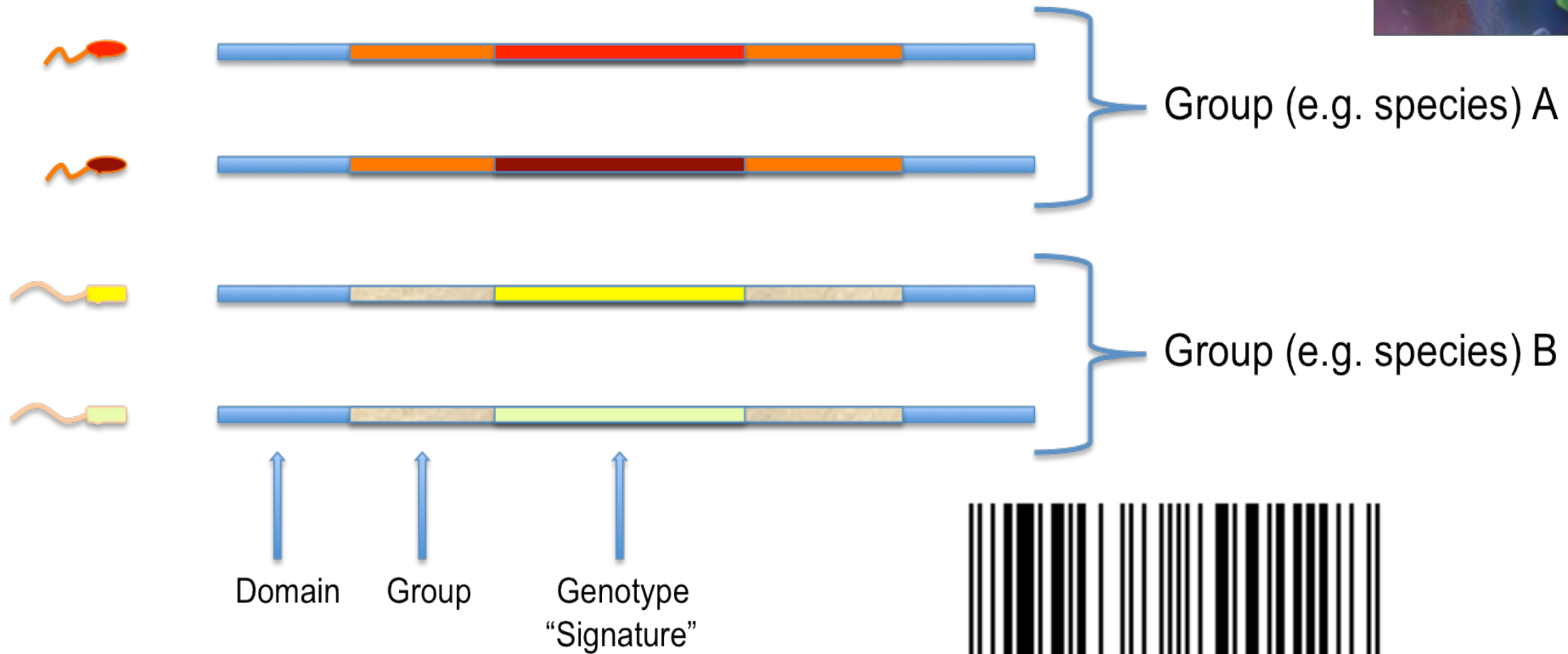


2a. E.g. 16S rRNA gene



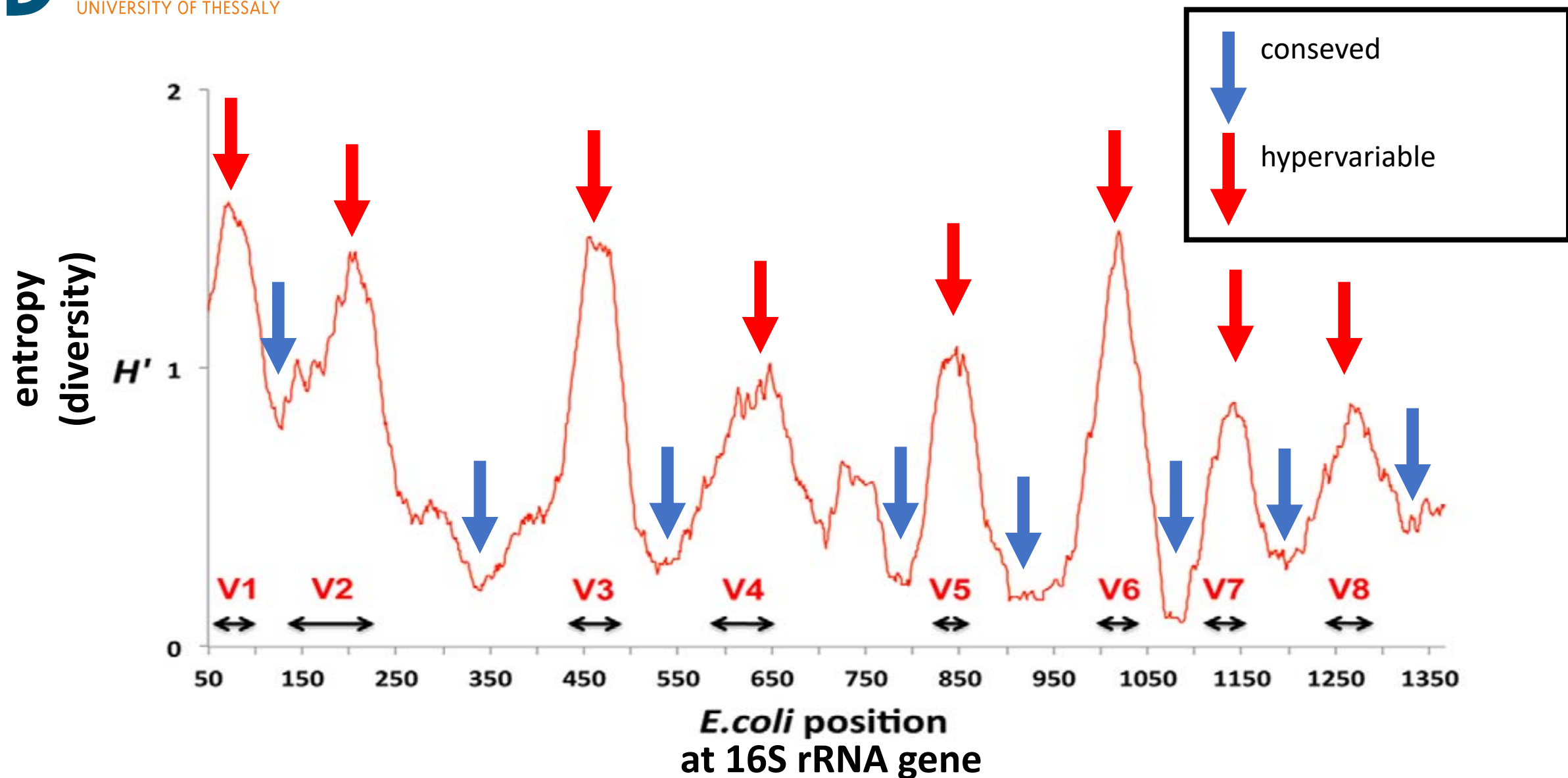
Microorganisms

Bacterial 16S rRNA strand fragment



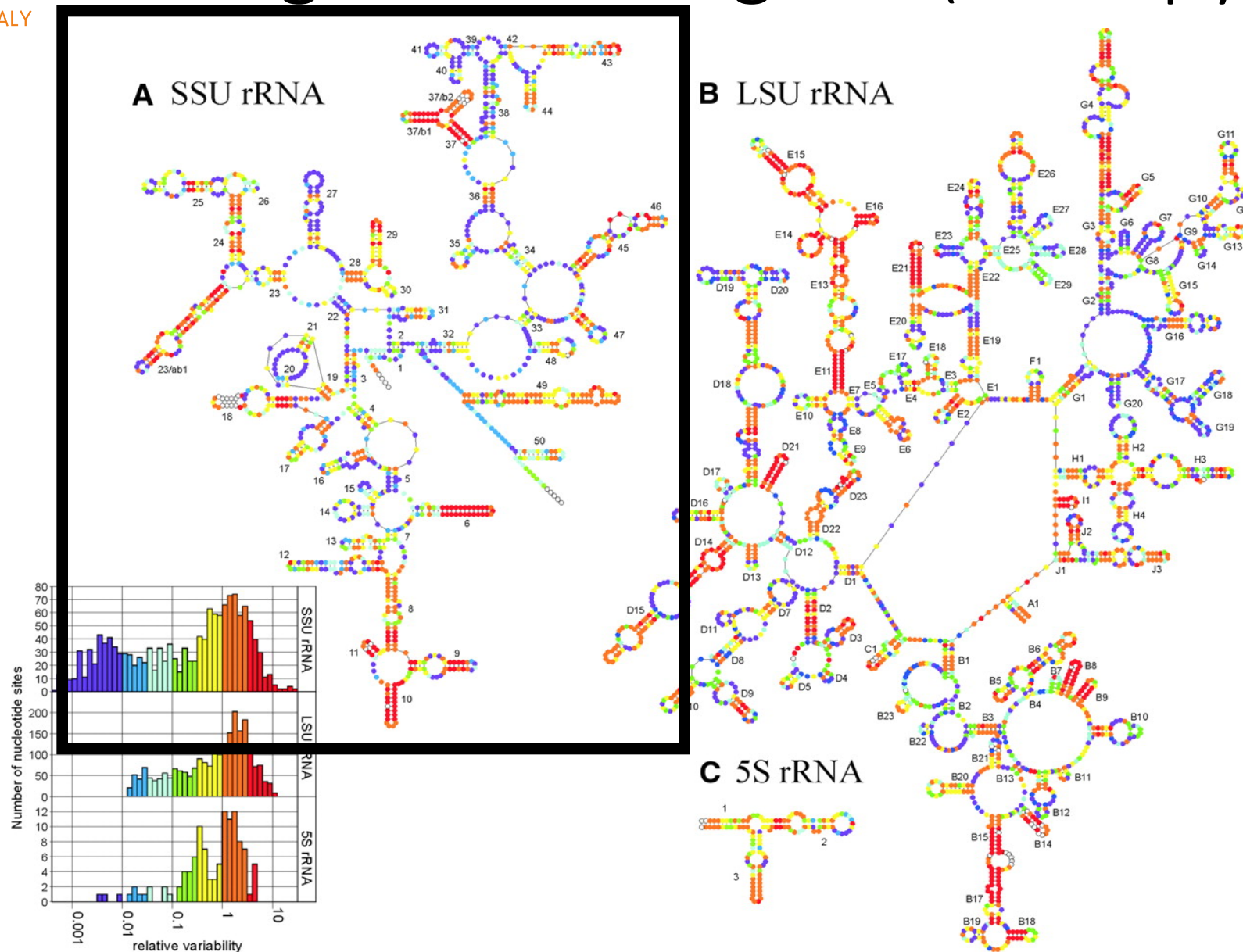


2a. E.g. 16S rRNA gene



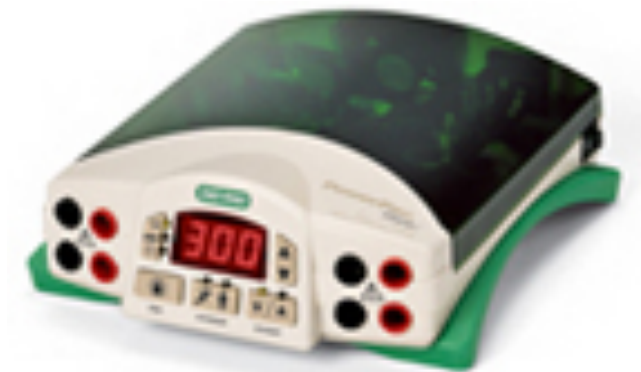
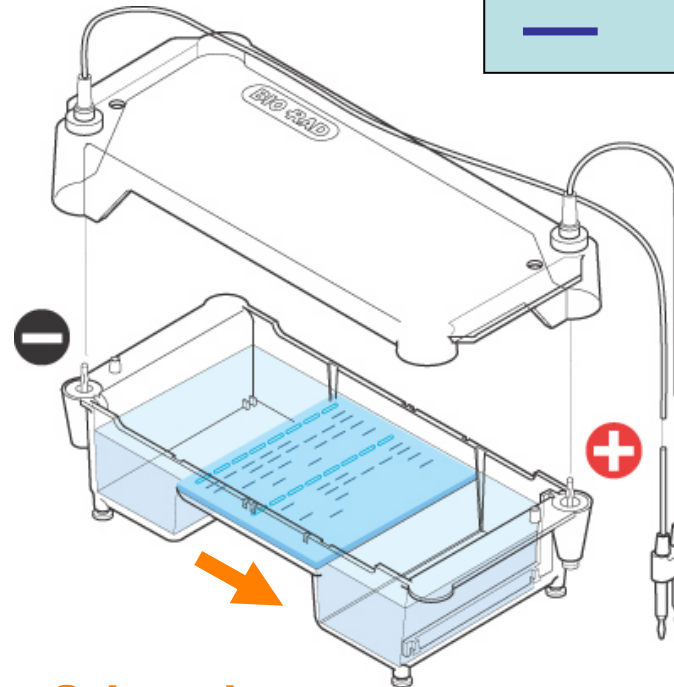
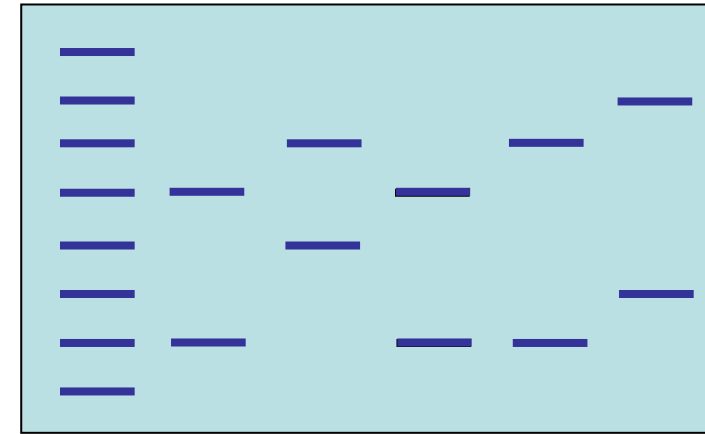


2a. E.g. 16S rRNA gene (entropy maps)



2b. PCR errors & prevention/correction

PCR Demo





2b. PCR biases & prevention/correction

- Prevention/correction:
 - ❖ mispriming -> stringency ↑: temperature/hybridization conditions & *
 - ❖ chimeras -> minimization of PCR cycles, high fidelity polymerases, relieving agents & *
 - ❖ amplification of homologous non-target genes (e.g. for 16S rRNA gene, mitochondrial or chloroplast homologue) -> selective extraction and selective environment sampling & *

*overall sequencing depth outweighs the errors (the higher the depth the more the sequences we can spare)



2b. PCR biases (mispriming)

Dealing with mispriming in data analysis:

- Contrasting each amplicon sequence with curated databases.



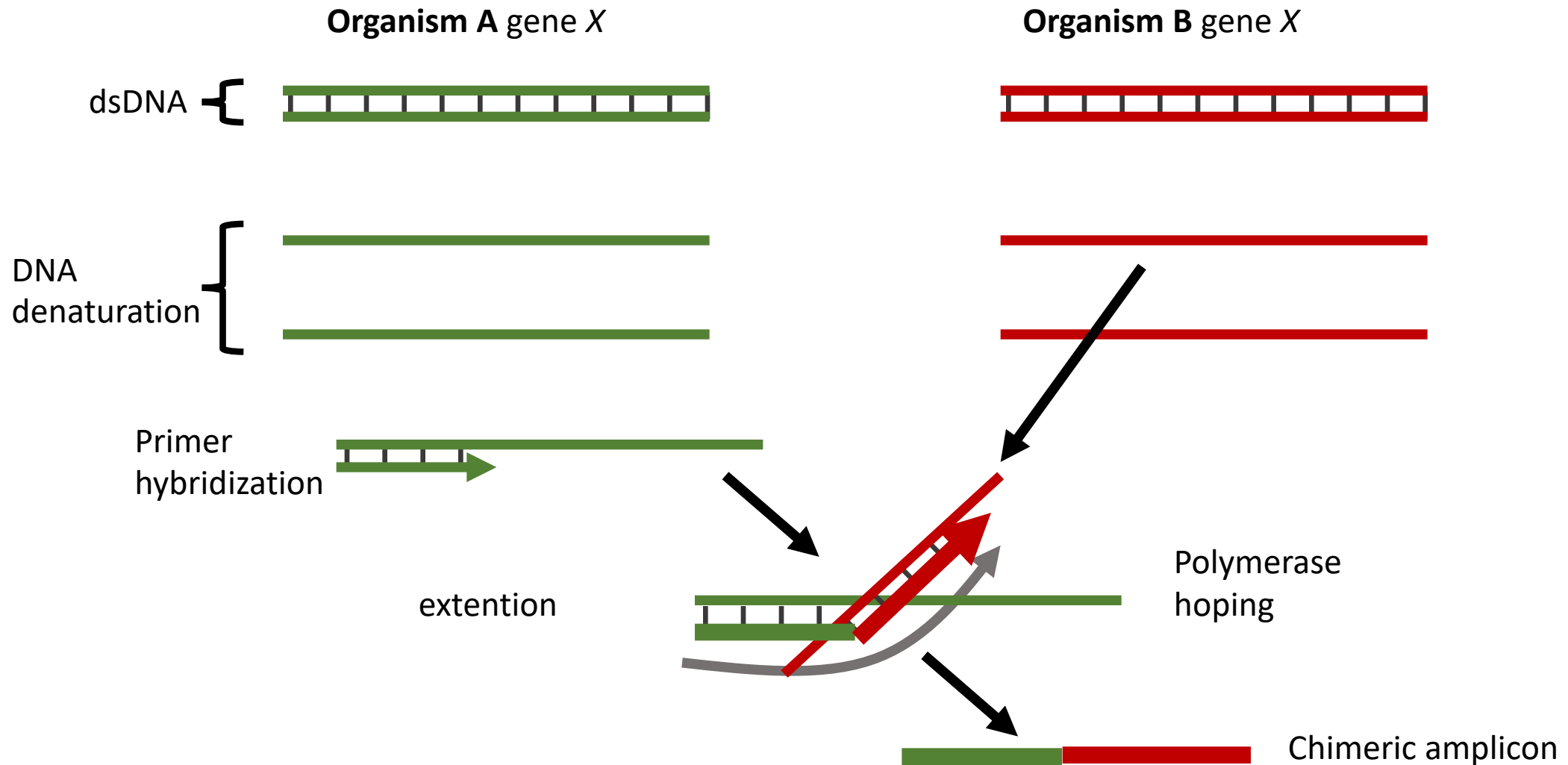
2b. PCR errors & prevention/correction

- Prevention/correction:

- ❖ mispriming -> temperature/hybridization conditions & *
- ❖ chimeras -> minimization of PCR cycles, high fidelity polymerases, relieving agents & *
- ❖ amplification of homologous non-target genes (e.g. for 16S rRNA gene, mitochondrial or chloroplast homologue) -> selective extraction and selective environment sampling & *

*overall sequencing depth outweighs the errors (the higher the depth the more the sequences we can spare)

2b. PCR errors (Chimeras)





2b. PCR errors & prevention/correction

- Prevention/correction:

- ❖ mispriming -> temperature/hybridization conditions & *
- ❖ chimeras -> minimization of PCR cycles, high fidelity polymerases, relieving agents & *
- ❖ amplification of homologous non-target genes (e.g. for 16S rRNA gene, mitochondrial or chloroplast homologue) -> selective extraction and selective environment sampling & *

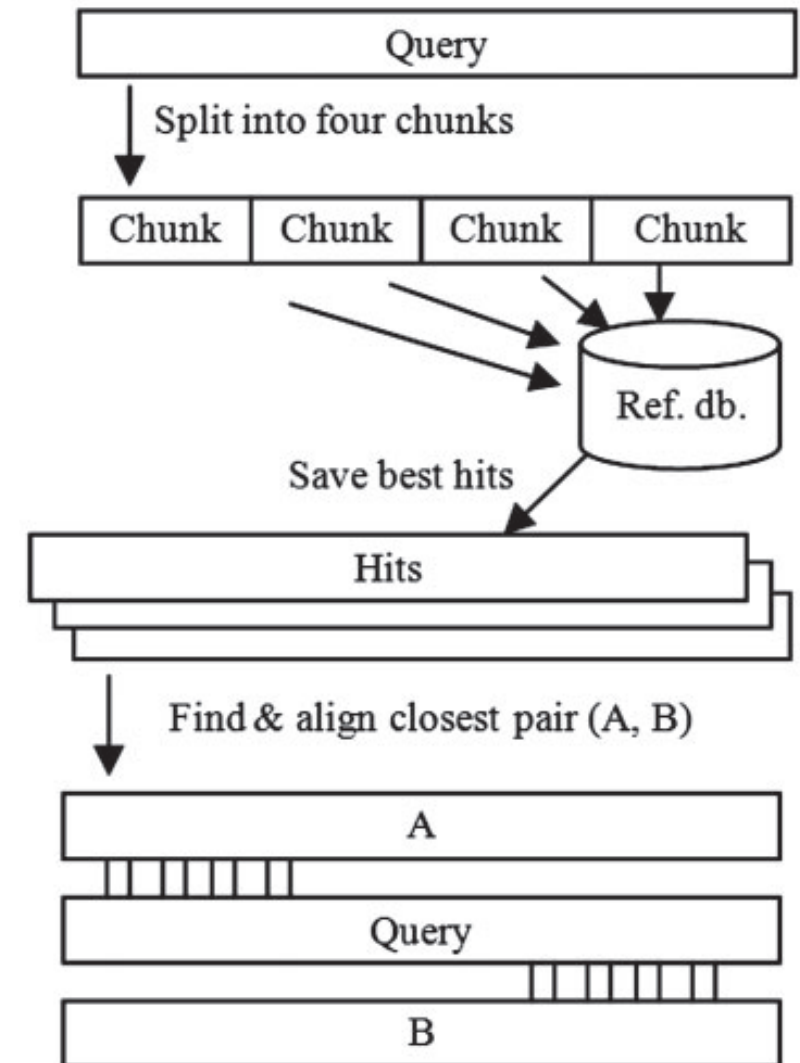
*overall sequencing depth outweighs the errors (the higher the depth the more the sequences we can spare)



2b. PCR errors (chimeras)

Dealing with chimeras in data analysis:

- Contrasting amplicon sequence portions with:
 - curated database
 - *de novo* built database using the exp. data (the most abundant sequences are the correct)
- A sequence containing portions of distant others in the database is chimeric





2b. PCR errors & prevention/correction

- Prevention/correction:

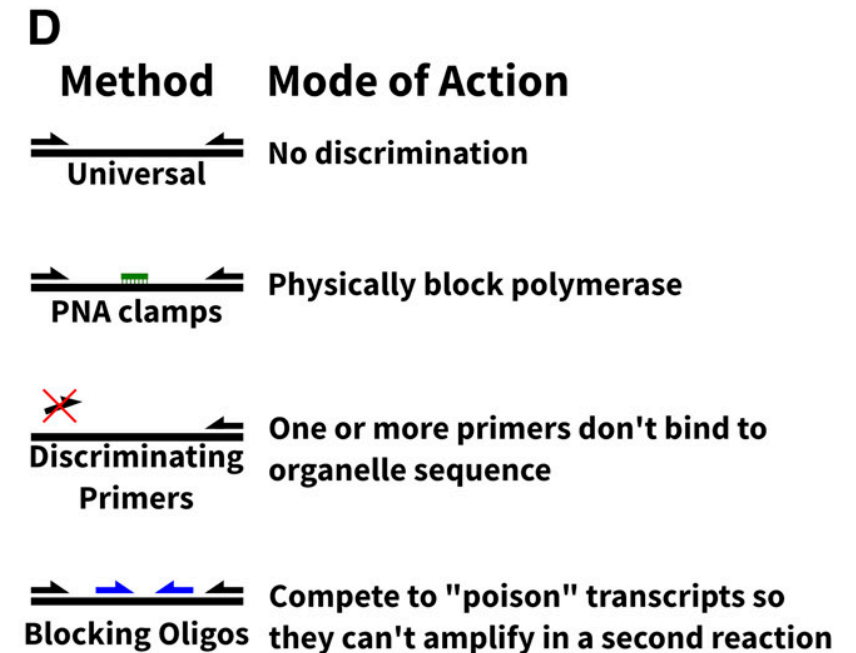
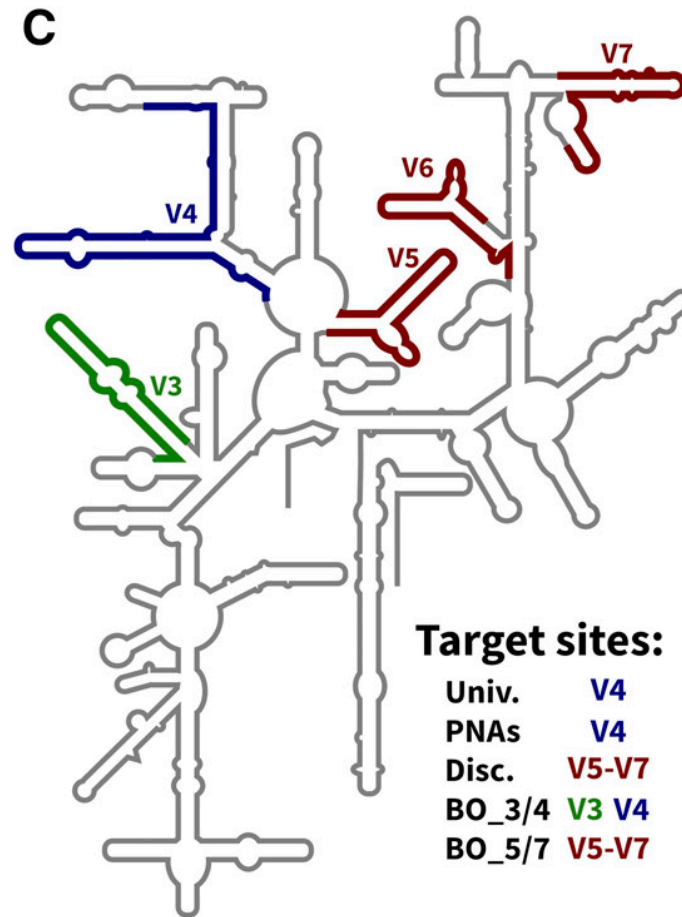
- ❖ mispriming -> temperature/hybridization conditions & *
- ❖ chimeras -> minimization of PCR cycles, high fidelity polymerases, relieving agents & *
- ❖ amplification of homologous non-target genes (e.g. for 16S rRNA gene, mitochondrial or chloroplast homologue) -> selective extraction and selective environment sampling & *

*overall sequencing depth outweighs the errors (the higher the depth the more the sequences we can spare)



2b. PCR biases (same gene, off-target org.)

- Strategies for endophyte community analysis:
 - ❖ Be selective in the bioinformatics
 - ❖ Use peptide nucleic acid (PNA) or lock nucleic acid (LNA) blockers in the lab
 - ❖ Use selective primers in the lab
 - ❖ Use blocking oligos in the lab



Microbial diversity with NGS-PMGA

3. Sequencing

The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

- 1st generation sequencing (single DNA fragment per reaction, ~700-1000 bp long)
 - Dideoxy termination method (**Sanger**)
- 2nd generation sequencing (massively parallel, max of ~300-600 bp total)
 - Sequencing by ligation (e.g. ABI-SOLiD)
 - Sequencing by synthesis (e.g. *Pyrosequencing*, **Illumina**, *Ion-torrent*)
- 3rd generation sequencing (massively parallel and very long reads, max of ~10kbp-250kbp)
 - Sequencing by synthesis (e.g. single molecule real time – **PacBio, Nanopore**)
 - Synthetic long reads (e.g. Moleculo, 10X)

The 3 generations of sequencing technologies

All three generations have been
proposed for PMGA sequencing



1st generation: Dideoxy termination (Sanger) method



Ingredients:

1) dsDNA



1st generation: Dideoxy termination (Sanger) method

...CAATACGTAACCTTCCCTTGCTAACTTCAGTCAGCATGGAAGCCCATTAGTCGGAAAGC...

...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...

Ingredients:

1) dsDNA

Process:

1) Denature dsDNA (95°C)



1st generation: Dideoxy termination (Sanger) method

...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...

Ingredients:

1) dsDNA

Process:

1) Denature dsDNA (95°C)



1st generation: Dideoxy termination (Sanger) method

5' CGTAACTTTCCCTT
| | | | | | | | | |
...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...

Ingredients:

- 1) dsDNA
- 2) Primer

Process:

- 1) Denature dsDNA (95°C)
- 2) Primer hybridization (e.g. 50 °C)



1st generation: Dideoxy termination (Sanger) method



Ingredients:

- 1) dsDNA
- 2) Primer
- 3) Polymerase

Process:

- 1) Denature dsDNA (95°C)
- 2) Primer hybridization (e.g. 50 °C)



1st generation: Dideoxy termination (Sanger) method

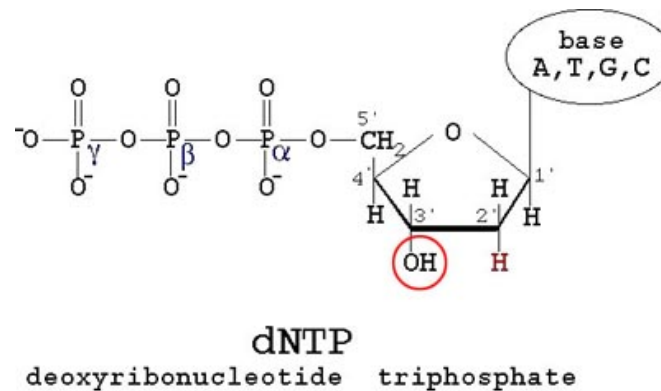


Ingredients:

- 1) dsDNA
- 2) Primer
- 3) Polymerase
- 4) dNTPs

Process:

- 1) Denature dsDNA (95°C)
- 2) Primer hybridization (e.g. 50 °C)
- 3) Extension/elongation (e.g. 72 °C)





1st generation: Dideoxy termination (Sanger) method

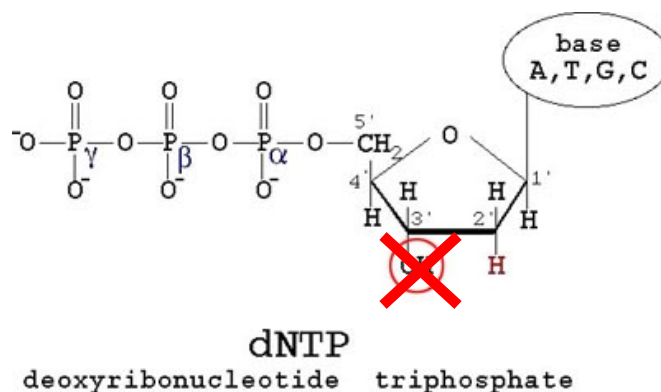


Ingredients:

- 1) dsDNA
- 2) Primer
- 3) Polymerase
- 4) dNTPs
- 5) ddNTPs (e.g. ddATP)

Process:

- 1) Denature dsDNA (95°C)
- 2) Primer hybridization (e.g. 50 °C)
- 3) Extension/elongation (e.g. 72 °C)
- 4) Termination



dNTP
deoxyribonucleotide triphosphate



1st generation: Dideoxy termination (Sanger) method

5' CGTAACTTTCCCTTTGCTAACTTC^A

...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...

Ingredients:

- 1) dsDNA
- 2) Primer
- 3) Polymerase
- 4) dNTPs
- 5) ddNTPs (e.g. ddATP)

Process:

- 1) Denature dsDNA (95°C)
- 2) Primer hybridization (e.g. 50 °C)
- 3) Extension/elongation (e.g. 72 °C)
- 4) Termination
- 5) Denature dsDNA (95°C)



1st generation: Dideoxy termination (Sanger) method

5' CGTAACTTCCCTTGGCTAACTTC^A

...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...

5' CGTAACTTCCCTTGGCTA^A

...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...

5' CGTAACTTCCCTTGGCTAACTTCAGTCAGCATGGA^A

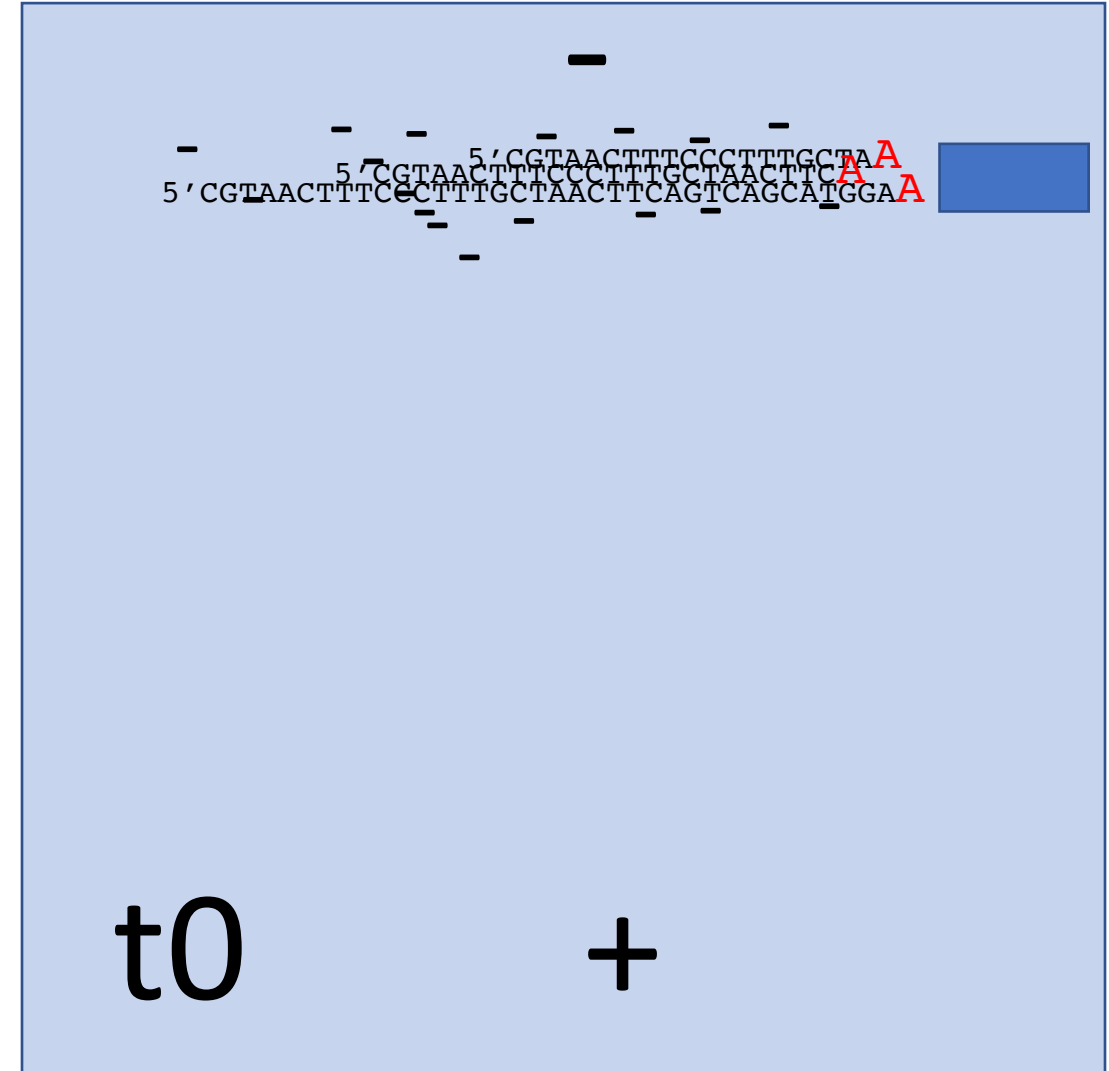
...GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGC...



1st generation: Dideoxy termination (Sanger) method

Diagnostics process:

- 1) Add fragments to a matrix that separates them according to size (e.g. polyacrylamide gel)
- 2) Apply current (negatively charged DNA moves toward positive electrode)
- 3) Identify terminator base according to location on gel

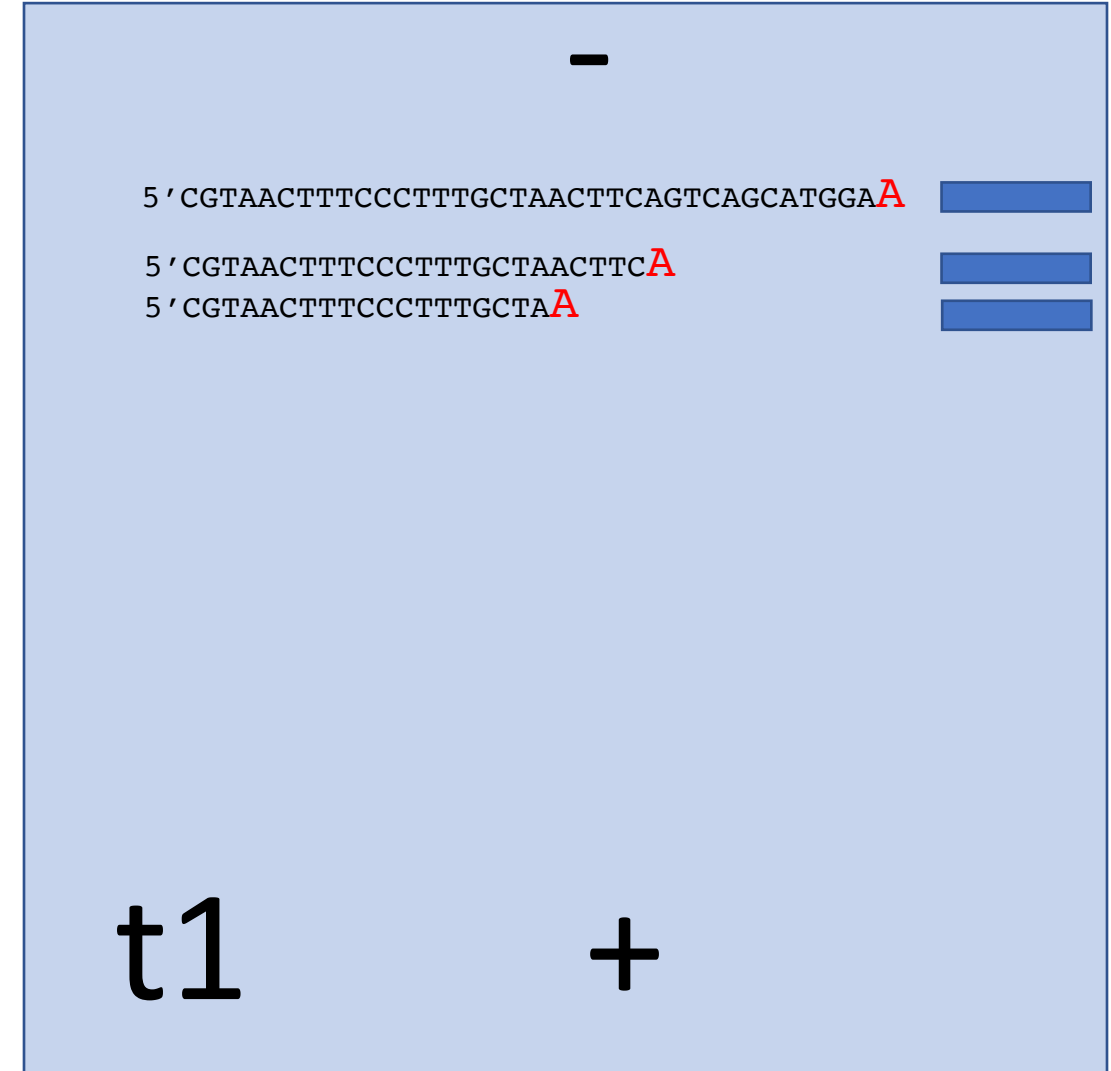




1st generation: Dideoxy termination (Sanger) method

Diagnostics process:

- 1) Add fragments to a matrix that separates them according to size (e.g. polyacrylamide gel)
- 2) Apply current (negatively charged DNA moves toward positive electrode)
- 3) Identify terminator base according to location on gel

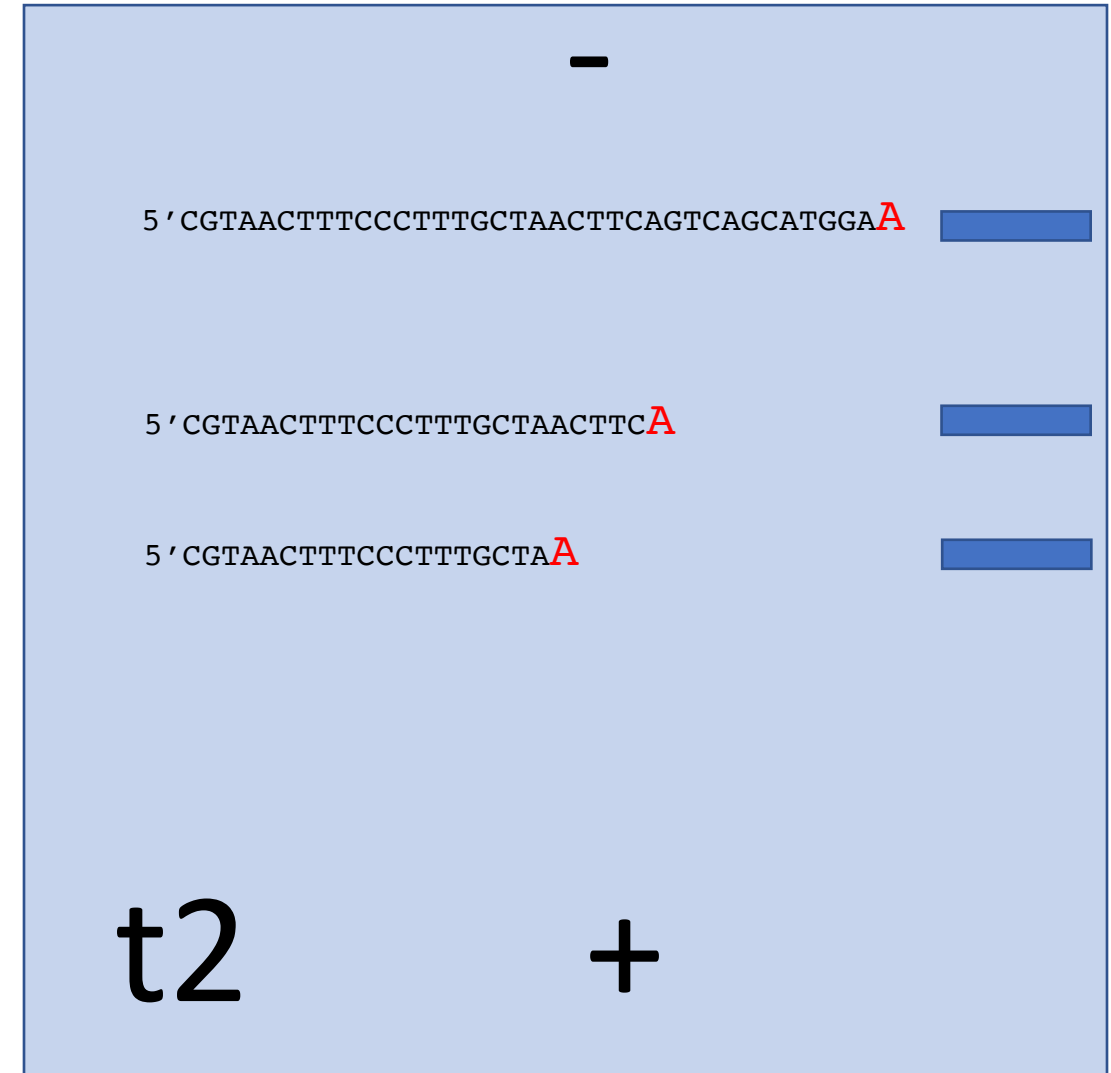




1st generation: Dideoxy termination (Sanger) method

Diagnostics process:

- 1) Add fragments to a matrix that separates them according to size (e.g. polyacrylamide gel)
- 2) Apply current (negatively charged DNA moves toward positive electrode)
- 3) Identify terminator base according to location on gel

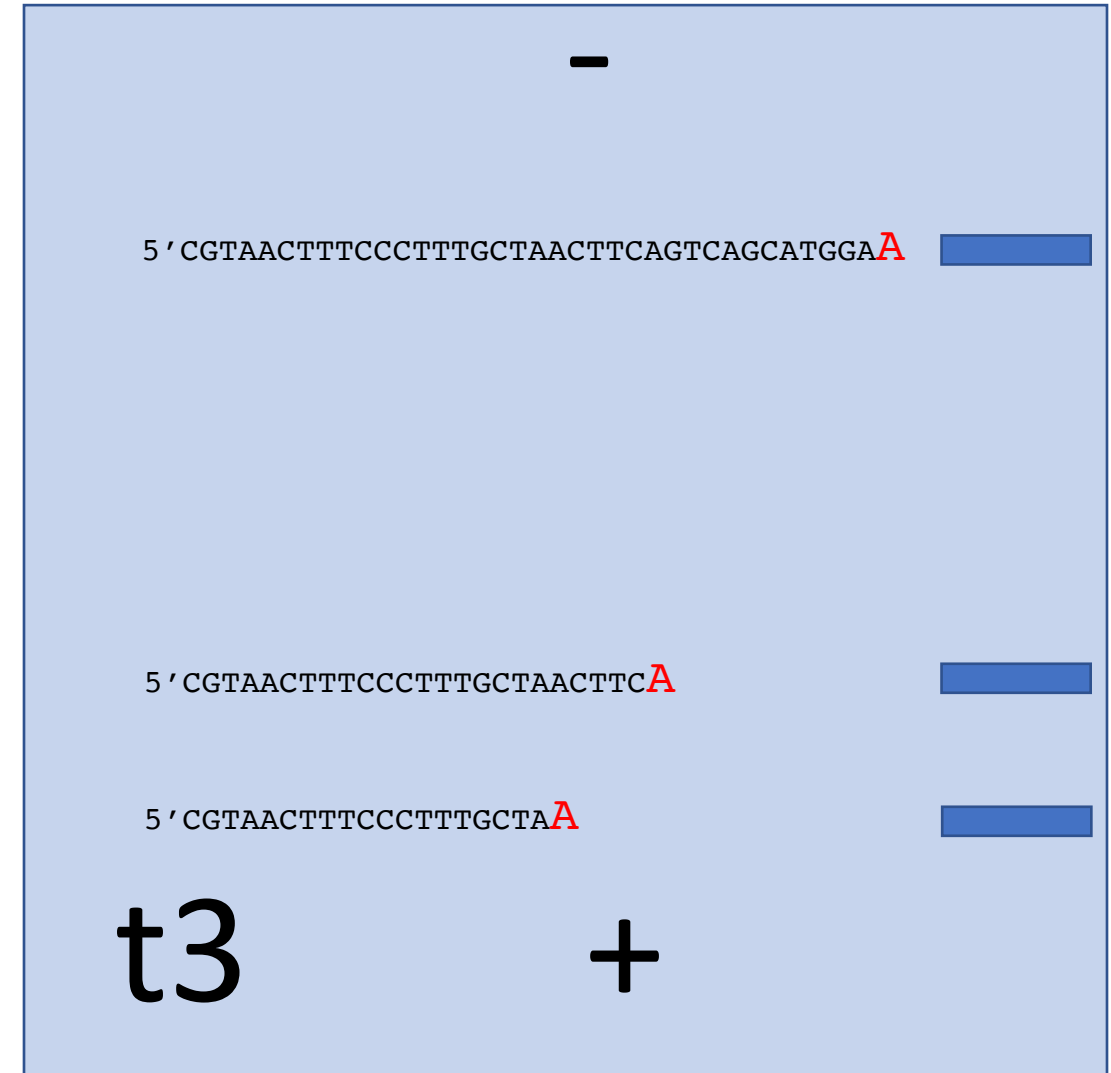




1st generation: Dideoxy termination (Sanger) method

Diagnostics process:

- 1) Add fragments to a matrix that separates them according to size (e.g. polyacrylamide gel)
- 2) Apply current (negatively charged DNA moves toward positive electrode)
- 3) Identify terminator base according to location on gel





1st generation: Dideoxy termination (Sanger) method

Diagnostics process:

- 1) Add fragments to a matrix that separates them according to size (e.g. polyacrylamide gel)
- 2) Apply current (negatively charged DNA moves toward positive electrode)
- 3) Identify terminator base according to location on gel

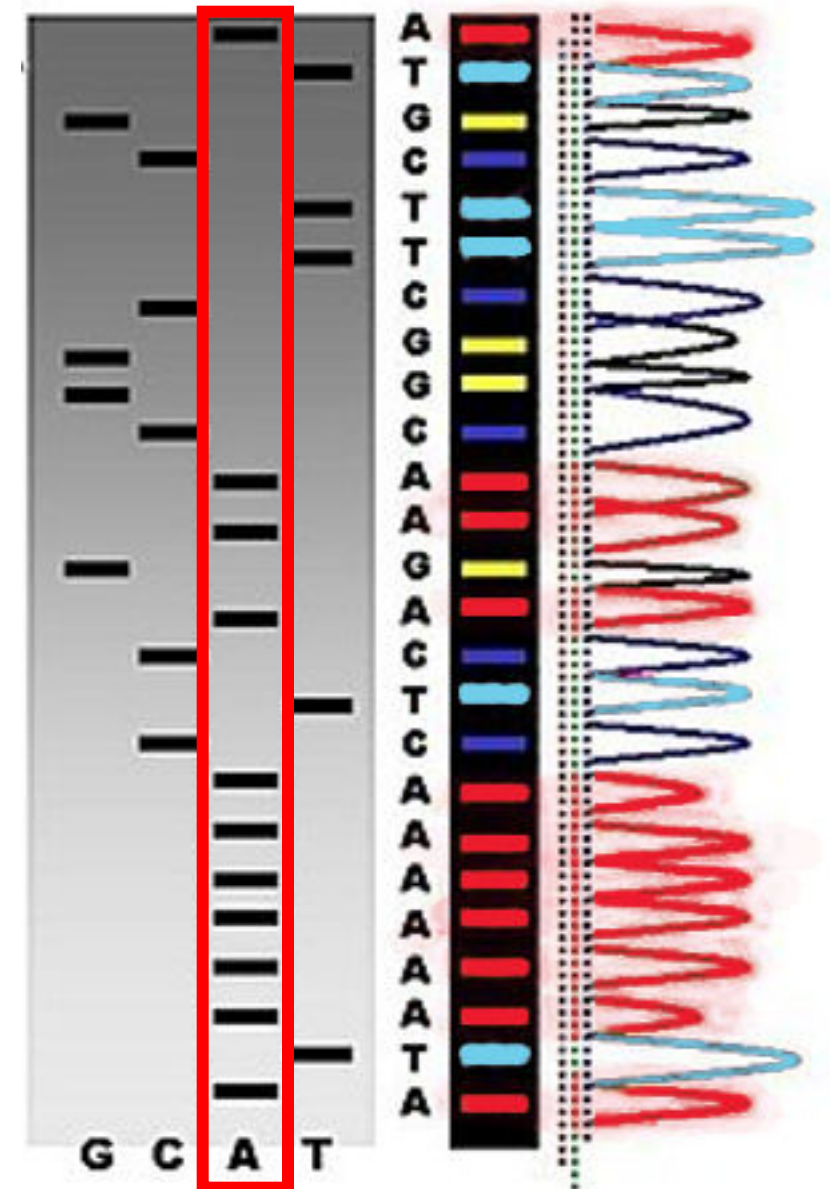




1st generation: Dideoxy termination (Sanger) method

Diagnostics process:

- 1) Add fragments to a matrix that separates them according to size (e.g. polyacrylamide gel)
- 2) Apply current (negatively charged DNA moves toward positive electrode)
- 3) Identify terminator base according to location on gel

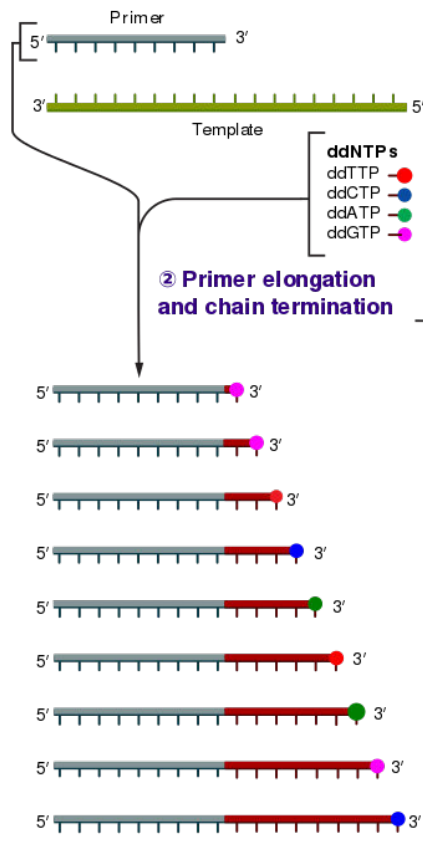




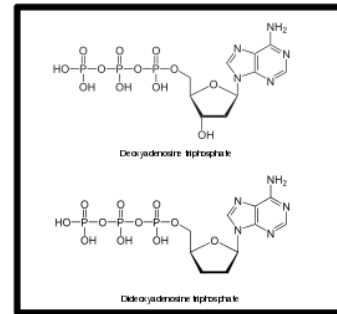
1st generation: Dideoxy termination (Sanger) method (Summary slide)

① Reaction mixture

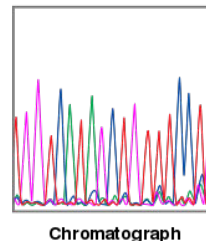
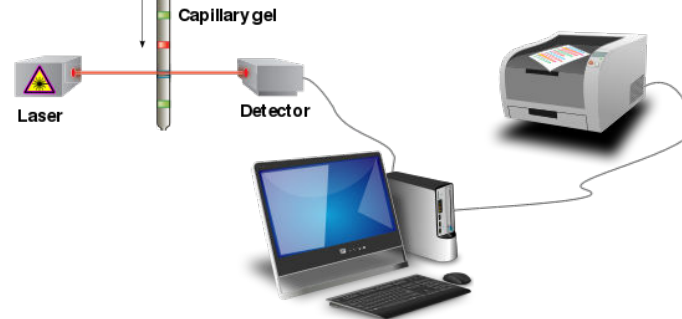
- ▶ Primer and DNA template
- ▶ DNA polymerase
- ▶ ddNTPs with flourochromes ▶ dNTPs (dATP, dCTP, dGTP, and dTTP)



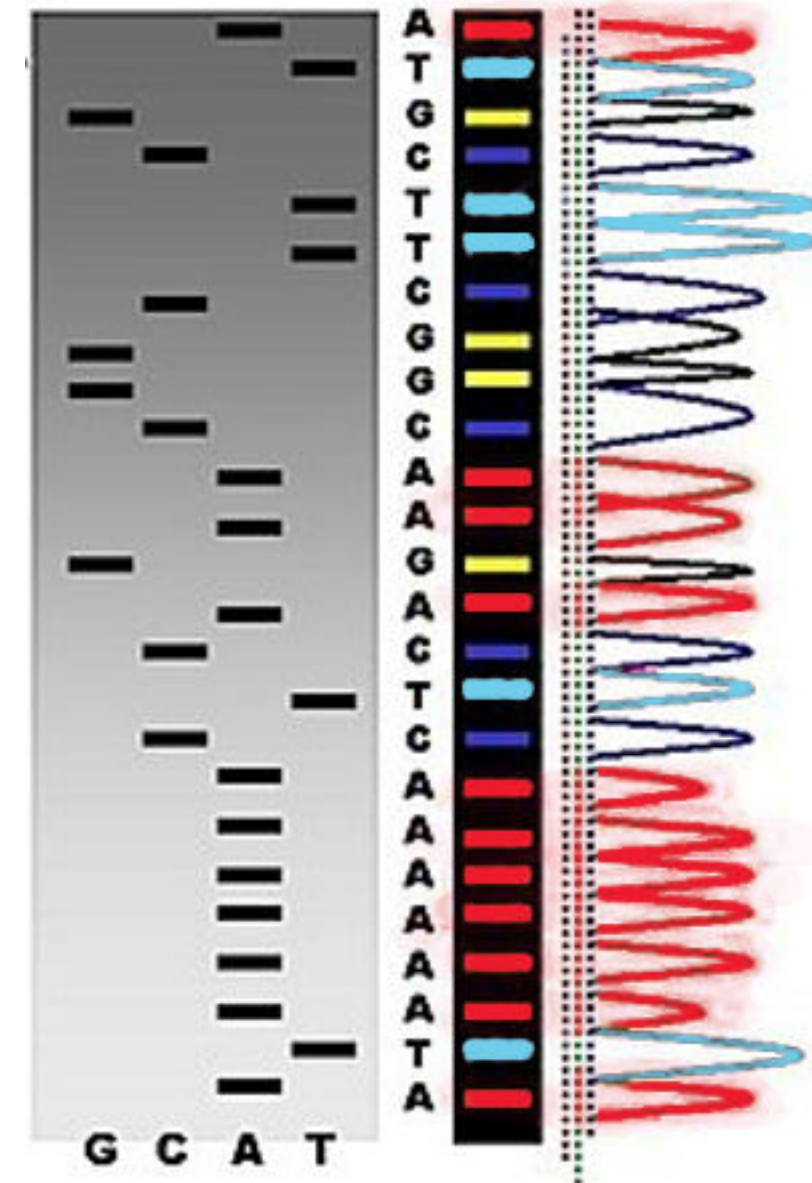
② Primer elongation and chain termination



③ Capillary gel electrophoresis separation of DNA fragments

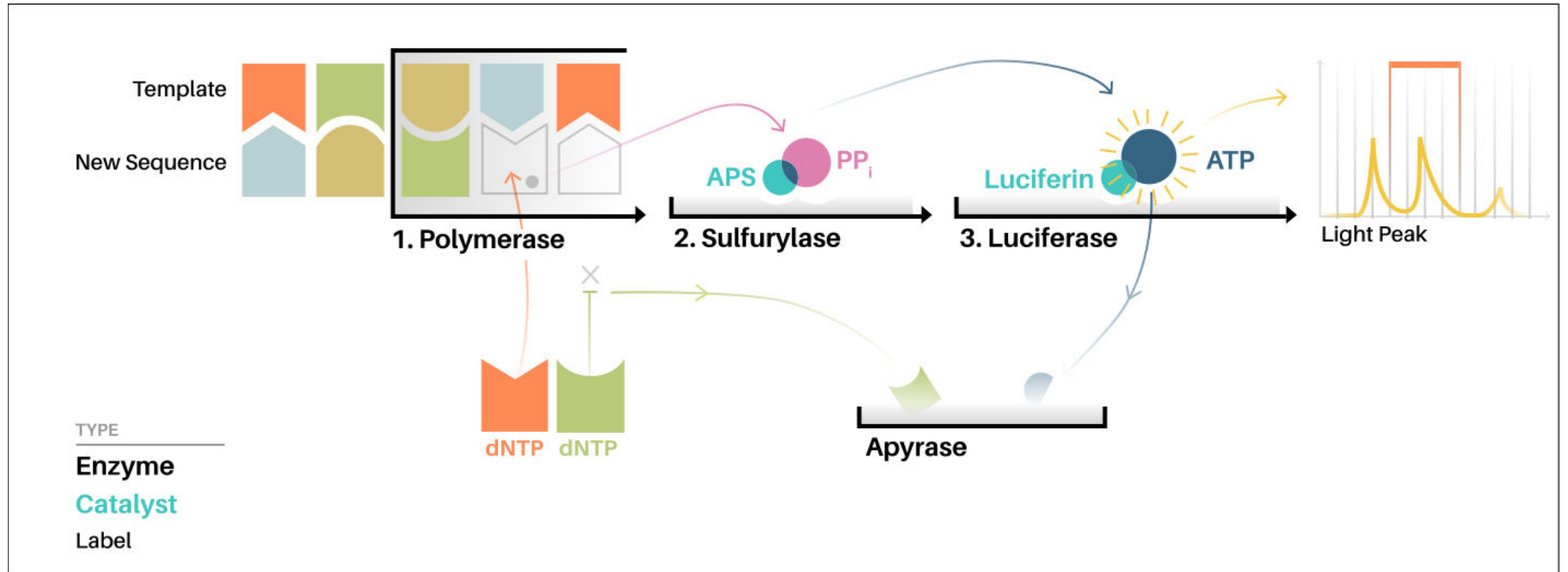


④ Laser detection of flourochromes and computational sequence analysis





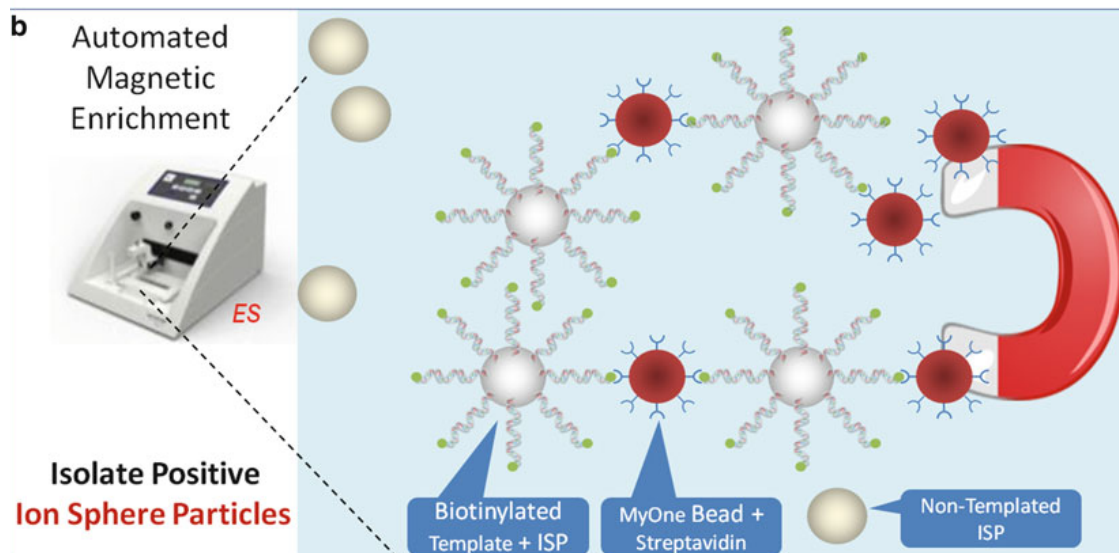
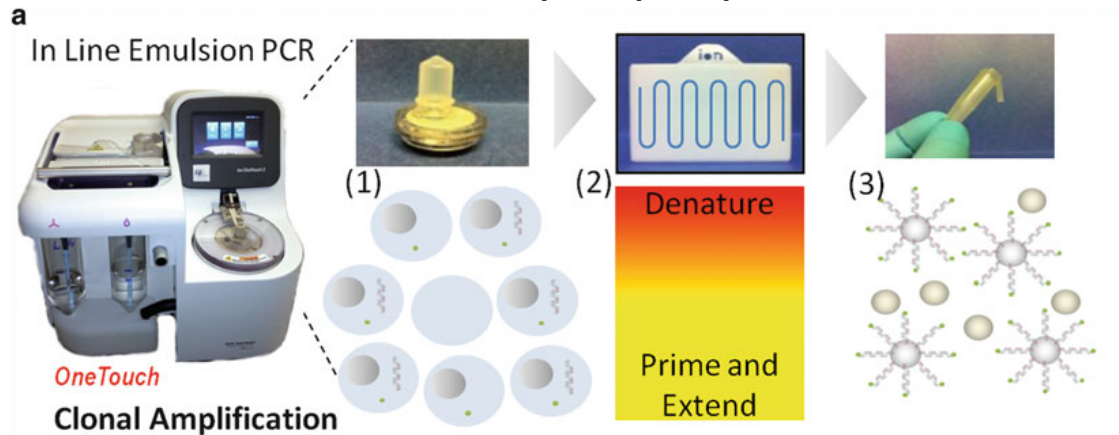
3. 2nd generation Pyrosequencing





3. 2nd generation Ion Torrent/Proton

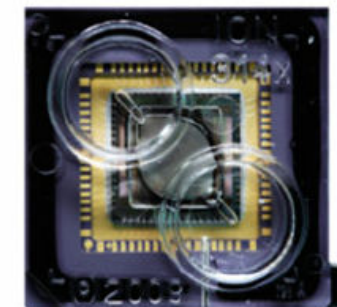
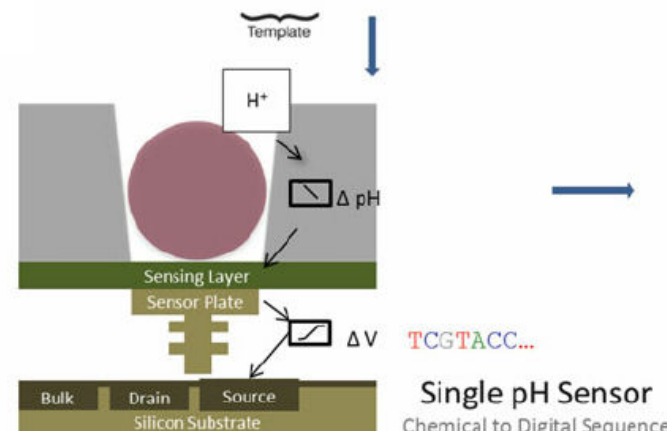
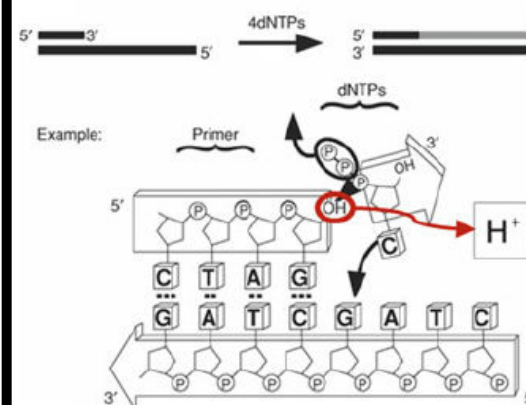
Sample prep



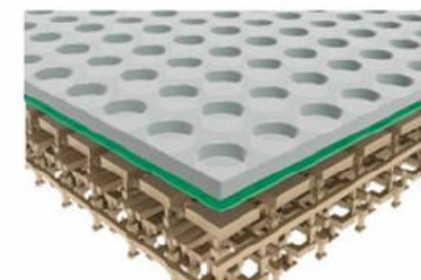
Sequencing

Principle and Elements of Semiconductor Sequencing

Simple Natural Chemistry of Sequencing-by-Synthesis with H^+ release detection



Sequencing Chip
Semiconductor Packaging

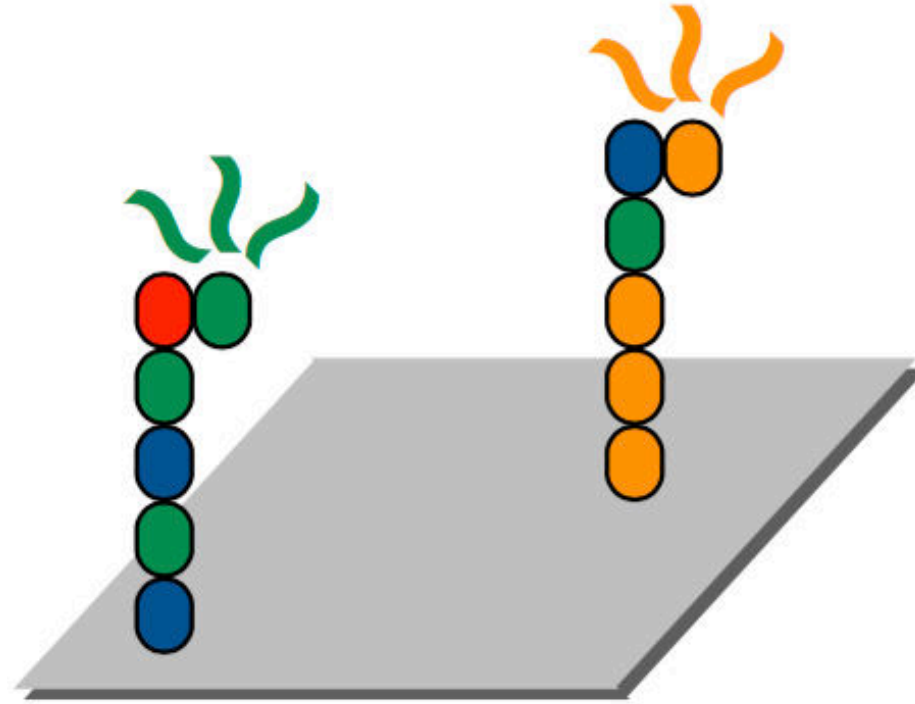


Millions of pH Sensors
Semiconductor Design



3. 2nd generation Illumina (sequencing)

