

Bring Flow Cell

Housekeeping

- Me
- Names
- Apology
- Test Wifi and IBIEM Computing Environment

Morning Bootcamp Topics

- Microbiome Analysis Overview
- Computing Environment: RStudio
- Reproducible Research
 - Git
- Introduction to R

Microbiome Analysis Overview

Josh Granek
August 8, 2018

Metagenomics

	What	Information	Analogy	Target Size	Cost
Amplicon	Marker Gene	Who is Present	Name	100bp - 1kb	Low
Shotgun Metagenome	Genomes	What Genes are Present	CV	100kb - 100Mb	High
Shotgun Metatranscriptome	All RNA	What Genes are Expressed	Twitter Feed	100kb - 100Mb	High

granek JA - PubMed - NCP

www.ncbi.nlm.nih.gov/pubmed/?term=granek+JA

NCBI Resources How To Sign In to NCBI

PubMed granek JA Search Help

Article types Summary - 20 per page - Sort by Most Recent - Send to: Filters: Manage Filters

Clinical Trial

Review

Customize ...

Text availability

Abstract

Free full text

Full text

PubMed Commons

Reader comments

Trending articles

Publication dates

6 years

10 years

Custom range...

Species

Humans

Other Animals

[Clear all](#)

[Show additional filters](#)

Search results

Items: 17

Evidence for distinct brain networks in the control of rule-based motor behavior.
1. Granek JA, Sergio LE.
J Neurophysiol. 2015 Aug;114(2):1298-309. doi: 10.1152/jn.00230.2014. Epub 2015 Jul 1.
PMID: 26133796
[Similar articles](#)

Rapid mapping of insertional mutations to probe cell wall regulation in *Cryptococcus neoformans*.
2. Eshet SK, Granek JA, Alspaugh JA.
Fungal Genet Biol. 2015 Dec;82:5-21. doi: 10.1016/j.fgb.2015.08.013. Epub 2015 Jun 23.
PMID: 26112690
[Similar articles](#)

Integrating chemical mutagenesis and whole-genome sequencing as a platform for forward and reverse genetic analysis of Chlamydia.
3. Kokes M, Dunn JD, Granek JA, Nguyen BD, Barker JR, Valdivie RH, Bestides RJ.
Cell Host Microbe. 2015 May 13;17(5):716-25. doi: 10.1016/j.chom.2015.03.014. Epub 2015 Apr 23.
PMID: 25920970 Free PMC Article
[Similar articles](#)

Antifungal drug resistance evoked via RNAi-dependent epimutations.
4. Calo S, Shantz-Wall C, Lee SC, Bestides RJ, Nicolas FE, Granek JA, Mieczkowski P, Torres-Martinez S, Ruiz-Vazquez RM, Cardenas ME, Heltman J.
Nature. 2014 Sep 26;513(7518):555-6. doi: 10.1038/nature13575. Epub 2014 Jul 27.
PMID: 25079329 Free PMC Article
[Similar articles](#)

Decoupled visually-guided reaching in optic ataxia: differences in motor control between canonical and non-canonical orientations in space.
5. Granek JA, Pisella L, Stemberger J, Vighetto A, Rossetti Y, Sergio LE.
PLoS One. 2013 Dec 31;8(12):e80136. doi: 10.1371/journal.pone.0080136. eCollection 2013.
PMID: 24302036 Free PMC Article
[Similar articles](#)

The genetic architecture of biofilm formation in a clinical isolate of *Saccharomyces cerevisiae*.
6. Granek JA, Murray D, Kayikçi Ö, Magwene PM.
Genetics. 2013 Feb;193(2):587-800. doi: 10.1534/genetics.112.142087. Epub 2012 Nov 19.
PMID: 23172850 Free PMC Article
[Similar articles](#)

The role of the caudal superior parietal lobe in updating hand location in peripheral vision: further evidence from optic ataxia.
7. Granek JA, Pisella L, Blangsted A, Rossetti Y, Sergio LE.
PLoS One. 2012 Jul 10:e48619. doi: 10.1371/journal.pone.0048619. Epub 2012 Jul 6.
PMID: 23071689 Free PMC Article
[Similar articles](#)

Pleiotropic signalling pathways orchestrate yeast development.
8. Granek JA, Kayikçi Ö, Magwene PM.
Curr Opin Microbiol. 2011 Dec;14(6):676-81. doi: 10.1016/j.mib.2011.08.004. Epub 2011 Sep 28. Review.
PMID: 21962291 Free PMC Article
[Similar articles](#)

Find related data Database: PubMed Find items

Search details granek JA [Author]

Search See more...

Recent Activity Turn Off Clear

granek JA (17) PubMed

granek J (20) PubMed

Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample PubMed

Scott Harrison (71) PubMed

See more...

Amplicon Sequencing

PCR amplify and sequence a marker gene

	Marker Gene
Bacteria	16s rRNA
Fungi	18s or ITS rRNA
Archaea	16s rRNA
Protozoa	18s rRNA
Viruses	?????

Metagenomics

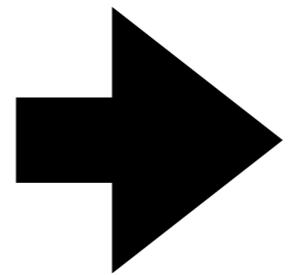
	What	Information	Analogy	Target Size	Cost
Amplicon	Marker Gene	Who is Present	Name	100bp - 1kb	Low
Shotgun Metagenome	Genomes	What Genes are Present	CV	100kb - 100Mb	High
Shotgun Metatranscriptome	All RNA	What Genes are Expressed	Twitter Feed	100kb - 100Mb	High

Metagenomics

	What	Information	Analogy	Target Size	Cost	Discovery?
Amplicon	Marker Gene	Who is Present	Name	100bp - 1kb	Low	+/-
Shotgun Metagenome	Genomes	What Genes are Present	CV	100kb - 100Mb	High	++
Shotgun Metatranscriptome	All RNA	What Genes are Expressed	Twitter Feed	100kb - 100Mb	High	++

I 6s Amplicon Analysis Overview

Big Picture



1. What is present?
2. How much?
3. Are there differences between treatments, host species, ...?
4. What are the differences?

Molecular Biology

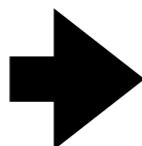


<http://www.geograph.org.uk/photo/2847164>

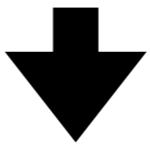
<https://commons.wikimedia.org/wiki/File:Pcr.jpg>

https://commons.wikimedia.org/wiki/File:Illumina_MiSeq_sequencer.jpg

Bioinformatic Analysis



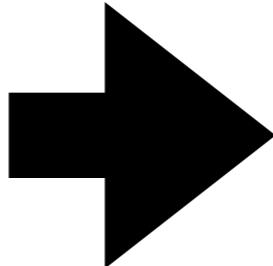
```
@M00698:36:000000000-AFBEL:1:1101:14738:1412 1:N:0:0
TTACGCTAACAGGCAGGTAGCCTGGCAGGGTCAGGAAATCAATTAACTCATCGGAAGTGGTATCTGTTCCATCAAGCGTGCAGCATCGTCAAAACGCC
+
ABBBABBBBAFFGGGGGGGGGGHGGHGGCG2GF3FFGHHHHHGGFGHEHHGGGEHHHHAGGHHGHHFFDHFHHGEGGGG@F@H?GHH/GBEFGGG
@M00698:36:000000000-AFBEL:1:1101:16483:1412 1:N:0:0
CTGCCAGTTGAACGACGGCGAGCAGTTATAAGCCAGCAGTTGCCGGATATTCGCGTGGATAGCTTGCAAAGCGACGCCAGTTCCAGATCCGGCG
+
AAABBFFFFFFFGGGGGGGGGGGHHHHHHHHGHGHGHGHHHHHGGGGHHHHGGGGGGHHHHFFHHHHHGHHGGGGGGGGHHHHHHHHGGG
@M00698:36:000000000-AFBEL:1:1101:15928:1413 1:N:0:0
GTAAAGTCCTGAGTGATACCGCAACTTTACCCCCAGTCCCACTTCGAACCCGAAACATATCGGAAAGAGGCCGTGCCTGATTTAAAGCCGTAGGT
+
```



	Sample 1	Sample 2	...	Sample N
Bacteria 1				
Bacteria 2				
...				
Bacteria N				

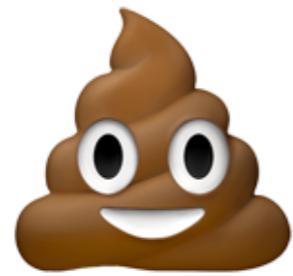
Statistical Analysis

	Sample 1	Sample 2	...	Sample N
Bacteria 1				
Bacteria 2				
...				
Bacteria N				



1. What is present?
2. How much?
3. Are there differences between treatments, host species, ...?
4. What are the differences?

Caveat



<http://www.lebleu.com>

Sequencing Details

DNA Sequencing Technologies (Abridged)

1st Generation	2nd Generation	3rd Generation
Chemical (Maxim-Gilbert)	Pyrosequencing (454)	Single molecule real time (PacBio)
Chain Termination (Sanger)	Chain Termination (Illumina)	Nanopore sequencing (Oxford Nanopore)
Pyrosequencing	Sequencing by ligation (SOLiD sequencing)	
	Ion semiconductor (Ion Torrent)	

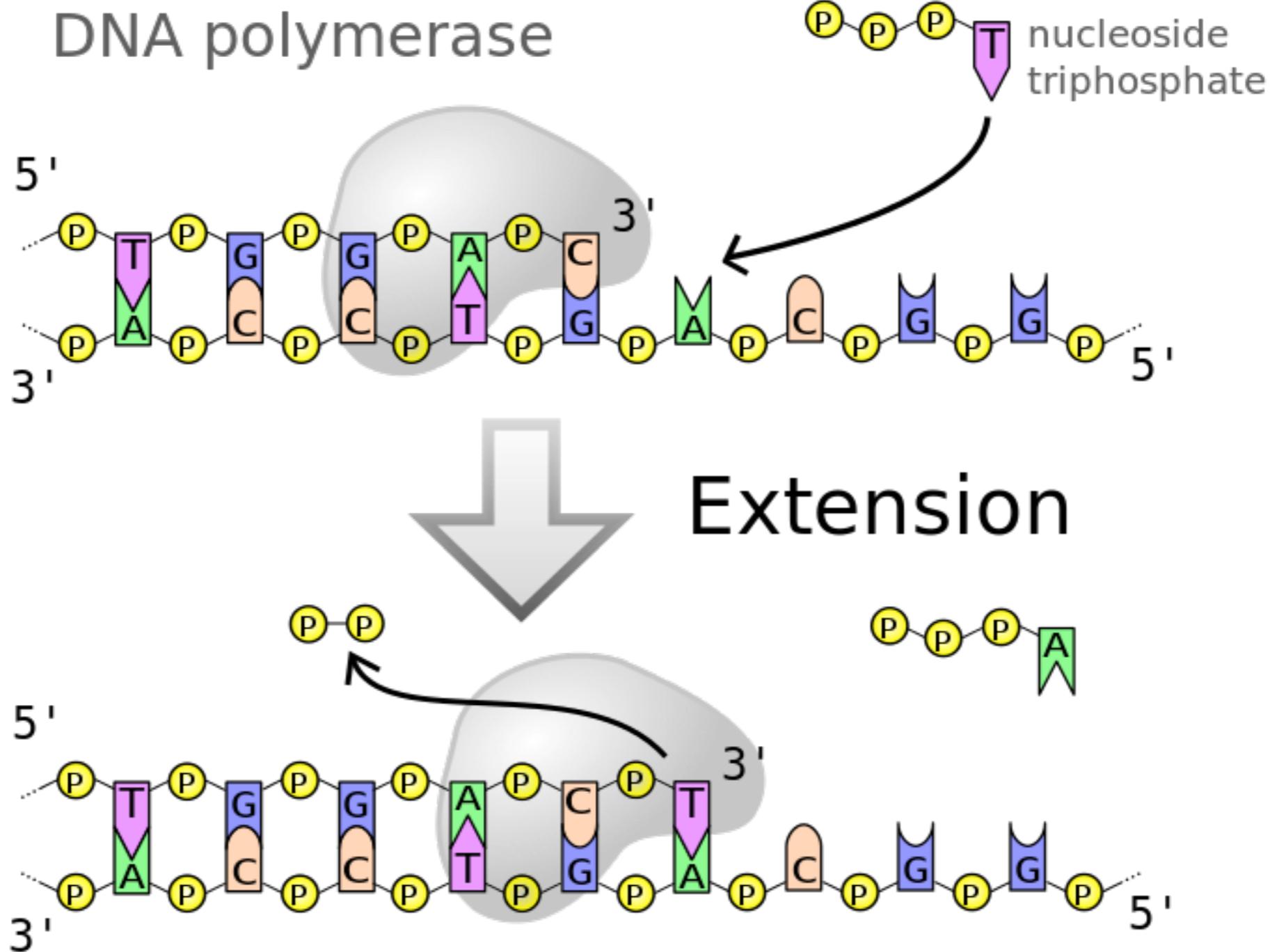
Sequencing by Synthesis

1st Generation	2nd Generation	3rd Generation
Chemical (Maxim-Gilbert)	Pyrosequencing (454)	Single molecule real time (PacBio)
Chain Termination (Sanger)	Chain Termination (Illumina)	Nanopore sequencing (Oxford Nanopore)
Pyrosequencing	Sequencing by ligation (SOLiD sequencing)	
	Ion semiconductor (Ion Torrent)	

DNA Synthesis

- What are the minimum components for DNA Replication?

