

# **Introduction to DNA Sequencing**

# **(Relevant) Trivia**

How many base pairs (bp) are there in a human genome?

How much did it cost to sequence the first human genome?

How long did it take to sequence the first human genome?

When was the first human genome sequence complete?

Whose genome was it?

# **(Relevant) Trivia**

How many base pairs (bp) are there in a human genome?

**~3 billion (haploid)**

How much did it cost to sequence the first human genome?

**~\$2.7 billion**

How long did it take to sequence the first human genome?

**~13 years**

When was the first human genome sequence complete?

**2000-2003**

Whose genome was it?

**Several people's, but actually mostly a dude from Buffalo**

# Overview

- Prologue: **Assembly**
- The Past: **Sanger**
- The Present: **Next-Gen (454, Illumina, ...)**
- The Future: ? (**Nanopore, MinION, Single-molecule**)

# Overview

- Prologue: **Assembly**
- The Past: Sanger
- The Present: Next-Gen (454, Illumina, ...)
- The Future: ? (Nanopore, MinION, Single-molecule)

Method	Read Length
Sanger	
454	
Illumina	
Ion Torrent	

Method	Read Length
Sanger	600-1000 bp
454	
Illumina	
Ion Torrent	

Method	Read Length
Sanger	600-1000 bp
454	300-500 bp
Illumina	
Ion Torrent	



Method	Read Length
Sanger	600-1000 bp
454	300-500 bp
Illumina	~100 bp
Ion Torrent	

Method	Read Length
<b>Sanger</b>	600-1000 bp
<b>454</b>	300-500 bp
<b>Illumina</b>	~100 bp
<b>Ion Torrent</b>	~200 bp

### **But...**

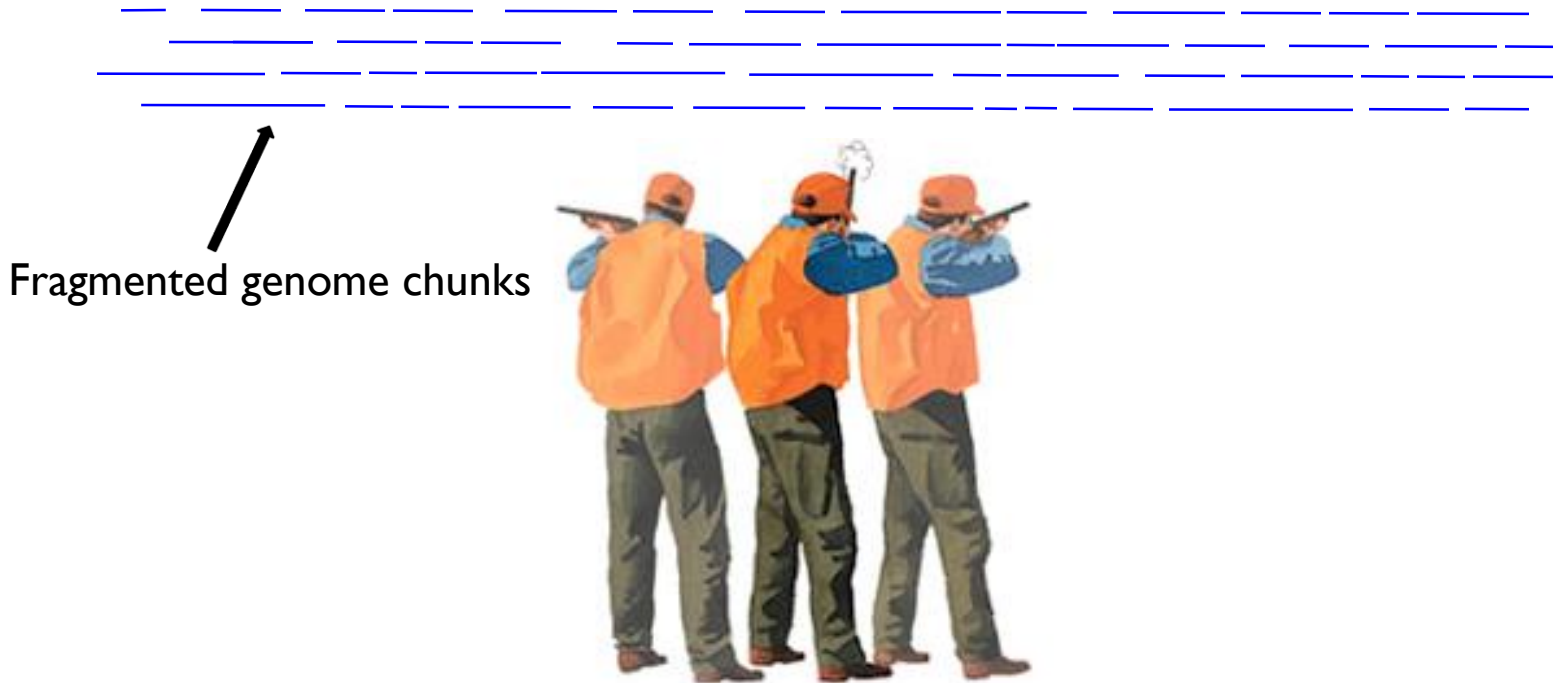
- Phage Genome: 30,000 to 500,000 bp
- Bacteria: Several million bp
- Human: 3 billion bp

# Shotgun Genome Sequencing

Fragmented genome  
Example: *Erwinia carotovora* plasmids

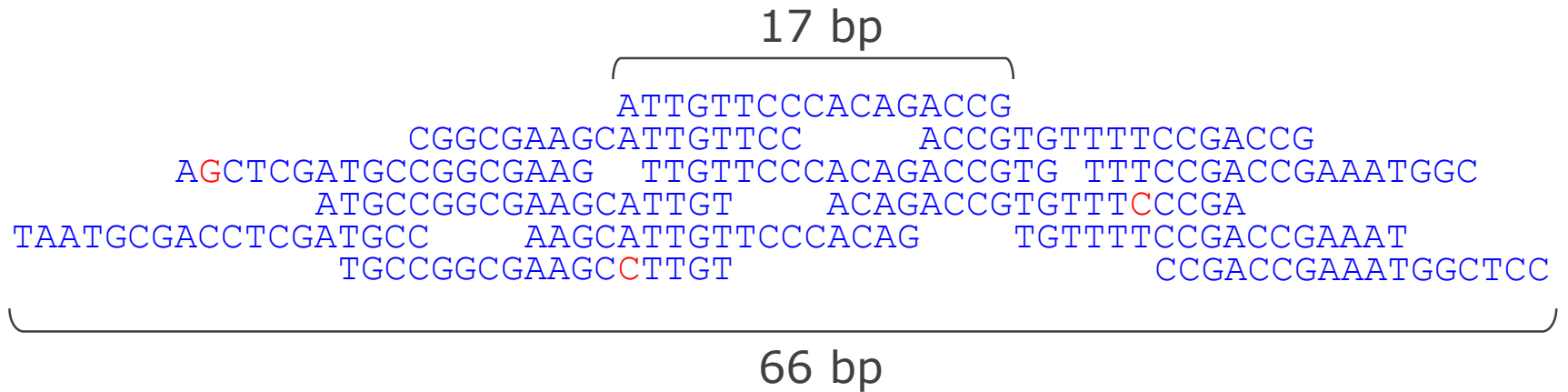


# Shotgun Genome Sequencing



NOT REALLY DONE BY DUCK HUNTERS  
Hydroshearing, sonication, enzymatic shearing

# Some Assembly Required



# Some Assembly Required

Consensus:

```
TAATGCGACCTCGATGCCGGCGAAGCATTGTTCCACAGACCGTGTTTTCCGACCGAAATGGCTCC
      ATTGTTCCACAGACCG
      CGGCGAAGCATTGTTCC      ACCGTGTTTTCCGACCG
AGCTCGATGCCGGCGAAG  TTGTTCCACAGACCGTG  TTTCCGACCGAAATGGC
      ATGCCGGCGAAGCATTGT      ACAGACCGTGTTTCCGA
TAATGCGACCTCGATGCC      AAGCATTGTTCCACAG      TGTTTTCCGACCGAAAT
      TGCCGGCGAAGCCTTGT      CCGACCGAAATGGCTCC
```

# Some Assembly Required

Consensus:

```
TAATGCGACCTCGATGCCGGCGAAGCATTGTTCCACAGACCGTGTTTTCCGACCGAAATGGCTCC
      ATGTTCCACAGACCG
      CGGCGAAGCATTGTTCC      ACCGTGTTTTCCGACCG
AGCTCGATGCCGGCGAAG  TTGTTCCACAGACCGTG  TTTCCGACCGAAATGGC
      ATGCCGGCGAAGCATTGT      ACAGACCGTGTTTCCGA
TAATGCGACCTCGATGCC      AAGCATTGTTCCACAG      TGTTTTCCGACCGAAAT
      TGCCGGCGAAGCCTTGT      CCGACCGAAATGGCTCC
```

**Coverage:** # of reads underlying the consensus

# Some Assembly Required

Consensus:

```
TAATGCGACCTCGATGCCGGCGAAGCATTGTTCCACAGACCGTGTTTTCCGACCGAAATGGCTCC
      ATTGTTCCACAGACCG
      CGGCGAAGCATTGTTCC      ACCGTGTTTTCCGACCG
AGCTCGATGCCGGCGAAG  TTGTTCCACAGACCGTG  TTTCCGACCGAAATGGC
      ATGCCGGCGAAGCATTGT      ACAGACCGTGTTTCCGA
TAATGCGACCTCGATGCC      AAGCATTGTTCCACAG      TGTTTTCCGACCGAAAT
      TGCCGGCGAAGCCTTGT      CCGACCGAAATGGCTCC
```

6x coverage  
100% identity

**Coverage:** # of reads underlying the consensus



# Assembly

Consensus:

TAATGCGACCTCGATGCCGGCGAAGCATTGTTCCACAGACCGTGTTTTCCGACCGAAATGGCTCC

ATTGTTCCACAGACCG  
CGGCGAAGCATTGTTCC ACCGTGTTTTCCGACCG  
AGCTCGATGCCGGCGAAG TTGTTCCACAGACCGTG TTTCCGACCGAAATGGC  
ATGCCGGCGAAGCATTGT ACAGACCGTGTTTCCGA  
TAATGCGACCTCGATGCC AAGCATTGTTCCACAG TGTTCGACCGAAAT  
TGCCGGCGAAGCCTTGT CCGACCGAAATGGCTCC

5x coverage  
80% identity

**Coverage:** # of reads underlying the consensus

# Assembly

Consensus:

TAATGCGACCTCGATGCCGGCGAAGCATTGTTCCACAGACCGTGTTTTCCGACCGAAATGGCTCC



ATTGTTCCACAGACCG  
CGGCGAAGCATTGTTCC ACCGTGTTTTCCGACCG  
AGCTCGATGCCGGCGAAG TTGTTCCACAGACCGTG TTTCCGACCGAAATGGC  
ATGCCGGCGAAGCATTGT ACAGACCGTGTTTCCGA  
TAATGCGACCTCGATGCC AAGCATTGTTCCACAG TGTTTTCCGACCGAAAT  
TGCCGGCGAAGCCTTGT CCGACCGAAATGGCTCC

2x coverage  
50% identity

**Coverage:** # of reads underlying the consensus

# Assembly

Consensus:

TAATGCGACCTCGATGCCGGCGAAGCATTGTTCCACAGACCGTGTTTTCCGACCGAAATGGCTCC

ATTGTTCCACAGACCG  
CGGCGAAGCATTGTTCC ACCGTGTTTTCCGACCG  
AGCTCGATGCCGGCGAAG TTGTTCCACAGACCGTG TTTCCGACCGAAATGGC  
ATGCCGGCGAAGCATTGT ACAGACCGTGTTTCCGA  
TAATGCGACCTCGATGCC AAGCATTGTTCCACAG TGTTCGACCGAAAT  
TGCCGGCGAAGCCTTGT CCGACCGAAATGGCTCC

1x coverage

**Coverage:** # of reads underlying the consensus

File Navigate Info Color Dim Misc

564\_454.fasta.screen.ace.12

Search for String:  Compl. Cont:  Compare Cont:  Find Main Min:  Err/10kb:  2.37

Contig6

Sone Tags Pos:

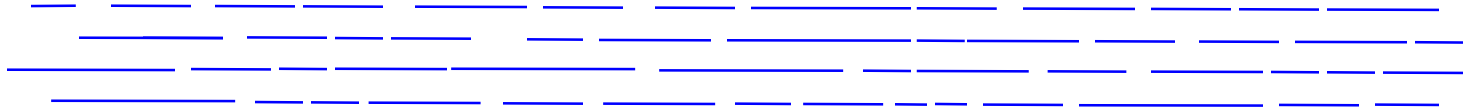
Help

clear

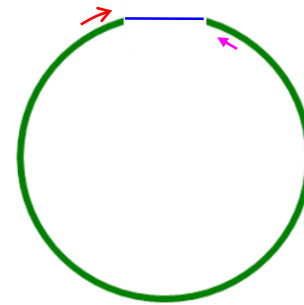
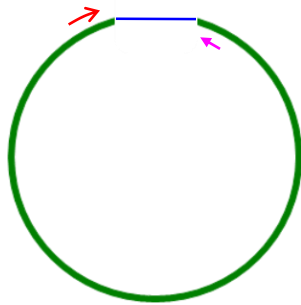
	1090	1100	1110	1120	1130	1140	1150	1160	1170	1180	1190	1200	1210	1220	1230	1240	1250	1260	1270	1280	1290	
CONS	56	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AVESM	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AI417	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
B083E	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AYM49	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
BH13F	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AS0HM	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AD04N	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
B6GFW	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AI130	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AI30J	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
CFRFR	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AI126	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
ADQNT	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AI0NP2	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
B8FC5	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
BY68Y	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
AI672L	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC	T	T
B8KDF	1	GGGAC	CACCA	CAACCA	CAAC	CGT	CGAC	CCCCG	CTCCCG	CTACCA	CAACCA	CAACCC	*ACCCA	TGGCT	GGAC	TC	GGAC	TA	TGAAG	GAAC		

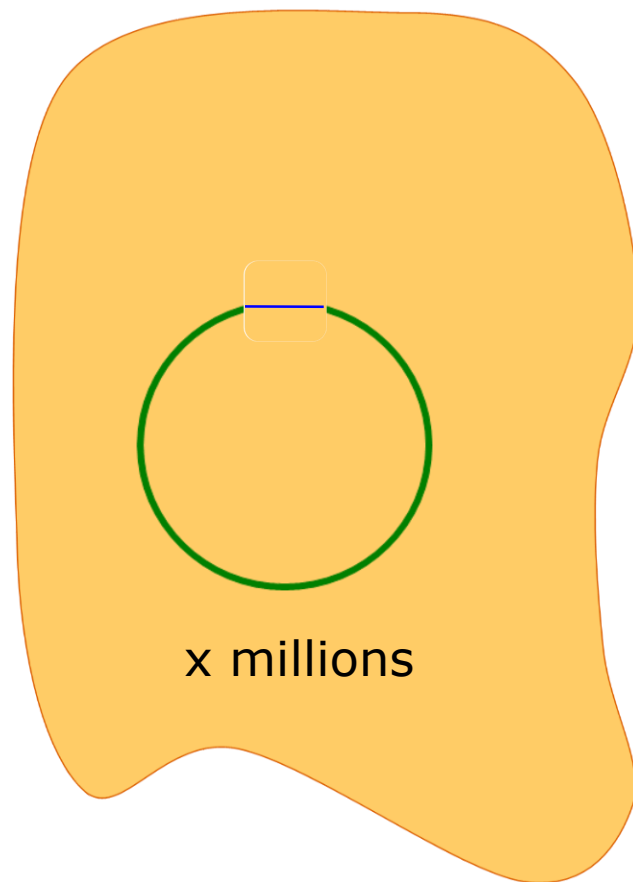
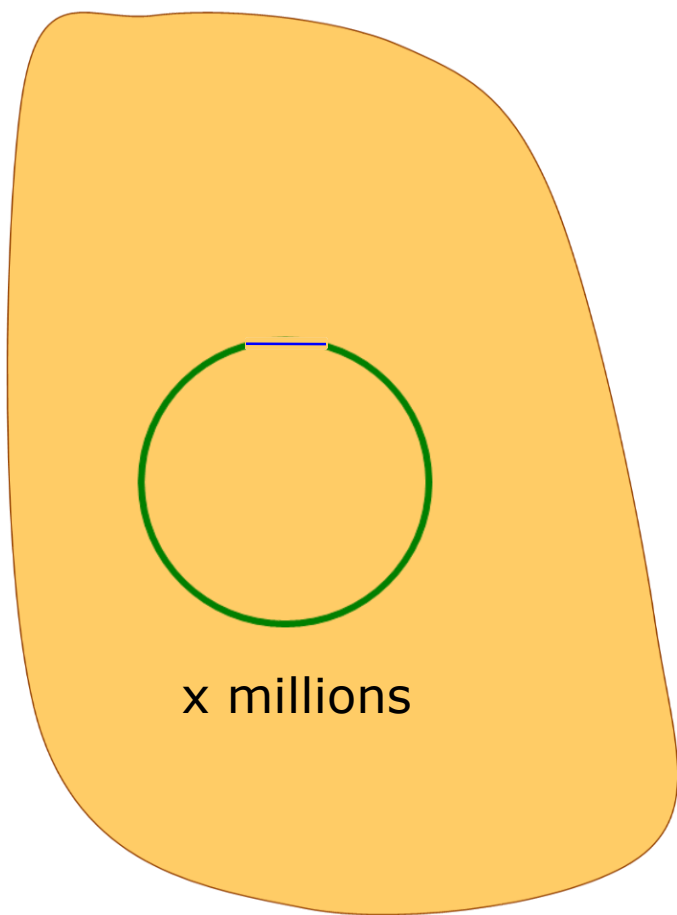
# Overview