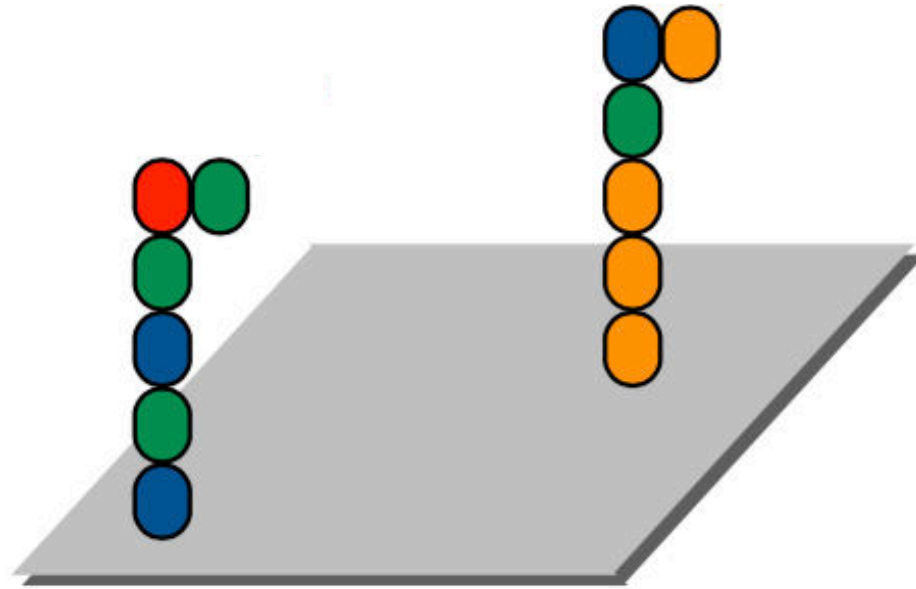
Besides bases,
polymerase and
primers





3. 2nd generation Illumina (sequencing)

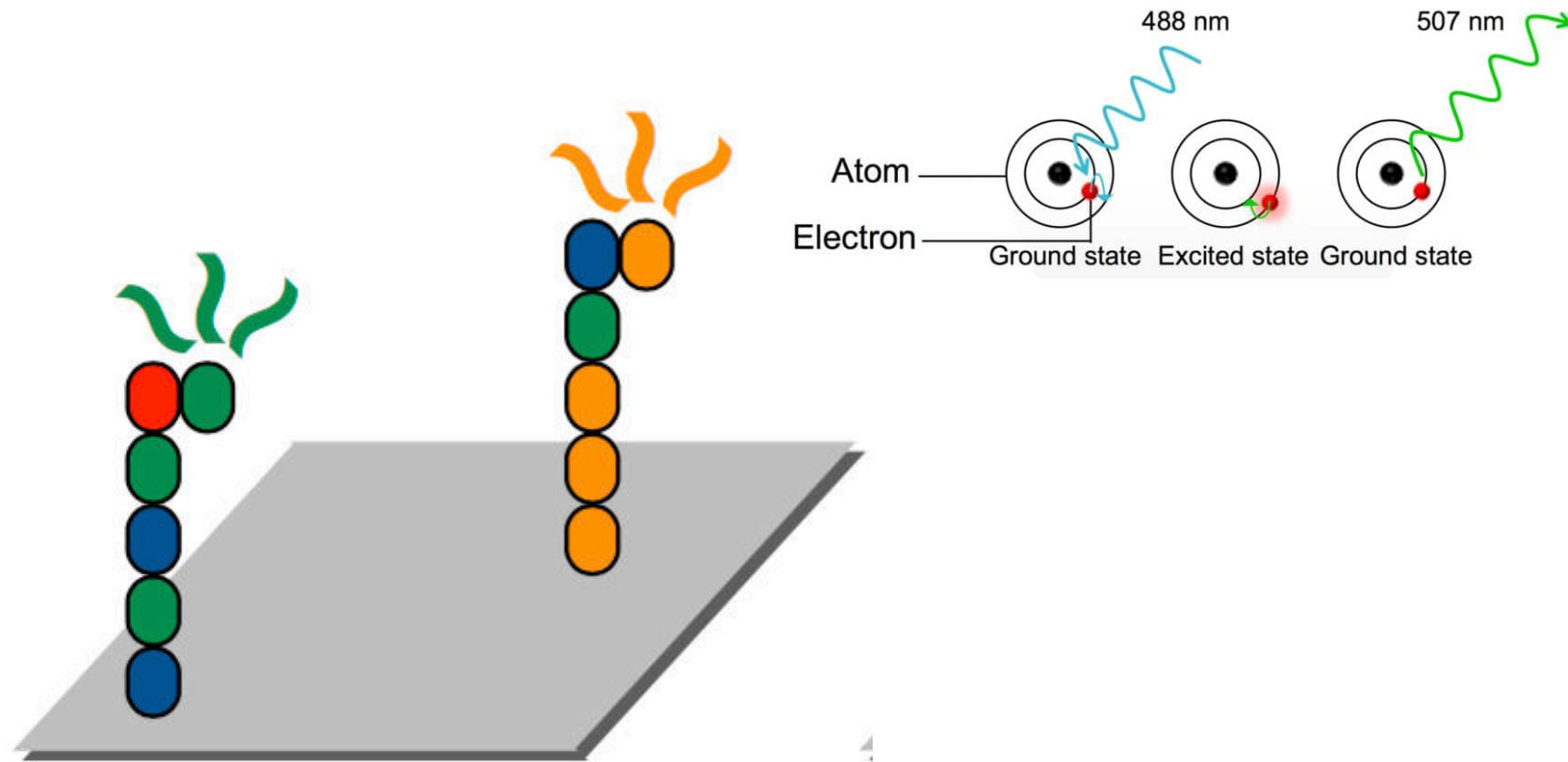
Important
ingredients are





3. 2nd generation Illumina (sequencing)

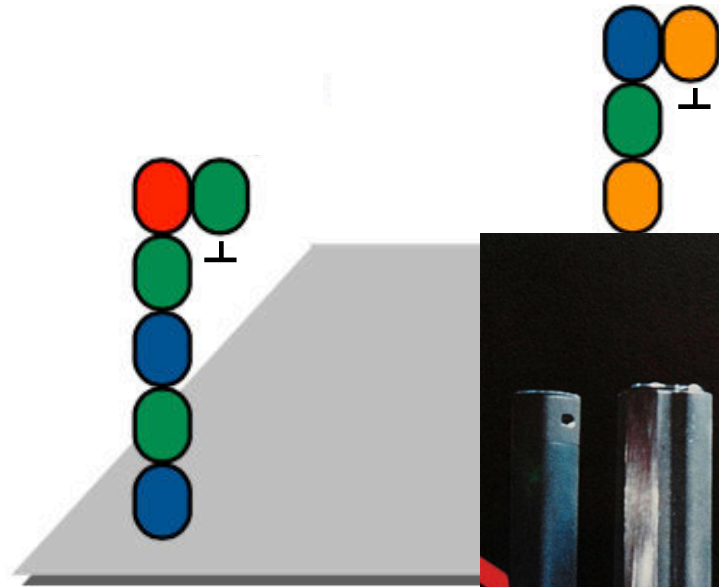
Important
ingredients are
the fluorophore





3. 2nd generation Illumina (sequencing)

Important
ingredients are
the fluorophore &
terminator!

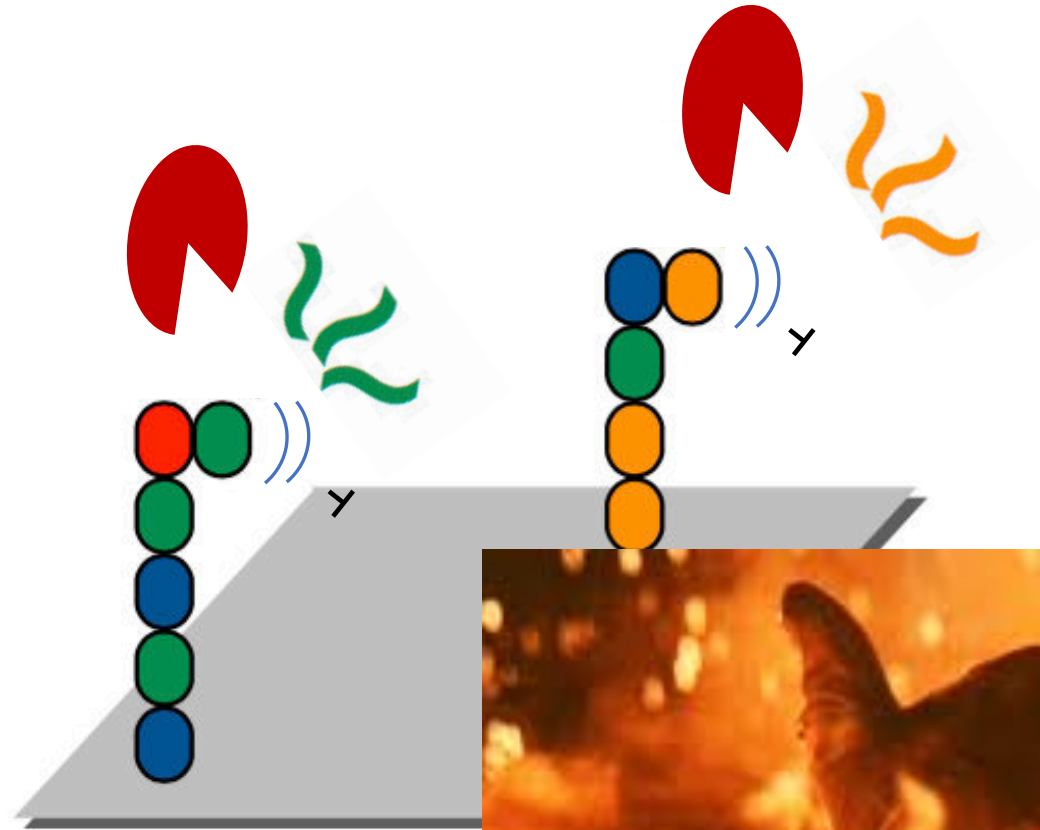




3. 2nd generation Illumina (sequencing)

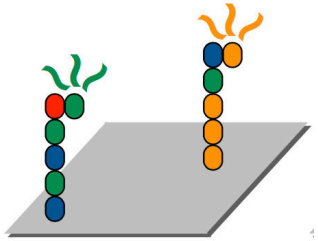
Important ingredients are the fluorophore & terminator!

... a cleaving enzyme ...
AND then, the terminator is...
NO MORE!





3. 2nd generation Illumina (sequencing)

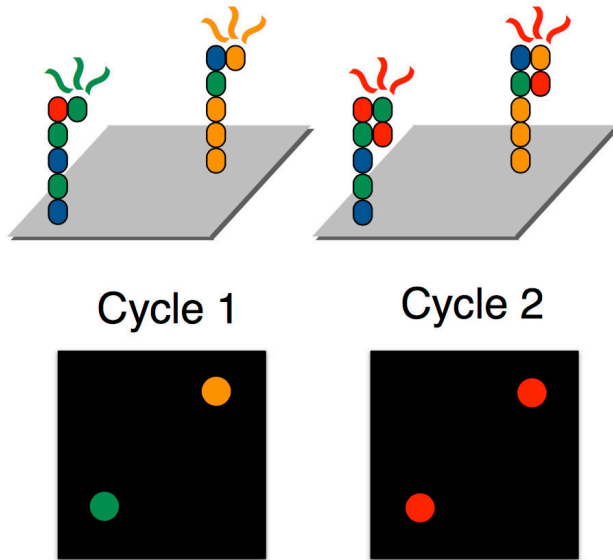


Cycle 1



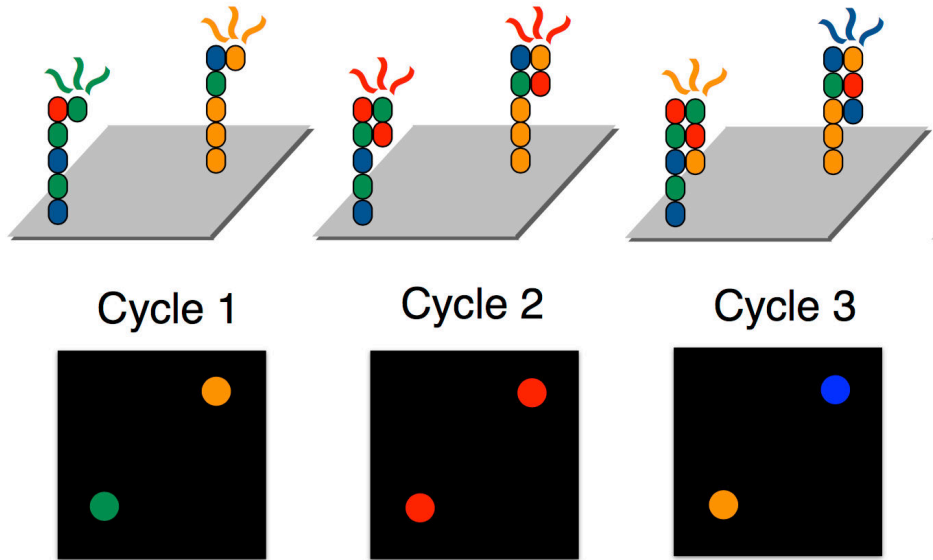


3. 2nd generation Illumina (sequencing)



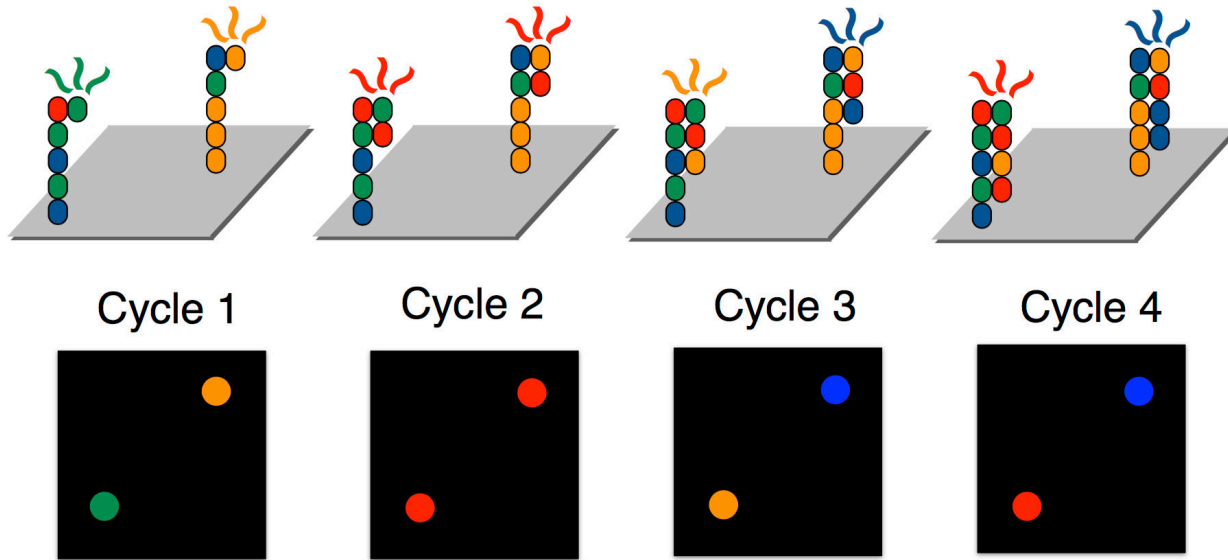


3. 2nd generation Illumina (sequencing)



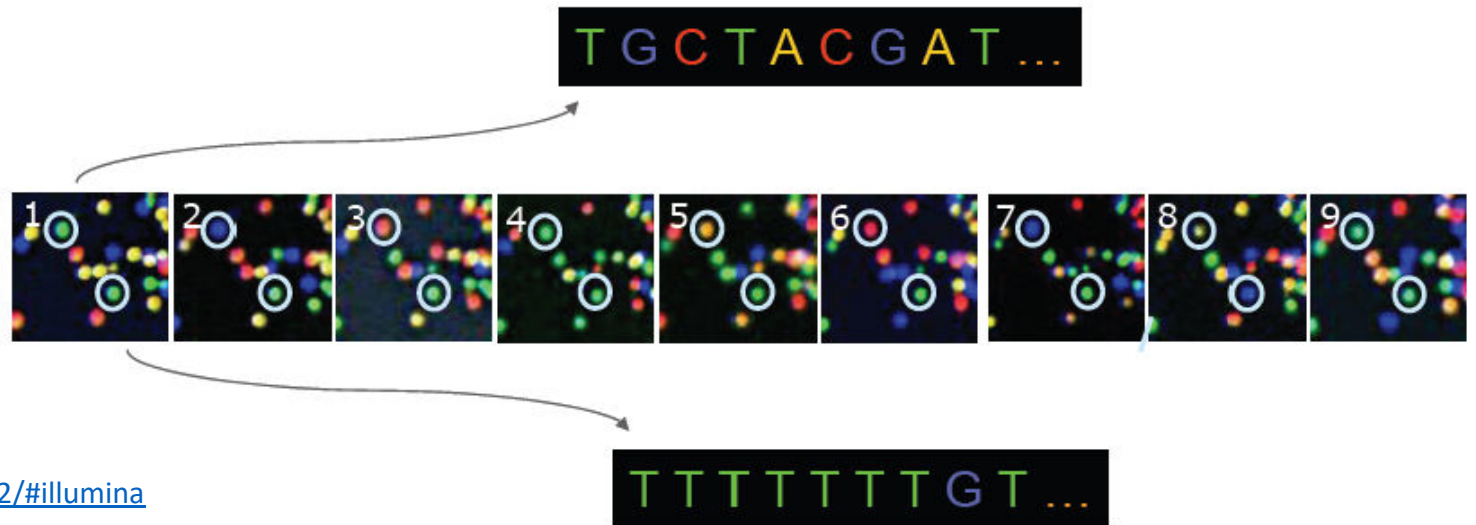
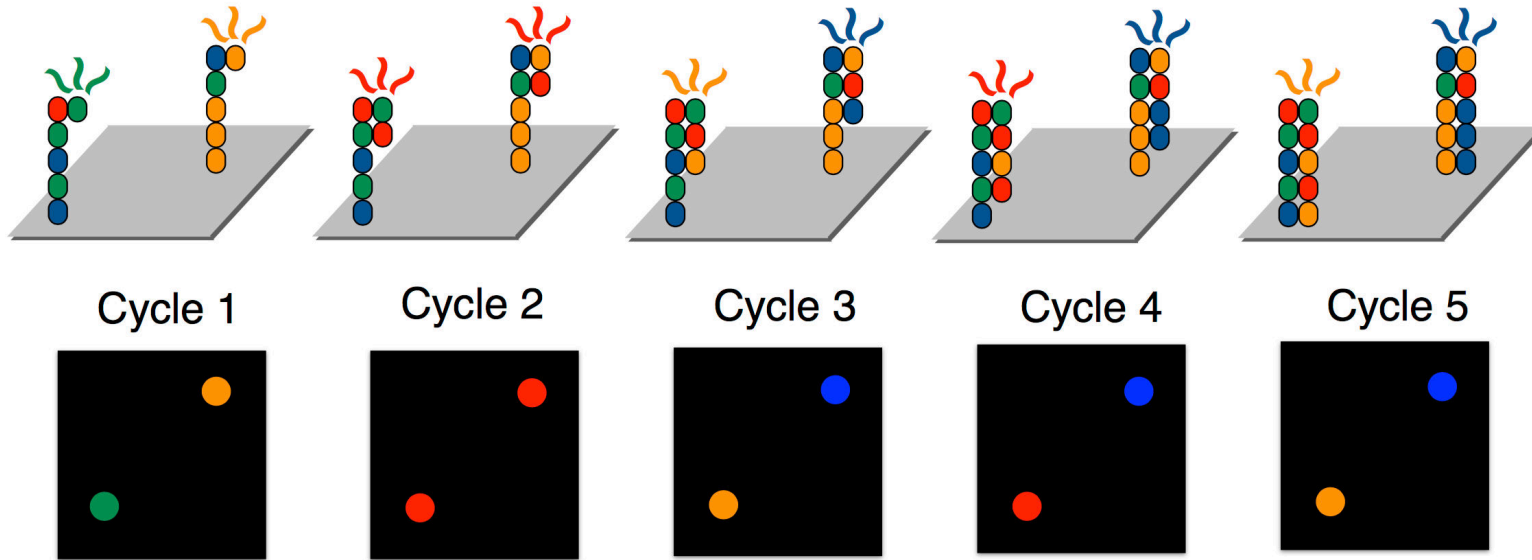



3. 2nd generation Illumina (sequencing)





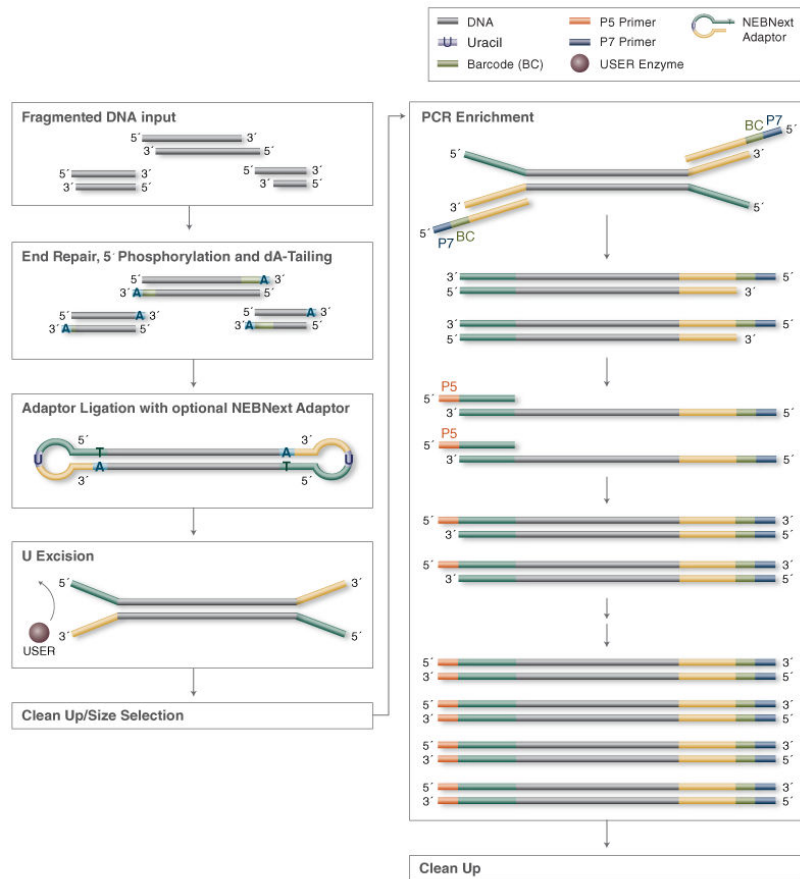
3. 2nd generation Illumina (sequencing)



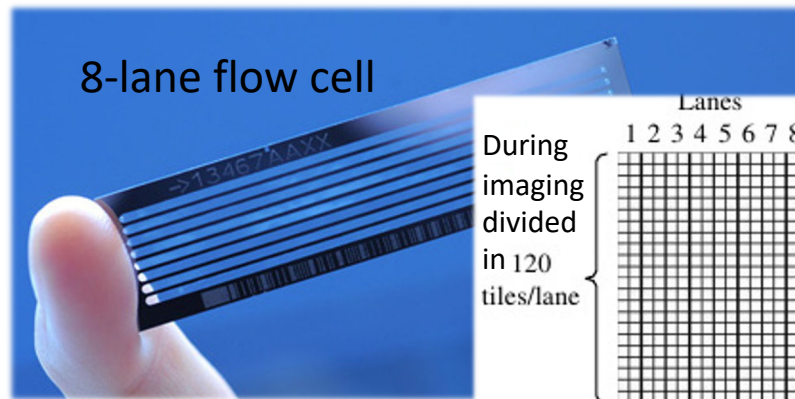
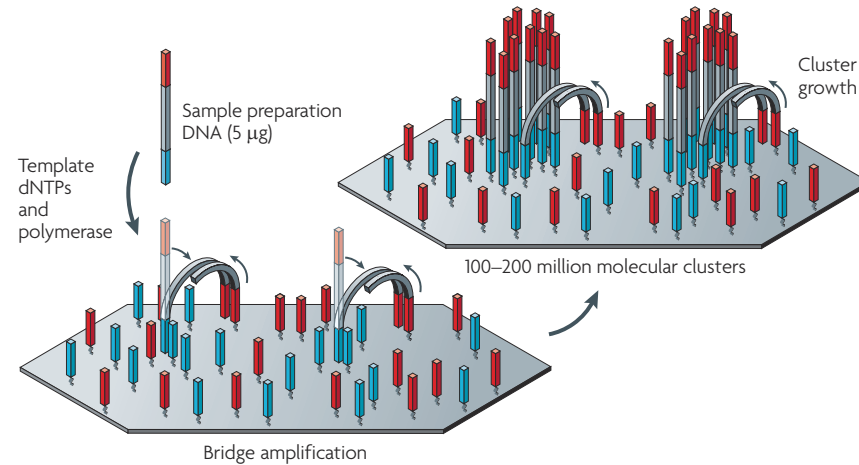