DNA Synthesis

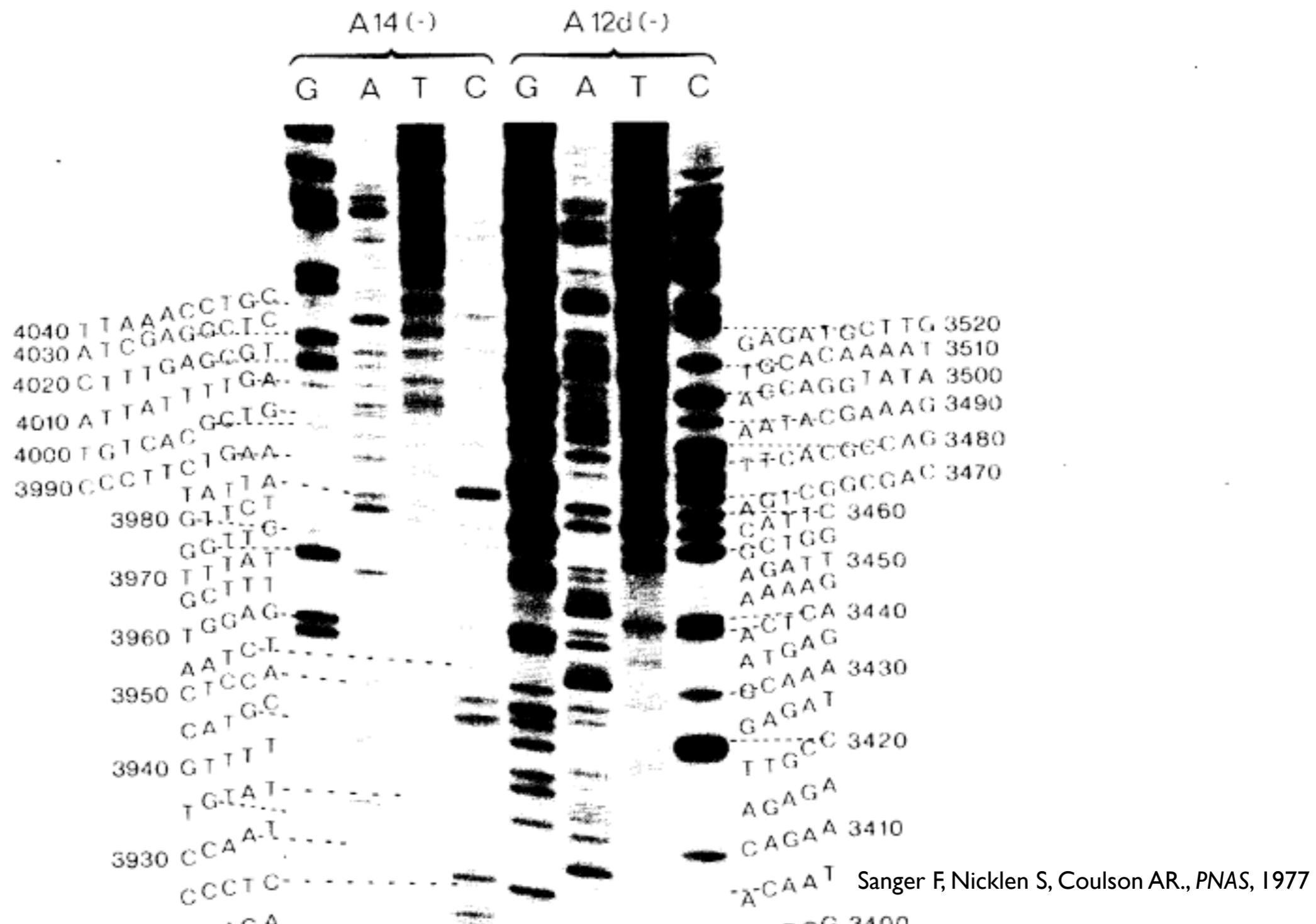


1st Generation	2nd Generation	3rd Generation
Chemical (Maxim-Gilbert)	Pyrosequencing (454)	Single molecule real time (PacBio)
Chain Termination (Sanger)	Chain Termination (Illumina)	Nanopore sequencing (Oxford Nanopore)
Pyrosequencing	Sequencing by ligation (SOLiD sequencing)	
	Ion semiconductor (Ion Torrent)	

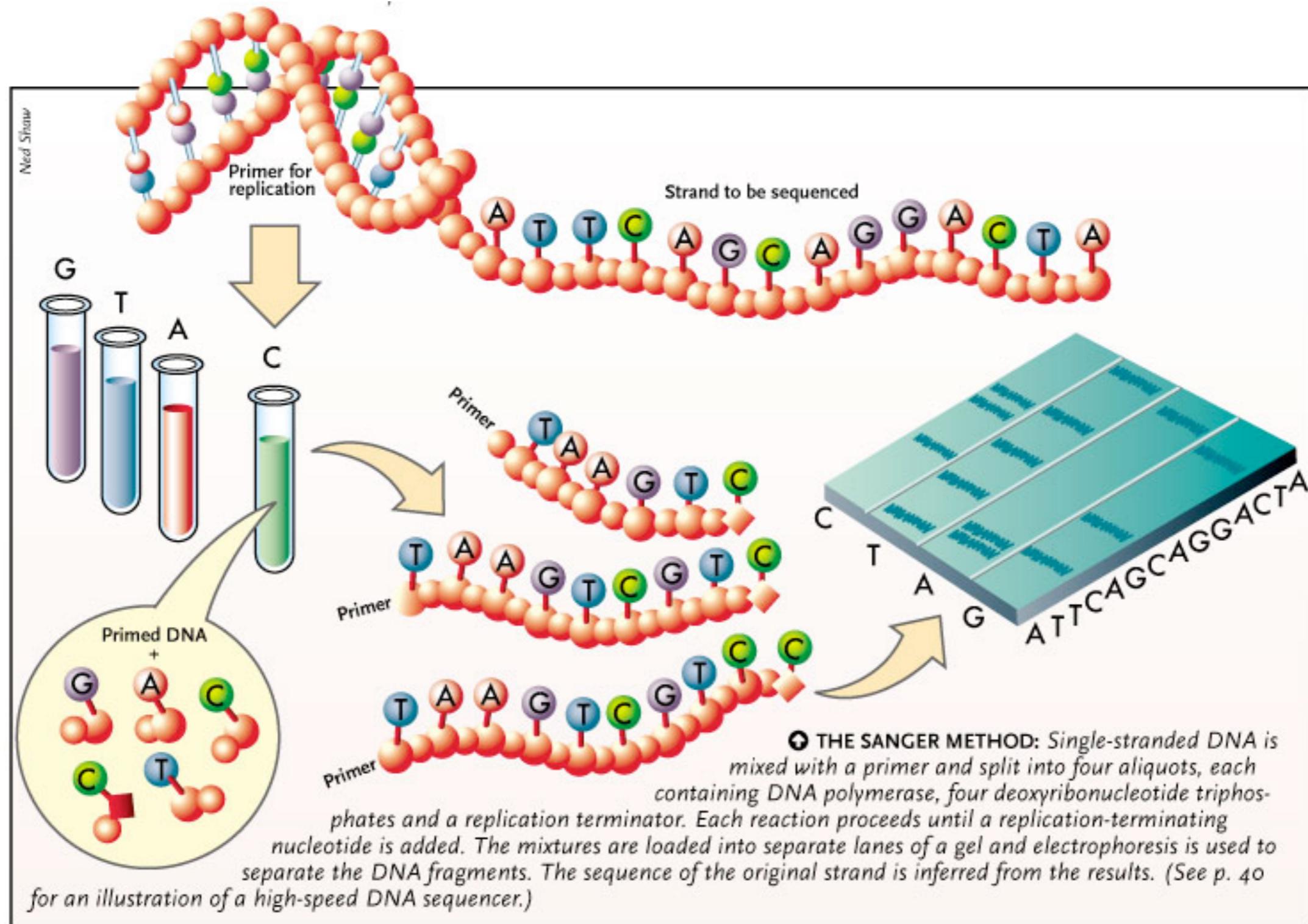
Sanger Sequencing

5464 Biochemistry: Sanger *et al.*

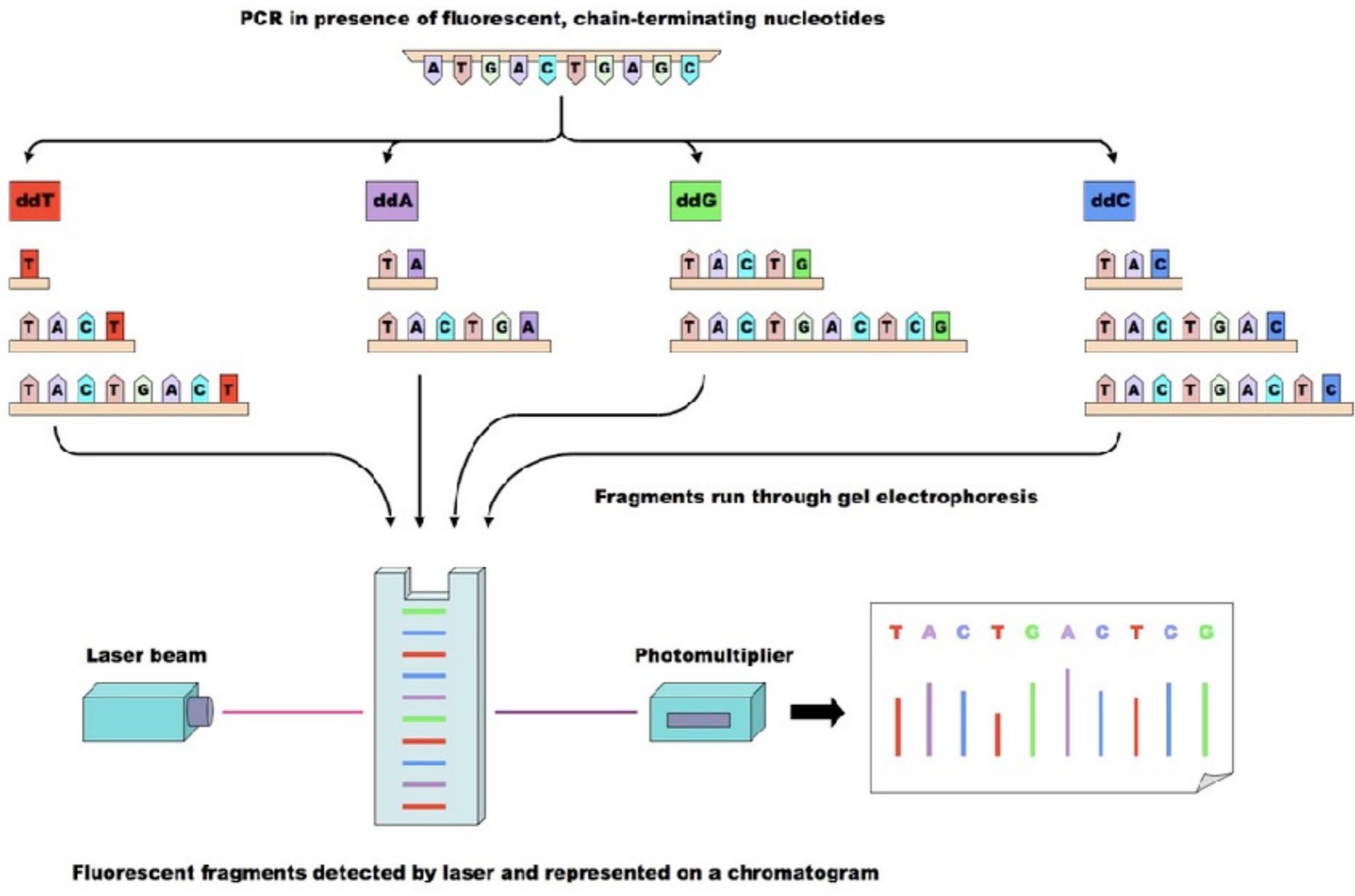
Proc. Natl. Acad. Sci. USA 74 (1977)