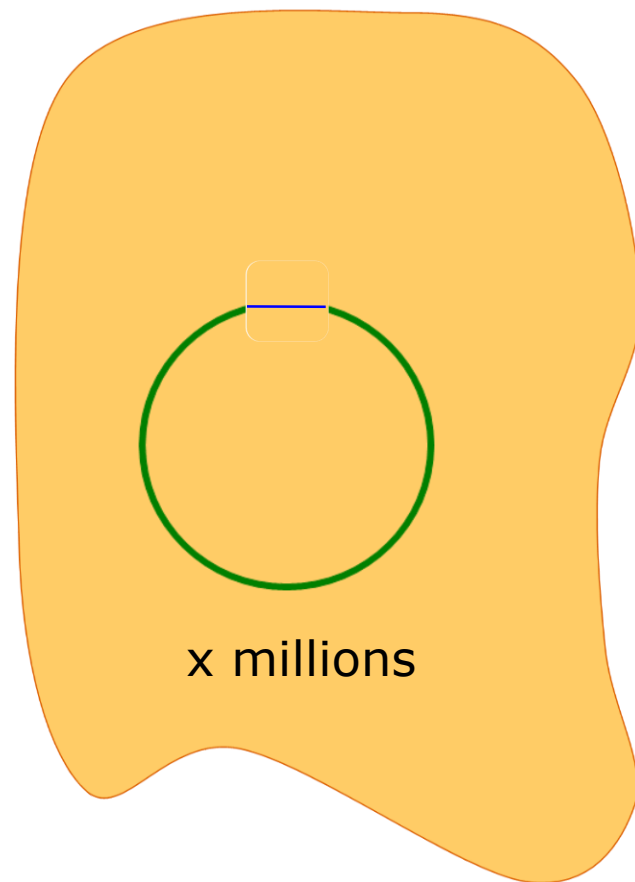
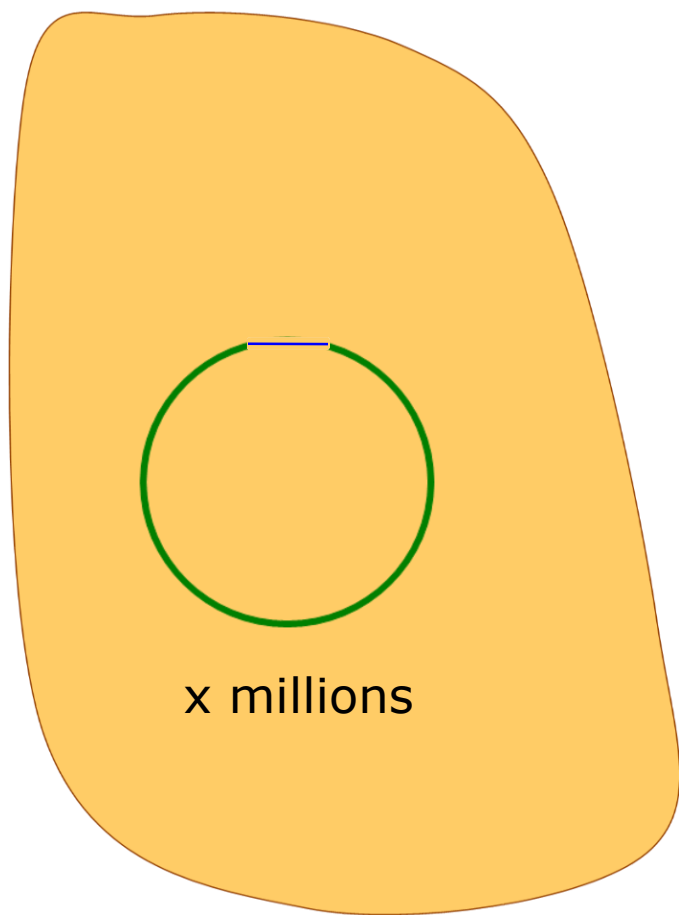
- Prologue: Assembly
- The Past: **Sanger**
- The Present: Next-Gen (454, Illumina, ...)
- The Future: ? (Nanopore, MinION, Single-molecule)



Fragments were cloned:









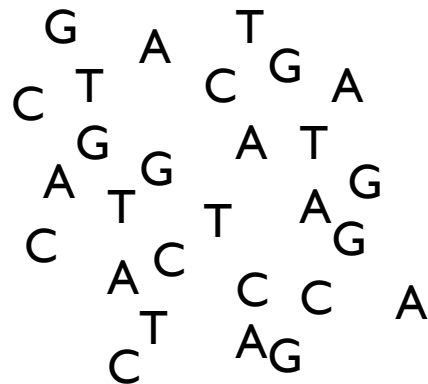
# Sanger Sequencing Reactions

For given template DNA, it's like PCR except:

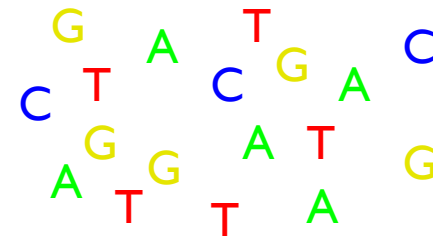
Uses only a single primer and polymerase to make new ssDNA pieces.

Includes regular nucleotides (A, C, G, T) for extension, but also includes dideoxy nucleotides.

## Regular Nucleotides

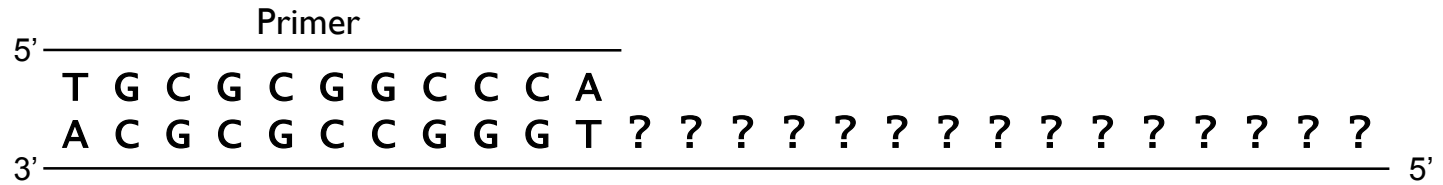


## Dideoxy Nucleotides



1. Labeled
2. Terminators

# Sanger Sequencing



# Sanger Sequencing



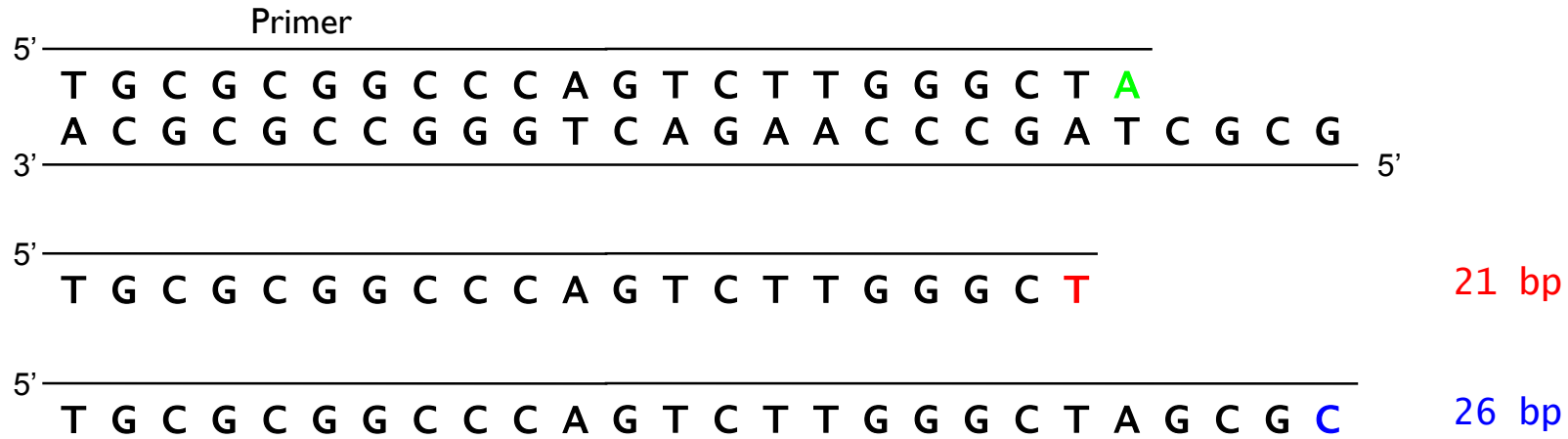
# Sanger Sequencing

5' Primer \_\_\_\_\_  
T G C G C G G C C C A G T C T T G G G C T A G C G **C**  
A C G C G C C G G G T C A G A A C C C G A T C G C G 5'

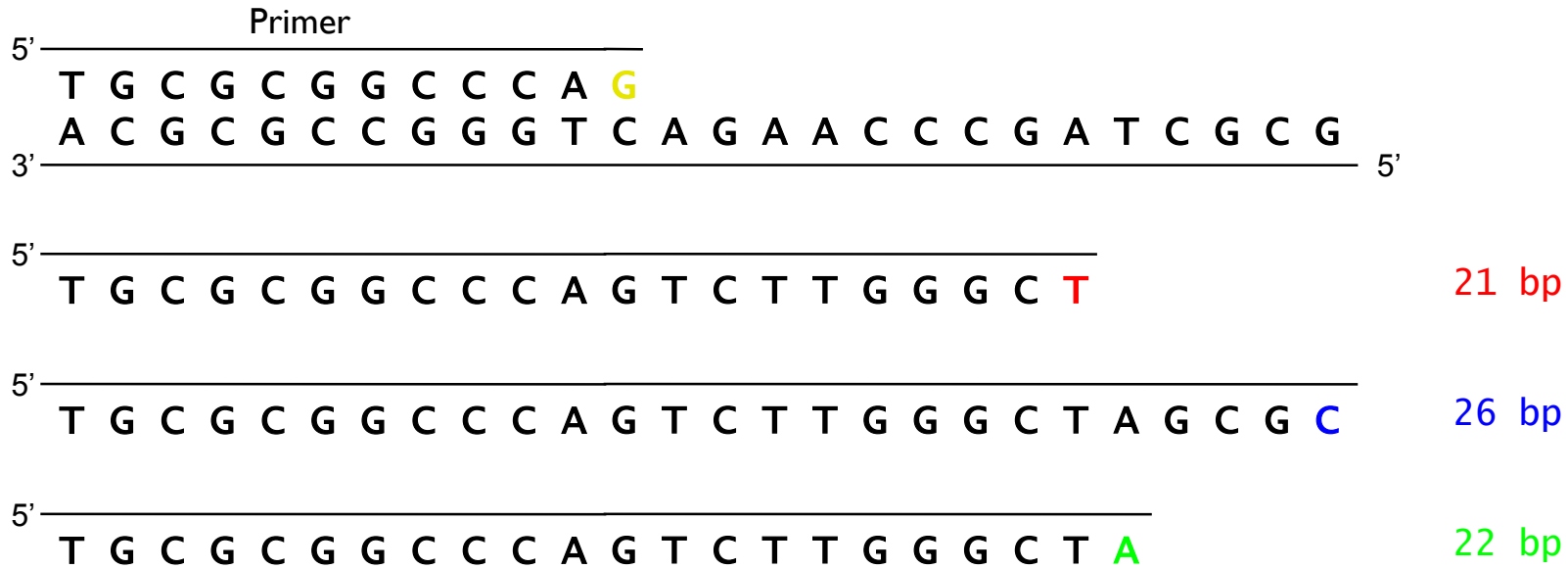
5' \_\_\_\_\_  
T G C G C G G C C C A G T C T T G G G C **T**