 DEPARTMENT OF
Biochemistry & Biotechnology
 UNIVERSITY OF THESSALY

3. 2nd generation Illumina

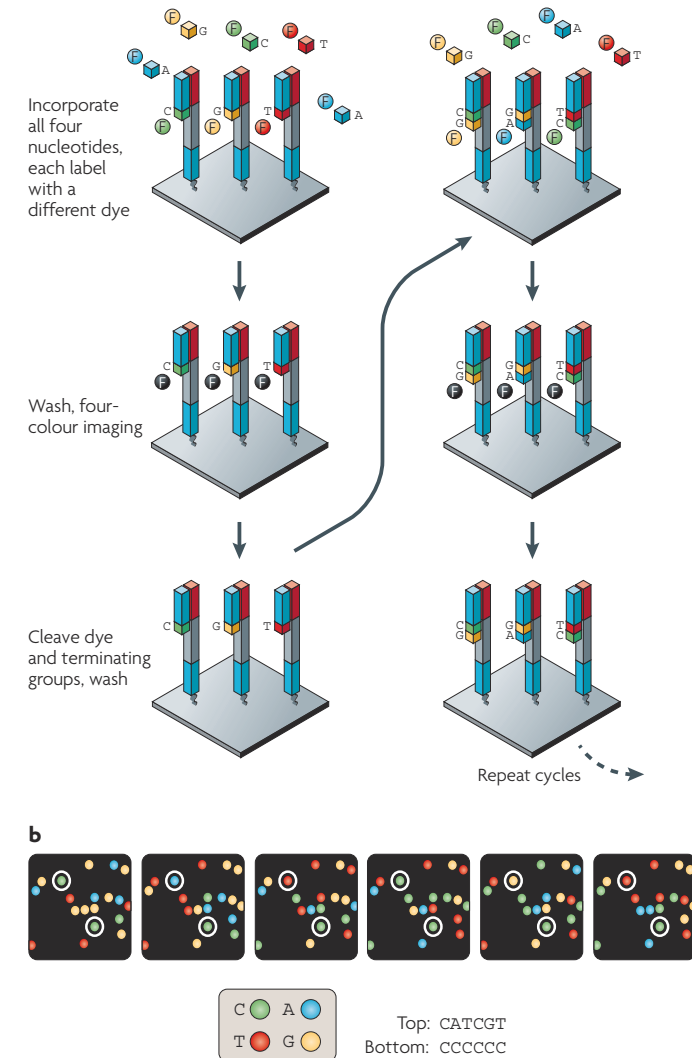
Library prep



Flow cell attachment and bridge amplification



Sequencing data generation



<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

<https://www.neb.com/applications/library-preparation-for-next-generation-sequencing/illumina-library-preparation>

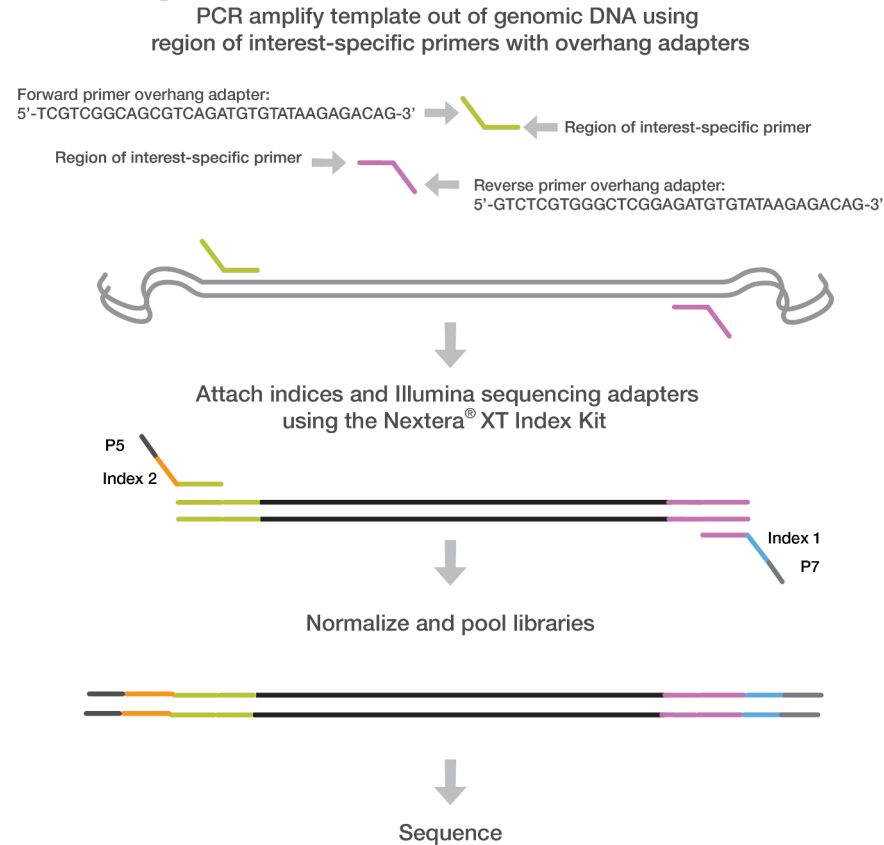
Kawashima, E., *et al.* (1998). Method of nucleic acid amplification, Google Patents

Metzker, M.L. (2009). Sequencing technologies — the next generation. *Nat Rev Genet* 11, 31 (modified)



3. 2nd generation Illumina (library prep)

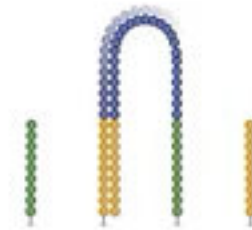
Figure 1 16S V3 and V4 Amplicon Workflow



User-defined forward and reverse primers that are complementary upstream and downstream of the region of interest are designed with overhang adapters, and used to amplify templates from genomic DNA. A subsequent limited-cycle amplification step is performed to add multiplexing indices and Illumina sequencing adapters. Libraries are normalized and pooled, and sequenced on the MiSeq system using v3 reagents.

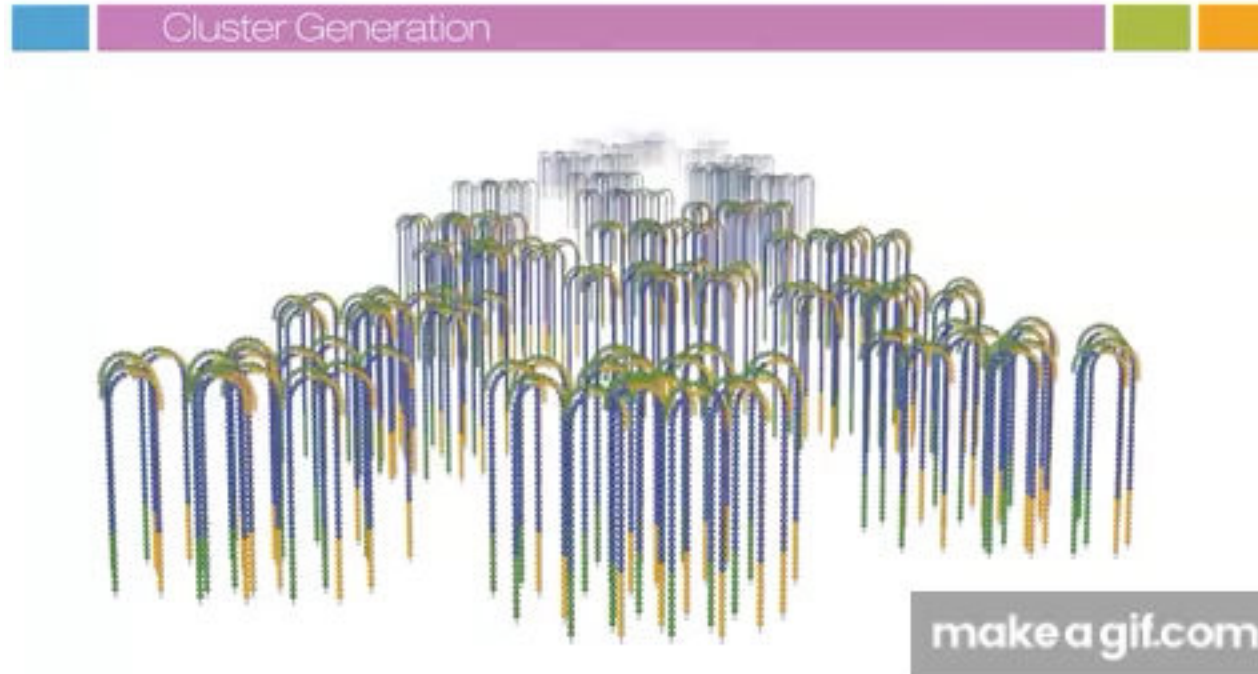


3. 2nd generation Illumina (cluster prep)





3. 2nd generation Illumina (cluster prep)

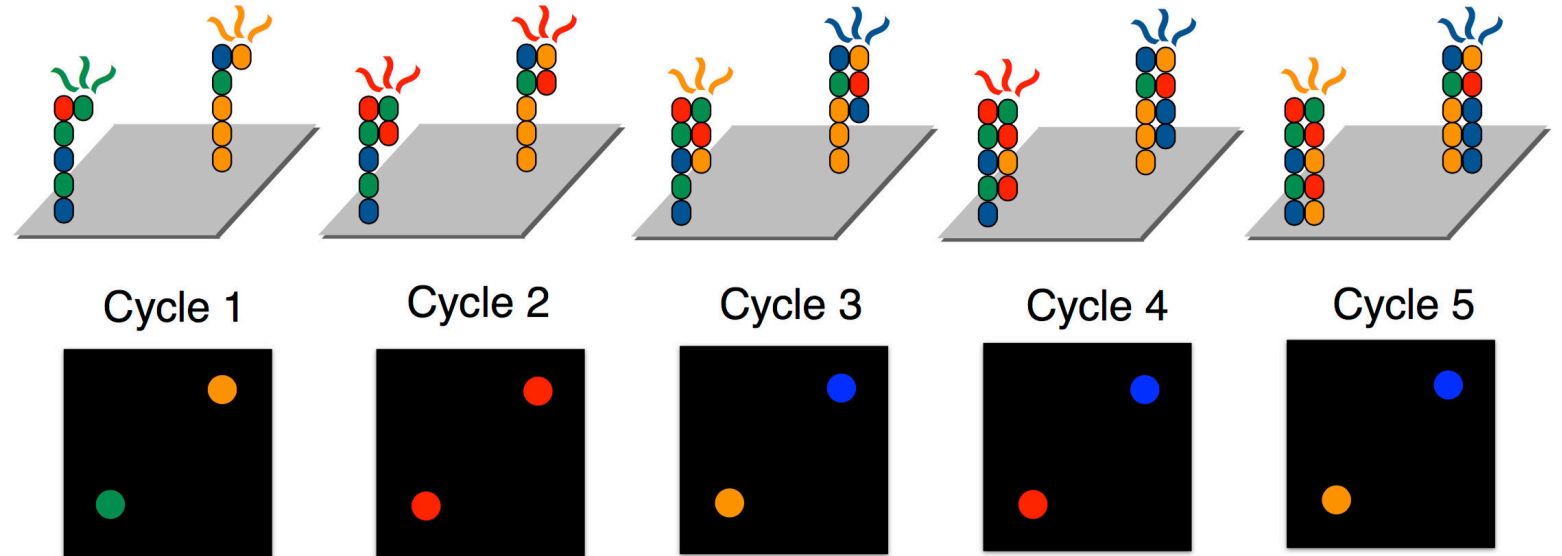




3. 2nd generation Illumina (sequencing)

Basic ingredients:

- 1) Single stranded template DNA attached to a slide (flow-cell lane)
- 2) Four nucleic acid bases with reversible terminators and cleavable dyes (each having its own colour: e.g. Adenine emits at yellow when excited, G at blue, T at green, C at red)
- 3) DNA polymerase



Process:

- 1) The 4 bases are added together with the polymerase on the slide with the attached template ssDNA fragments
- 2) A base with its dye and terminator is added in a complementary fashion against each ssDNA (5' -> 3') to synthesize dsDNA and the reaction is terminated
- 3) The dye is excited and an image is acquired
- 4) The dye and the terminator are enzymatically removed
- 5) The cycle is repeated (go to step 1)



3. 2nd generation Illumina (sequencing)



3. Illumina reads output (the fastq format)

[illegible]

```
@M02542:38:000000000-ABF47:1:1101:17854:1061 1:N:0:1
GTTATGCATTGAAAGGGAAACGATTGAAGTCAGTCGTACCTTCGGGTAATCAGCCTTTCGNGTGCCNNNCTGNNGAGCGTGAGGAGTTAAA
+
CCCCCGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG# :DF##### :C## :DFGGGGGGGGGGGG
@M02542:38:000000000-ABF47:1:1101:18832:1061 1:N:0:1
TACTGCGGTTGAAAGGGAAACGATTGAAGTCAGTCGTACCTGCGGGTAATCAGCCTTTTGNGTGTNNNNATGNNAAGCGTGAGGAGTTAAT
+
CCCCCGGGFGGFGGGGGGGGGEGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG# :DGF##### :D## :@FGGGGGGGGGGG
@M02542:38:000000000-ABF47:1:1101:17854:1061 2:N:0:1
TACGTCAACATCCTTAACGNATTTANCNNGTNCTAACCTTTTGAACGNNANNTNTGATNACCGTANGCAAGCTTTCGGNACCAGAGCAACT
+
CCCCCGGGGGGGGGGGGGGG# =CFG# :## :# :CDFGGGGGGGGGG# :## :# :FG# :CFGFF# :CDGGGGGGGE# :AFGGFGGGGG
@M02542:38:000000000-ABF47:1:1101:18832:1061 2:N:0:1
TACGTCAACATCCTTAACGNATTTANCNNGTNCTAACCTATTGAACGNNANNTNTGATNCCATAANCAAGCTTTCAGCNCCAGAGCAACTT
+
CCCCCGGGGGGFGGGGGGG# =CFFF# :## :# :CFGFGGFGGFGDF# :## :# : ,C# :CFFF# :DFGGGGGGGG# :AFGCFGGDGG
.....
```



3. Error probability (Phred Q values)

$$Q = -10 \times \log_{10}(P_{err}) \quad \Leftrightarrow \quad P_{err} = 10^{Q/-10}$$

Representation of quality scores

Table 1 ASCII Characters Encoding Q-scores 0-40

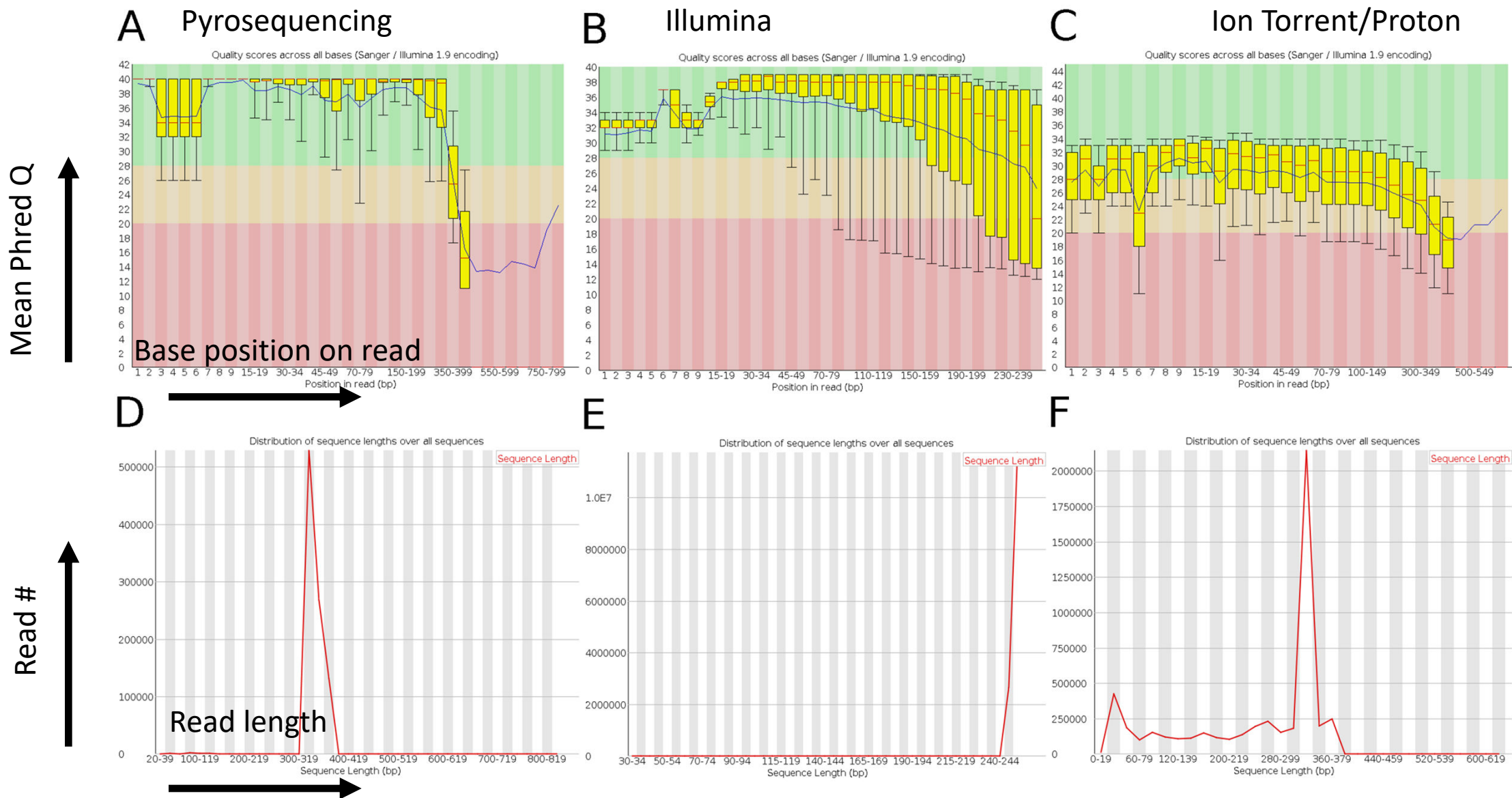
Symbol	ASCII Code	Q-Score	Symbol	ASCII Code	Q-Score	Symbol	ASCII Code	Q-Score
!	33	0	/	47	14	=	61	28
"	34	1	0	48	15	>	62	29
#	35	2	1	49	16	?	63	30
\$	36	3	2	50	17	@	64	31
%	37	4	3	51	18	A	65	32
&	38	5	4	52	19	B	66	33
'	39	6	5	53	20	C	67	34
(40	7	6	54	21	D	68	35
)	41	8	7	55	22	E	69	36
*	42	9	8	56	23	F	70	37
+	43	10	9	57	24	G	71	38
,	44	11	:	58	25	H	72	39
-	45	12	;	59	26	I	73	40
.	46	13	<	60	27			

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%



3. 2nd generation error probability





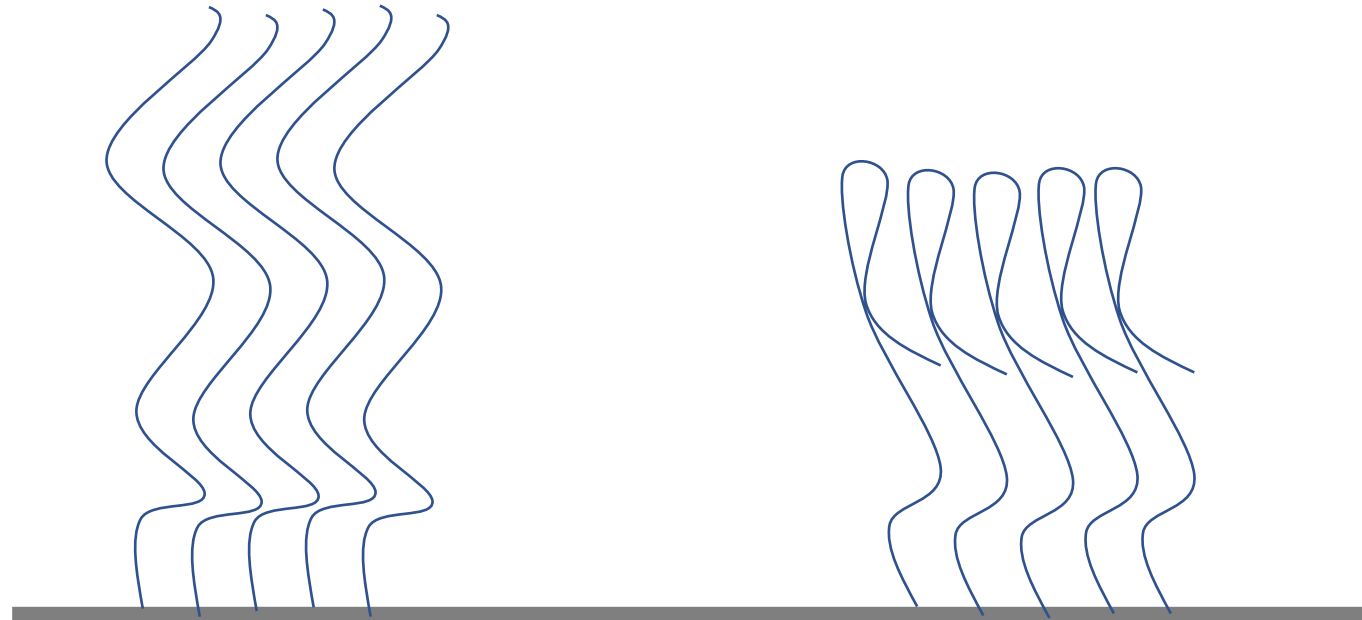
3. 2nd gen Inherent sequencing error types

- Substitutions in Illumina read-ends (due to phasing): can be dealt with using paired-end sequencing overlaps
- Homopolymers in pyrosequencing and Ion Torrent/Proton reads



Error sources (inverted repeats)