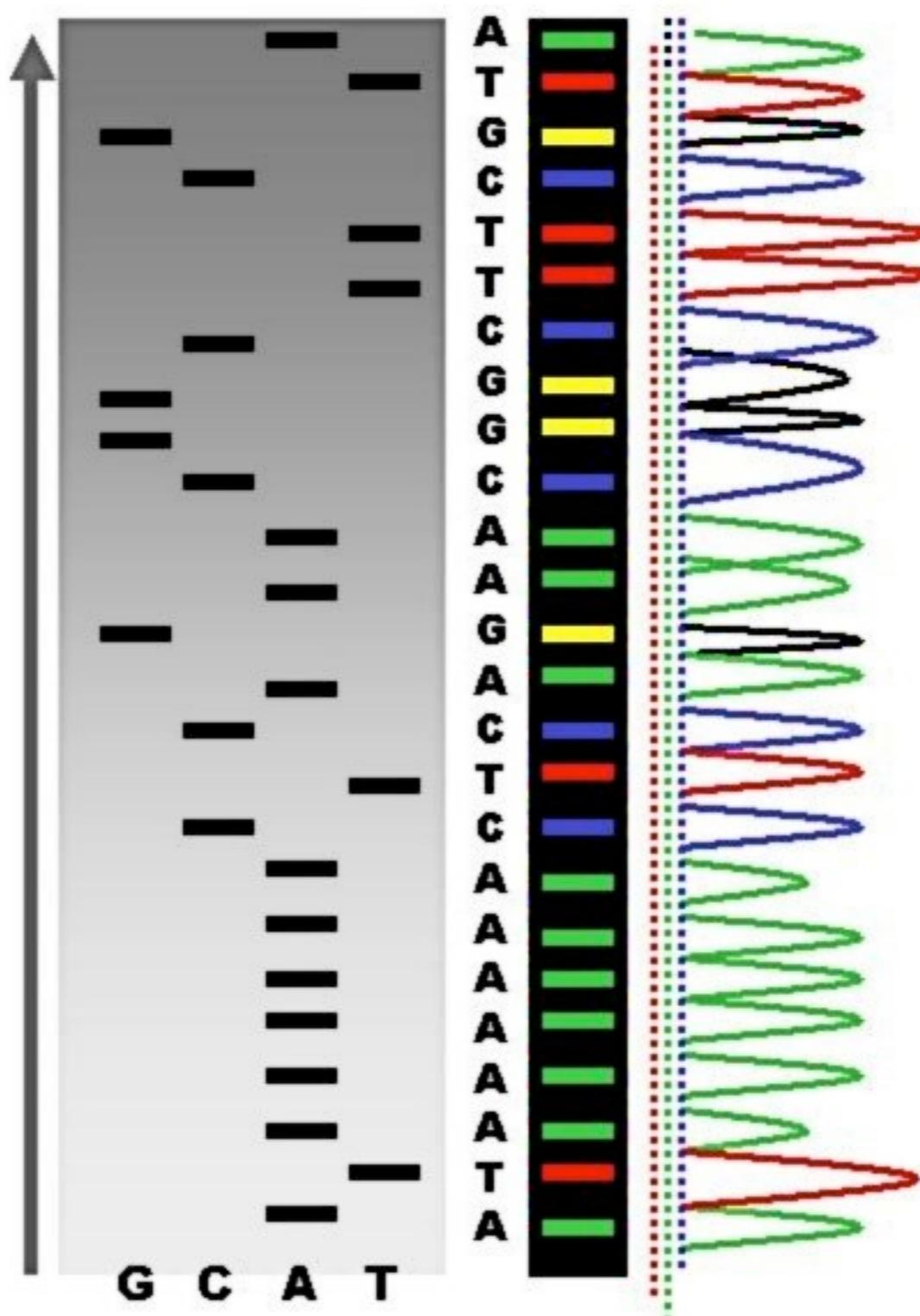
Sanger Sequencing



Dye-terminator



Radiolabel vs. Dye



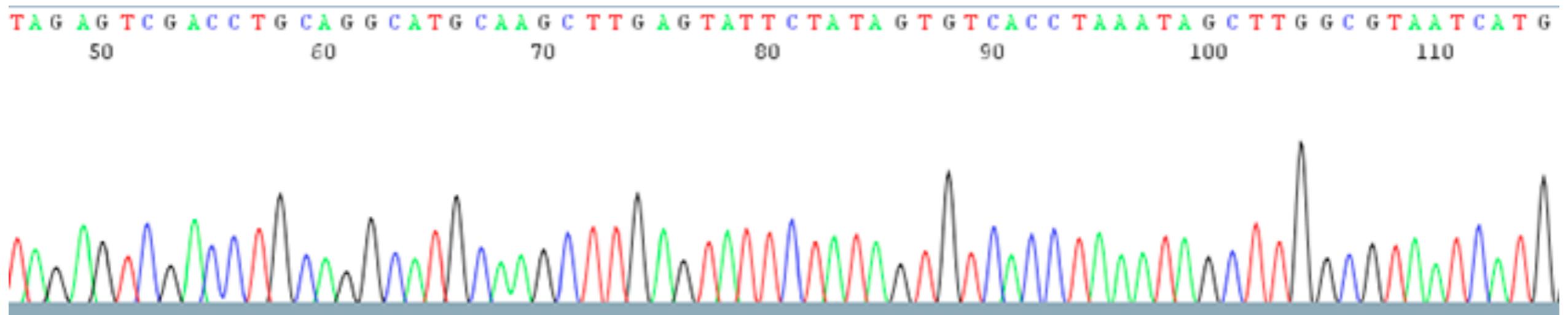
https://upload.wikimedia.org/wikipedia/commons/3/3d/Radioactive_Fluorescent_Seq.jpg

High-Throughput

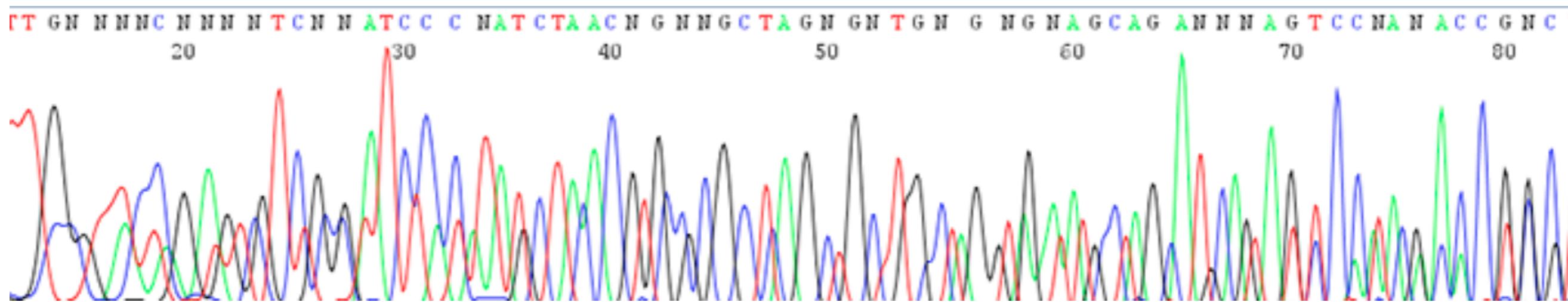
How?

- Separate
- Detect

Dye-terminator Sanger Sequencing



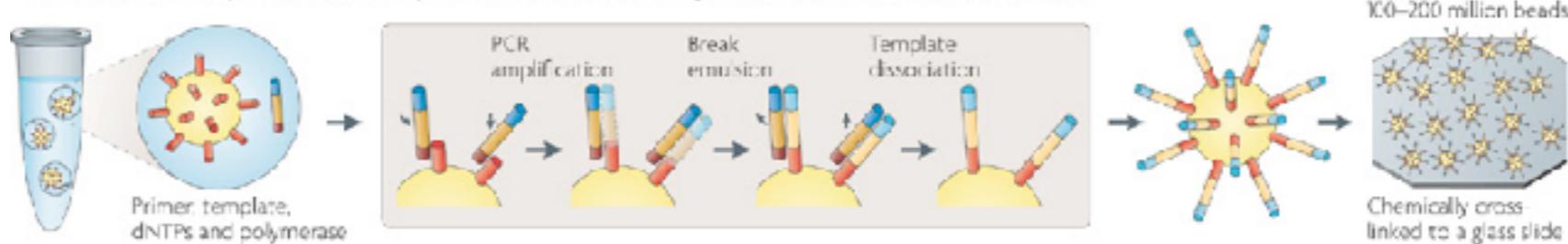
Bad Sanger Sequence



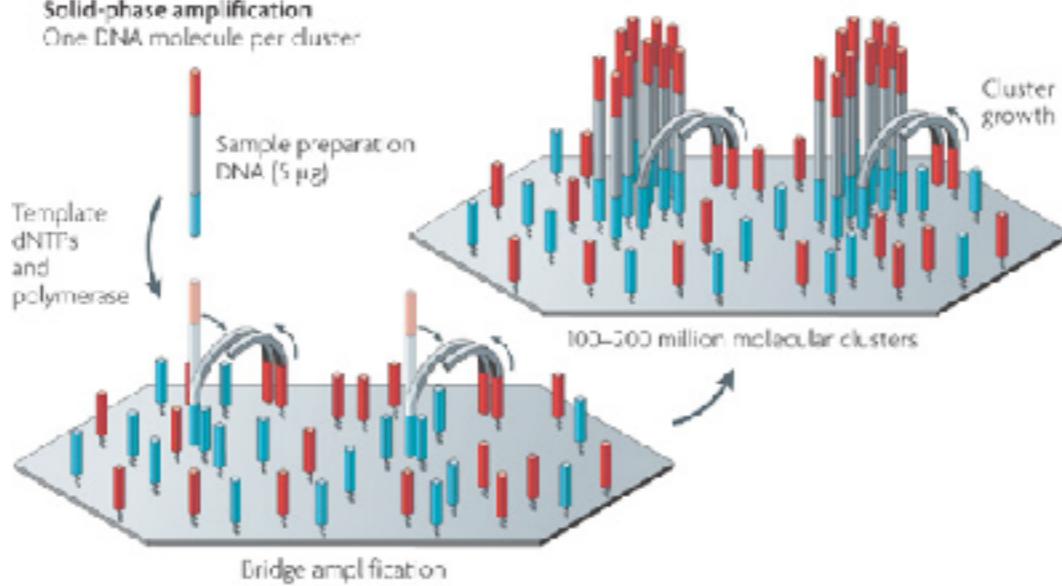
Template immobilization

a Roche/454, Life/APG, Polonator
Emulsion PCR

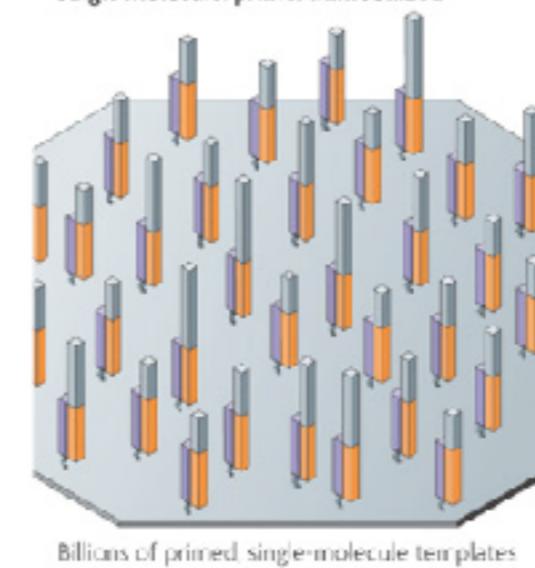
One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion



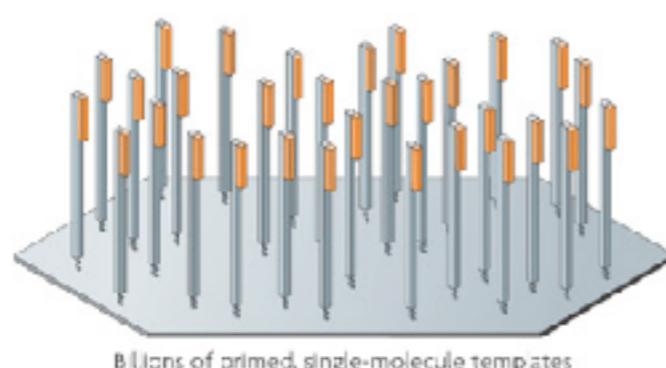
b Illumina/Solexa
Solid-phase amplification
One DNA molecule per cluster



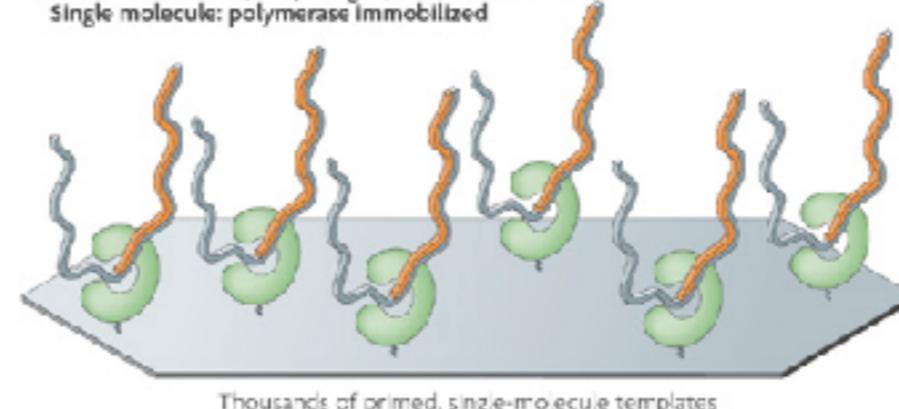
c Helicos BioSciences: one-pass sequencing
Single molecule: primer immobilized



d Helicos BioSciences: two-pass sequencing
Single molecule: template immobilized

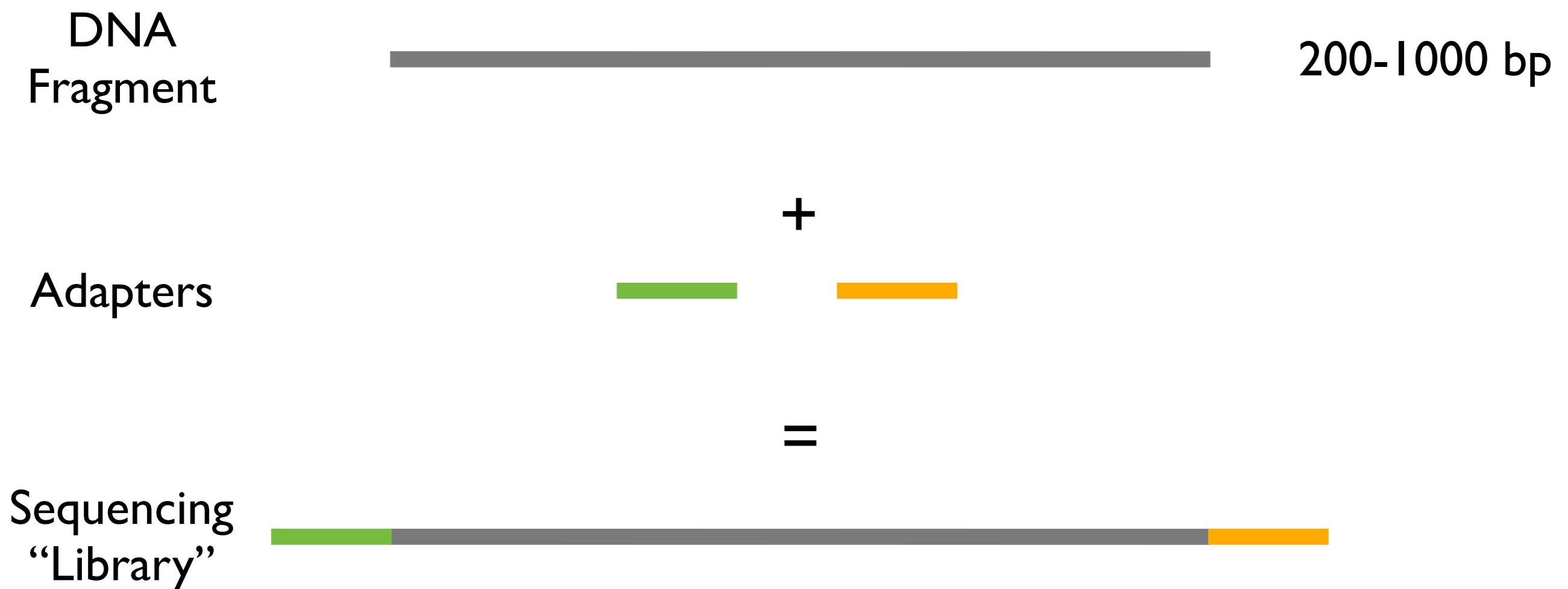


e Pacific Biosciences, Life/Visigen, LI-COR Biosciences
Single molecule: polymerase immobilized

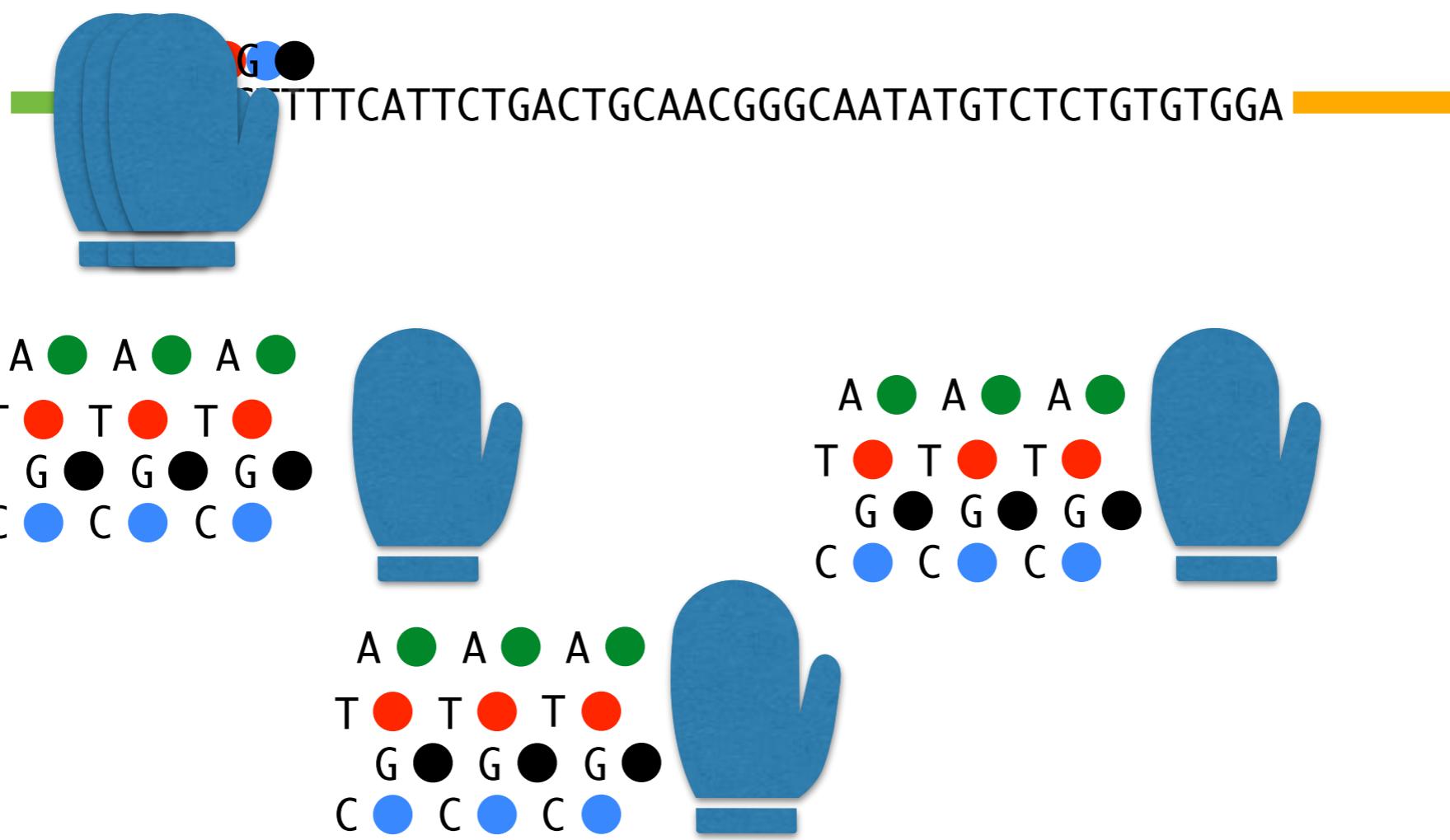


1st Generation	2nd Generation	3rd Generation
Chemical (Maxim-Gilbert)	Pyrosequencing (454)	Single molecule real time (PacBio)
Chain Termination (Sanger)	Chain Termination (Illumina)	Nanopore sequencing (Oxford Nanopore)
Pyrosequencing	Sequencing by ligation (SOLiD sequencing)	
	Ion semiconductor (Ion Torrent)	