21 bp

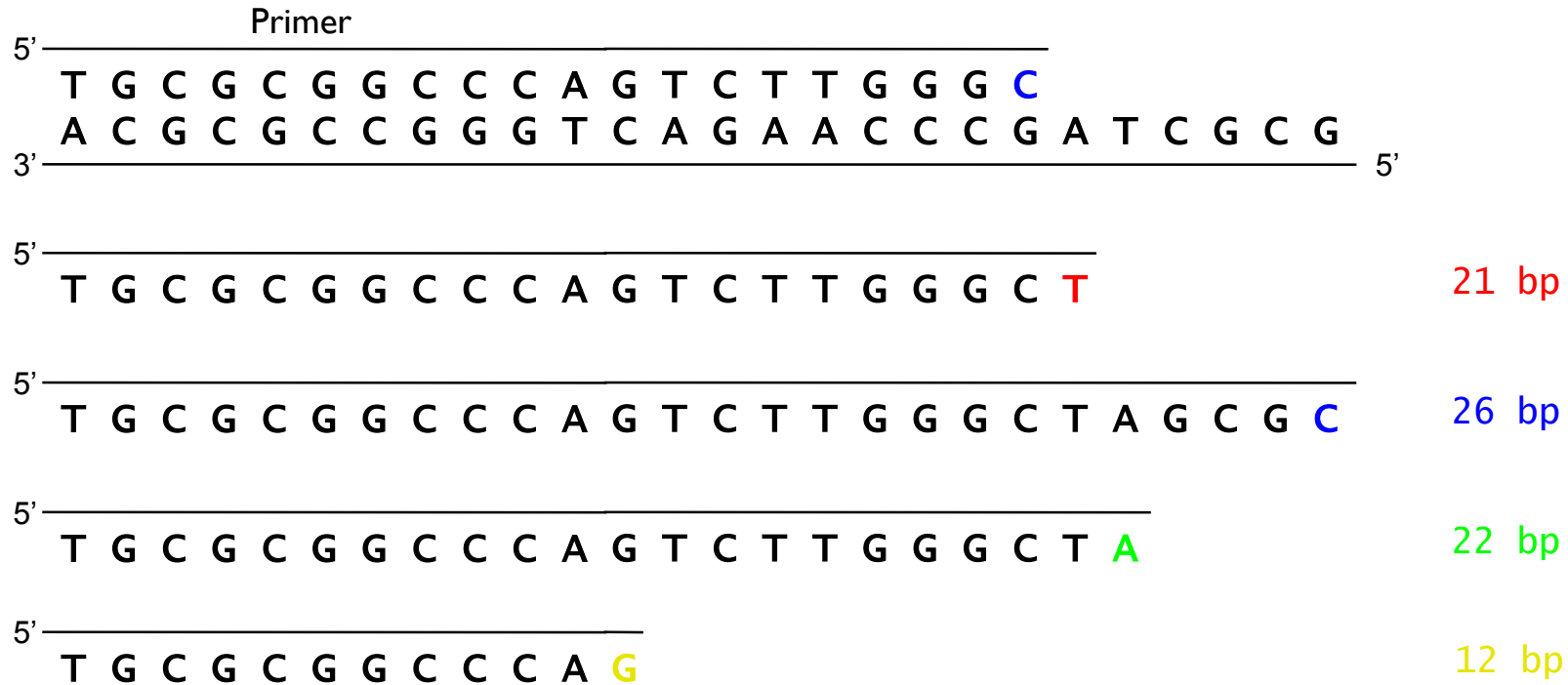
# Sanger Sequencing



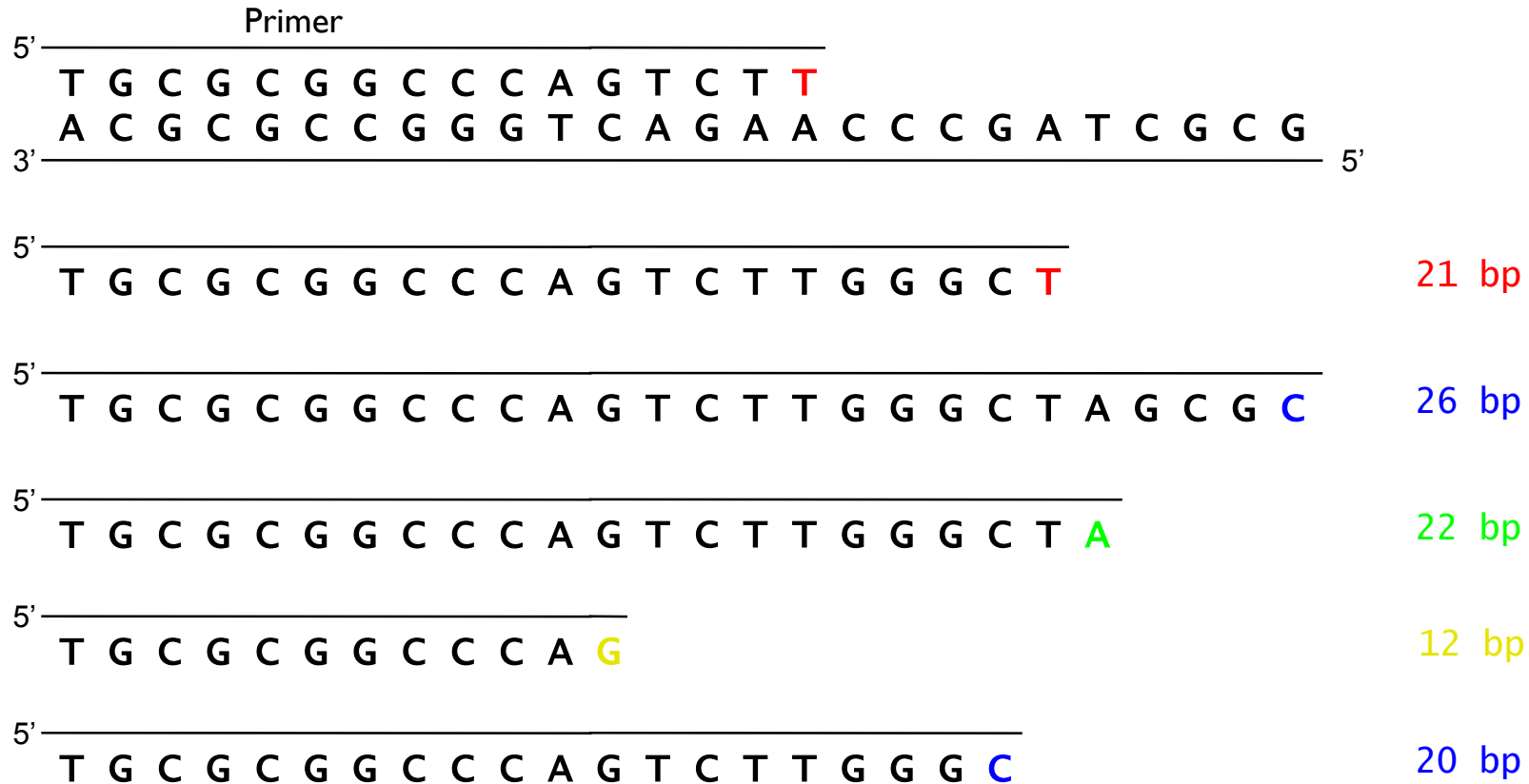
# Sanger Sequencing



# Sanger Sequencing

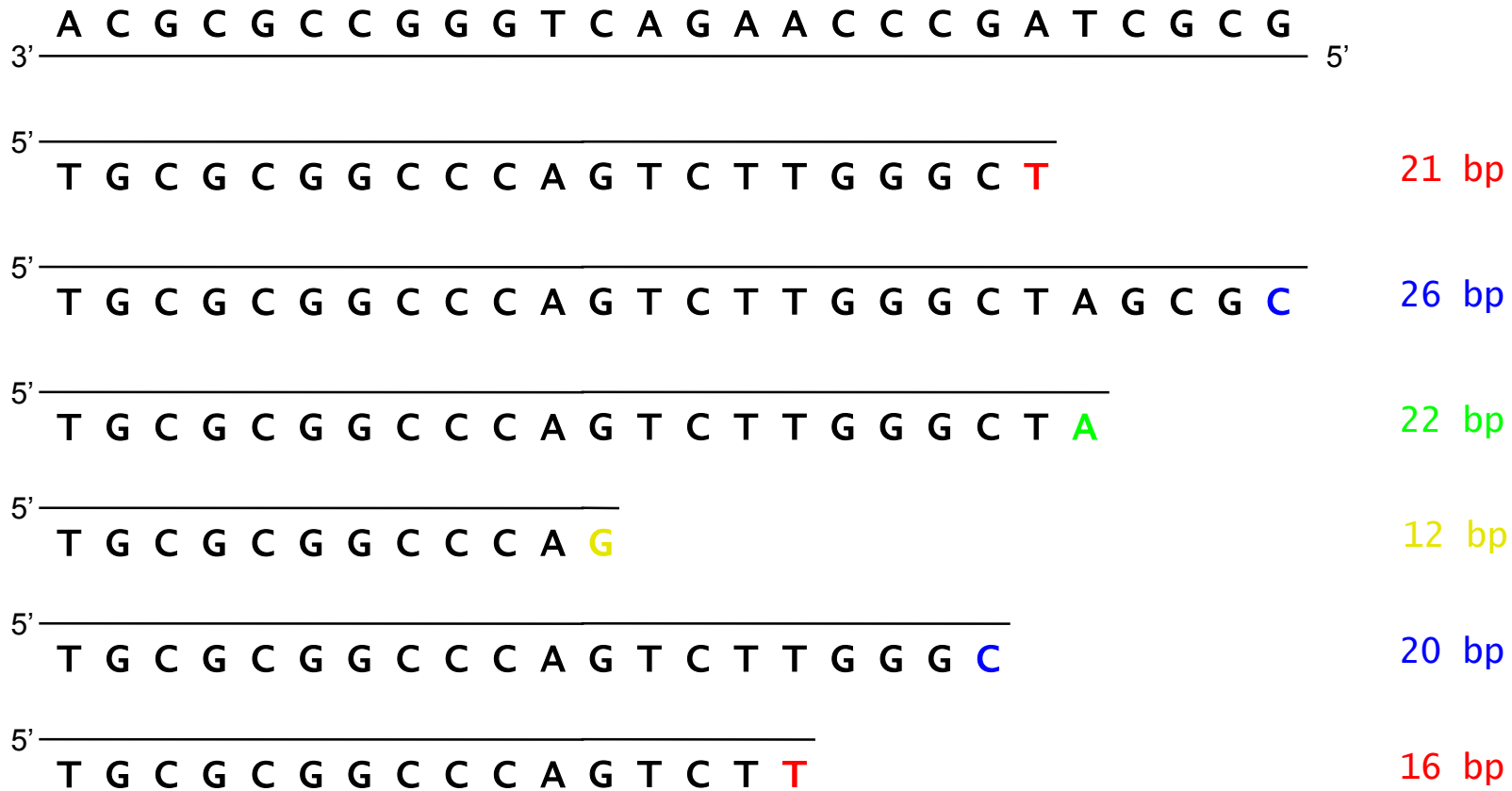


# Sanger Sequencing

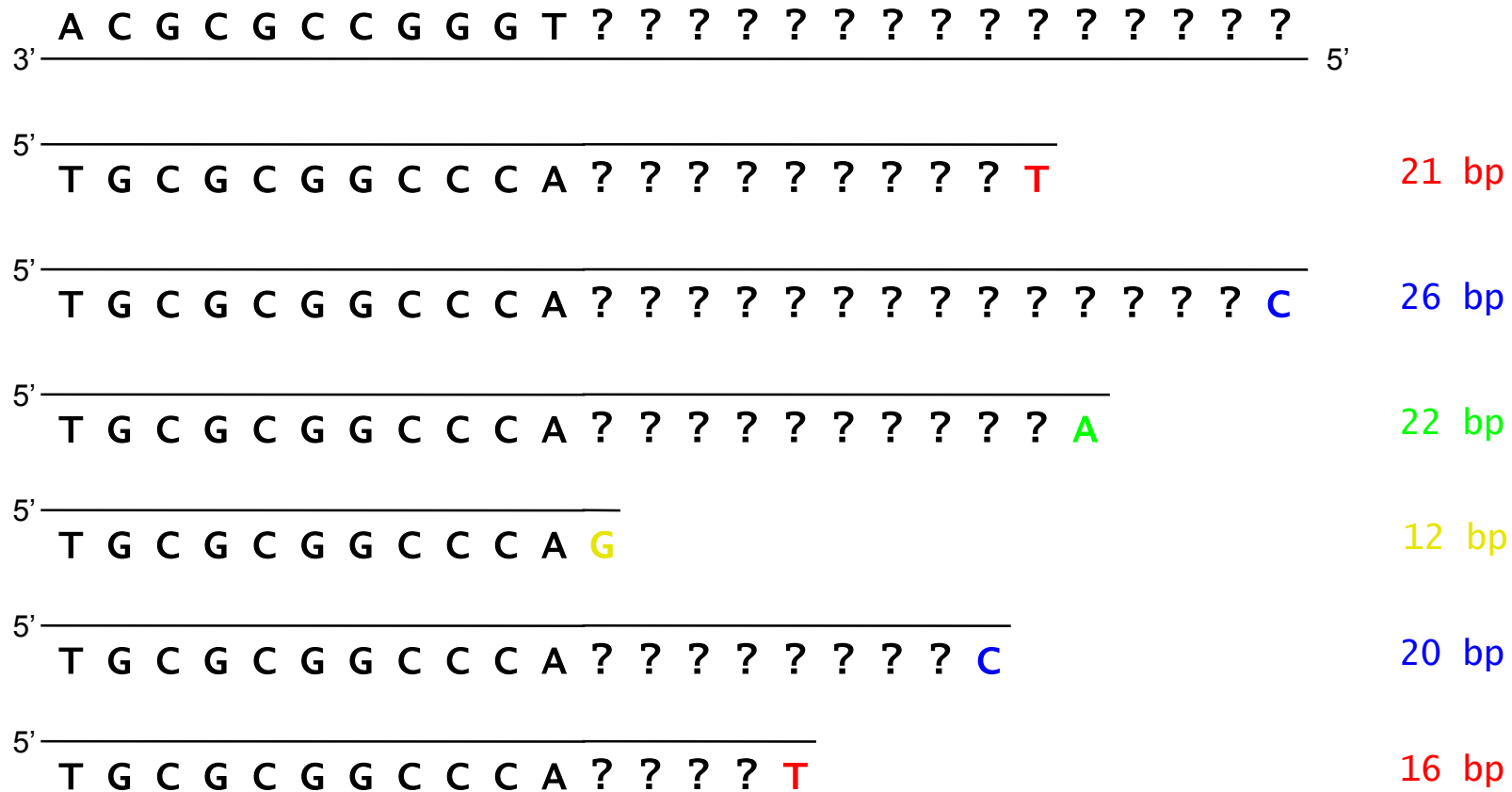




# Sanger Sequencing



# Sanger Sequencing

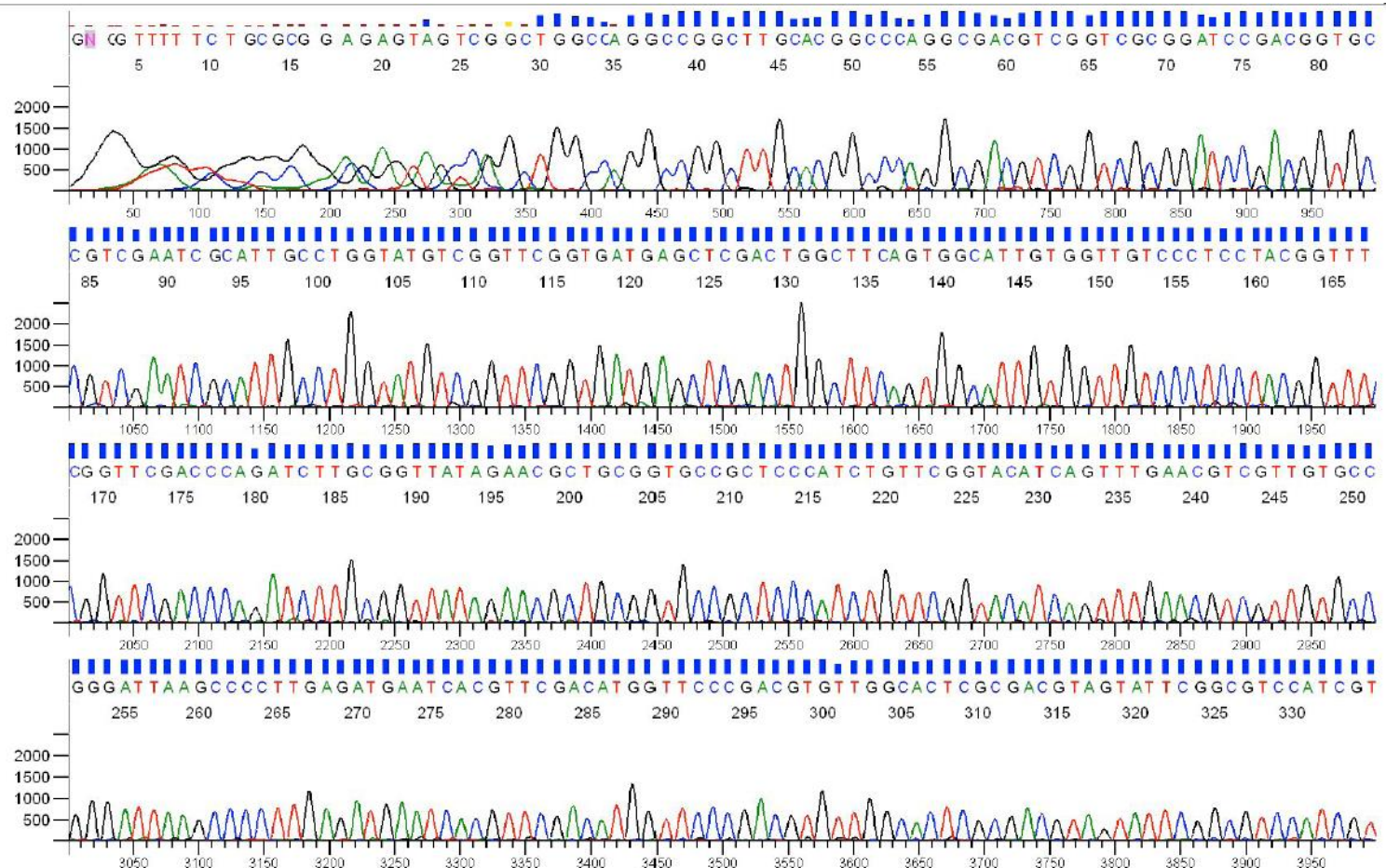


# Sanger Sequencing



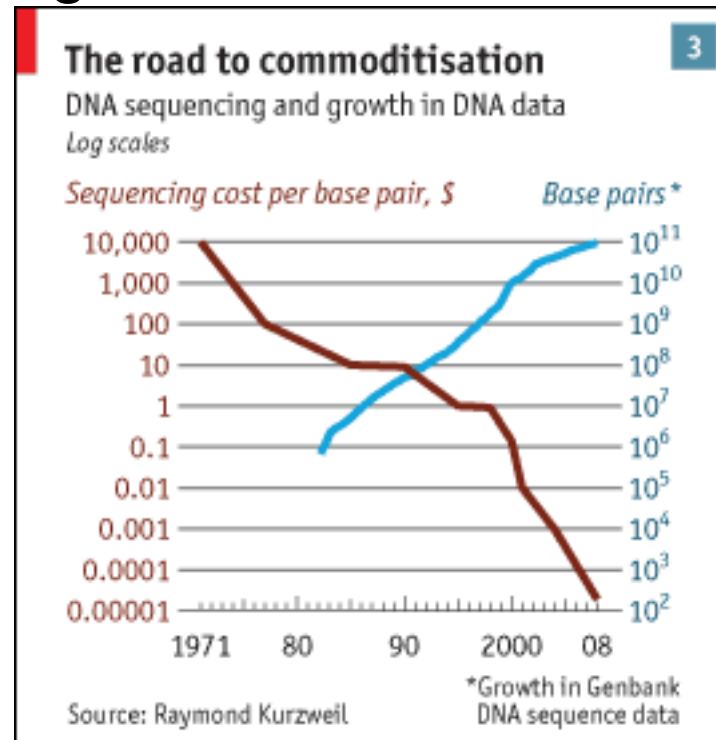
# Sanger Sequencing Output

Each sequencing reaction gives us a **chromatogram**, usually ~600-1000 bp:



# Sanger Throughput Limitations

- Must have 1 colony picked for every 2 reactions
- Must do 1 DNA prep for every 2 reactions
- Must have 1 PCR tube for each reaction
- Must have 1 gel lane for each reaction

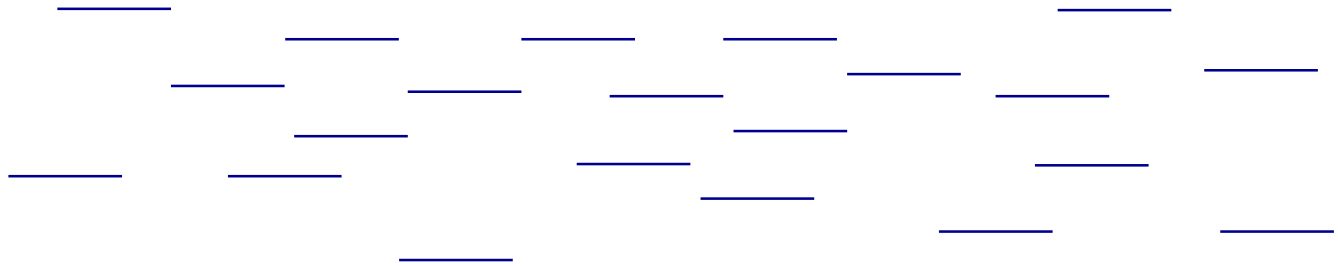


from *The Economist*

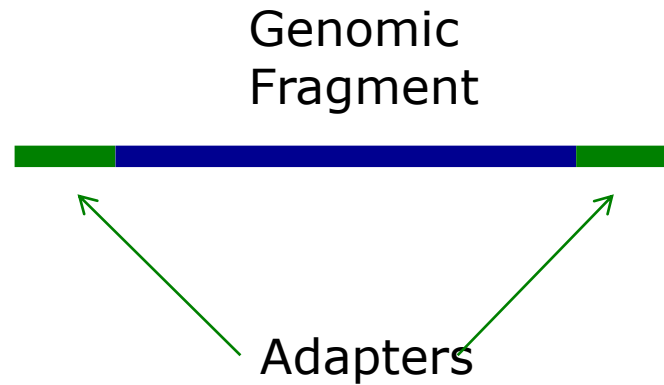
# Overview

- Prologue: Assembly
- The Past: Sanger
- The Present: **Next-Gen (454, Illumina, ...)**
- The Future: ? (Nanopore, MinION, Single-molecule)