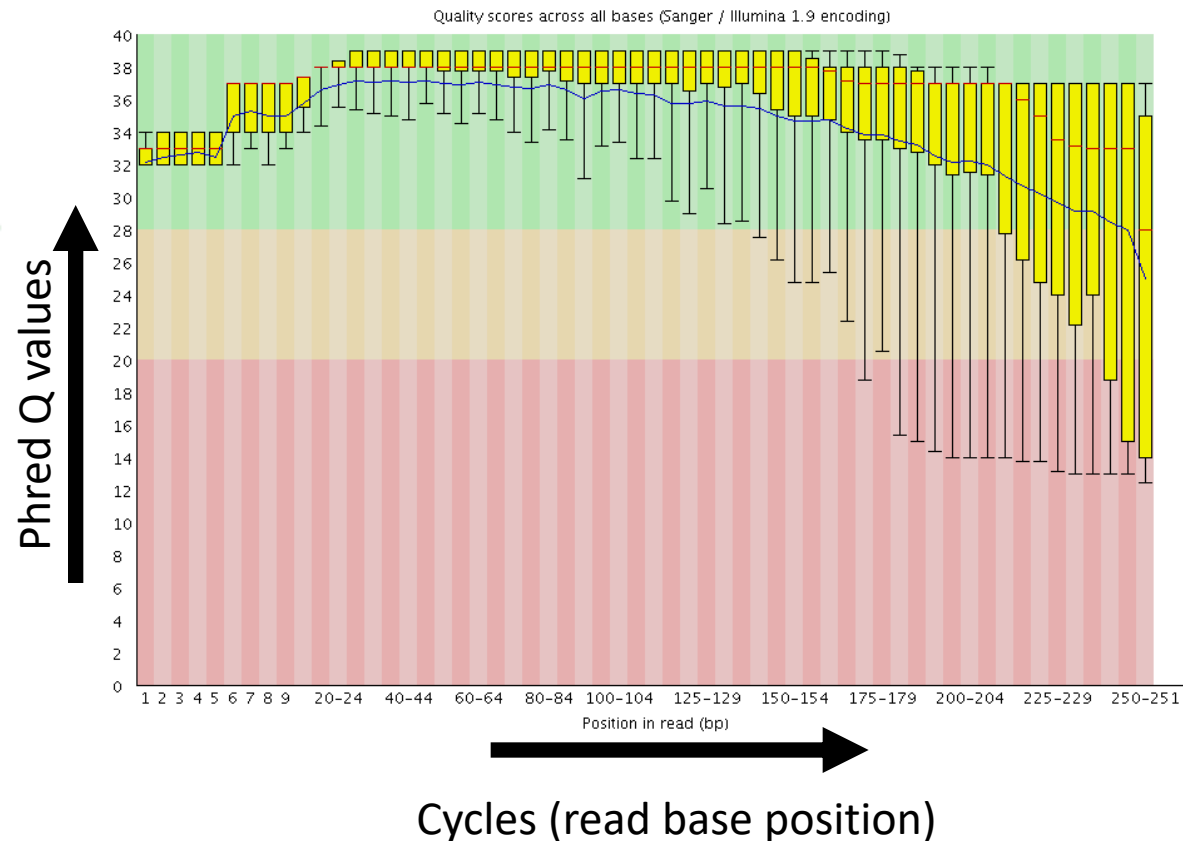
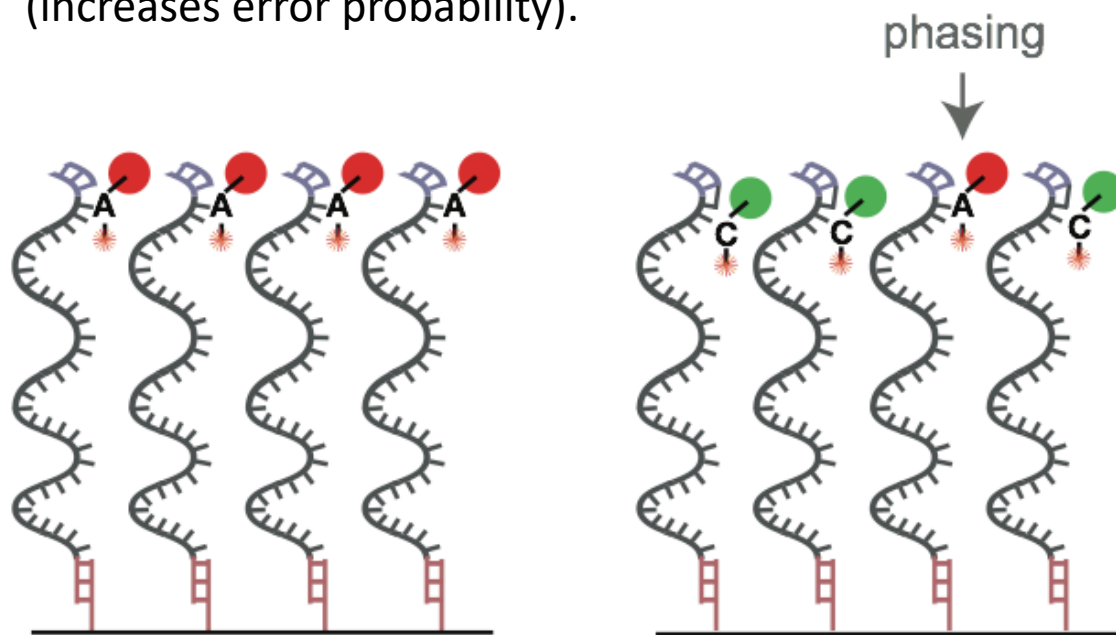
Difficulties with template replication due to secondary structures (e.g. palindromic repeats).





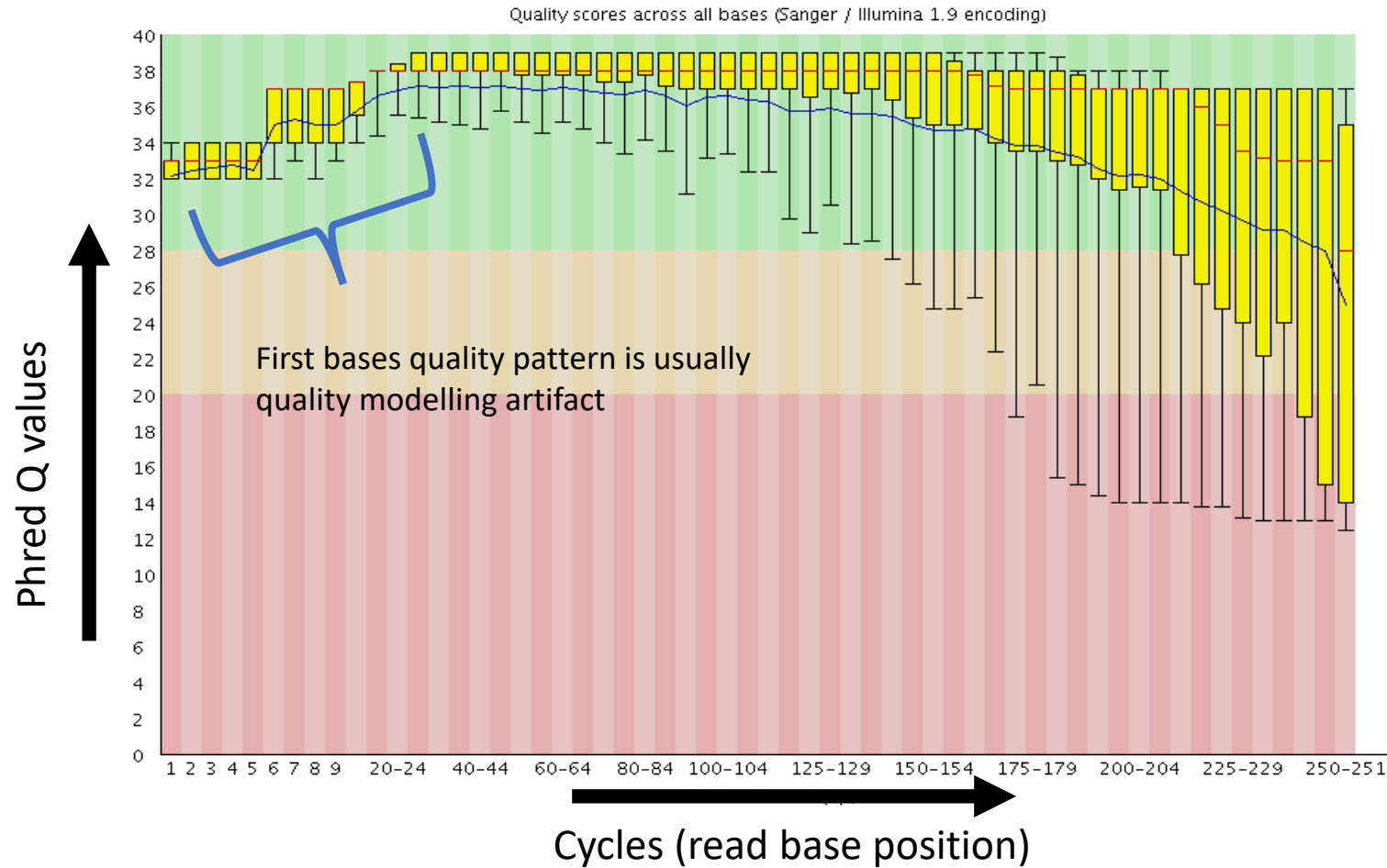
3. Error source in Illumina (phasing)

Fluorescence phasing due to imperfect function of enzymes. Occasional lack of terminator/fluorophore removal accumulates toward the read-end and reduces Phred Q (increases error probability).





3. Error source in Illumina (initial; basecalling algorithmic)



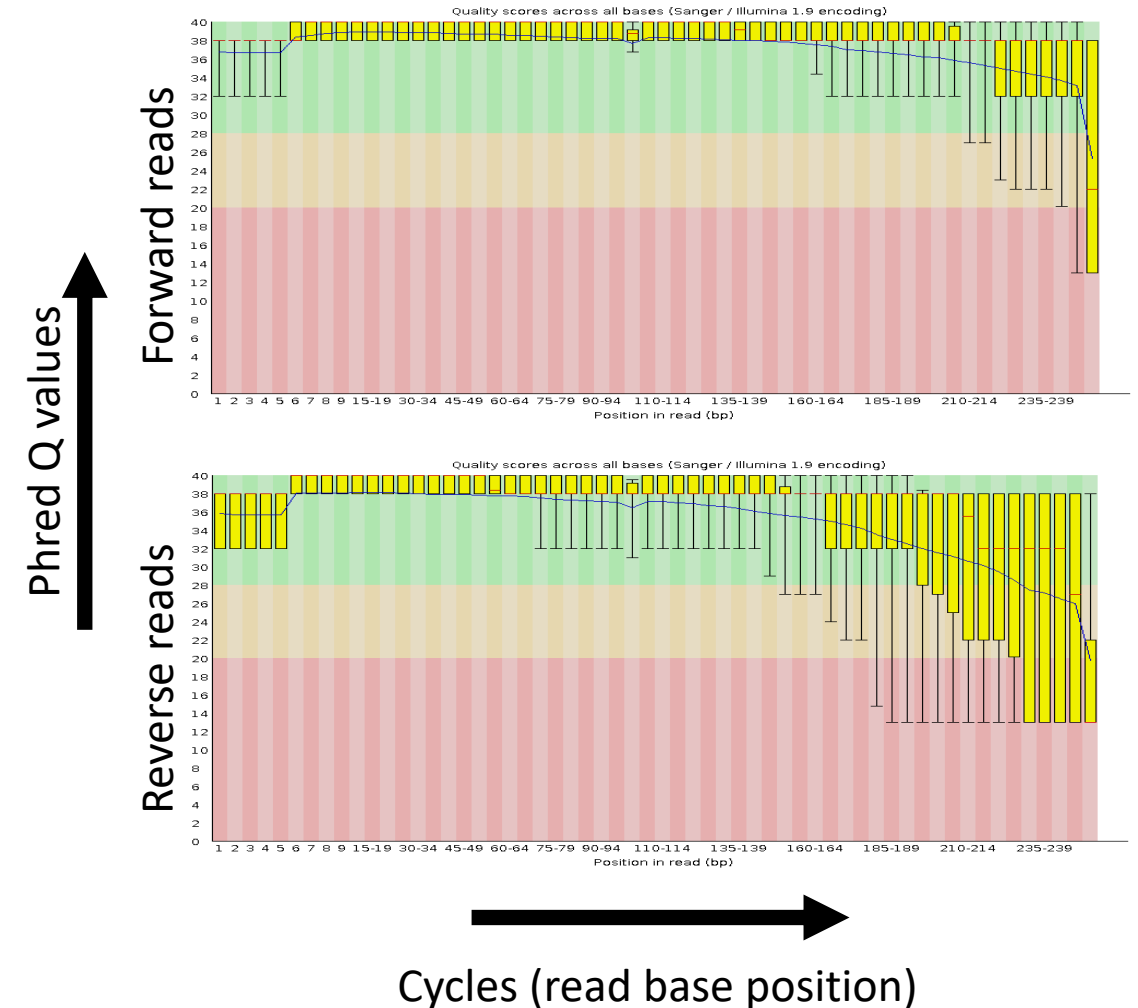


3. Error source in Illumina (for vs rev)

Besides their similar per-base error patterns, forward reads tend to have slightly higher per base quality score statistics than the reverse reads.

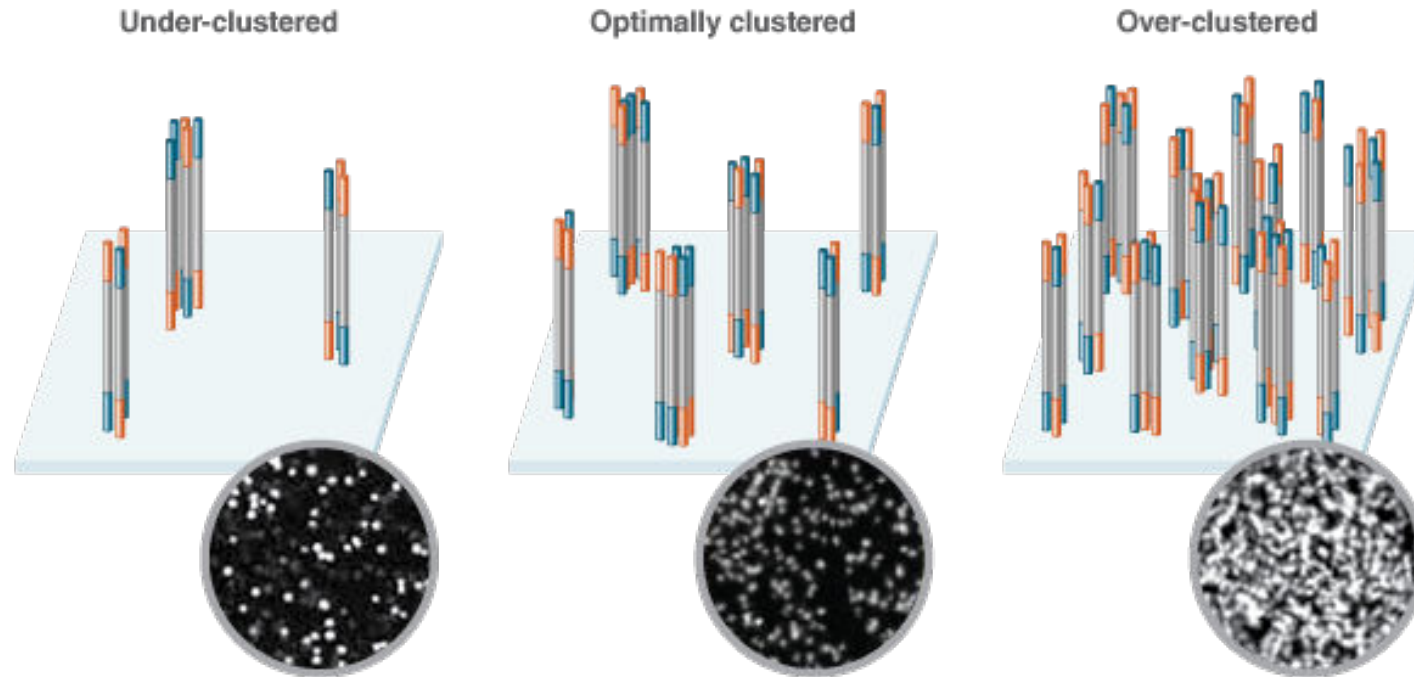
Possibly attributed to

- overclustering (extra bridge amplification for second strand)
- Possible DNA wear-out





3. Error source in Illumina (over-clustering)

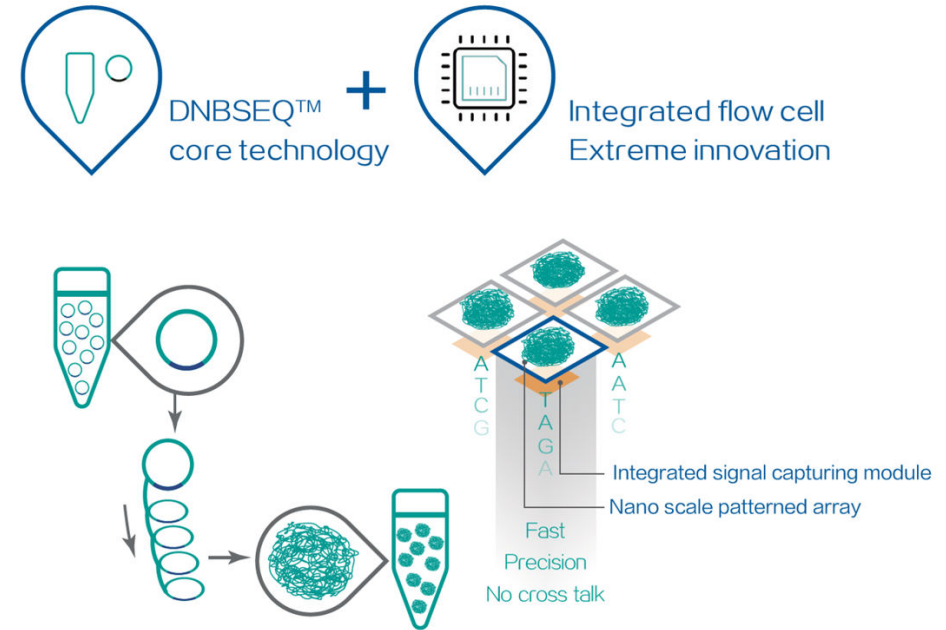




DNA nanoball

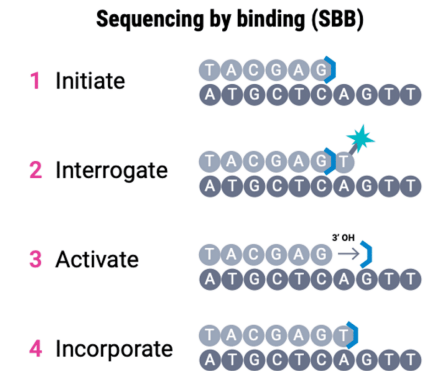
DNA nanoball

- Basic difference: rolling cycle amplification during library prep -> no secondary errors



DNA sequencing by binding similar with Illumina

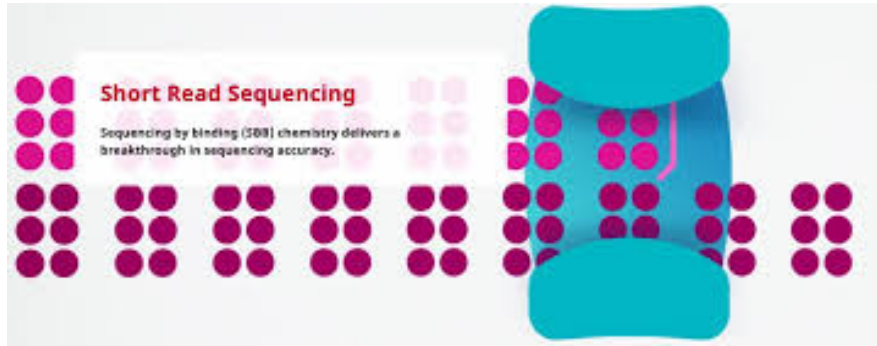
- Enzymatic template interrogation with fluorescent base matching prior polymerization.



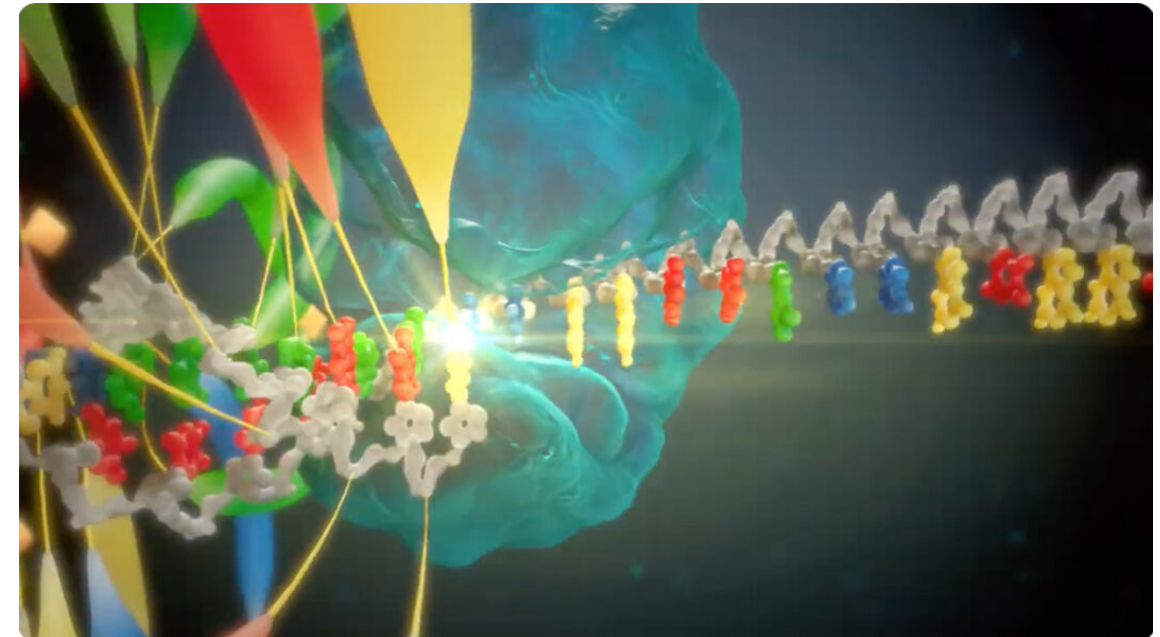


Several others

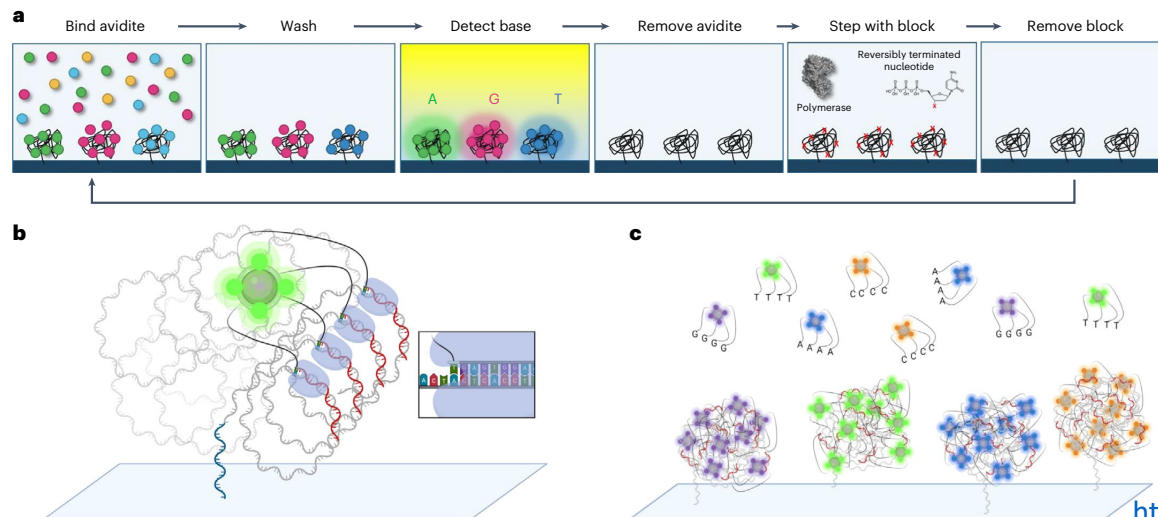
Pacific Biosciences... sequencing by binding (Onso instrument)