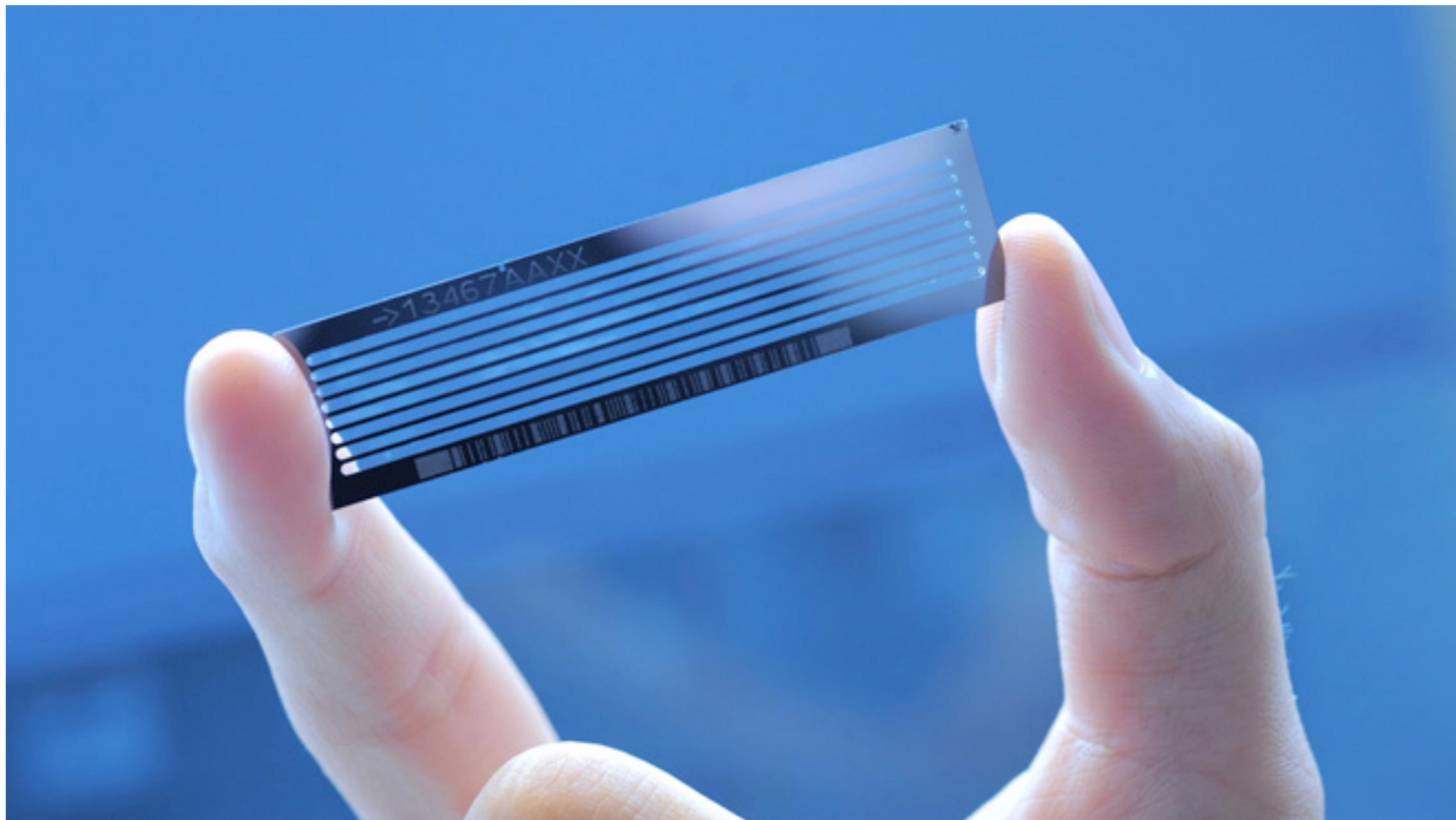
Library Preparation



Sequencing



A Flow Cell

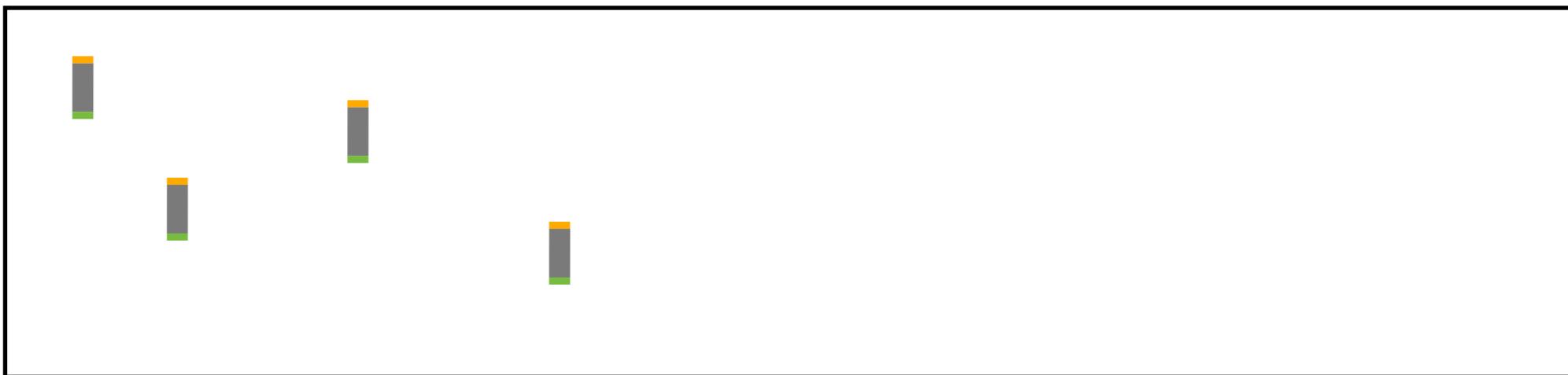


Show Flow Cell

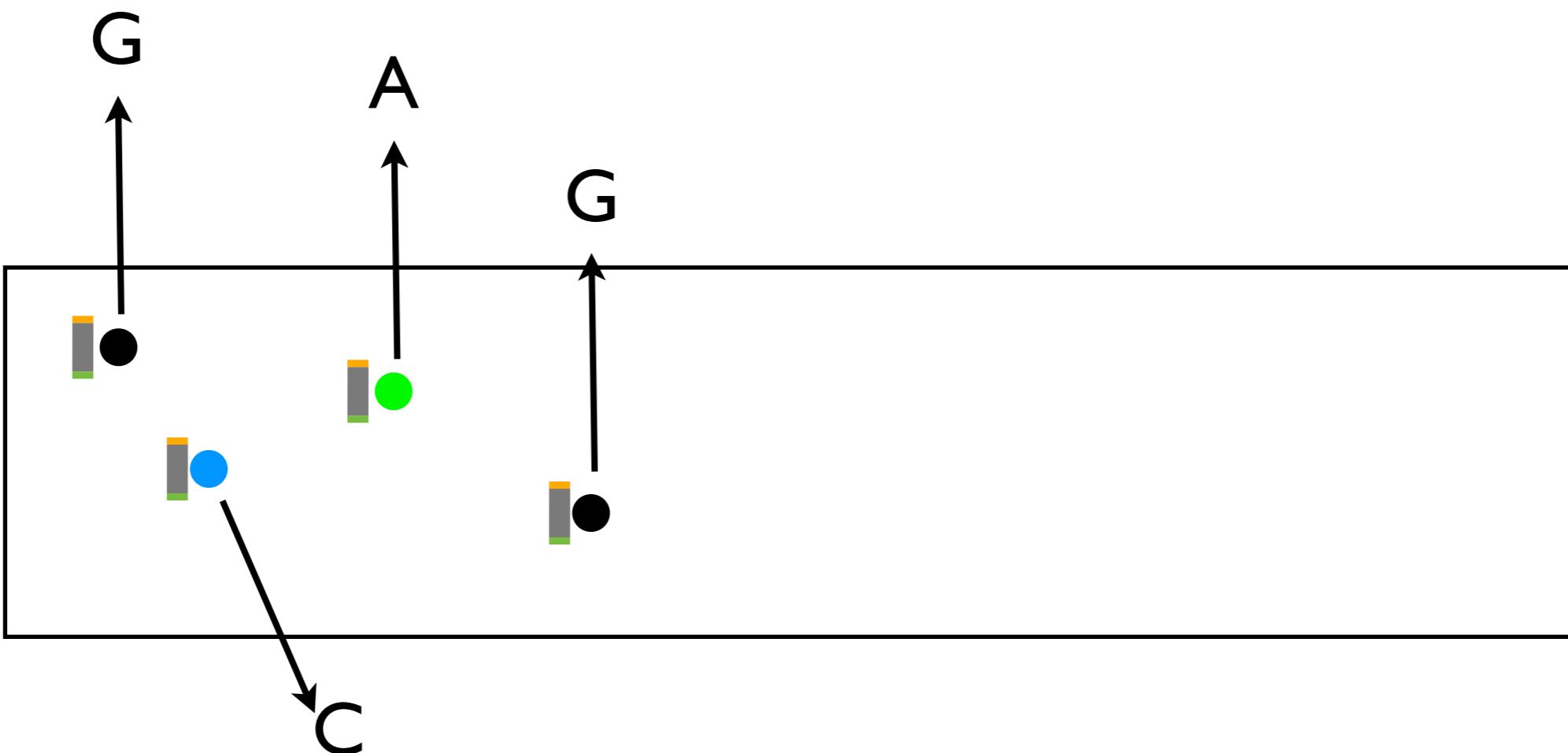
A Flow Cell



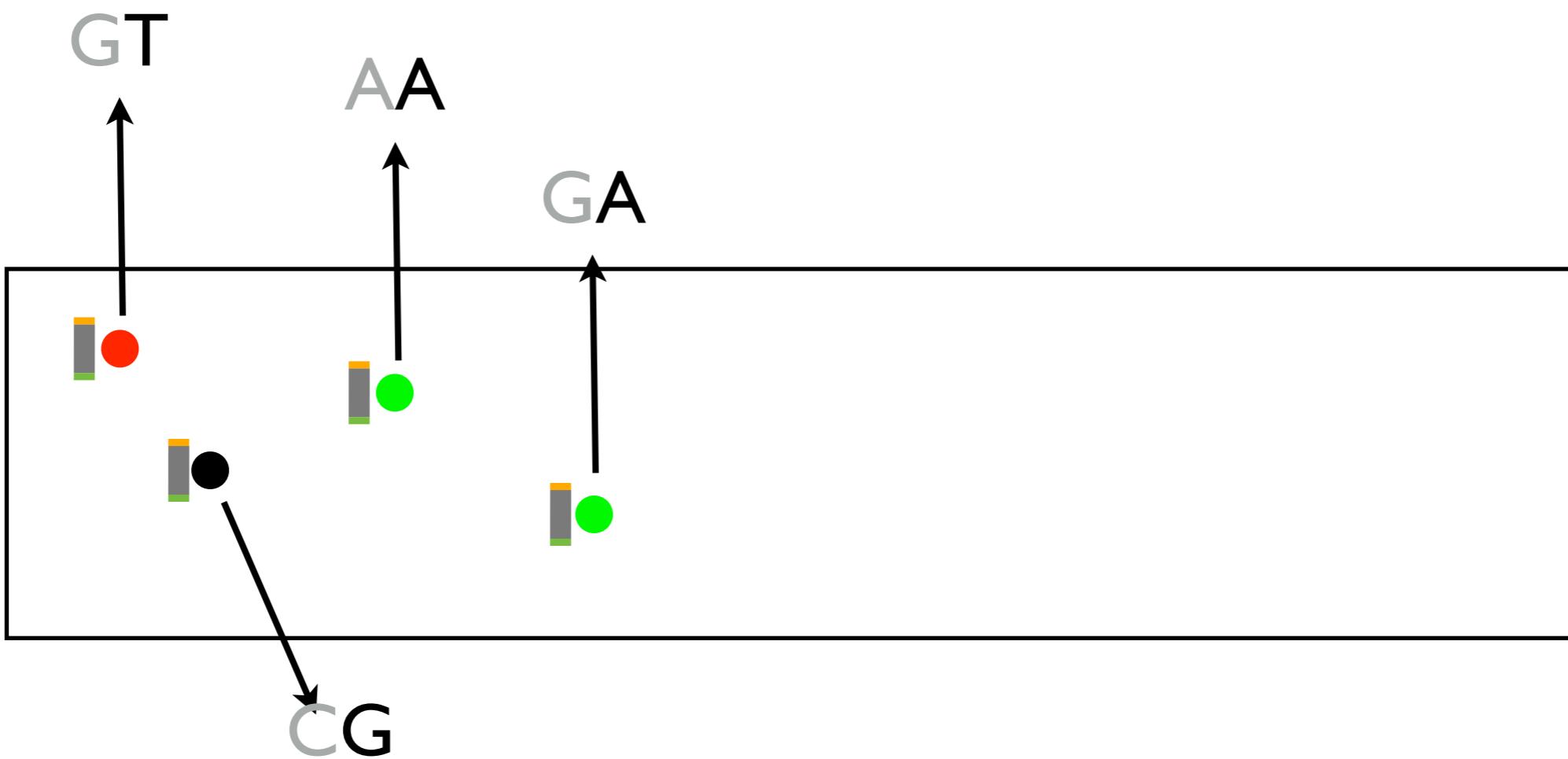
Bind Library



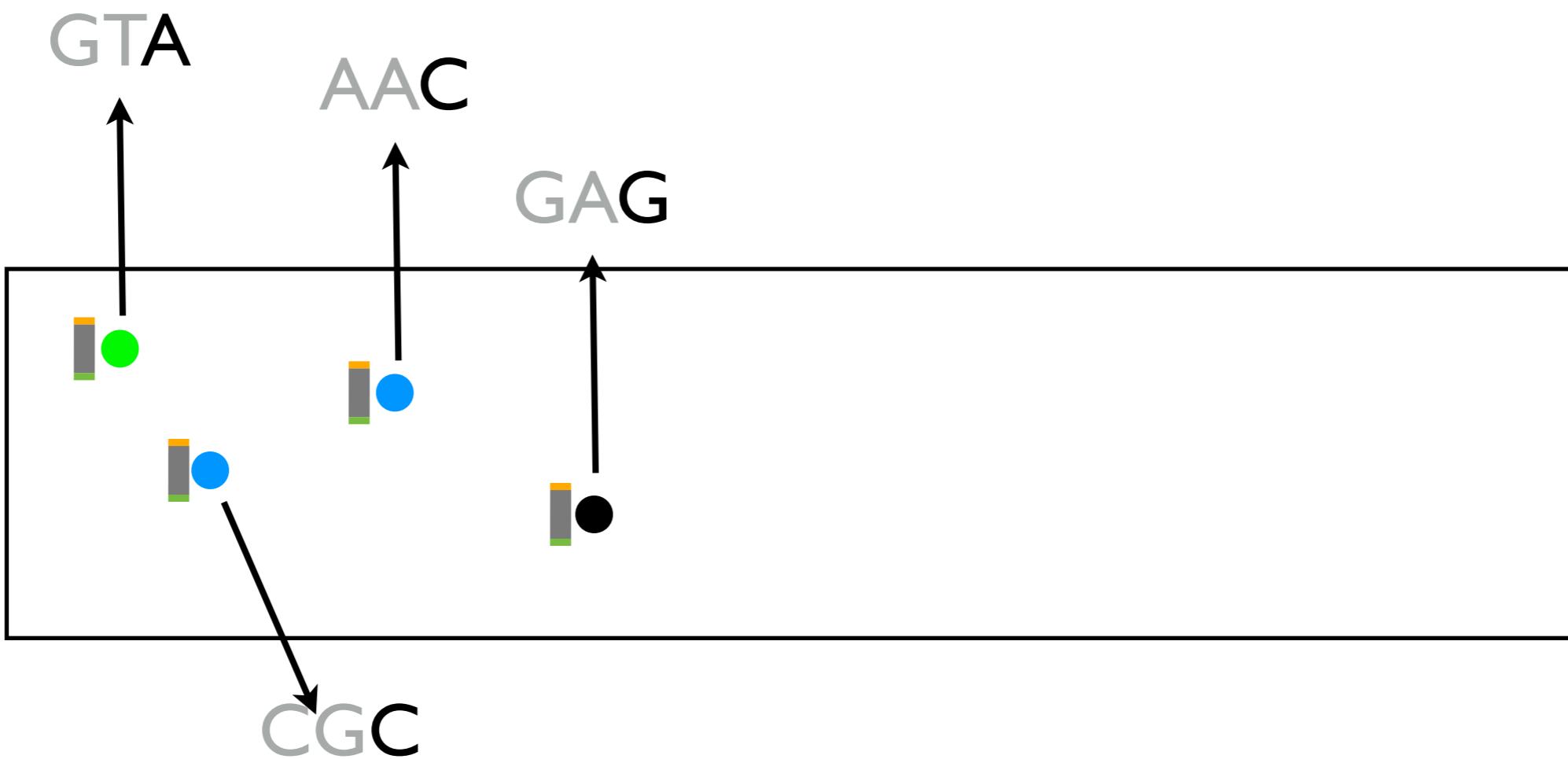
1st Cycle



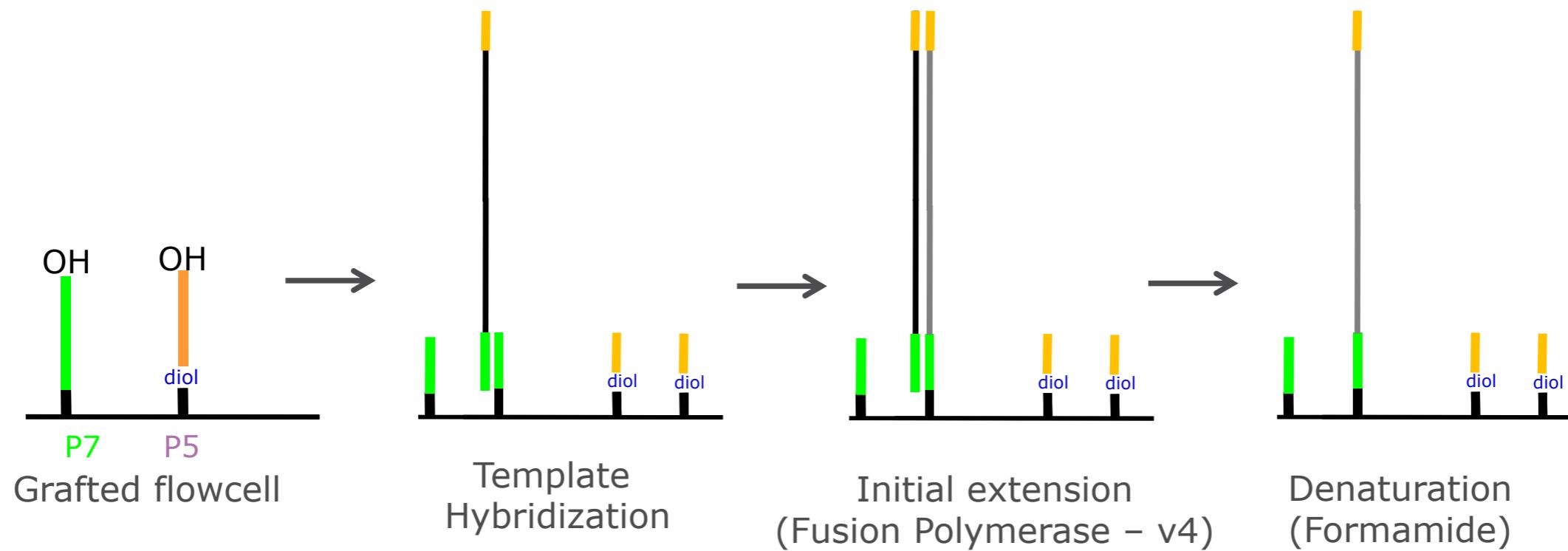
2nd Cycle



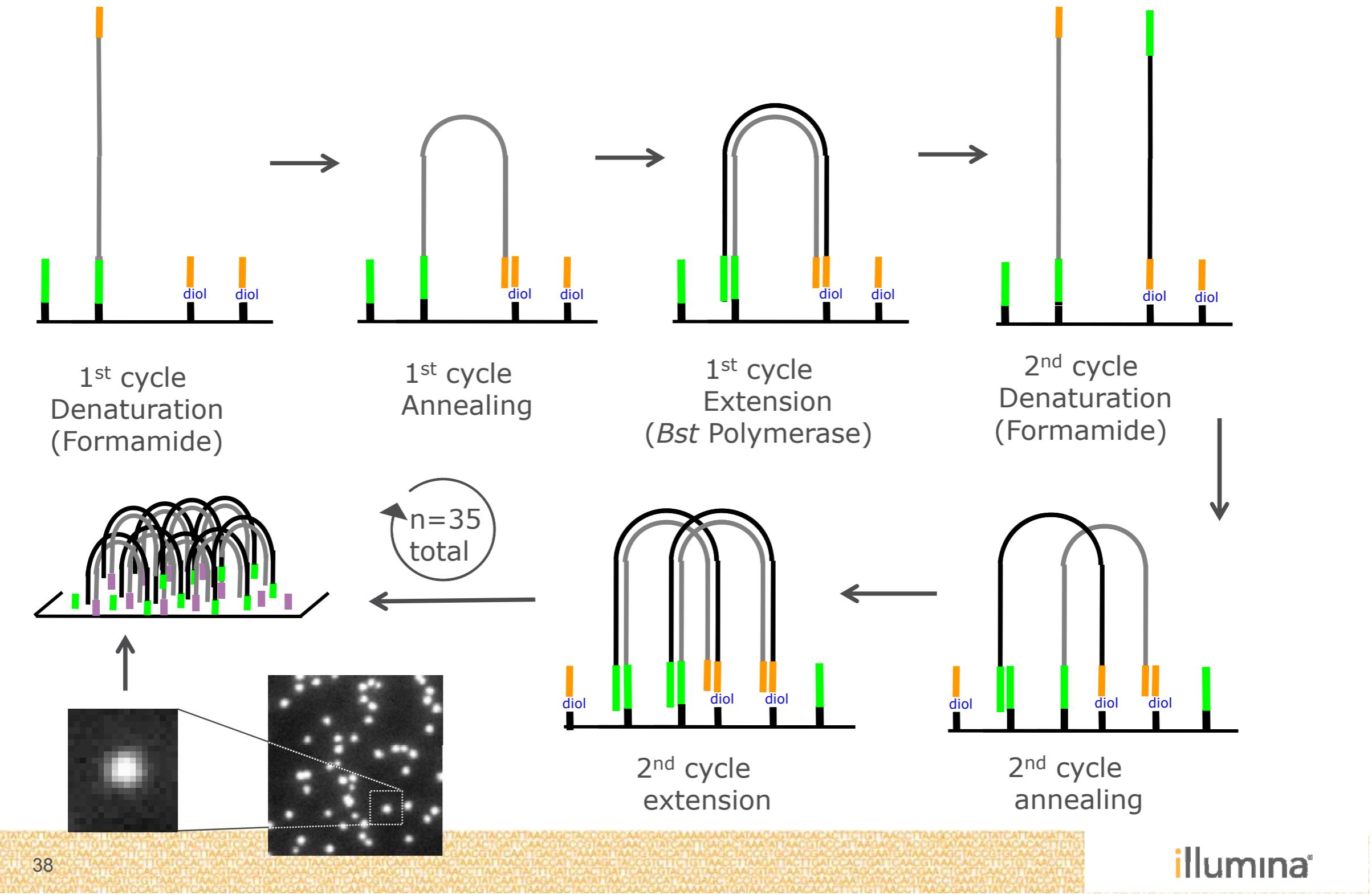
3rd Cycle



Cluster generation – hybridization and amplification



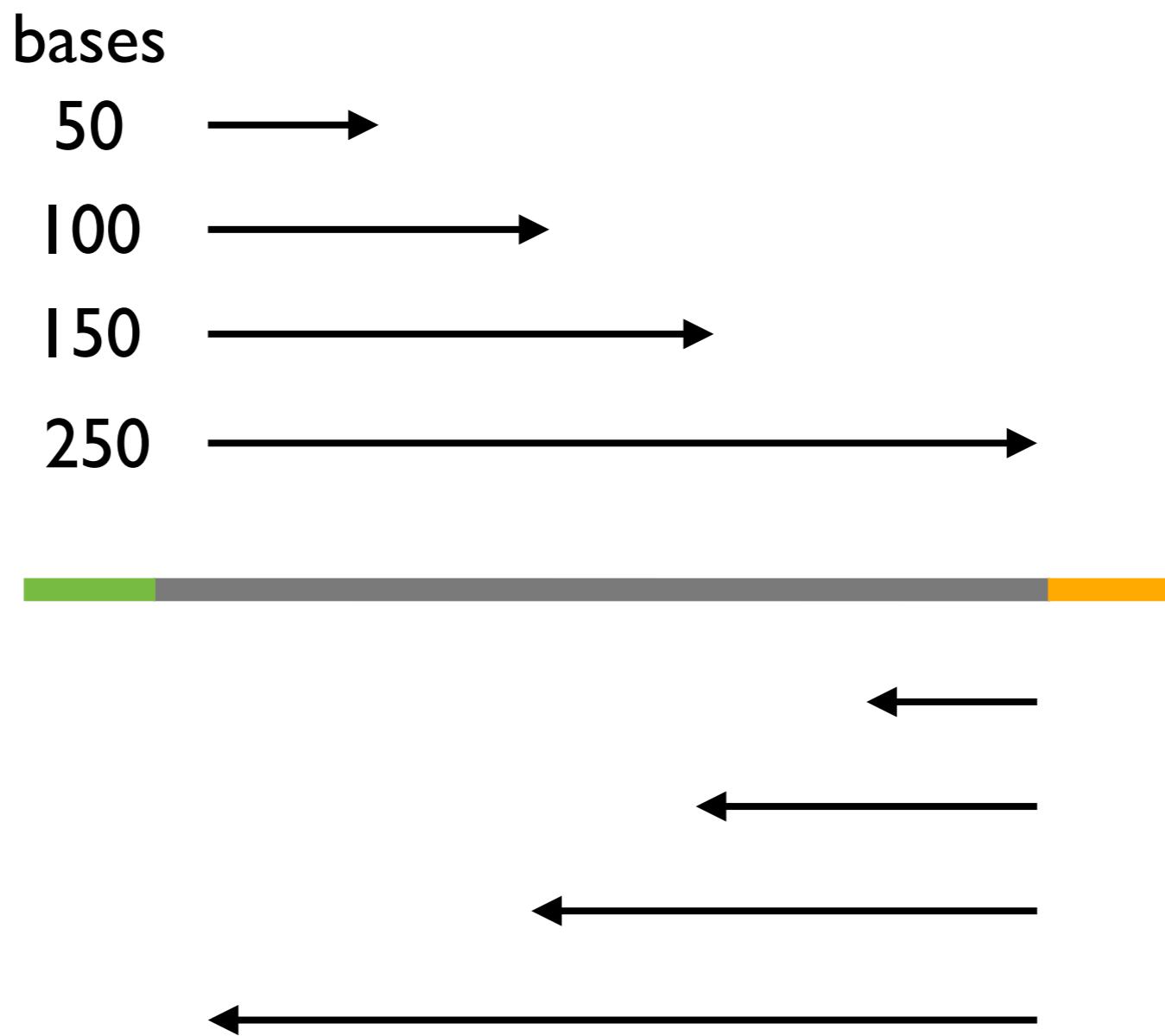
Cluster generation – hybridization and amplification



Paired End