# Shotgun sequencing by Ion Torrent Personal Genome Machine and 454

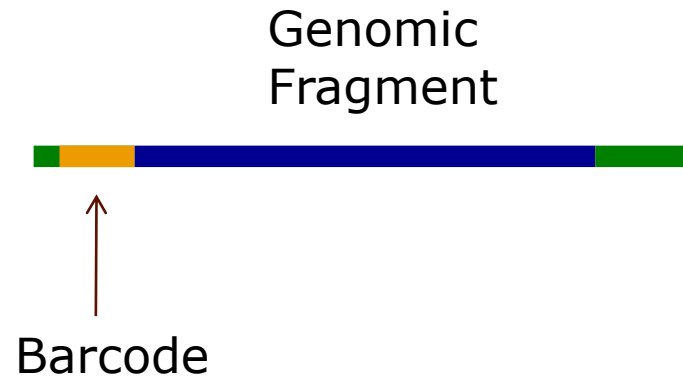


# Shotgun sequencing by PGM/454

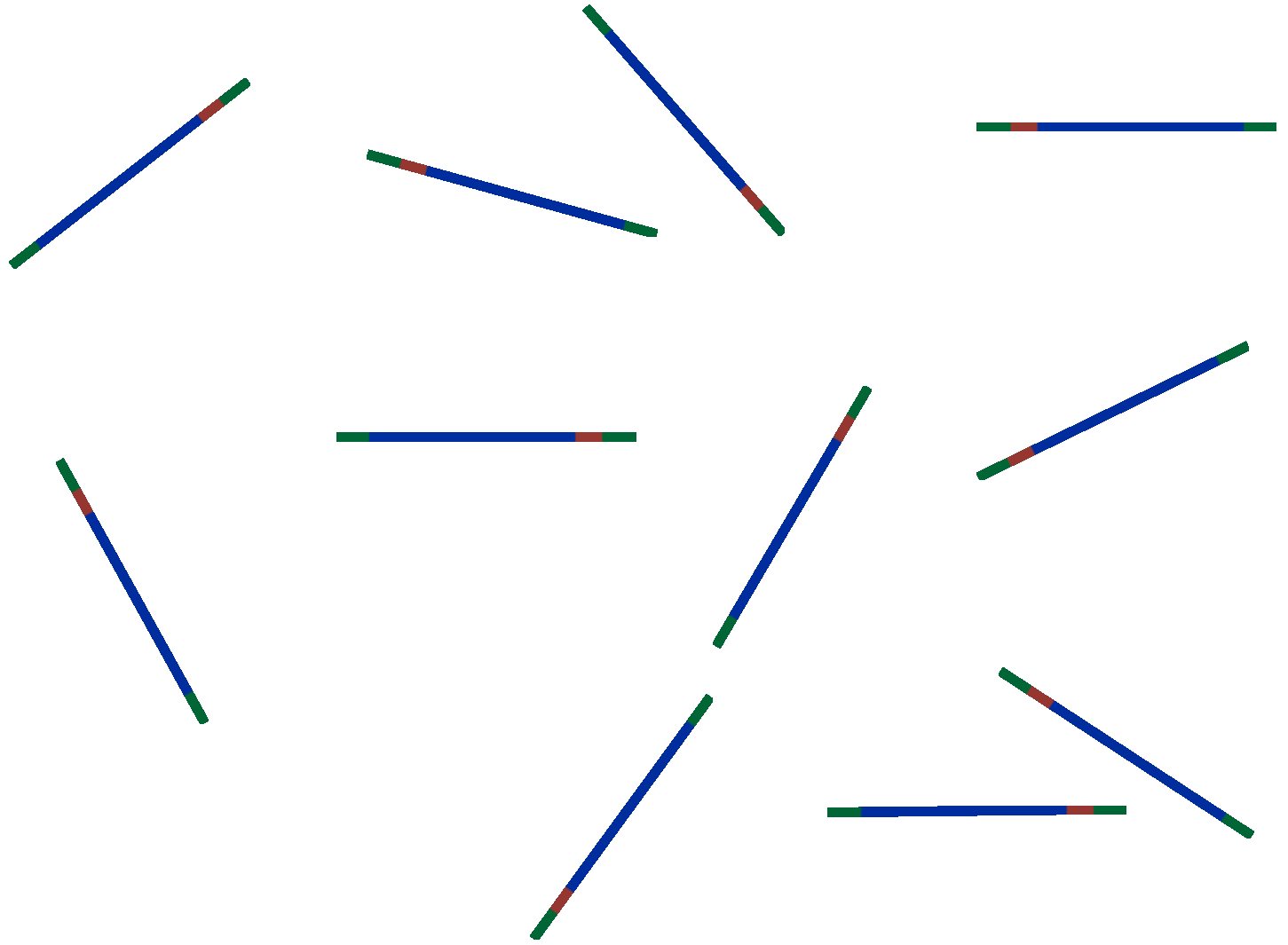




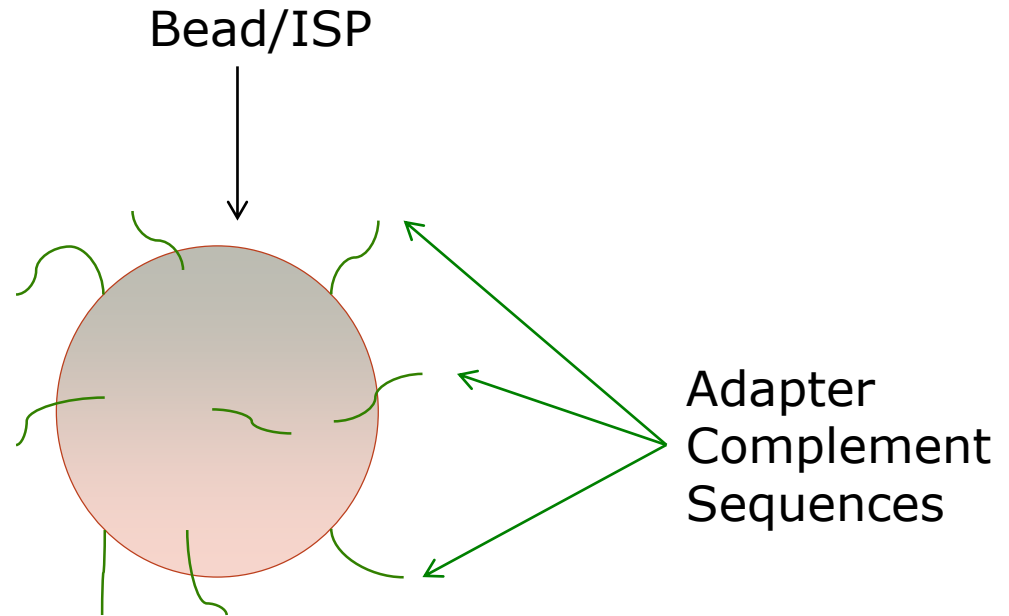
# Shotgun sequencing by PGM/454



# Shotgun sequencing by PGM/454

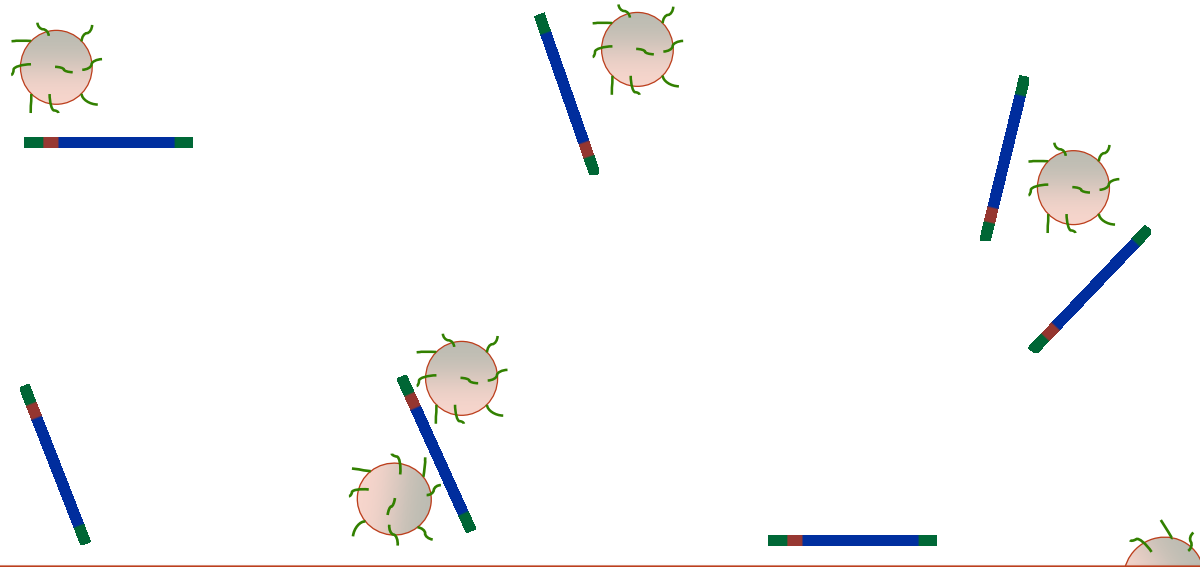


# Shotgun sequencing by PGM/454



The idea is that each bead should be amplified all over with a SINGLE library fragment.

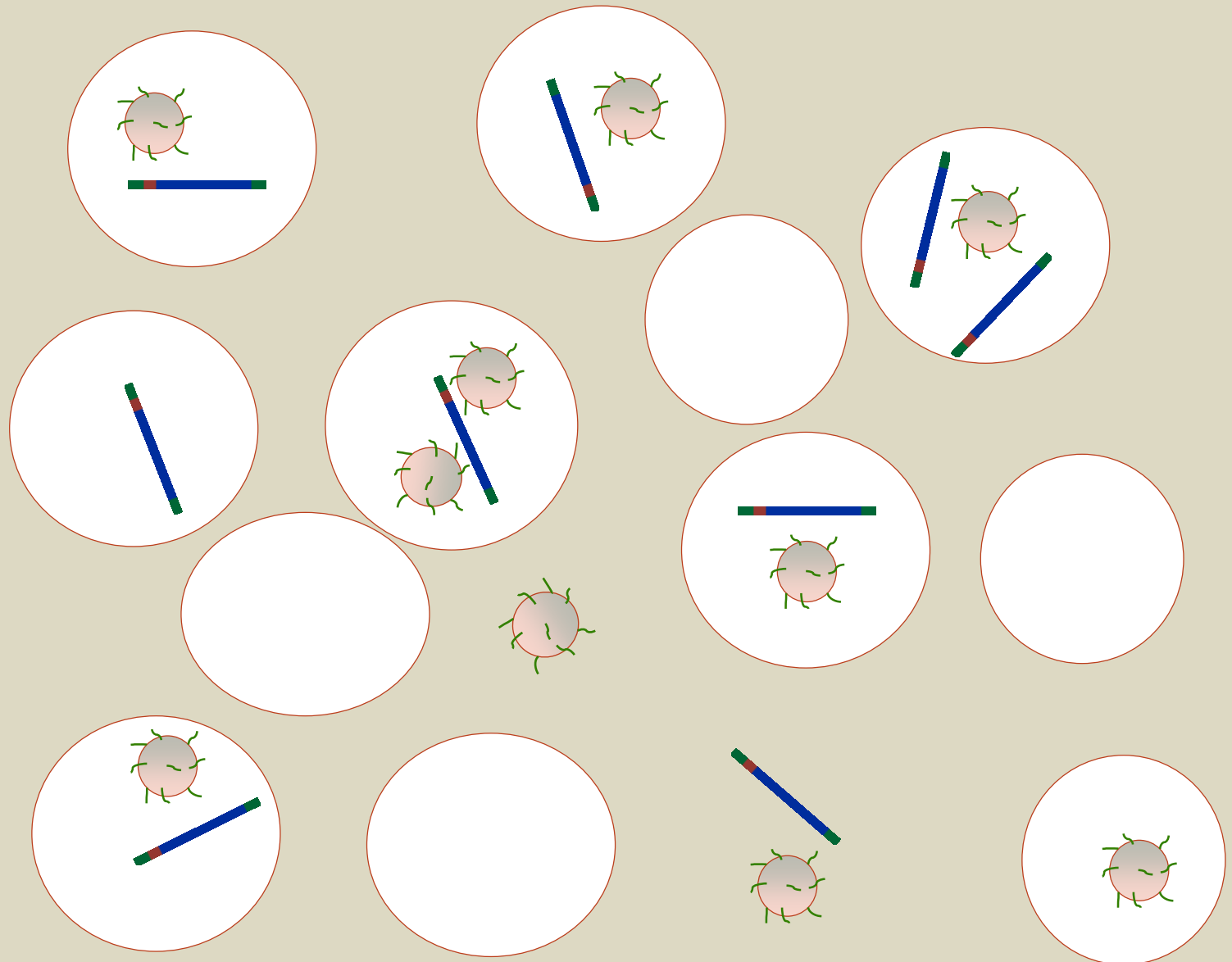
# Shotgun sequencing by PGM/454



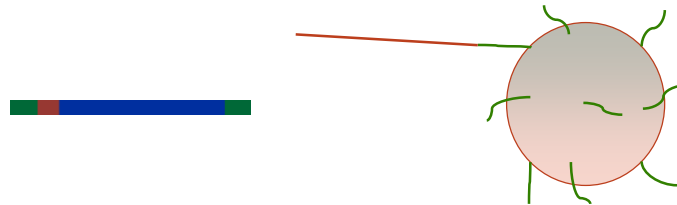
Problem: How do I do PCR to amplify the fragments without having to use 1 tube for each reaction?