Roche...Sequencing by expansion (SBX)



Element Biosciences... rolling cycle amplification and Avidites



<https://www.genengnews.com/topics/omics/the-best-of-ngs-instrument-companies-to-watch-in-2025/>

Arslan, S., et al. (2024). Nat Biotechnol 42, 132-138, doi: 10.1038/s41587-023-01750-7

<https://www.pacb.com/technology/sequencing-by-binding/>



3. Synthetic long reads

Barcode and Amplify

Ligation or direct amplification add LoopSeq-specific sequences. Amplification then adds one unique barcode to each molecule and creates thousands of copies of each molecule and its unique barcode



<https://www.elementbiosciences.com/products/loopseq>

For sFL16S see Callahan et al 2021, Microbiome <https://doi.org/10.1186/s40168-021-01072-3>

Or Jeong et al (2021). Scientific Reports 11, 1727, [10.1038/s41598-020-80826-9](https://doi.org/10.1038/s41598-020-80826-9)

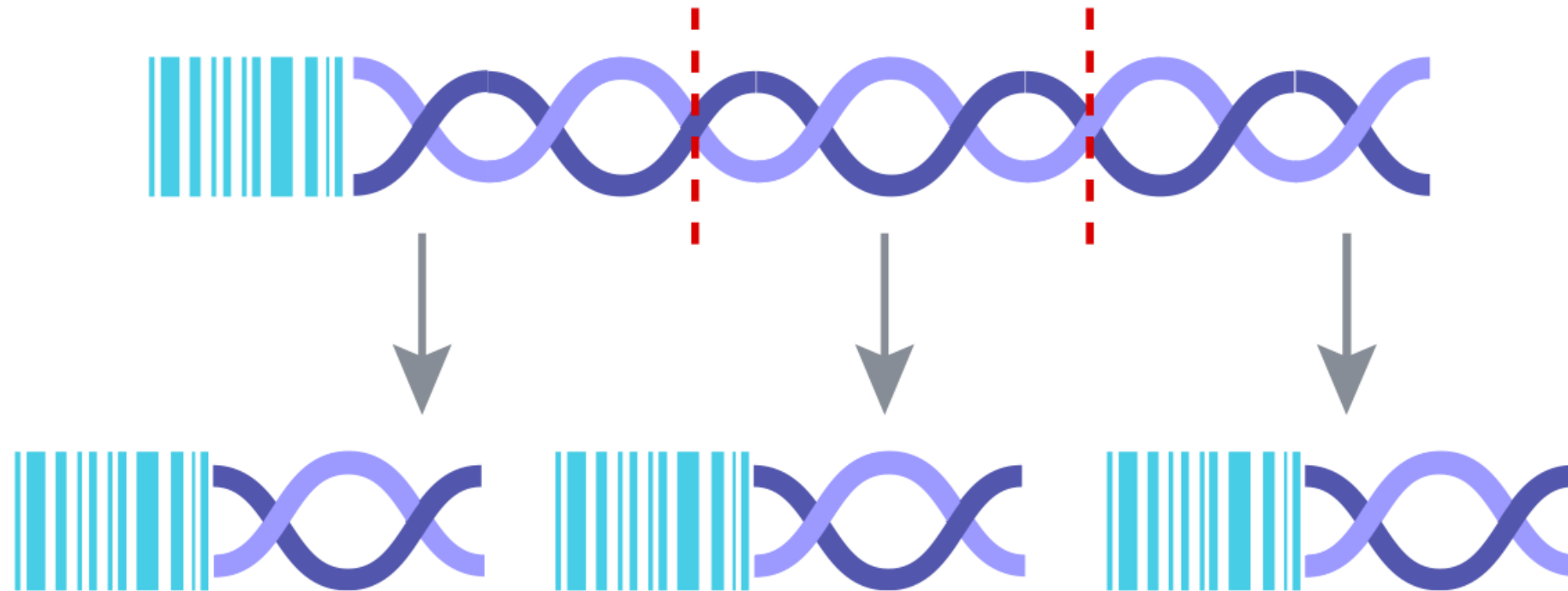
https://hpst.cz/sites/default/files/oldfiles/loopseq-16s-18s-microbiome-24-plex-user-manual_0.pdf



3. Synthetic long reads

Distribute and Prep

Enzymes digest each molecule into fragments and distribute copies of the barcode throughout the fragments, so each fragment has the same barcode as the source molecule.



<https://www.elementbiosciences.com/products/loopseq>

For sFL16S see Callahan et al 2021, Microbiome <https://doi.org/10.1186/s40168-021-01072-3>

Or Jeong et al (2021). Scientific Reports 11, 1727, [10.1038/s41598-020-80826-9](https://doi.org/10.1038/s41598-020-80826-9)

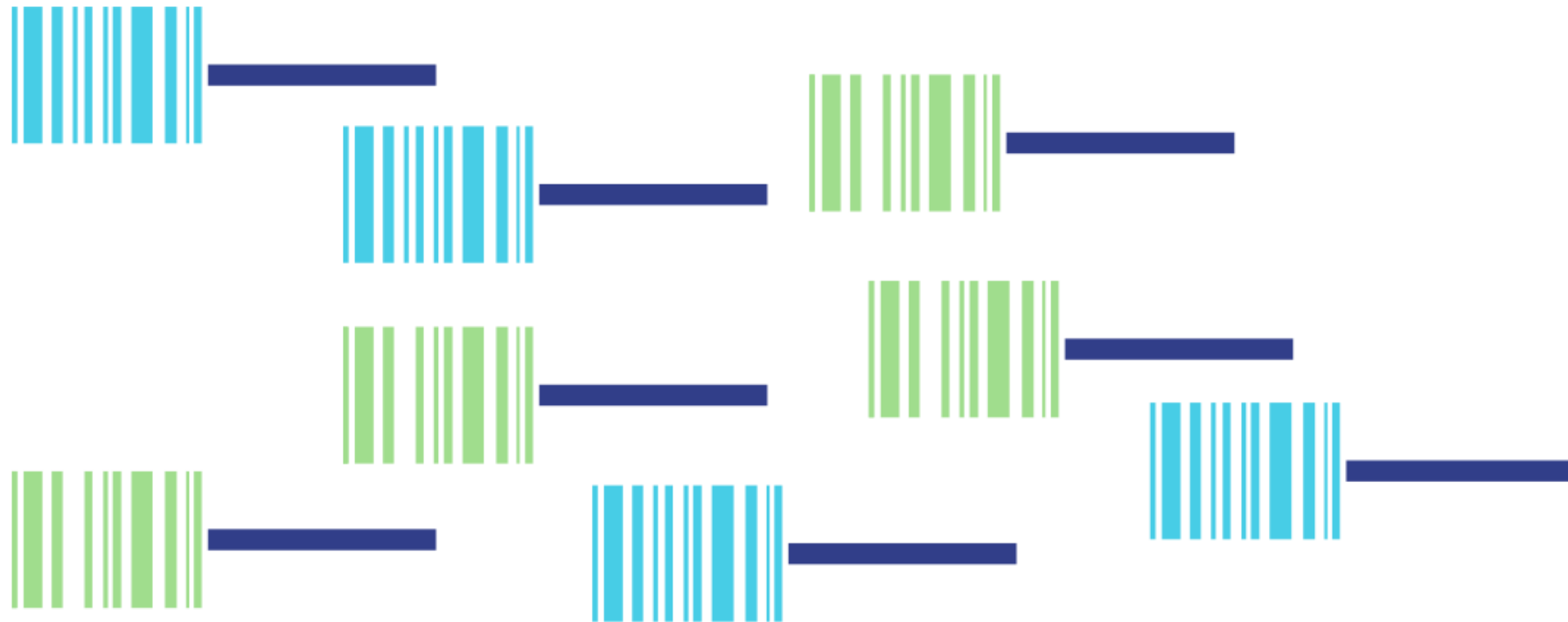
https://hpst.cz/sites/default/files/oldfiles/loopseq-16s-18s-microbiome-24-plex-user-manual_0.pdf



3. Synthetic long reads

Sequence

The AVITI System sequences the fragments as short reads. Subsequent demultiplexing groups the short reads by barcode.



<https://www.elementbiosciences.com/products/loopseq>

For sFL16S see Callahan et al 2021, Microbiome <https://doi.org/10.1186/s40168-021-01072-3>

Or Jeong et al (2021). Scientific Reports 11, 1727, 10.1038/s41598-020-80826-9

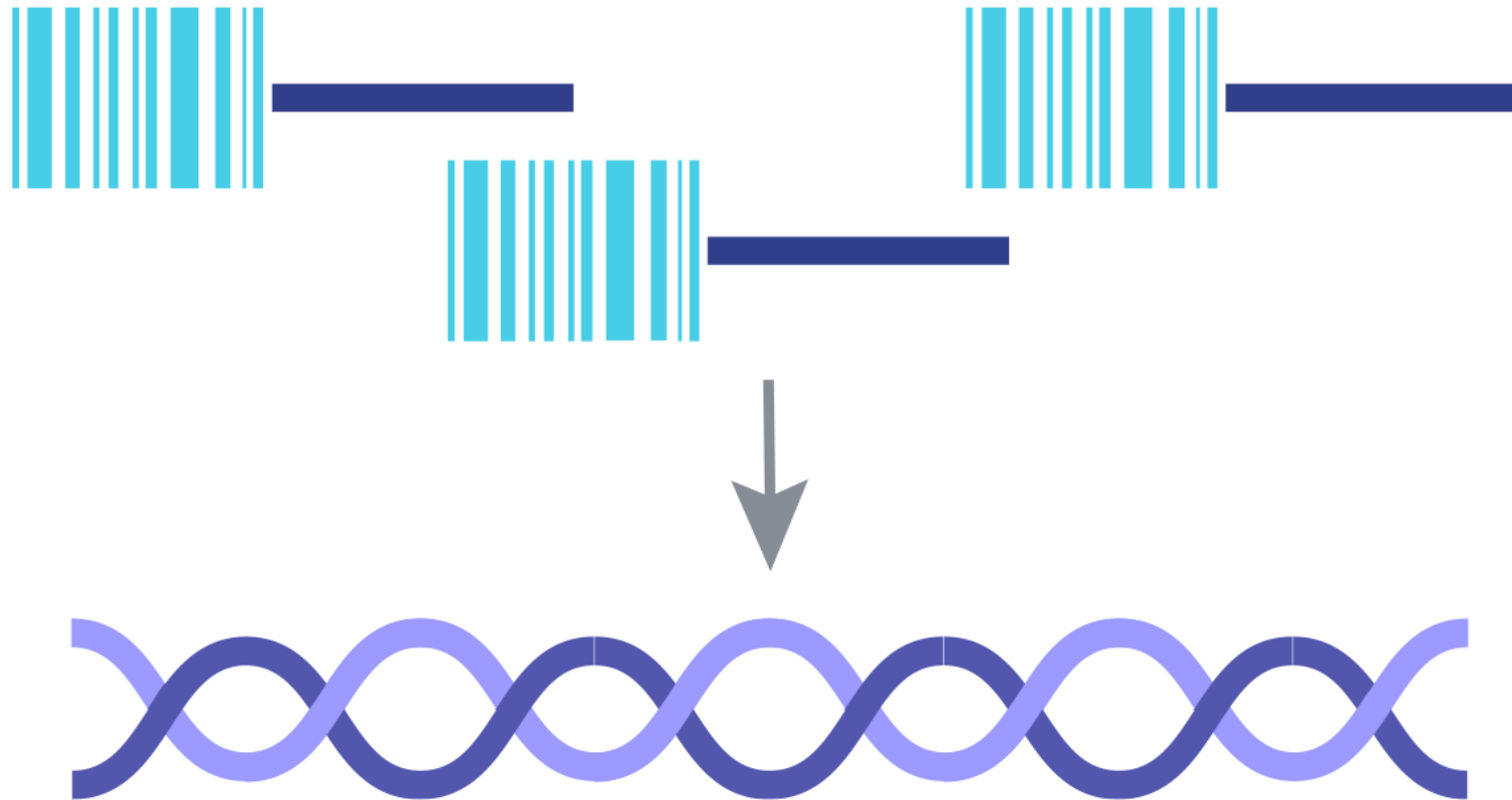
https://hpst.cz/sites/default/files/oldfiles/loopseq-16s-18s-microbiome-24-plex-user-manual_0.pdf



3. Synthetic long reads

Reassemble

De novo assembly reassembles overlapping short reads into the original full-length molecule.



<https://www.elementbiosciences.com/products/loopseq>

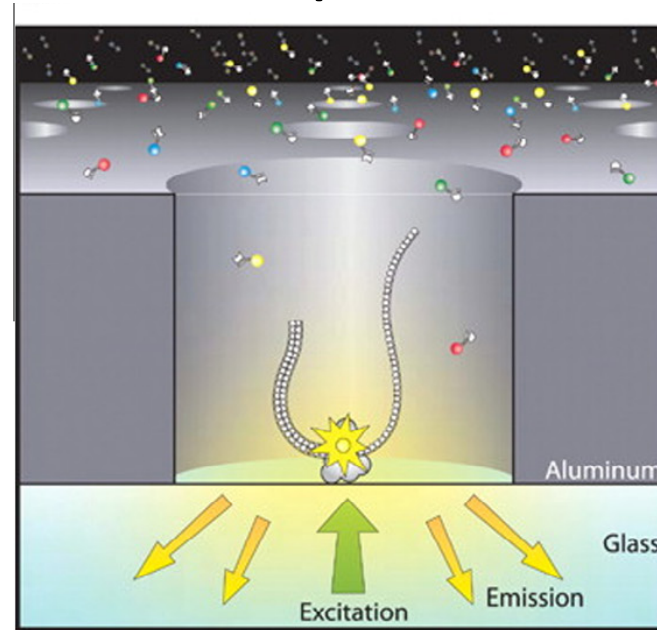
For sFL16S see Callahan et al 2021, Microbiome <https://doi.org/10.1186/s40168-021-01072-3>

Or Jeong et al (2021). Scientific Reports 11, 1727, [10.1038/s41598-020-80826-9](https://doi.org/10.1038/s41598-020-80826-9)

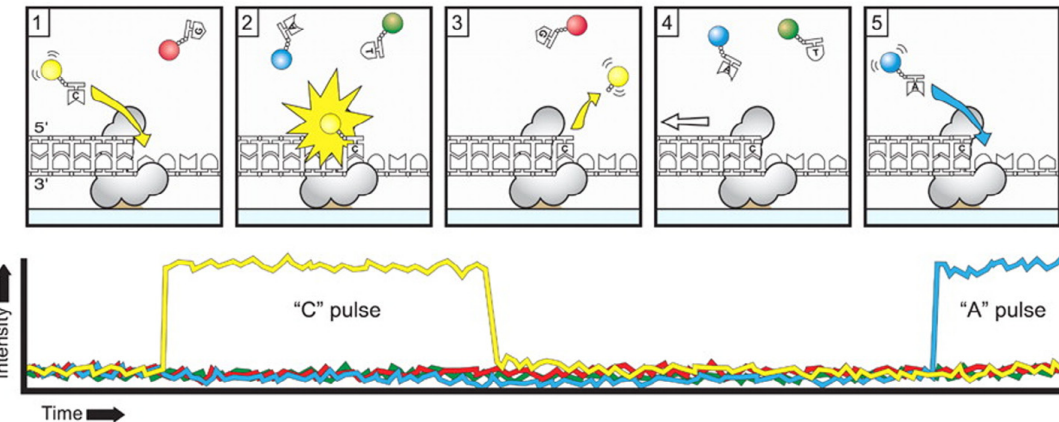
https://hpst.cz/sites/default/files/oldfiles/loopseq-16s-18s-microbiome-24-plex-user-manual_0.pdf



3rd generation: PacBio (single molecule real-time – SMRT – sequencing)



Zero-Mode Waveguide (ZMW)



Very long reads (currently of mean length of > 10 kbp)
but...
Random errors up to 10 % (mostly InDels)...

HiFi-reads!!!
Solution

Circular consensus sequencing (CCS)
+ overlapping-based correction

DNA fragment

ligate adaptors

sequence

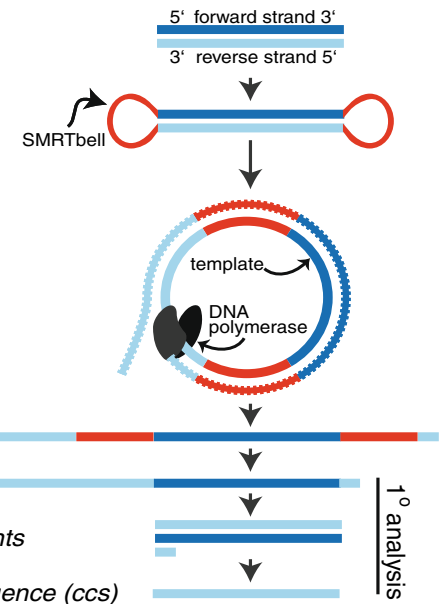
data analysis

raw long read

processed long read

single-molecule fragments

circular consensus sequence (ccs)



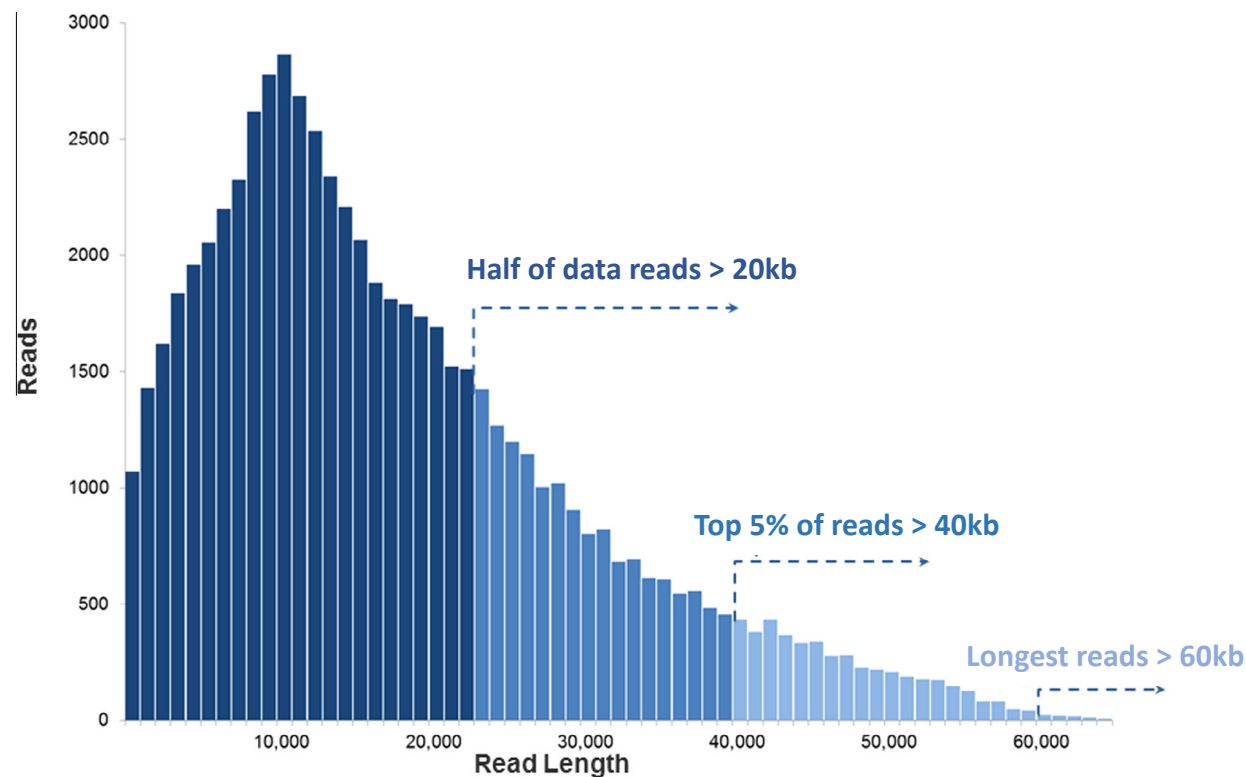
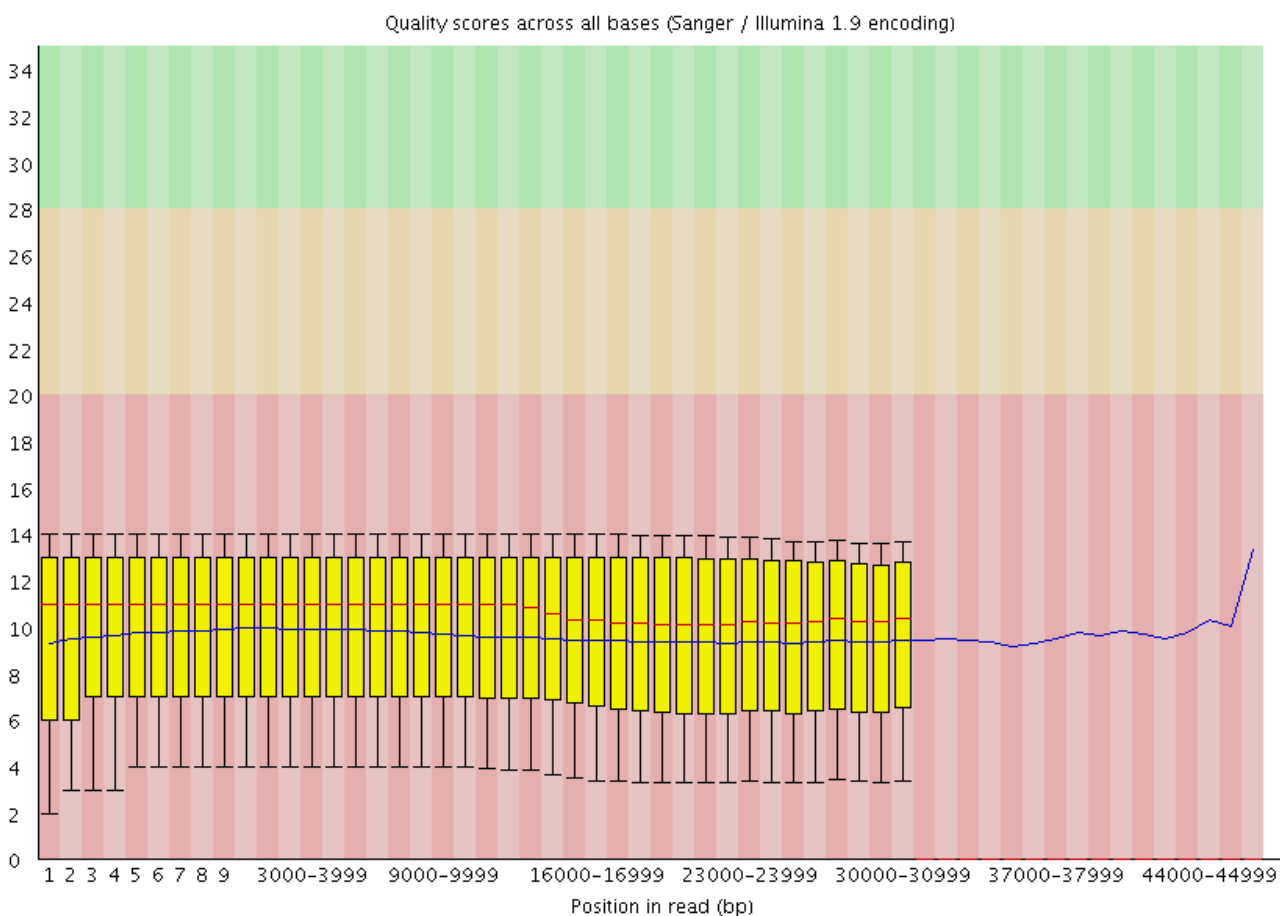
<https://www.youtube.com/watch?v=v8p4ph2MAvI>