Paired End



Paired End

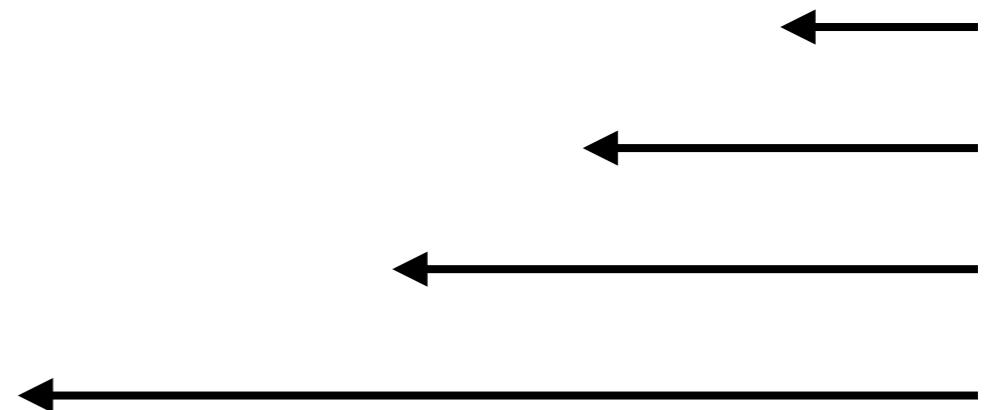
bases

50 →

100 →

150 →

250 →



Sequencing Library

Amplicon Library



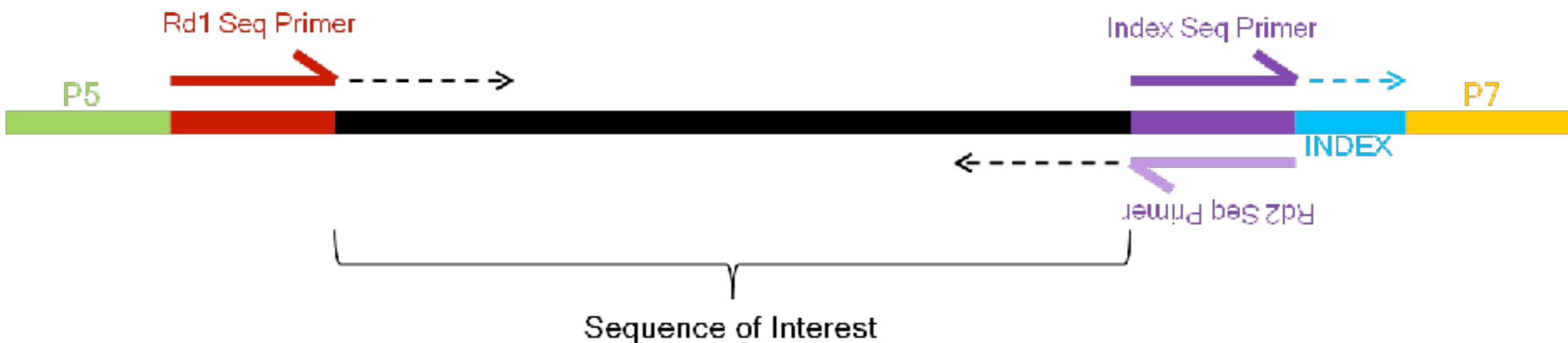
Shotgun Library



Stopped Here 8/9/2018

Multiplexing (Barcodes)

STRUCTURE DETAILS



MiSeq, NextSeq, and More Seqs

	MiSeq	NextSeq	HiSeq 4000	NovaSeq 6000
Maximum Output	15 Gb	120 Gb	750 Gb	3000 Gb
Maximum Reads per Run	25 million	400 million	2.5 billion	10 billion
Maximum Read Length	2 × 300 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp
Run Time	4-56 hours	15-29 hours	< 1–3.5 days	13-45 hours