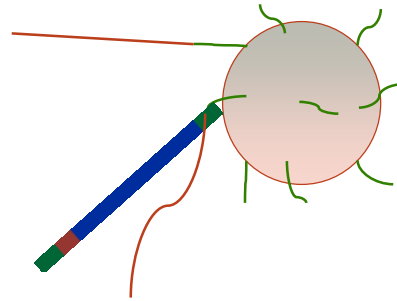
# Shotgun sequencing by PGM/454



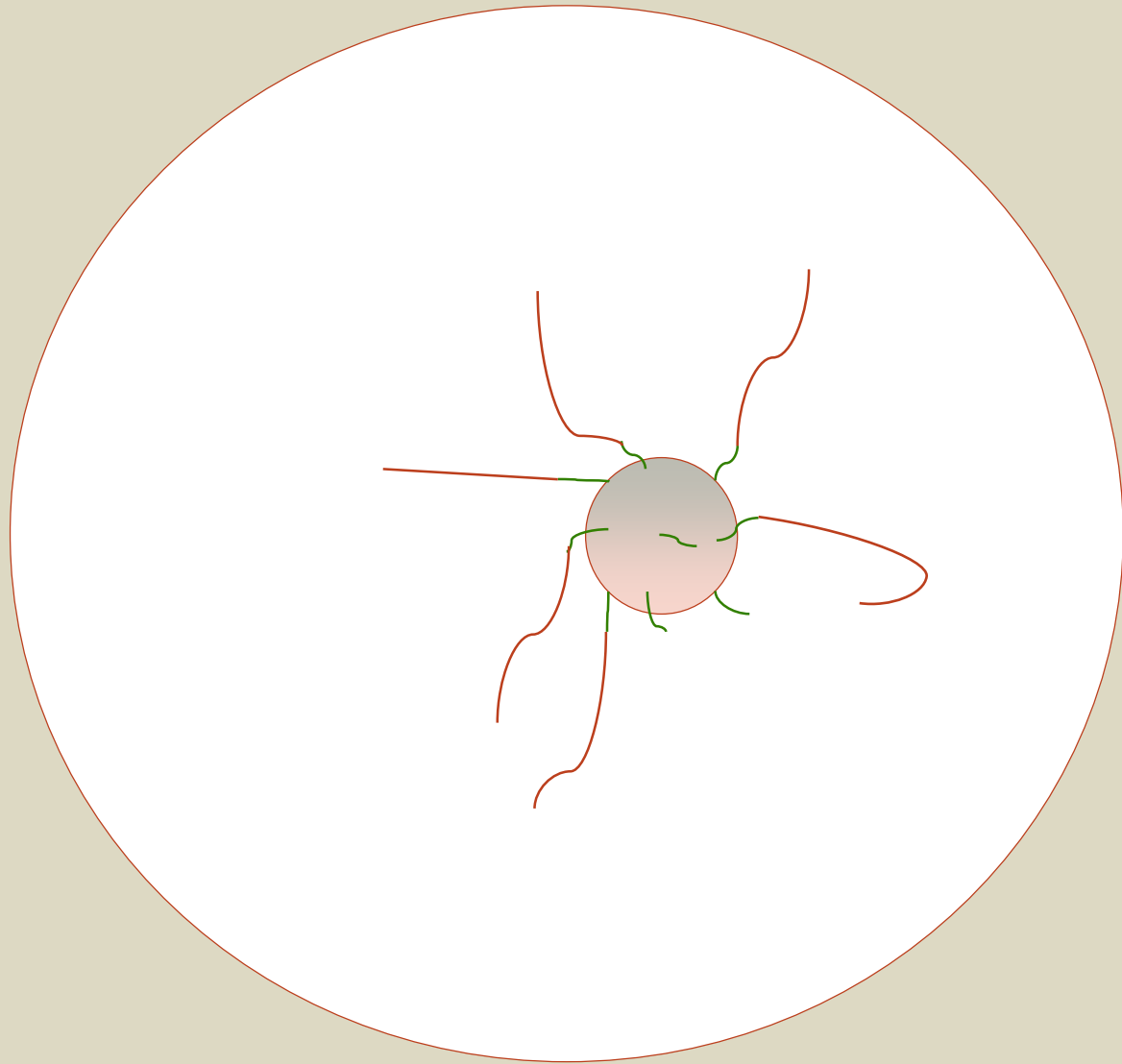
# Shotgun sequencing by PGM/454



# Shotgun sequencing by PGM/454

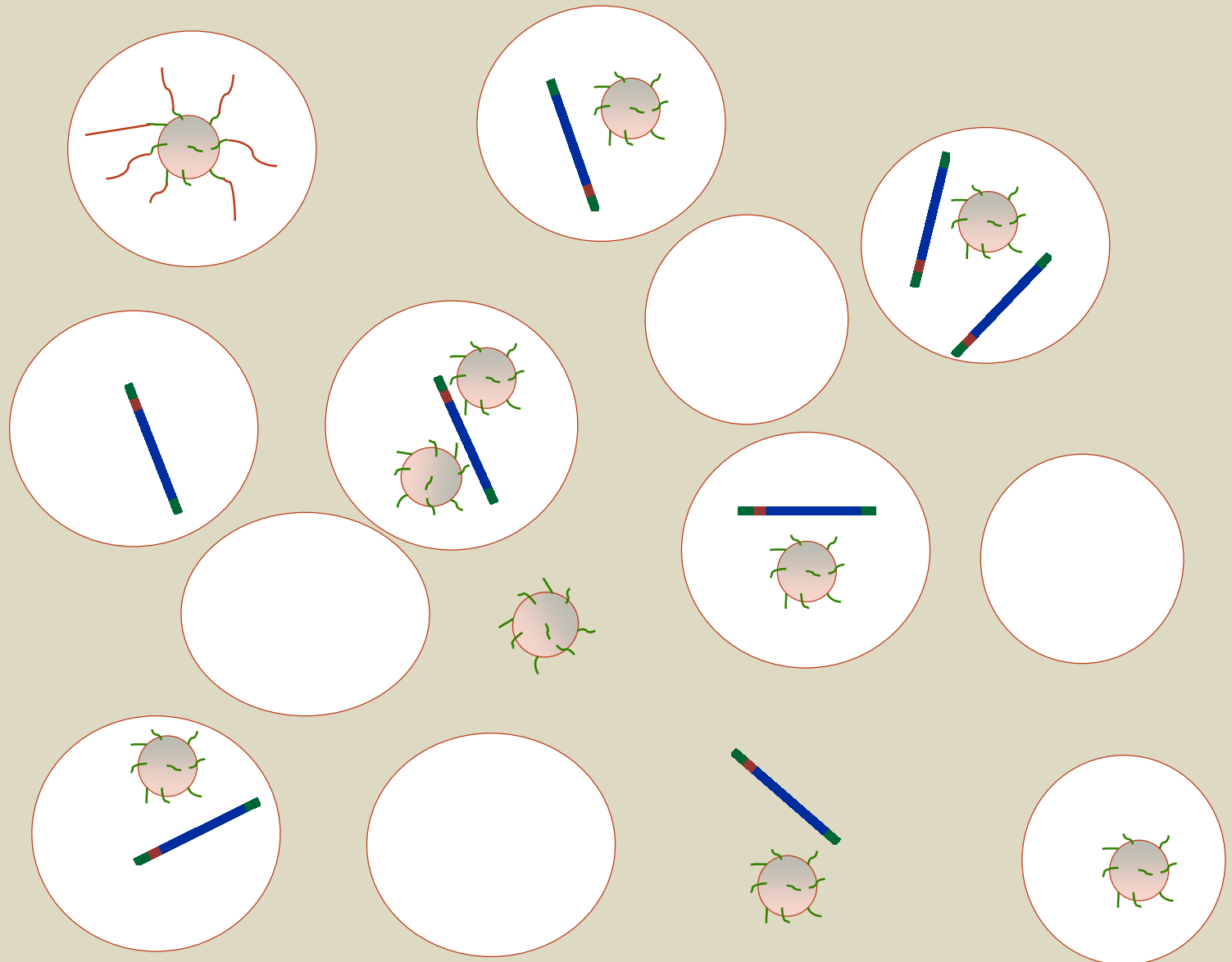


# Shotgun sequencing by PGM/454

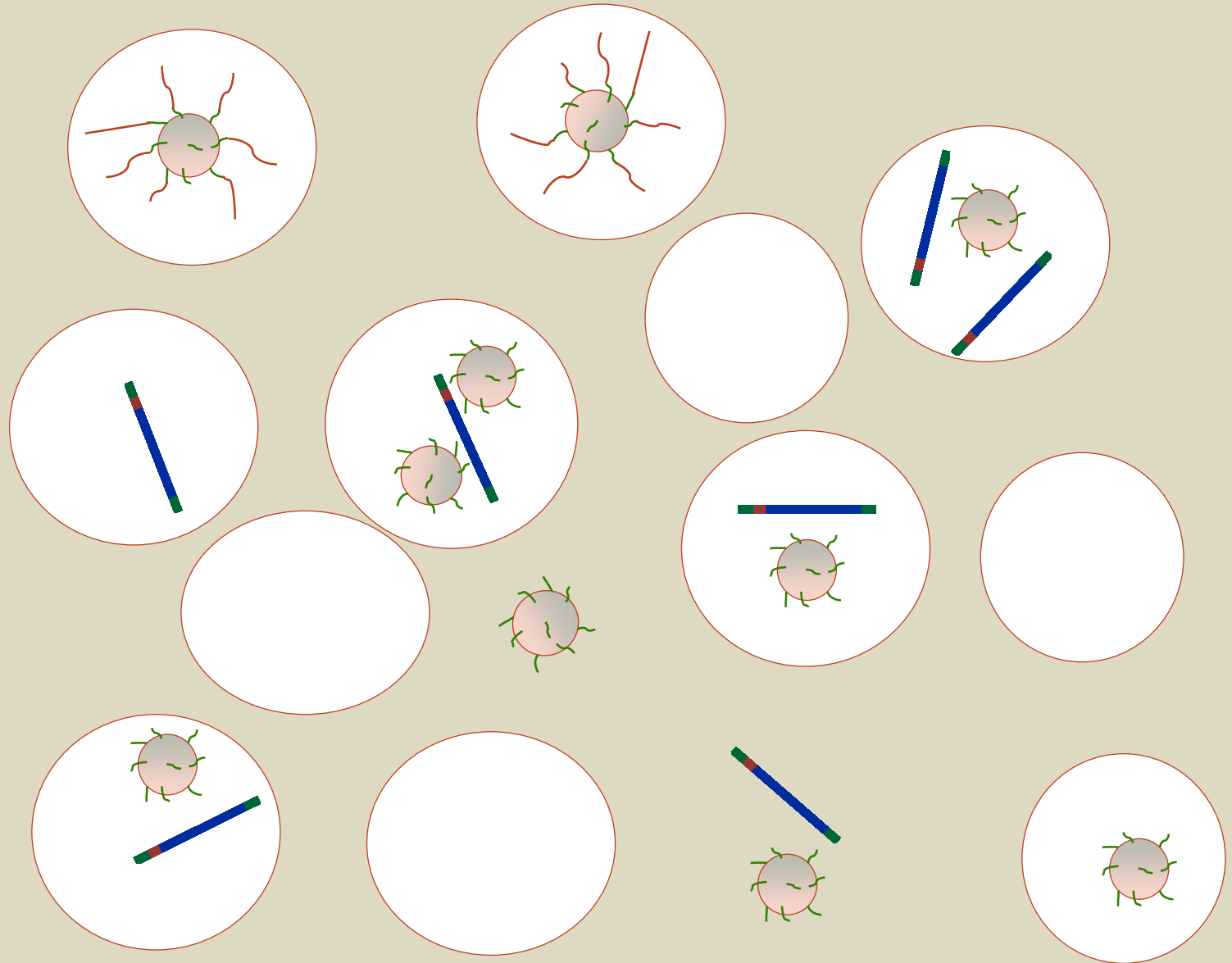




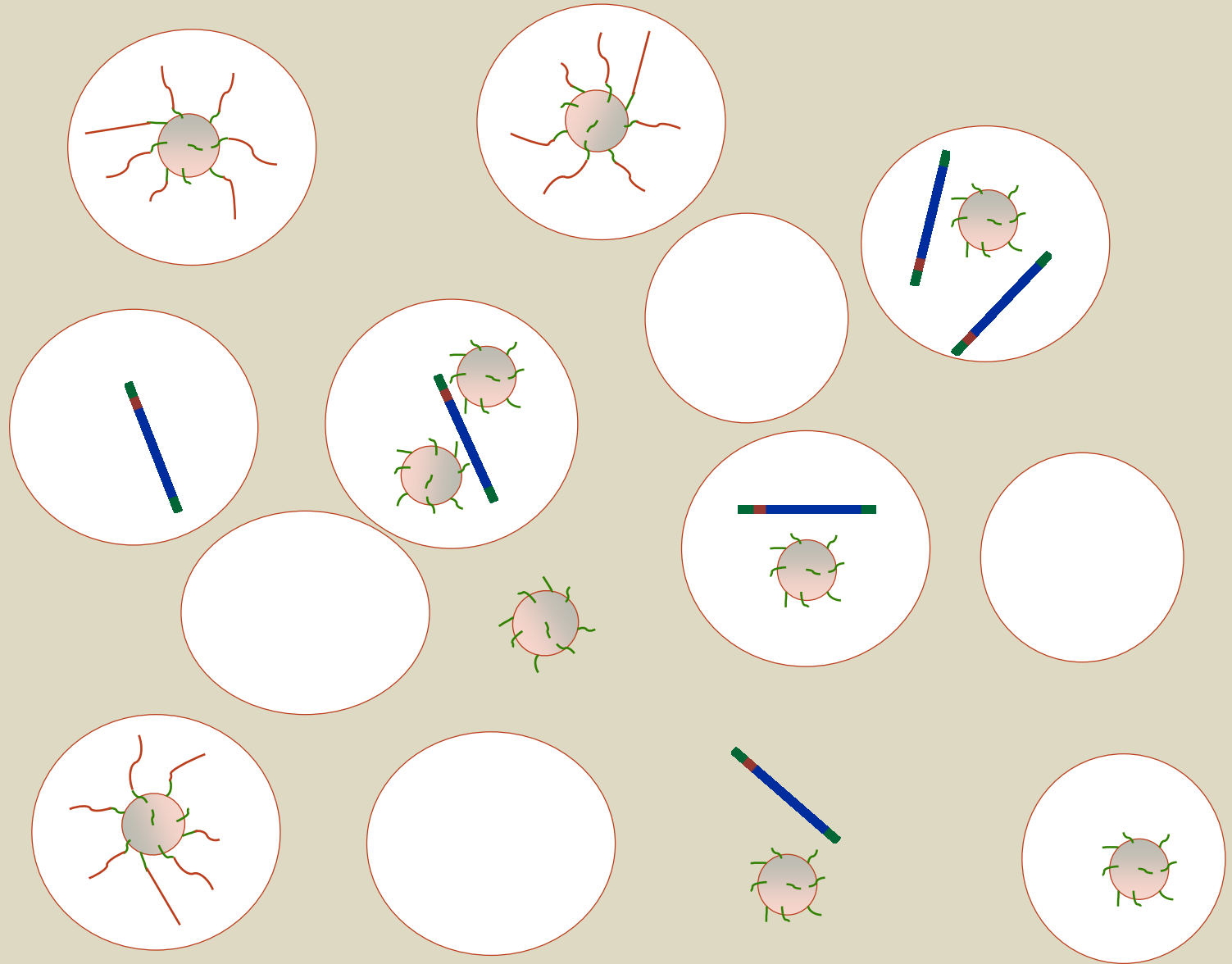
# Shotgun sequencing by PGM/454



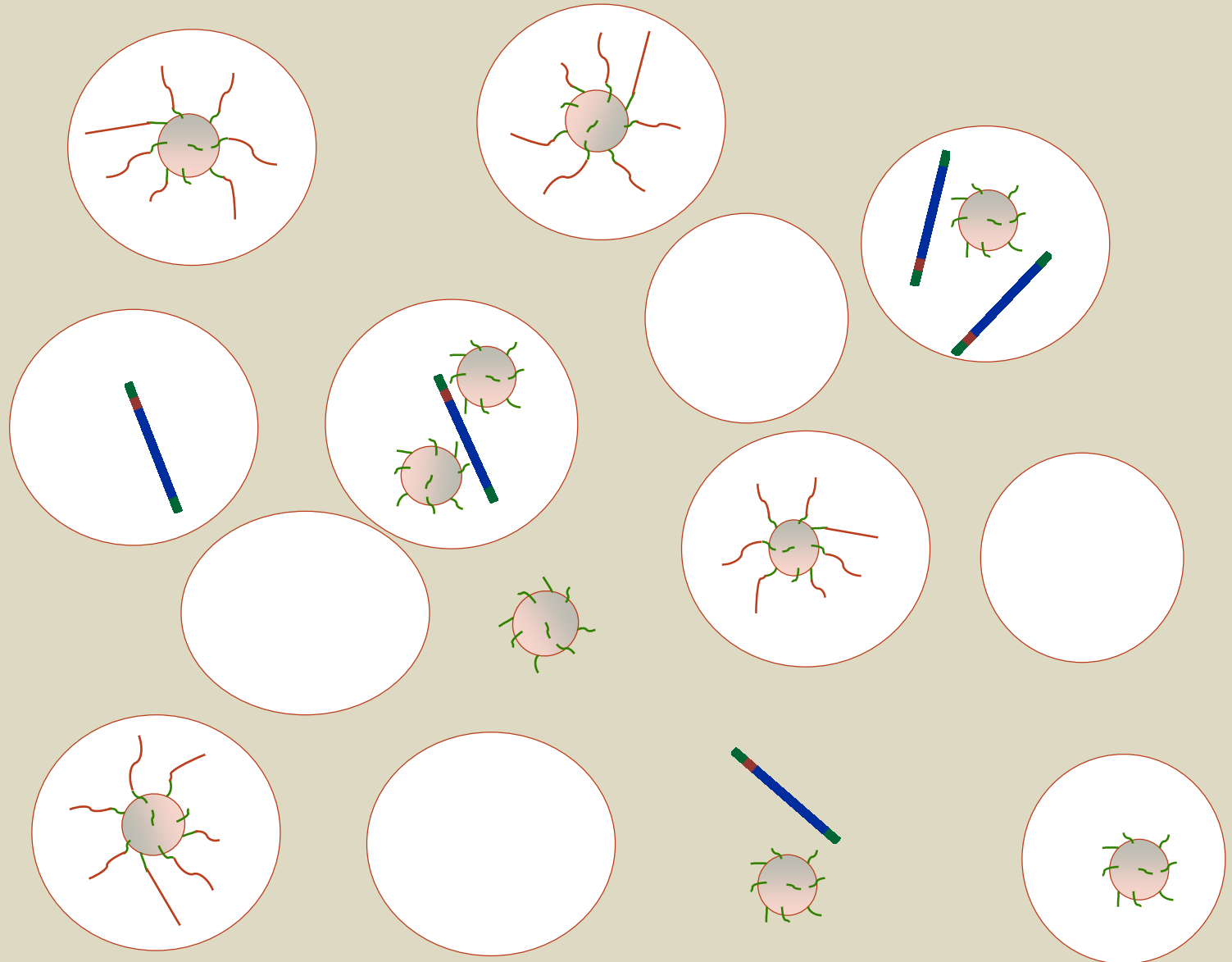
# Shotgun sequencing by PGM/454



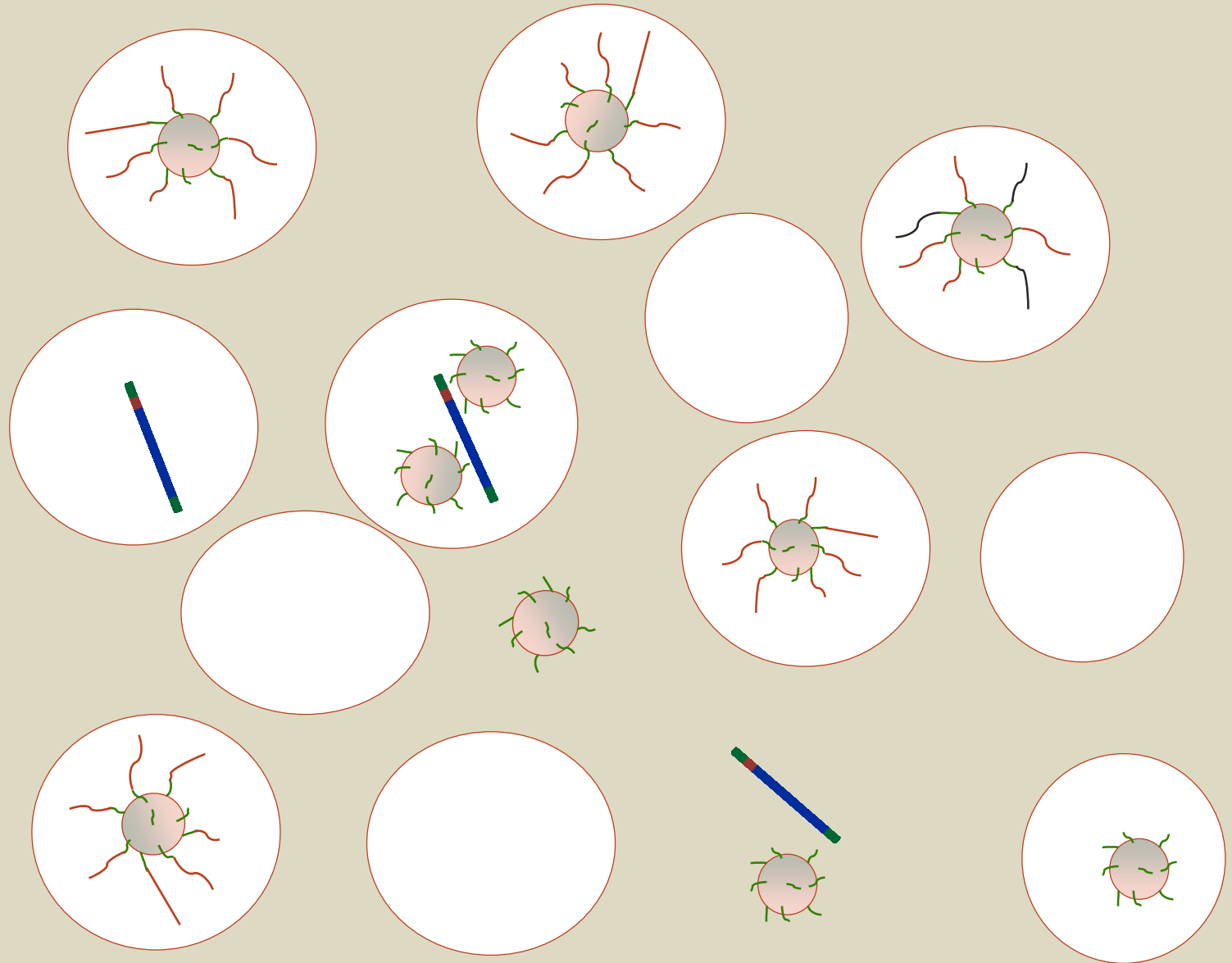
# Shotgun sequencing by PGM/454



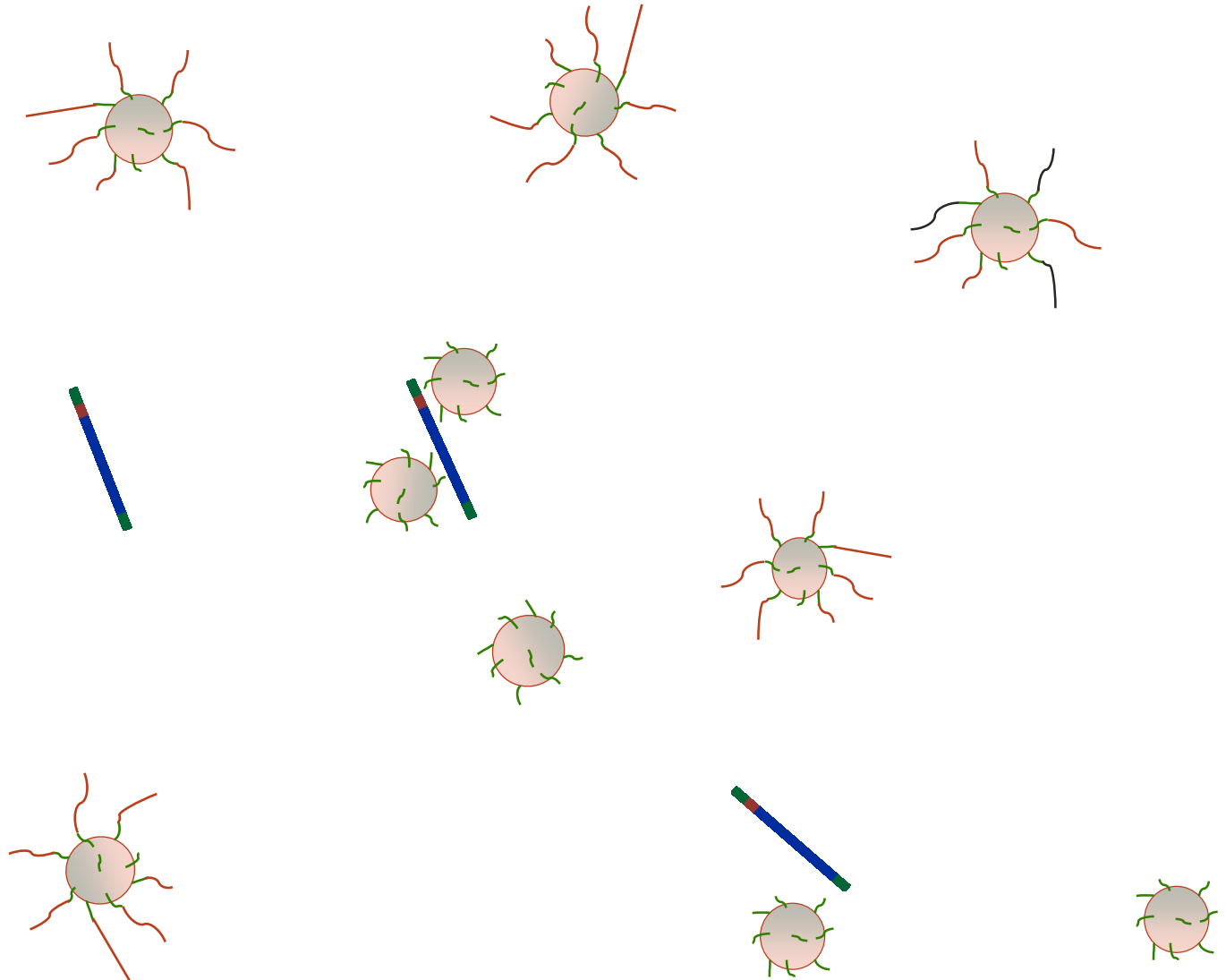
# Shotgun sequencing by PGM/454



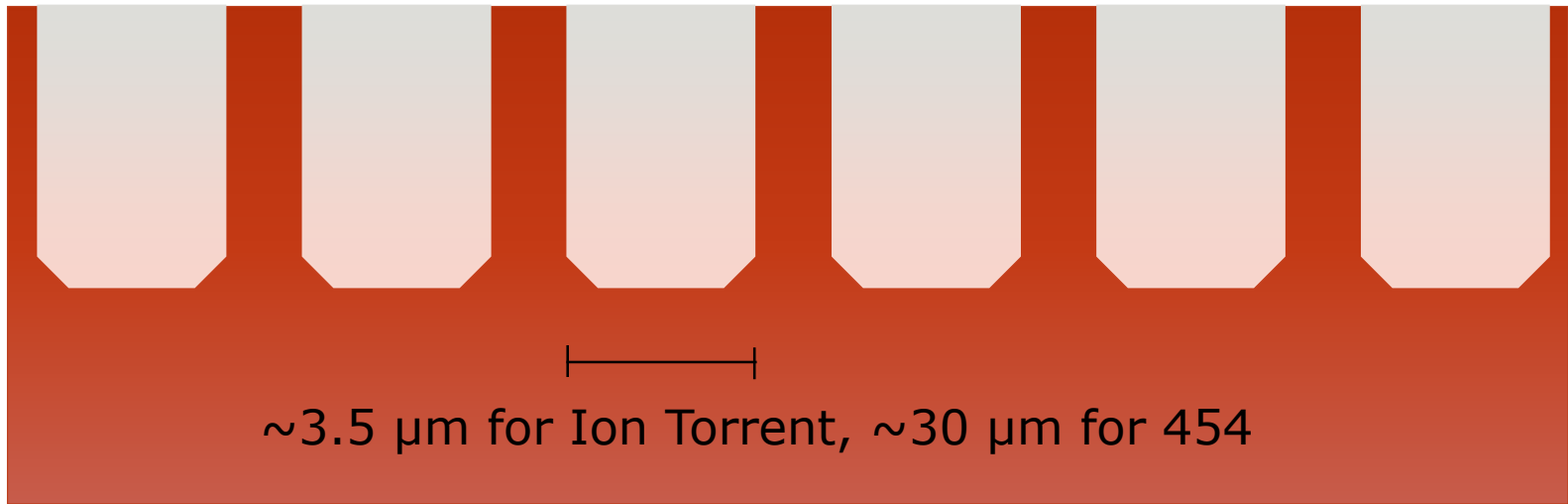
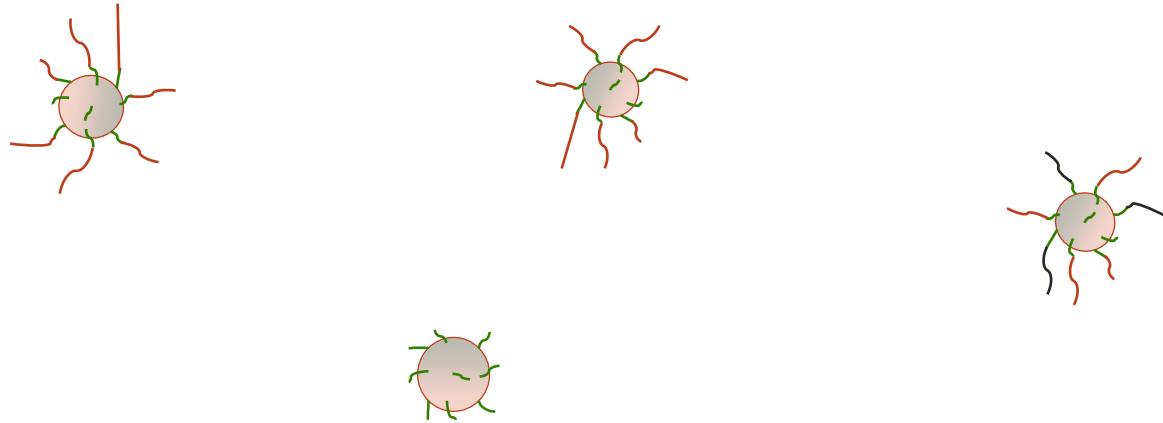
# Shotgun sequencing by PGM/454



# Shotgun sequencing by PGM/454

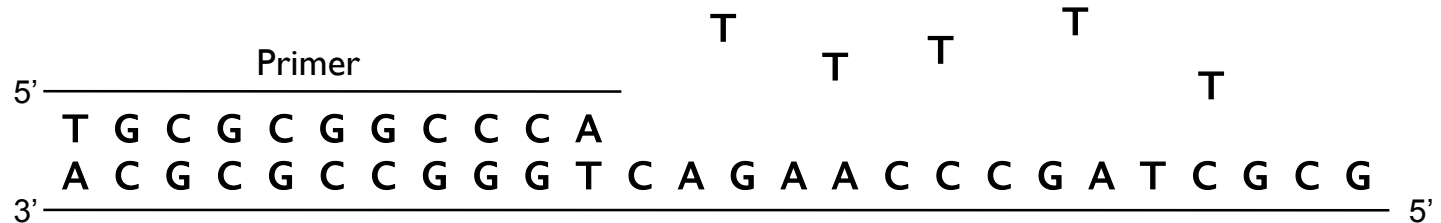


# Shotgun sequencing by PGM/454



# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:



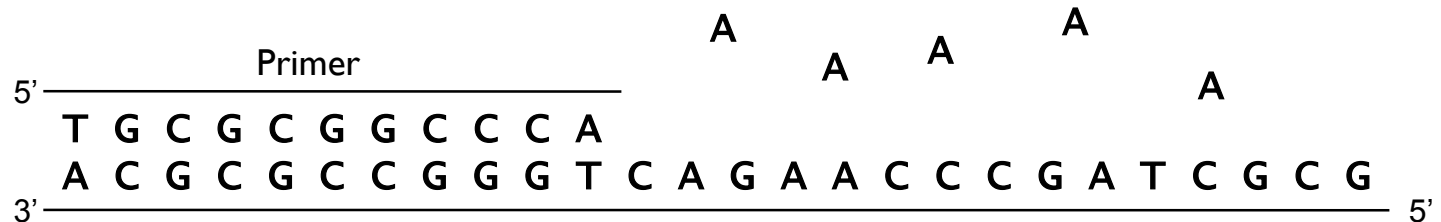
If that nucleotide is incorporated, enzymes turn by-products into light:





# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

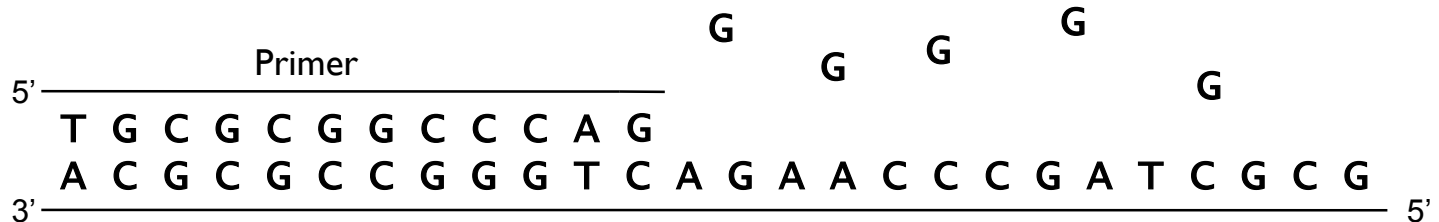


If that nucleotide is incorporated, enzymes turn by-products into light:

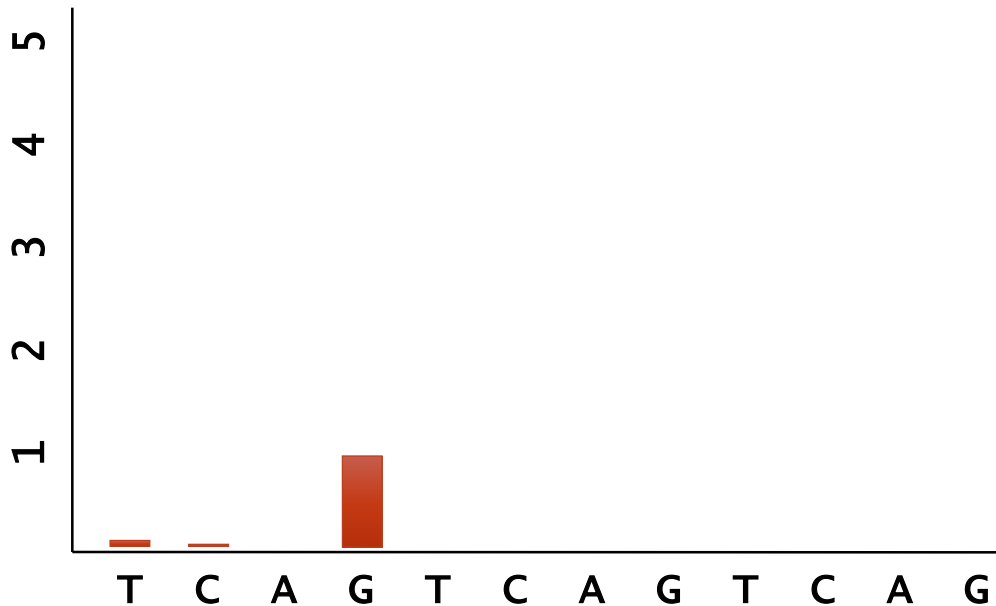


# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

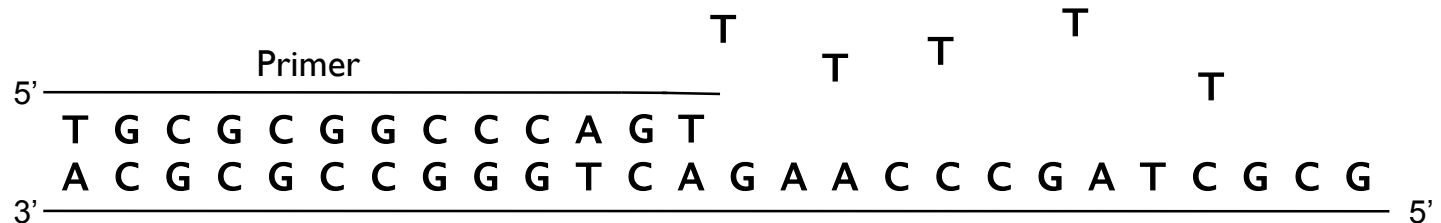


If that nucleotide is incorporated, enzymes turn by-products into light:

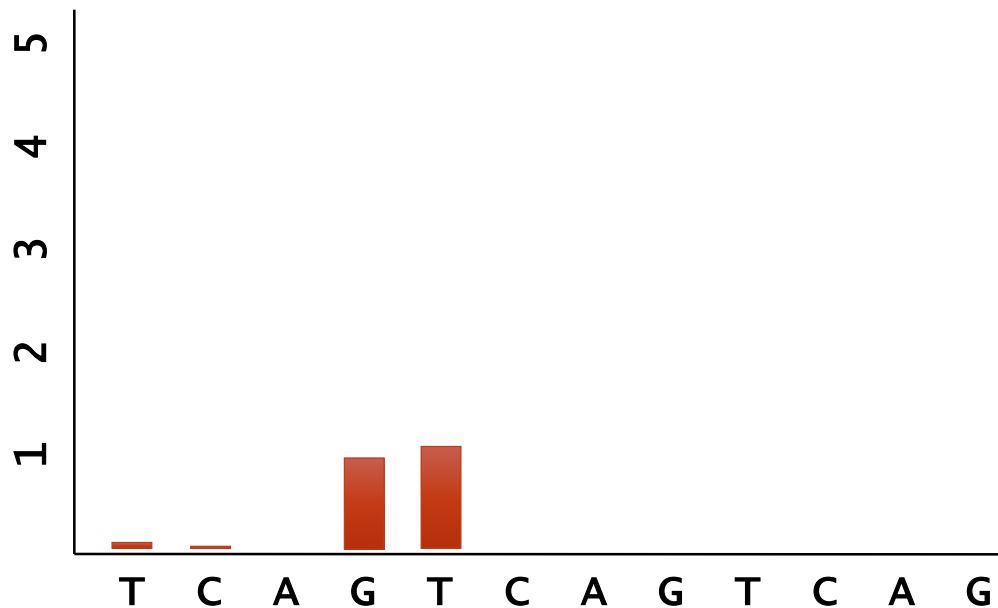


# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

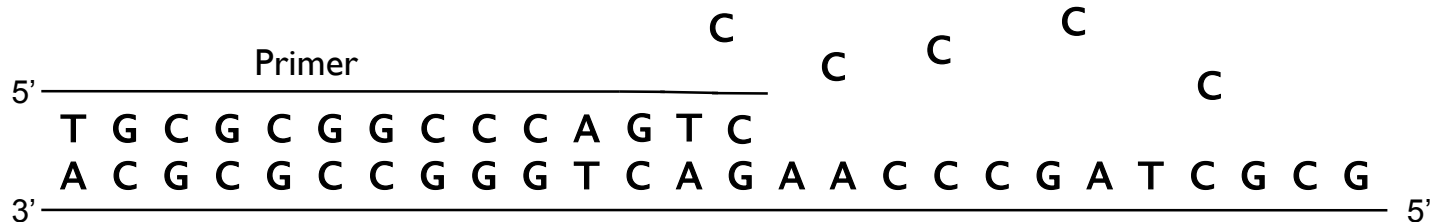


If that nucleotide is incorporated, enzymes turn by-products into light:

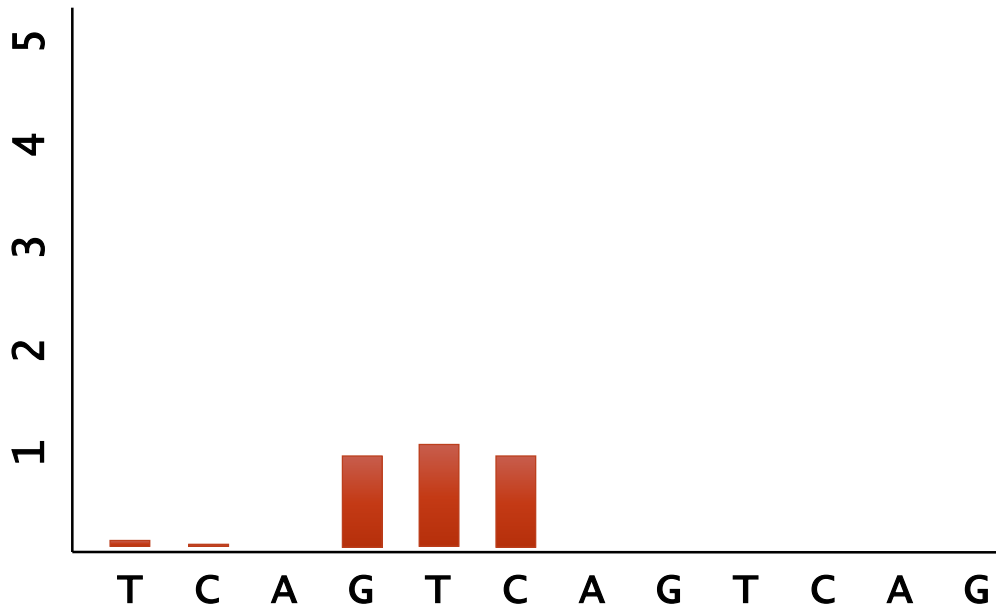


# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

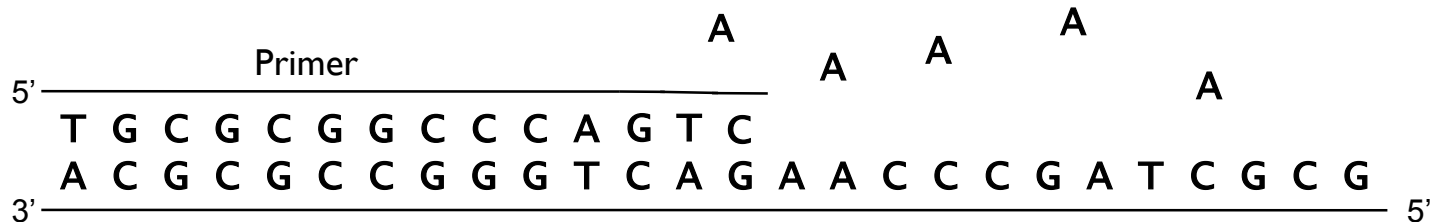


If that nucleotide is incorporated, enzymes turn by-products into light:

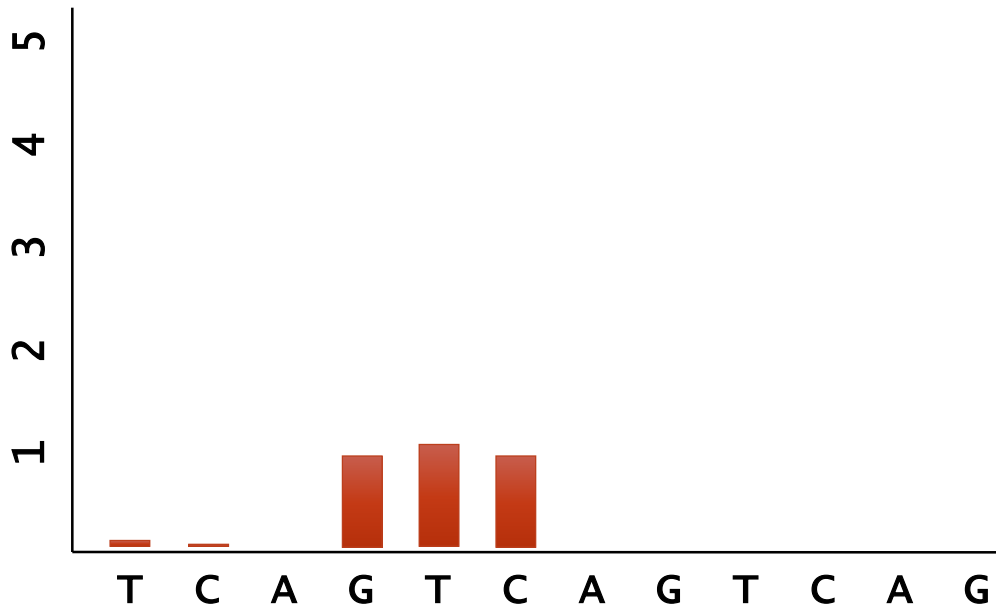


# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

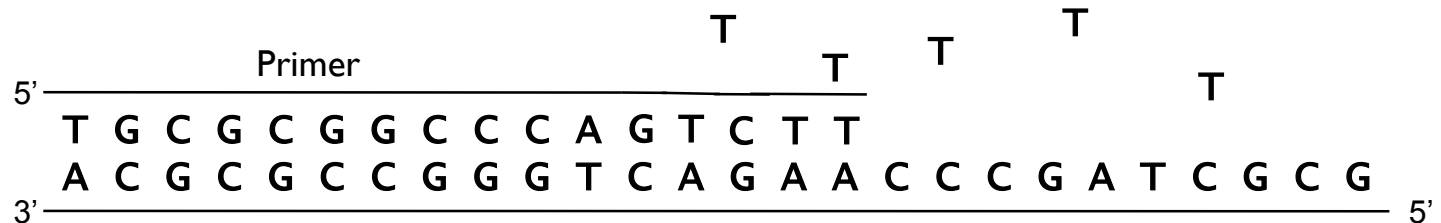


If that nucleotide is incorporated, enzymes turn by-products into light:

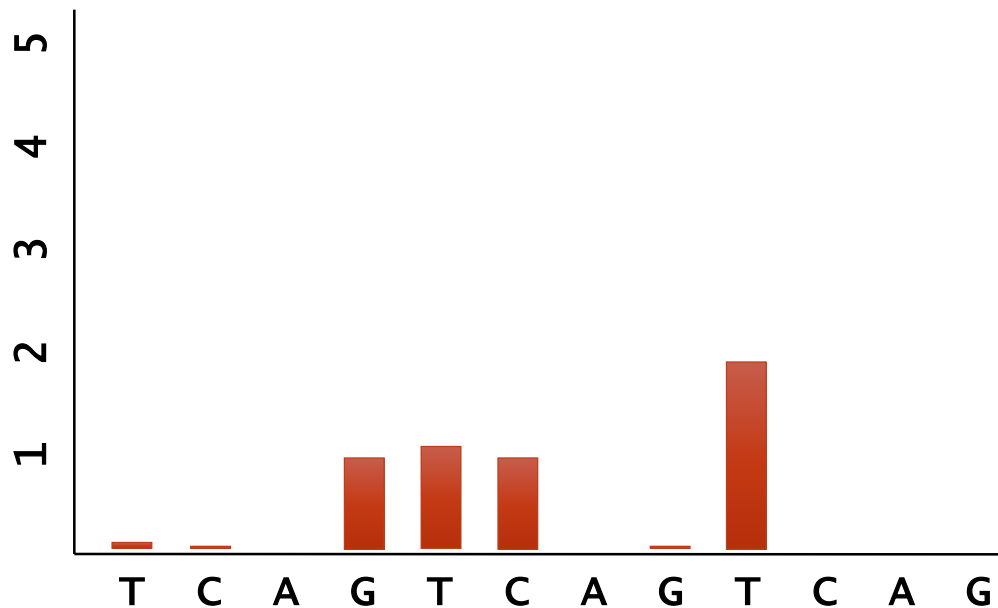


# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

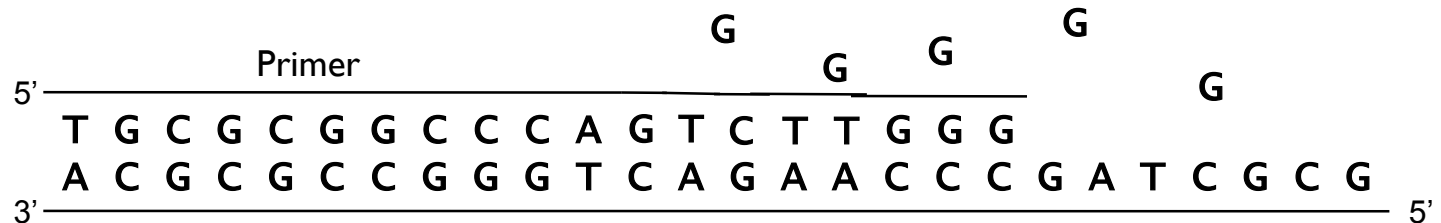


If that nucleotide is incorporated, enzymes turn by-products into light:

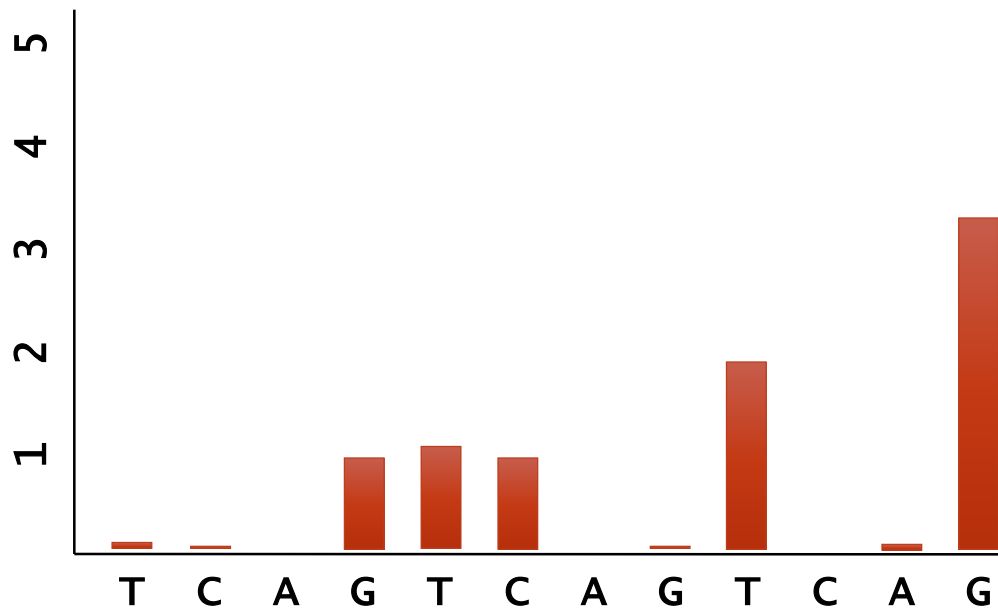


# Shotgun sequencing by PGM/454

Only give polymerase one nucleotide at a time:

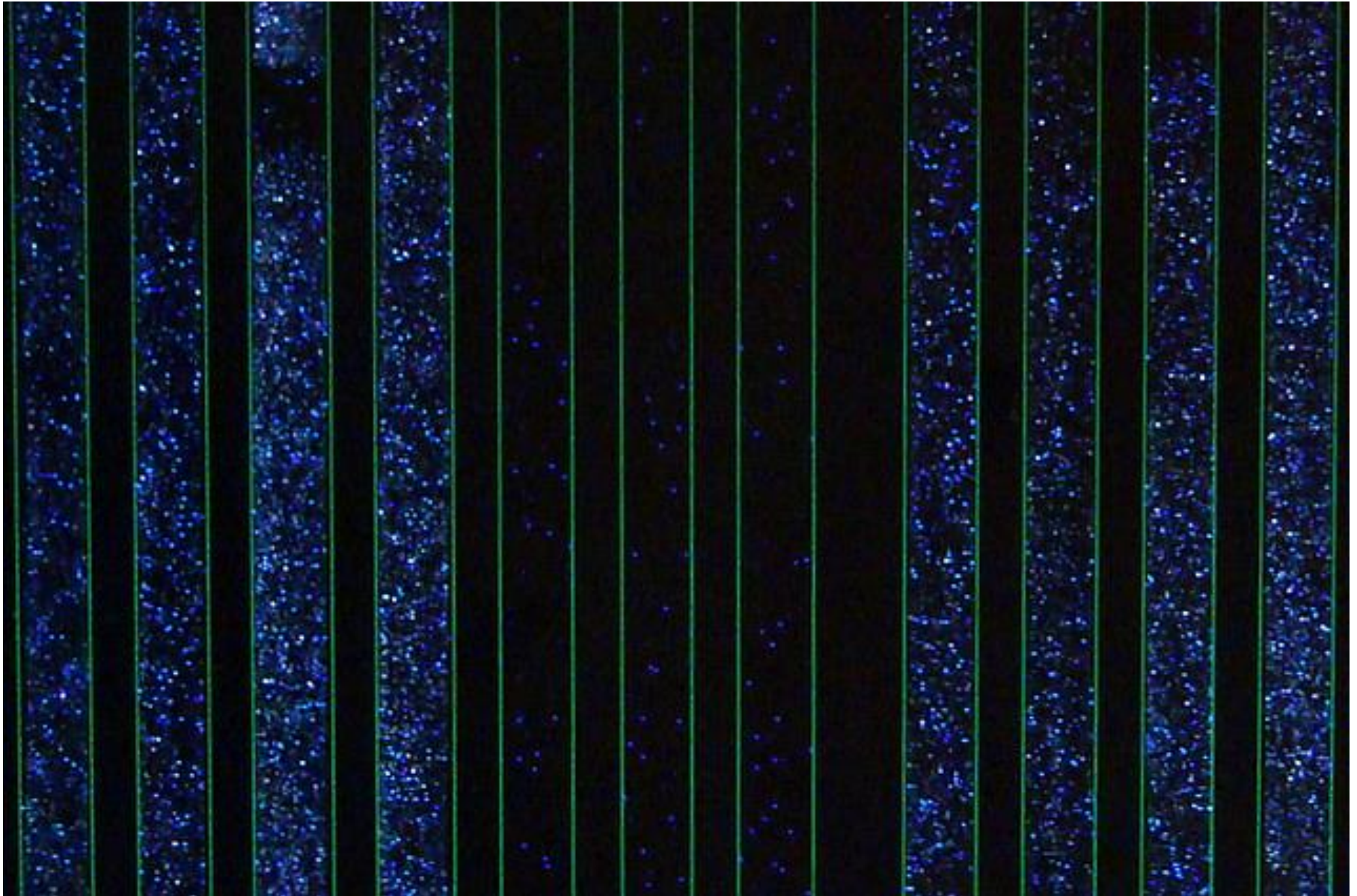


If that nucleotide is incorporated, enzymes turn by-products into light:



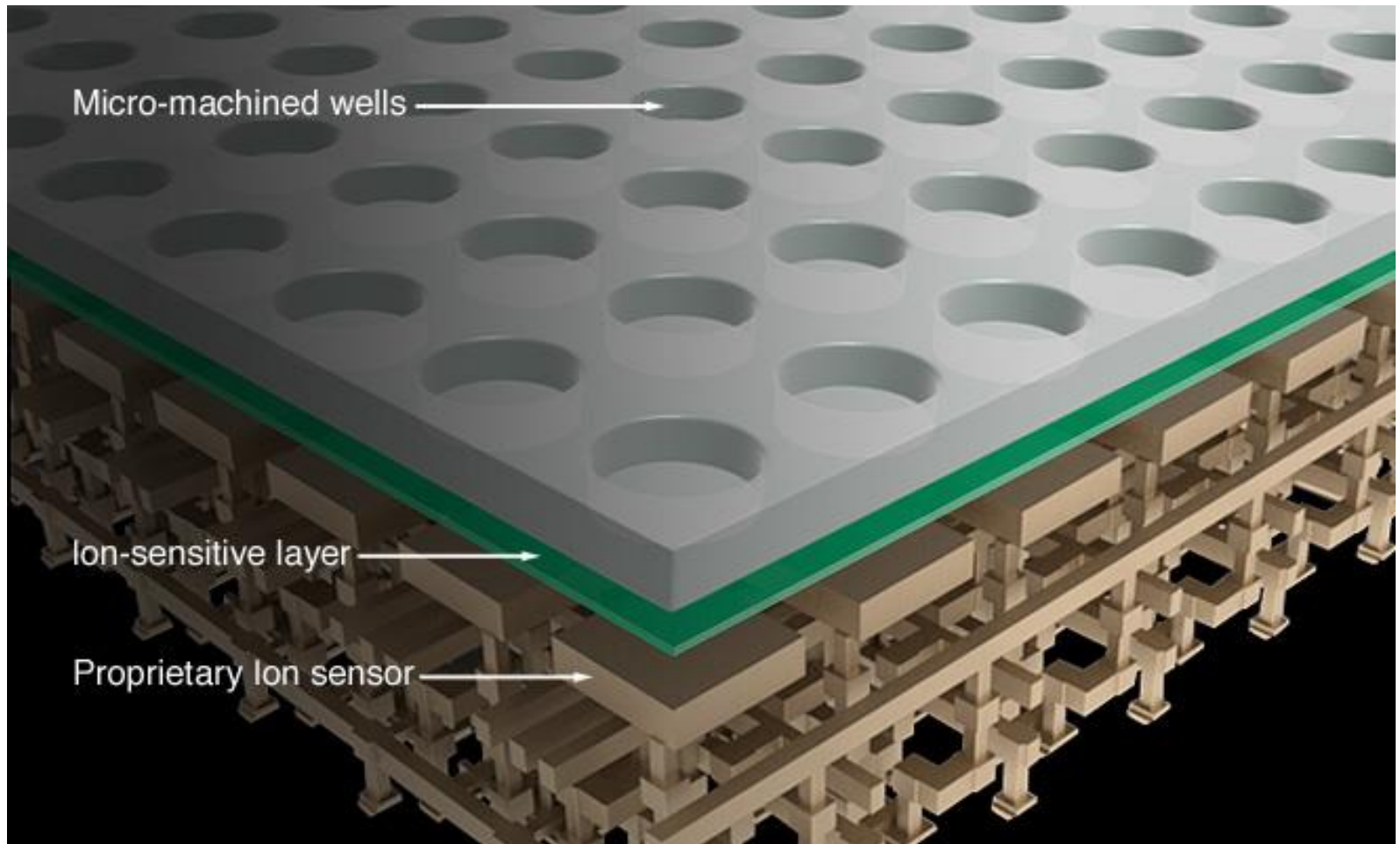
The real power of this method is that it can take place in millions of tiny wells in a single plate at once.

# Raw 454 data

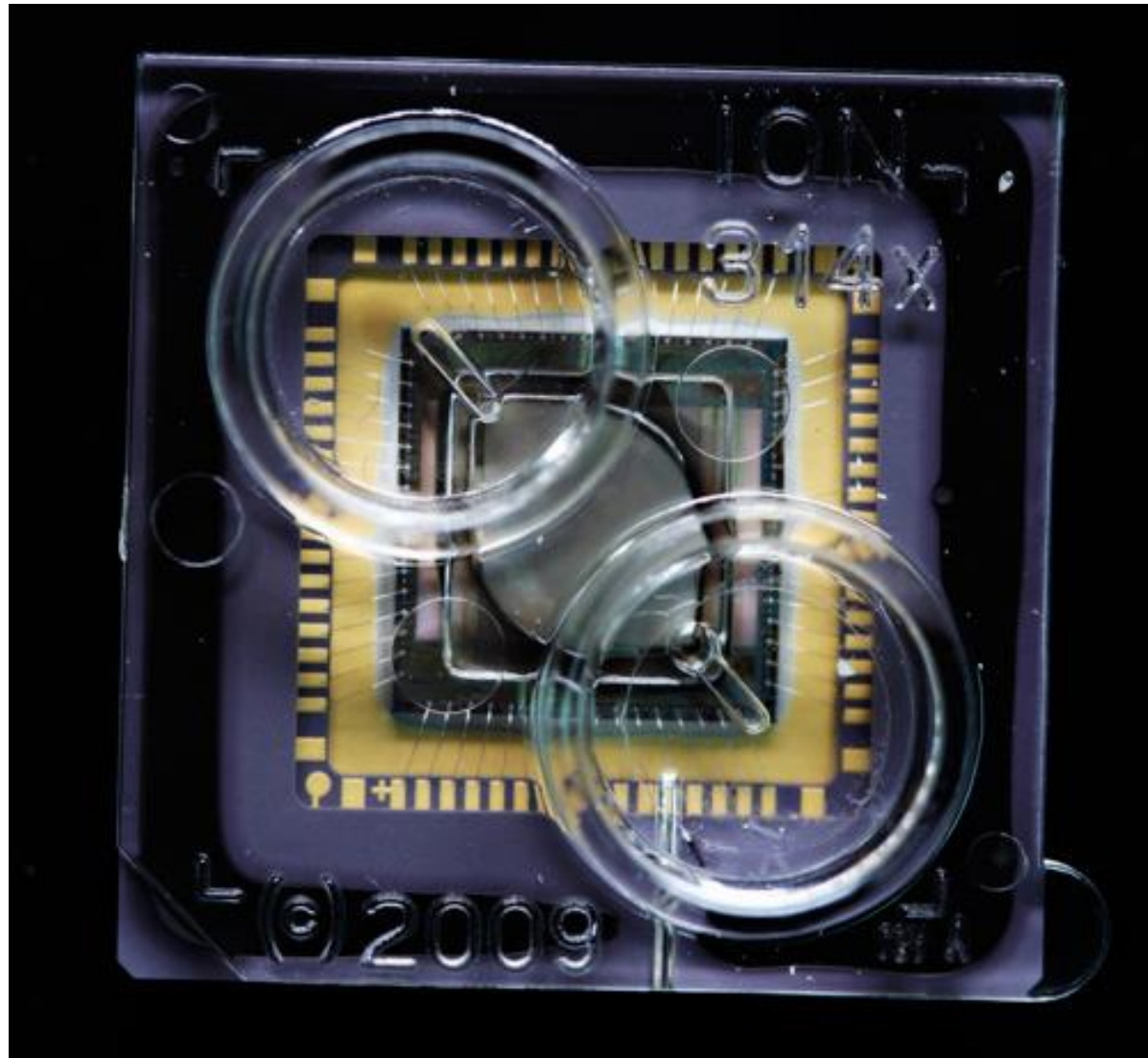




# Ion Torrent Sequencing

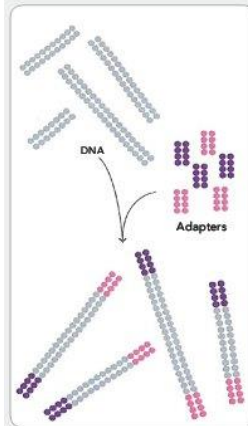


# Ion Torrent Sequencing



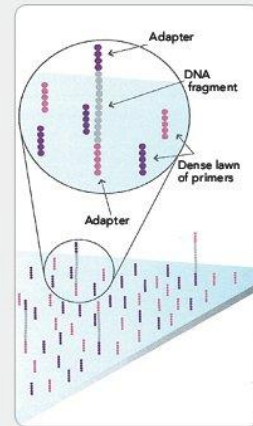
# Illumina Sequencing

1. PREPARE GENOMIC DNA SAMPLE



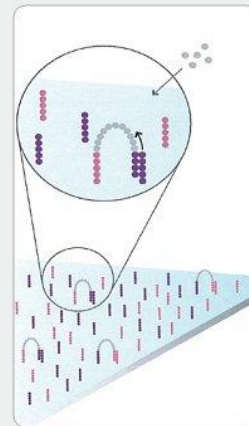
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



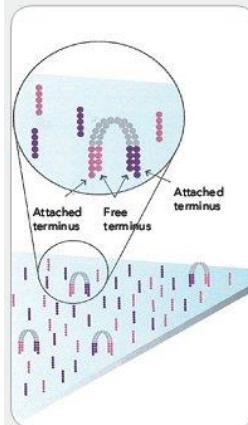
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

3. BRIDGE AMPLIFICATION



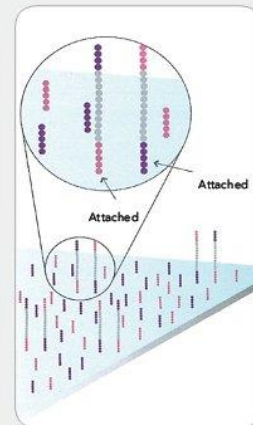
Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

4. FRAGMENTS BECOME DOUBLE STRANDED



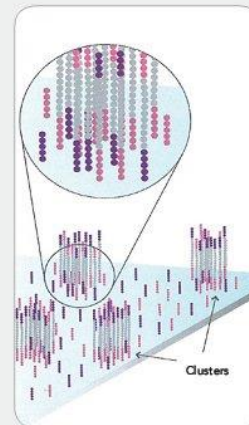
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



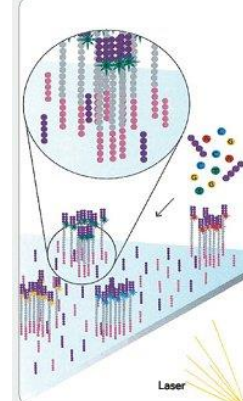
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

7. DETERMINE FIRST BASE



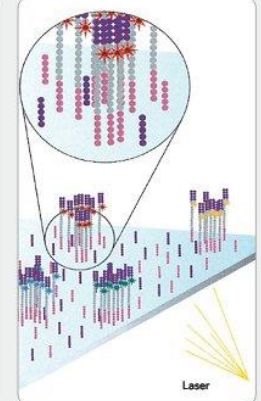
First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.