Rhoads, A. and Au, K.F. (2015). PacBio Sequencing and Its Applications. Genomics, Proteomics & Bioinformatics 13, 278-289

Metzker, M.L. (2009). Sequencing technologies — the next generation. Nat Rev Genet 11, 31



3rd generation: PacBio per-base error probabilities



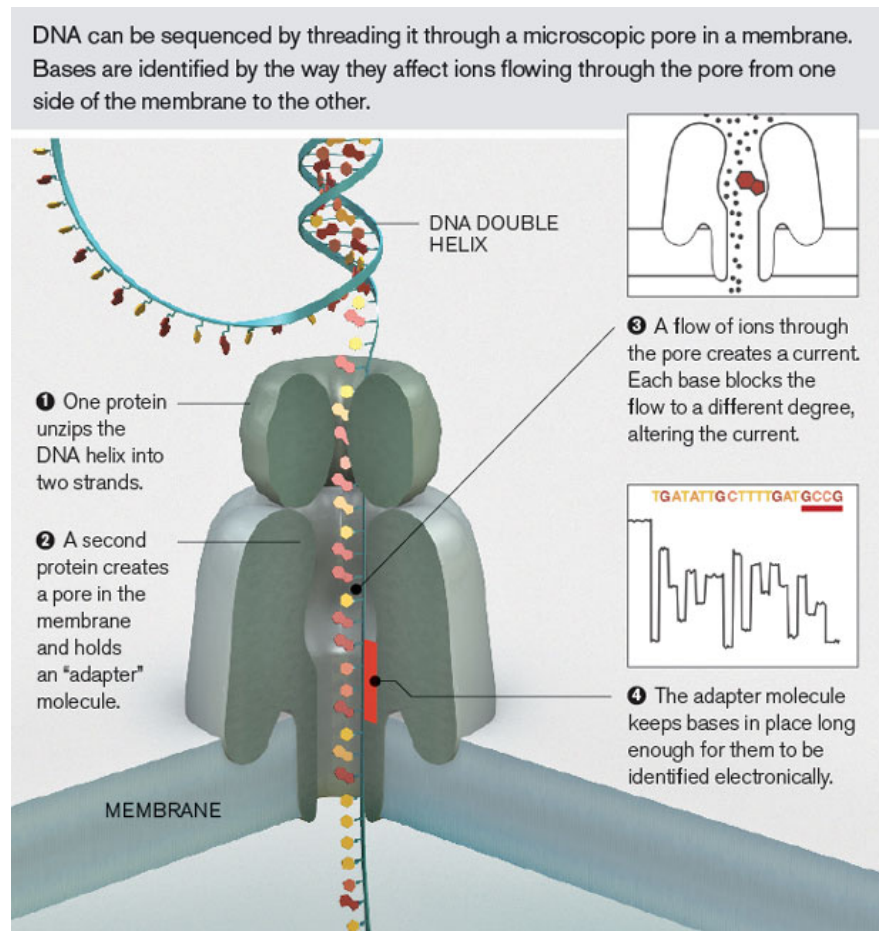


3rd generation: PacBio read error types

- Mainly Insertion/deletion
- Due to:
 - Polymerase monitoring sync issues (e.g. polymerization speed variation – usually deletion due to polymerase high speed)
 - Unincorporated nucleotide pass-bys (insertions)
- Corrected via the circular consensus sequencing process (multiple passes of the same sequences) and also computationally (overlapping-based correction)



3rd generation: Nanopore



- Higher error rates than PacBio
- Error types are usually InDels
- Minimal sample prep and benchtop convenience



<https://nanoporetech.com/products/minion>

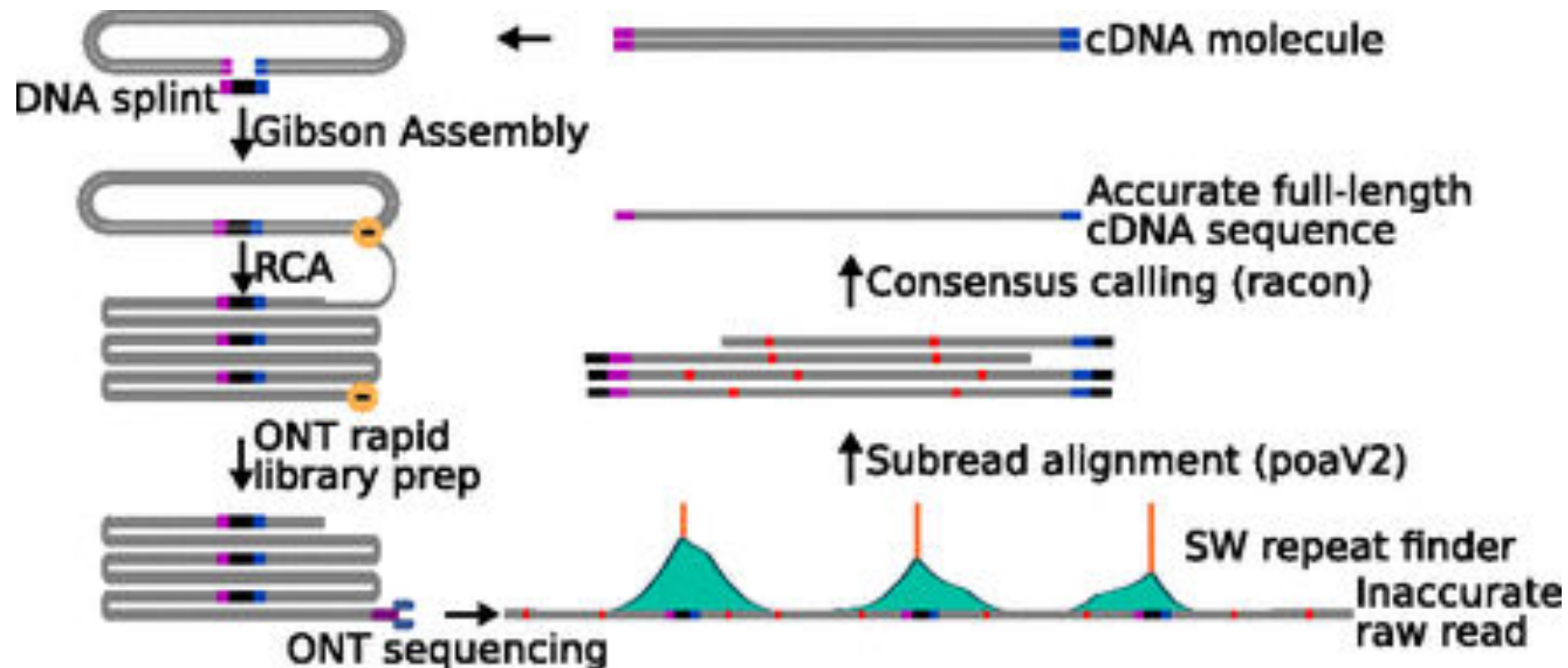
Lu, H., *et al.* (2016). Oxford Nanopore MinION Sequencing and Genome Assembly. *Genomics, Proteomics & Bioinformatics* 14, 265-279

O'Donnell, C.R., *et al.* (2013). Error analysis of idealized nanopore sequencing. *Electrophoresis* 34, 2137-2144



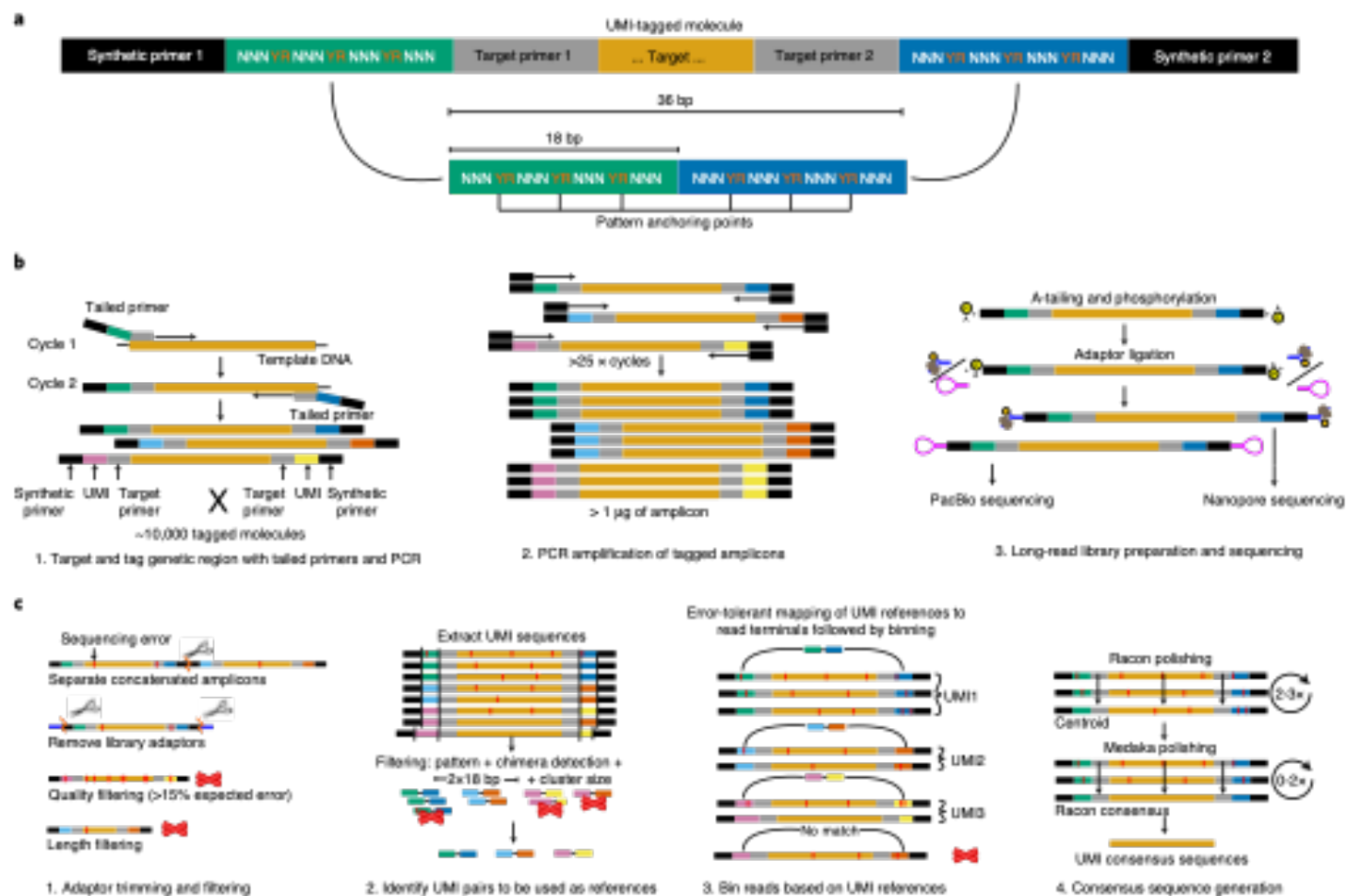
3rd generation: Nanopore

- Dramatic error rate reduction approach Rolling Circle Amplification (RCA)... similarly to PacBio





3rd generation: high quality amplicon sequencing with Nanopore/PacBio



The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

- 1st generation sequencing (single DNA fragment per reaction, ~700-1000 bp long)
 - Dideoxy termination method (**Sanger**)
- 2nd generation sequencing (massively parallel, max of ~300-600 bp total)
 - Sequencing by ligation (e.g. ABI-SOLiD)
 - Sequencing by synthesis (e.g. *Pyrosequencing*, **Illumina**, *Ion-torrent*)
- 3rd generation sequencing (massively parallel and very long reads, max of ~10kbp-250kbp)
 - Sequencing by synthesis (e.g. single molecule real time – **PacBio, Nanopore**)
 - Synthetic long reads (e.g. Moleculo, 10X... e.g. see sFL16S)



The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

- 1st generation sequencing (single DNA fragment per reaction, ~700-1000 bp long)
 - Dideoxy termination method (**Sanger**)
- 2nd generation sequencing (massively parallel, max of ~300-600 bp total)
 - Sequencing by ligation (e.g. ABI-SOLiD)
 - Sequencing by synthesis (e.g. *pyrosequencing*, **Illumina**, *Ion-torrent*)
- 3rd generation sequencing (massively parallel and very long reads, max of ~10kbp-250kbp)
 - Sequencing by synthesis (e.g. single molecule real time – **PacBio, Nanopore**)
 - Synthetic long reads (e.g. Molecule, 10X... e.g. see sFL16S)

Which method to select?



The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

Which method to select?

- Amplicon size (150-450bp most established cases; complete marker genes – 1000-1500bp – are always desirable): 1st, 2nd, 3rd generation
- 10s of thousands of reads per sample required: 2nd generation is more cost-effective
- Paired reads desirable: Illumina (2x250bp, 2x300bp)
- Lowest error rates desirable: Illumina, synthetic, highQ-long

The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

Which method to select?

- Amplicon size (150-450bp most established cases; complete marker genes – 1000-1500bp – are always desirable): 1st, 2nd, 3rd generation
- 10s of thousands of reads per sample required: 2nd generation is more cost-effective
- Paired reads desirable: Illumina (2x250bp, 2x300bp)
- Lowest error rates desirable: Illumina, synthetic, highQ-long



The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

Which method to select?

- Amplicon size (150-450bp most established cases; complete marker genes – 1000-1500bp – are always desirable): 1st, 2nd, 3rd generation
- 10s of thousands of reads per sample required: 2nd generation is more cost-effective
- Paired reads desirable: Illumina (2x250bp, 2x300bp)

- Lowest error rates desirable: Illumina, synthetic, highQ-long

Read 1

GGGCTCTCTGGATTAGATACCC TGGTAGAAATTGAAGAGTTTGATCATGGCTCAGATTGAACGCTCGGCAGGCCTAACACATGCAAGTCGAACGGAACAGGAAGAAGCTTGCTGACGGGTGAGTAATGTCTGGGAAACTGCCGTGATACAGGAAGAAGCTTAAACTYAAARRAATTGACGGC
CCCGAGAGACTAATCTATGGGACCATCTTAACTTCTCAAACAGTACCGAGTCTAACTTGCGAGCCGTCGGATTGTGTACGTTACGTTGCTTCTTCGAACGACTGCCCACTCATTACAGACCCTTGACGGACTATGCTCTTCTCGAATTGARTTTYYTTAACTGCC

The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

Which method to select?

- Amplicon size (150-450bp most established cases; complete marker genes – 1000-1500bp – are always desirable): 1st, 2nd, 3rd generation
- 10s of thousands of reads per sample required: 2nd generation is more cost-effective
- Paired reads desirable: Illumina (2x250bp, 2x300bp)
- Lowest error rates desirable: Illumina, synthetic, highQ-long



The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

Which method to select?

- Amplicon size (150-450bp most established cases; complete marker genes – 1000-1500bp – are always desirable): 1st, 2nd, 3rd generation (3rd generation is advantageous here)
- 10s of thousands of reads per sample required: 2nd generation is more cost-effective still
- Paired reads desirable: Illumina (2x250bp, 2x300bp)
- Lowest error rates desirable: Illumina, synthetic, highQ-long



Main traits of sequencing technologies

Three generations of sequencing (read size, analyzed fragment numbers)

- 1st: single fragment per reaction, ~700-1000 bp long (Sanger – dideoxy termination method)
- 2nd: millions of fragments per reaction up to 2x300 bp reads (Illumina – sequencing by synthesis)
- 3rd: 10's-100's of thousands of reads of up to 250kbp long (PacBio, Nanopore; reports for reads of [4Mbp for Nanopore](#))

Sequencing error types and positional bias per technology

- Sanger: low quality read beginnings/ends
- Illumina: substitutions, error prone read ends (phasing)
- PacBio, Nanopore: indels, randomly distributed (high error rates which can be corrected with improved base-calling algorithms, CCS)

Current main applications (practically every method is used for several tasks but mostly for...)

- Sanger: single locus analysis
- Illumina (extremely high throughput): genomics, **amplicon diversity**, ChIPseq, RNAseq...
- PacBio, Nanopore: genomics, meta-gnomics, RNAseq

The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

Which method to select?

Illumina? (low-cost/high-quality/high-depth)

synthetic? (high-cost/high-quality/high-depth)

PacBio/Nanopore? (high-cost/intermediate-to-high-quality/low-depth)

The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

How much can the read length of a single marker gene affect our study outcome?

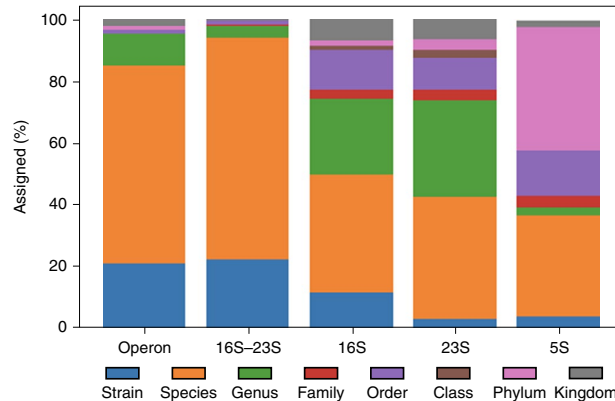
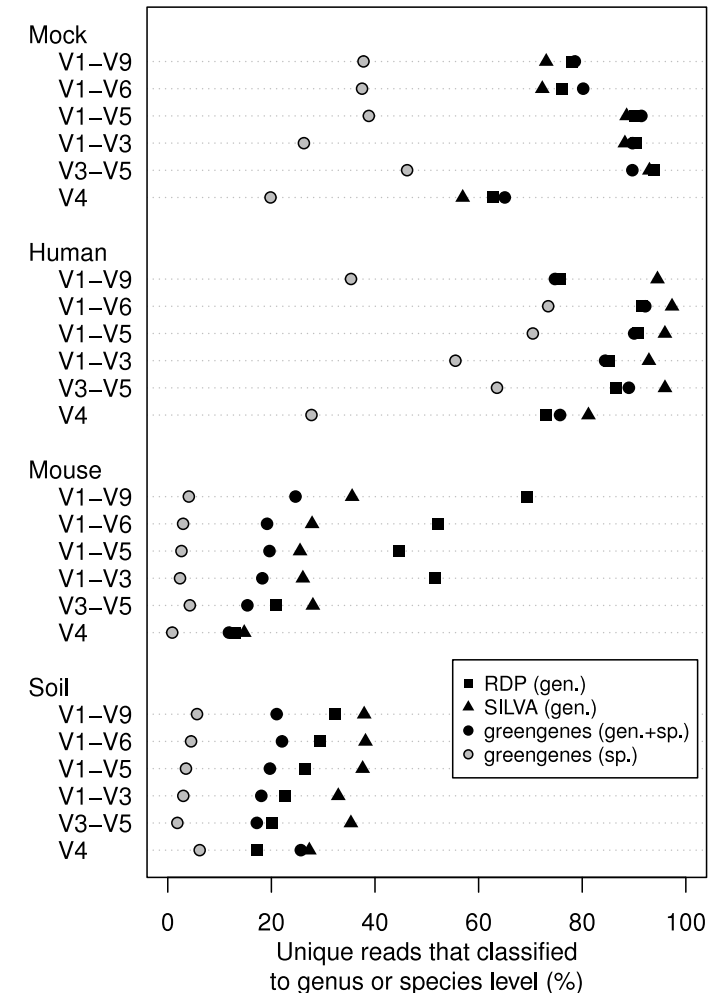
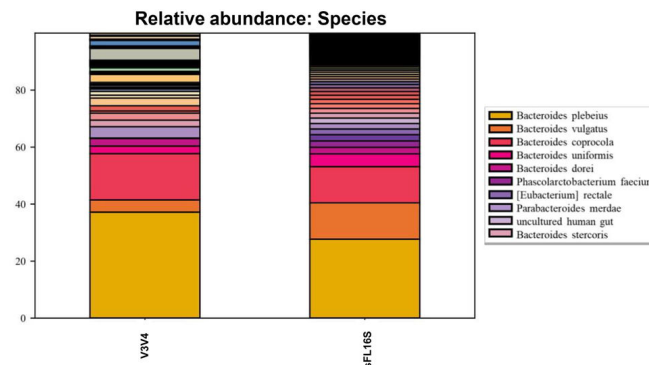
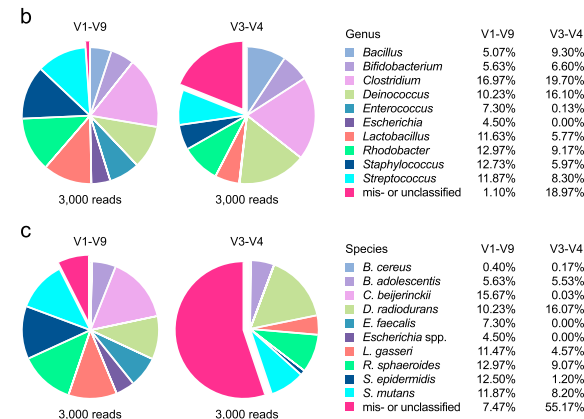


Fig. 3 | BLAST-based consensus taxonomic assignment against the Web of Life database for whole rRNA operons, using the combination of 16S and 23S rRNAs and individual rRNA genes. In the dataset, 253,089 operons were available and used for assignment. Of these, $n = 253,087$ had an annotatable 23S rRNA gene, $n = 253,088$ had an annotatable 16S rRNA gene and $n = 50,560$ had an annotatable 5S rRNA gene.



Karst, et al 2021. Nat Methods <https://doi.org/10.1038/s41592-020-01041-y>

Matsuo, et al 2021, BMC Microbiol, <https://doi.org/10.1186/s12866-021-02094-5>

Jeong et al 2021, Scientific Reports <https://doi.org/10.1038/s41598-020-80826-9>

Schloss et al 2016, Peer J, <https://doi.org/10.7717/peerj.1869>

Johnson et al 2019, Nat Commun <https://doi.org/10.1038/s41467-019-13036-1>



The 3 generations of sequencing technologies (size vs analyzed sequence variant numbers)

How much can the read length of a single marker gene affect our study outcome?

- Prokaryotes (genomic plasticity)?
- Environment?
- Reference database?
- Eukaryotes (makes sense)?

Can we resort to cost-effective hybrid solutions? E.g.:

- Generate habitat/experiment specific (pooled DNA template) long-read reference databases for our short-read sample data (AutoTax)