Cost*	\$1,787	\$4,695	\$19,206	\$35,538
Cost/Mbp*	\$0.119	\$0.039	\$0.026	\$0.012

* Duke Sequencing and Genomic Technologies Shared Resource, July 2018

Illumina Video

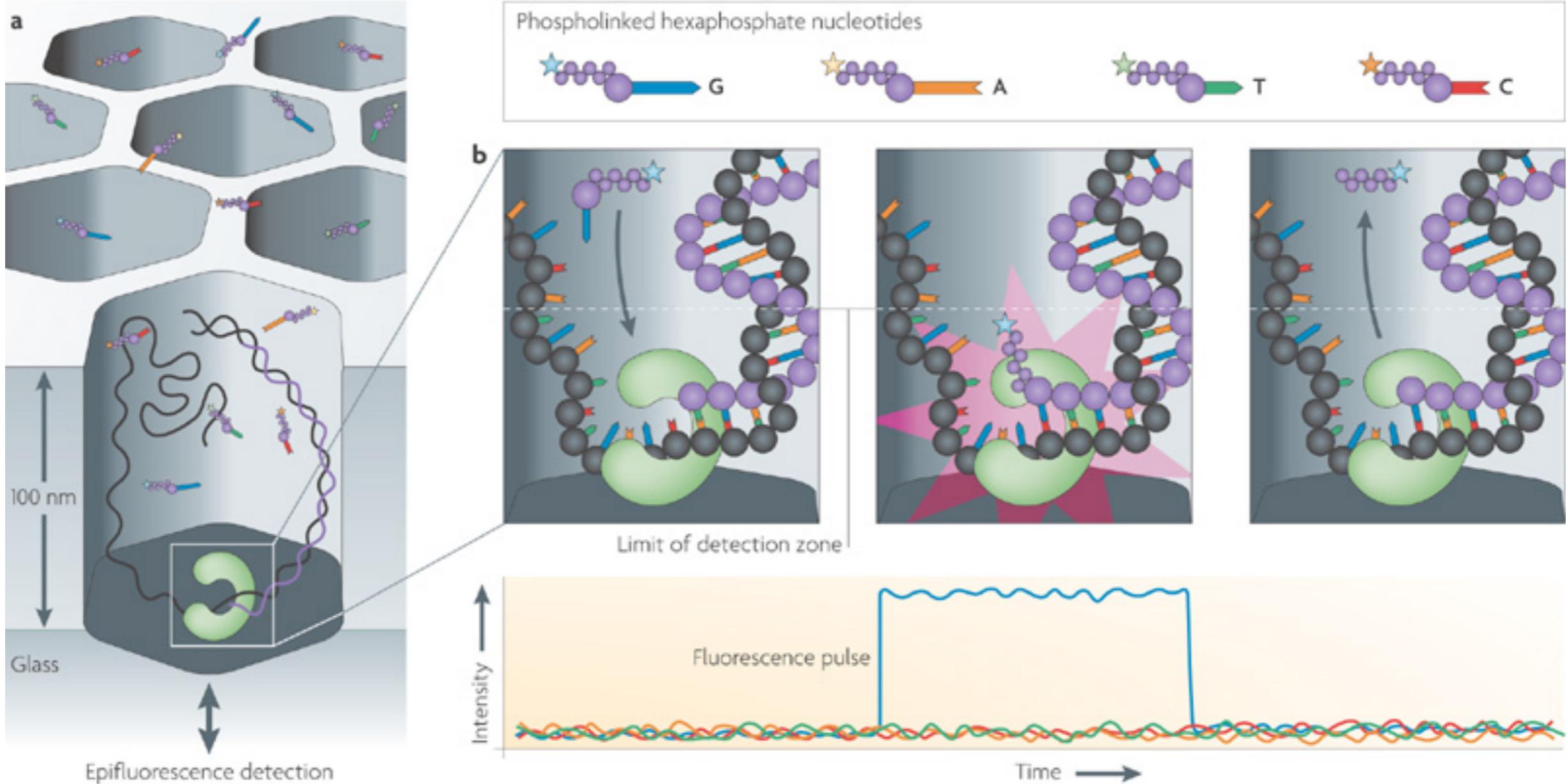
<https://www.youtube.com/watch?v=HMyCqWhwB8E>

Single Molecule Technologies

1st Generation	2nd Generation	3rd Generation
Chemical (Maxim-Gilbert)	Pyrosequencing (454)	Single molecule real time (PacBio)
Chain Termination (Sanger)	Chain Termination (Illumina)	Nanopore sequencing (Oxford Nanopore)
Pyrosequencing	Sequencing by ligation (SOLiD sequencing)	
	Ion semiconductor (Ion Torrent)	

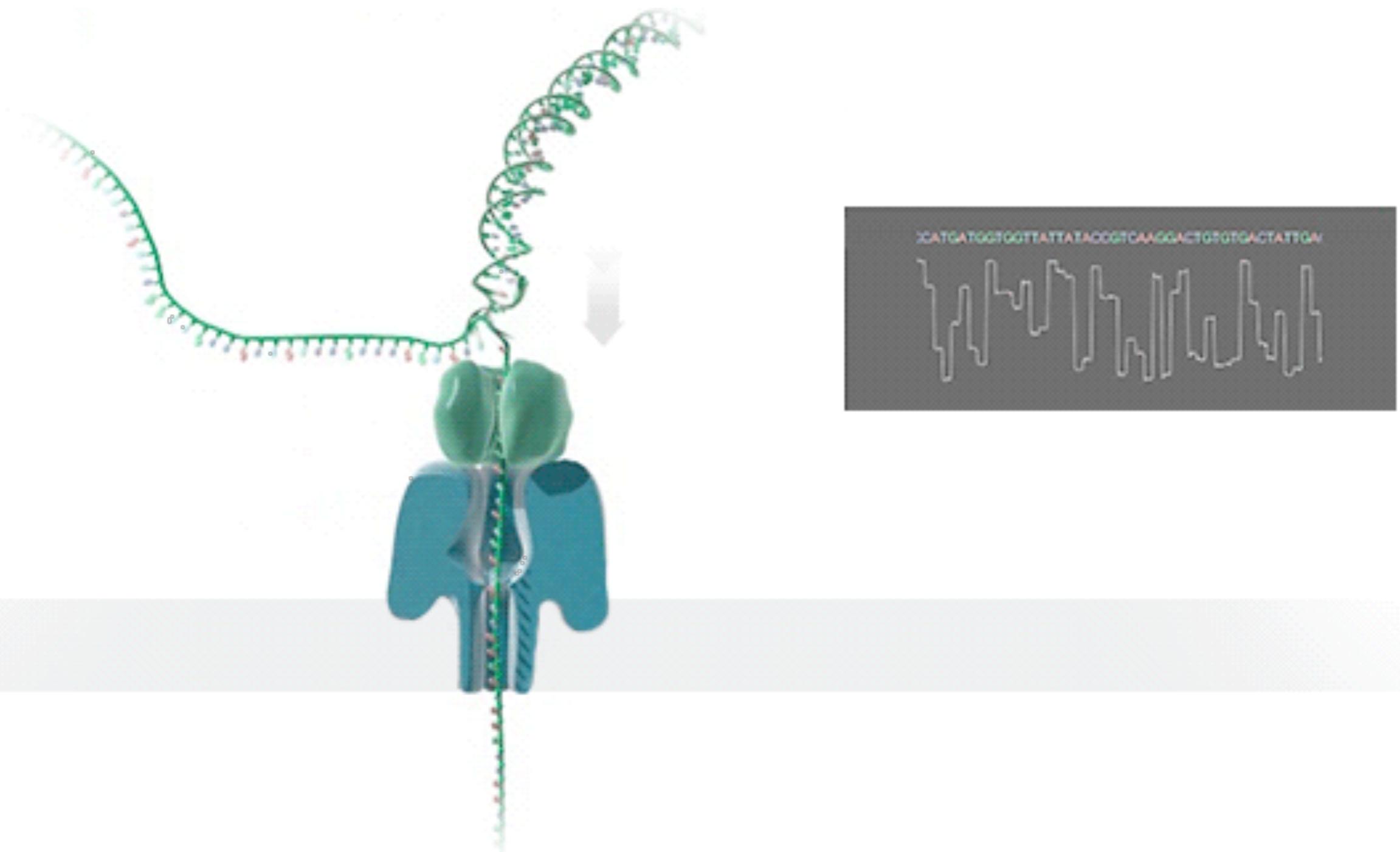
Pacific Biosciences

Pacific Biosciences — Real-time sequencing



1st Generation	2nd Generation	3rd Generation
Chemical (Maxim-Gilbert)	Pyrosequencing (454)	Single molecule real time (PacBio)
Chain Termination (Sanger)	Chain Termination (Illumina)	Nanopore sequencing (Oxford Nanopore)
Pyrosequencing	Sequencing by ligation (SOLID sequencing)	
	Ion semiconductor (Ion Torrent)	

Oxford Nanopore



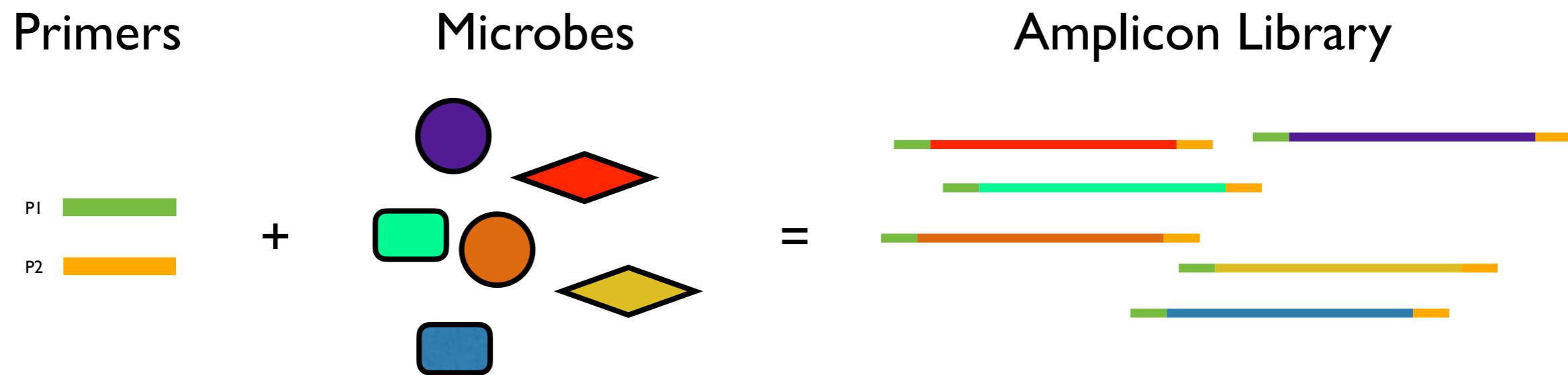
Sequencers



<https://www2.nanoporetech.com/images/product-page/MinION-Banner.jpg>
<http://www.gatc-biotech.com/en/gatc/sequencing-technologies/pacbio-rs-ii.html>
<http://www.dnavision.com/illumina.php>

Method	Read length	Accuracy	Reads per run	Max Output	Run Time	Cost (\$/Mb)	Pros	Cons
Sanger	400-900 bp	99.9%	1	900 bp	0.3-3 hours	\$2400	Long reads.	Expensive. Low Output
454	700 bp	99.9%	1 million	700 Mb	24 hours	\$10	Long reads Fast	Expensive. Homopolymer errors.
Ion Torrent	-400 bp	98%	up to 80 million	32 Gb	2 hours	\$1	Cheap equipment Fast.	Homopolymer errors. Short reads
Illumina	50-600 bp	99.9%	1 million to 3 billion	1500 Gb	1-11 days	\$0.05 to \$0.15	High yield Cheap	Equipment expense. Short reads
PacBio	>10kb ave. >40kb max	87%*	50,000	500 Mb	0.5-4 hour	\$0.13 to \$0.60	Very Long reads Fast Detects base modification	Homopolymer errors. Moderate Output. Equipment expense.

Amplicon Sequencing



FASTQ

- FASTA with Quality
- https://en.wikipedia.org/wiki/FASTQ_format

FASTQ: Filenames

“What’s in a name? that which we call a **FASTQ**

By any other name would smell as sweet”

— *Romeo and Juliet*, William Shakespeare

FASTQ: Filenames

Sample Name	Barcode	Lane #	Read #	File Part
1C_TAAGGCGA	L006	R1	001	.fastq.gz
1C_TAAGGCGA	L006	R1	002	.fastq.gz
1C_TAAGGCGA	L006	R1	003	.fastq.gz
1C_TAAGGCGA	L006	R1	004	.fastq.gz
1C_TAGGCATG	L003	R1	001	.fastq.gz
1C_TAGGCATG	L005	R1	001	.fastq.gz
2C_CGTACTAG	L003	R1	001	.fastq.gz

FASTQ: Format

- **FASTA Format**

>NC_007779.1 Escherichia coli str. K-12 substr. W3110, complete genome
AGCTTTCACTGACTGAAACGGCAATATGTCTCTGTGGATTAAAAAGAGTGTCTGATAGCAGCTTCTGA
GTTACCTGCCGTGAGTAAATTAAAATTATTGACTTAGTCATAACTTAACCAATATAGGCATAGCGCACAGAC
AGATAAAAATTACAGAGTACACAACATCCATGAAACGCATTAGCACCACCATACCACACCACATTAC
AACGGTGCGGGCTGACCGTACAGGAAACACAGAAAAAGCCGCACCTGACAGTGC
GGGCTTTTTGACCAAAGG
TAACGAGGTAA
ACAACCATGCGAGTGTGAAGTT
CGCGGTACATCAGTGGCAAATGCAGAACGTT
CTGCGTGTGCCG
ATATTCTGGAAAGCAATGCCAGGCAGGGCAGGTGGCCACCGTC
CTCTGCC
CCCCGCC
AAATC
ACCAACCAC
CTGGTGCG
GCGATGATTGAAAAAAC
CATTAGCGGCCAGGATGCTTACCCA
ATATCAGCGATGCC
AACGT
ATTTTGCC
AACTTT
GACGGGACTGCCGCC
CAGCCGGGTT
CCGCTGGCG
CAATTGAAA
ACTTC
GTGATCAGGA
ATTGCC
AAATAA
AACATGTC
CTGC
ATGGC
ATTAGTT
GTT
GGGG
CAGT
GCCCG
GATAGC
ATCAAC
GCT
GC
GT
GATT
TGC
CTGGCG
GAGAAA
AA

Header
Sequence

- **FASTQ Format**

@M00698:36:00000000-AFBEL:1:1101:16483:1412 1:N:0:0
CTGCCAGTTGAACGACGGCGAGCAGTTATAAGCCAGCAGTTGCCCGGATATT
CGCGTGGATAGCTTG
CAAAGCGACGCCAG
TCCAGATCCGGCG
+
AAABBFFFFFFFGGGGGGGGGHHHHHHHHGHGHGHHHHHGGGGGGHHHGHFFHHHHGHGGGGGGGGHHHHHHHHGGG

Quality
Score

Quality Scores

- Go to Notebook:

Wk4_Day4_AM/quality_scores.ipynb

FASTQ: Read Files

Combined_R1.fastq.gz

```
@M00698:36:00000000-AFBEL:1:1101:14738:1412 1:N:0:0
TTACGCTAACAGGCGGTAGCCTGGCAGGGTCAGGAAATCAATTAACTCATCGGAAGTGGTATCTGTTCCATCAAGCGTGCAGCATCGTCAAAACGCC
+
ABBBABBBBAFFFGGGGGGGGGGHGGHGGCG2GF3FFGHHHHHGGFGHEHHGGGEHHHAGGHGHHFFDHFHHHGEGGGG@F@H?GHH/GBEFGGG
@M00698:36:00000000-AFBEL:1:1101:16483:1412 1:N:0:0
CTGCCAGTTAACGACGGCGAGCAGTTAACGCCAGCAGTTGCCCGGATATTCGCGTGGATAGCTTGCAAAGCGACGCCAGTTCCAGATCCGGCG
+
AAABBFFFFFFGGGGGGGGGGHHHHHHHHGHGHGHGHHHHHGGGGHHHHGGGGGGHHHHGHFFHHHHGHGGGGGGGGHHHHHHHHGGG
```

Combined_R2.fastq.gz

```
@M00698:36:00000000-AFBEL:1:1101:14738:1412 2:N:0:0
GGAAGATCGGGCGACGGCTGAAATTCCCGTACCTCGATCTGGCAGTGGATCCATCATAAAAACGTTGAGCAATGGCAAACCGGTGACCAAAGCCTTGT
+
ABBABFFFFDBDGCG??FFGGGHGHFEG3EAEGGFHAE3GFBGGHGGHHCFGHFGBFHDFEGGHFHEFHHH3BFGF0GFEGGGGGHHA/FGHFHHH
@M00698:36:00000000-AFBEL:1:1101:16483:1412 2:N:0:0
GCTTCTTCCGTACTCATCGGGCATTGAGCAAGCGATCAGCCGTGGCCTGGCGTATGCCATATGCTGACCTGGTCTGGTGTGAAACCTCCACGCCGGAT
+
CCCCCFFFFBFFGGGGGGGGCECGHHHHHHHHGGHGGGGHGGCGCHGFHGGGGHHGGGGHHHHHHHHHEHGHHHGHHHGGGGGG
```

Combined_I1.fastq.gz

```
@M00698:36:00000000-AFBEL:1:1101:14738:1412 1:N:0:0
AGTTCC
+
CCCCDF
@M00698:36:00000000-AFBEL:1:1101:16483:1412 1:N:0:0
CCTGTC
+
A11>>1
```