8. IMAGE FIRST BASE



After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

9. DETERMINE SECOND BASE



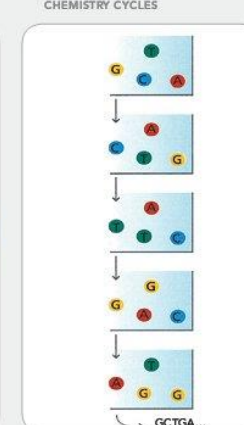
Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

10. IMAGE SECOND CHEMISTRY CYCLE



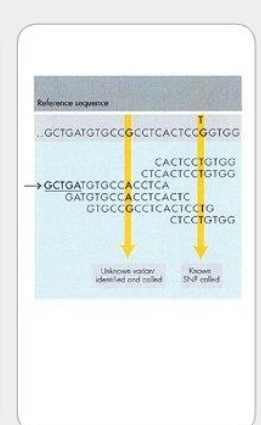
After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

12. ALIGN DATA



Align data, compare to a reference, and identify sequence differences.

# Next-Gen Sequencing

Take home message: Massively Parallel

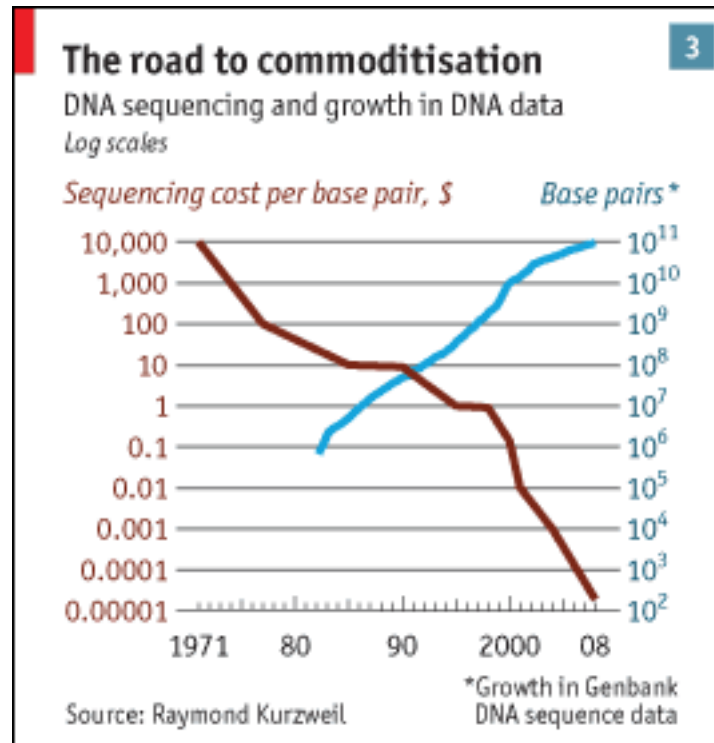
1,000 monkeys at 1,000 typewriters is nothing

We're talking 100,000 to 100 million concurrent reads

# Overview

- Prologue: Assembly
- The Past: Sanger
- The Present: Next-Gen (454, Illumina, ...)
- The Future: ? (**Nanopore, MinION, Single-molecule**)

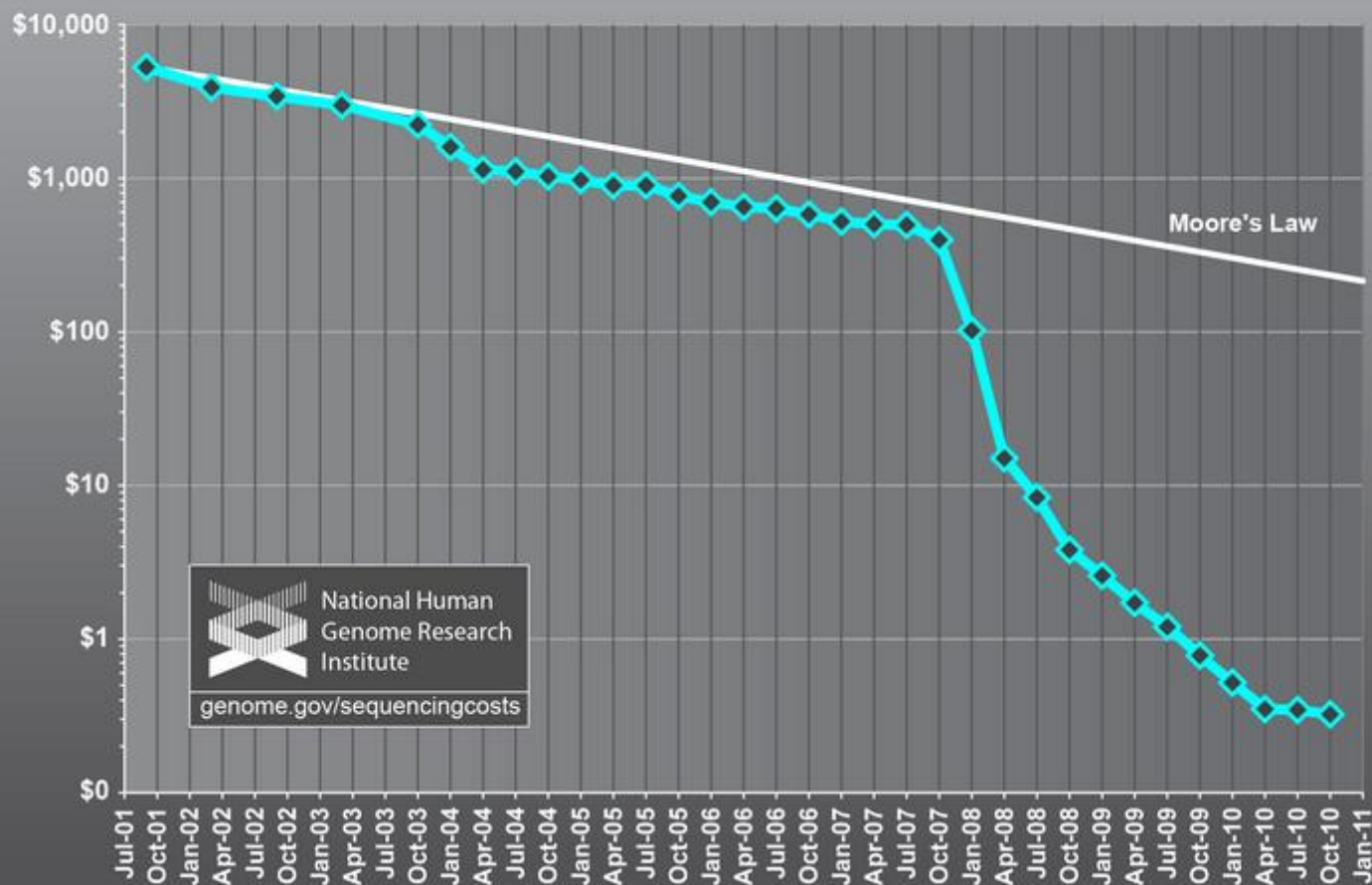
# DNA Sequencing over Time



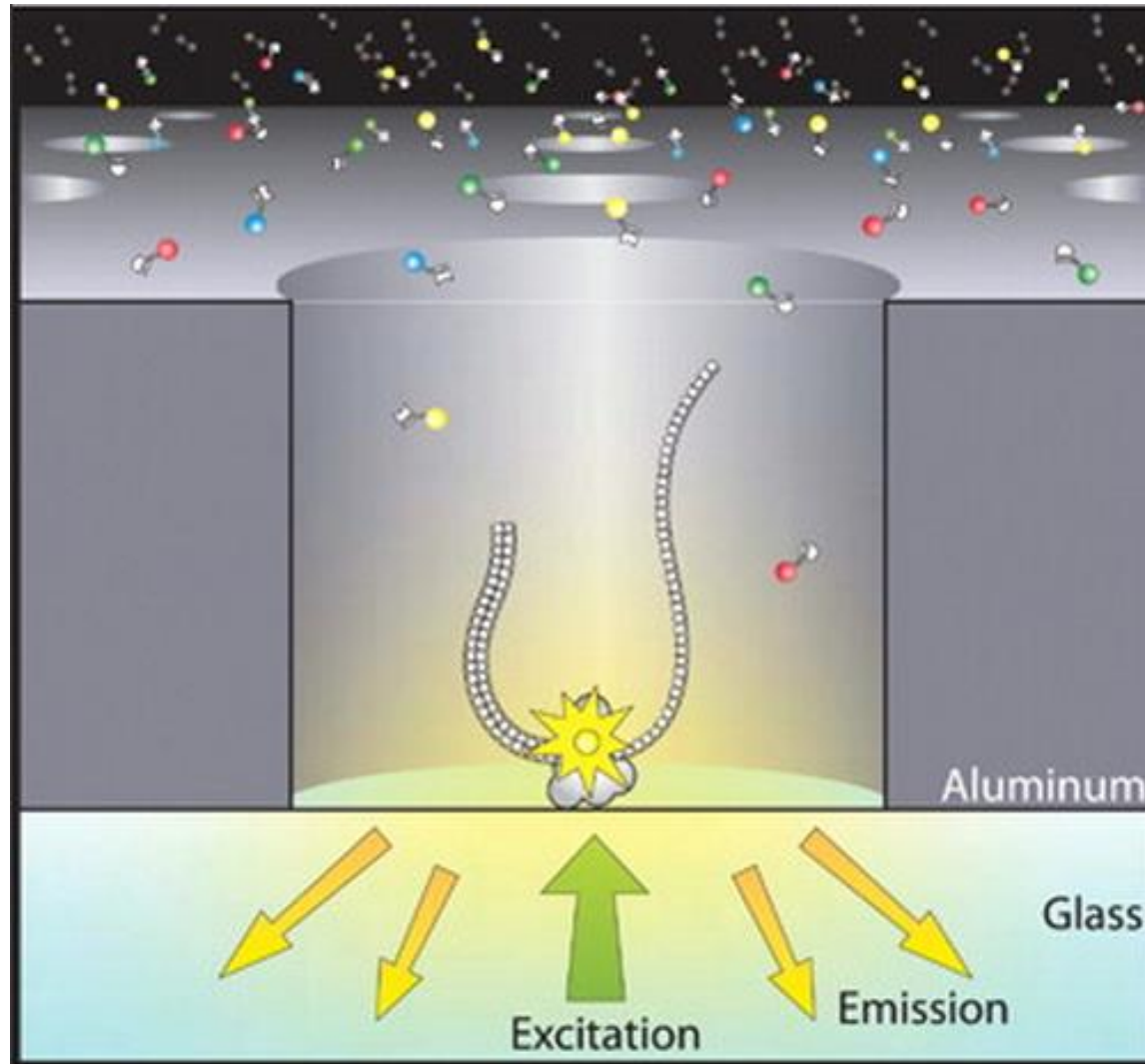
from *The Economist*



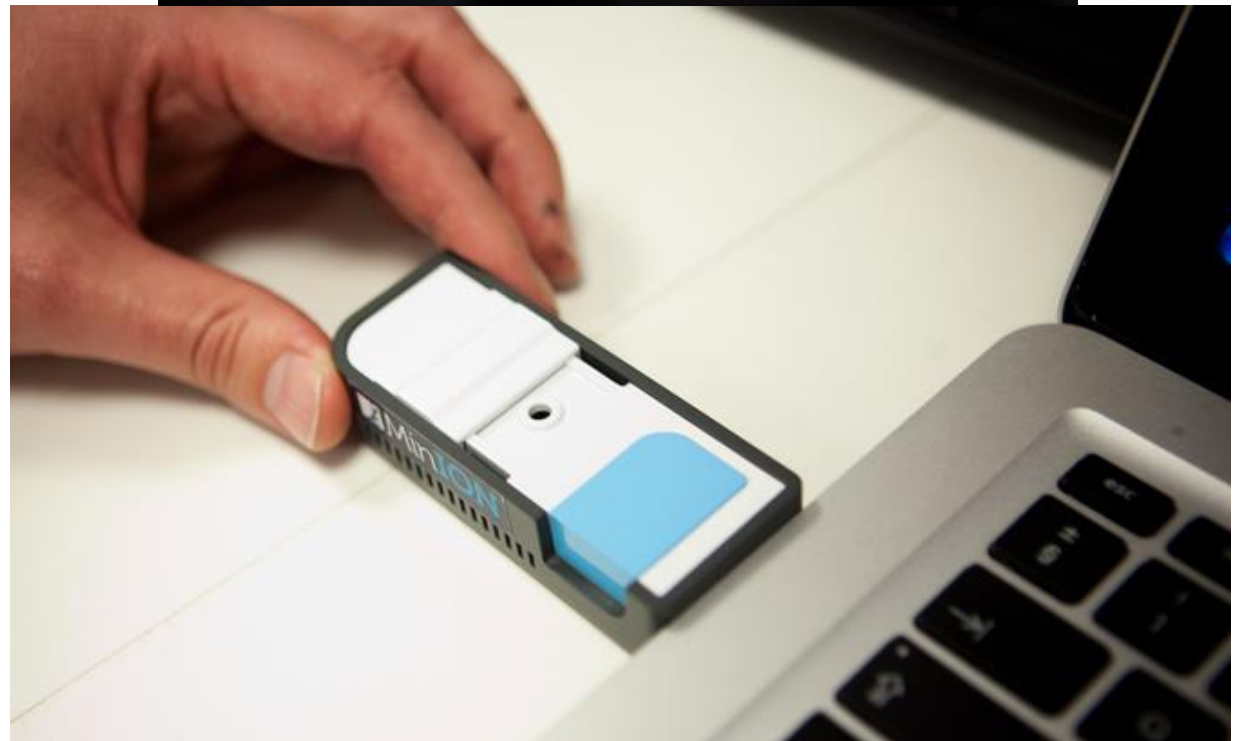
## Cost per Megabase of DNA Sequence



# Single Molecule Sequencing







“The MinION has been used to successfully read the genome of a lambda bacteriophage, which has 48,500-ish base pairs, twice during one pass. That's impressive, because reading 100,000 base pairs during a single DNA capture has never been managed before using traditional sequencing techniques.

The operational life of the MinION is only about six hours, but during that time it can read more than 150 million base pairs. That's somewhat short of the larger human chromosomes (which contain up to 250 million base pairs), but Oxford Nanopore has also introduced GridION -- a platform where multiple cartridges can be clustered together. The company reckon that a 20-node GridION setup can sequence a complete human genome in just 15 minutes.”

—*Wired*

# **(Relevant) Trivia**

How many base pairs (bp) are there in a human genome?

**~3 billion (haploid)**

How much did it cost to sequence the first human genome?

**~\$2.7 billion**

How long did it take to sequence the first human genome?

**~13 years**

When was the first human genome sequence complete?

**2000-2003**

Whose genome was it?

**Several people's, but actually mostly a dude from Buffalo**

# Final Thoughts

- DNA sequencing is becoming vastly faster and more affordable
- Generating data is no longer the bottleneck, understanding it is
- Bioinformatics types should be in high demand in the near future

# Comparing Different Technologies

## Sanger Sequencing

Advantages	Disadvantages
<ul style="list-style-type: none"><li>Lowest error rate</li><li>Long read length (~750 bp)</li><li>Can target a primer</li></ul>	<ul style="list-style-type: none"><li>High cost per base</li><li>Long time to generate data</li><li>Need for cloning</li><li>Amount of data per run</li></ul>

# Comparing Different Technologies

## 454 Sequencing

Advantages	Disadvantages
<p>Low error rate</p> <p>Medium read length (~400-600 bp)</p>	<p>Relatively high cost per base</p> <p>Must run at large scale</p> <p>Medium/high startup costs</p>

# Comparing Different Technologies

## Ion Torrent Sequencing

Advantages	Disadvantages
<ul style="list-style-type: none"><li>Low startup costs</li><li>Scalable (10 – 1000 Mb of data per run)</li><li>Medium/low cost per base</li><li>Low error rate</li><li>Fast runs (&lt;3 hours)</li></ul>	<ul style="list-style-type: none"><li>New, developing technology</li><li>Cost not as low as Illumina</li><li>Read lengths only ~100-200 bp so far</li></ul>

# Comparing Different Technologies

## Illumina Sequencing

Advantages	Disadvantages
Low error rate Lowest cost per base Tons of data	Must run at very large scale Short read length (50-75 bp) Runs take multiple days High startup costs De Novo assembly difficult



# Comparing Different Technologies

## PacBio Sequencing

Advantages	Disadvantages
<p>Can use single molecule as template</p> <p>Potential for very long reads (several kb+)</p>	<p>High error rate (~10-15%)</p> <p>Medium/high cost per base</p> <p>High startup costs</p>

# Sequencing illumina®

Platform	Reads per run	Read length (mode or average)	Bases per run (gigabases)
ABI Sanger	96	800	0.0000768
454	1 millions	700	0.7
IonTorrent	75 millions	200	15
SOLiD	3 billions	75	320
<b>Illumina</b>	<b>600 millions to 6 milliards</b>	<b>100 to 300</b>	<b>7.5 to 2 000</b>



MiniSeq System  
Up to 7.5 Gb



MiSeq Series +  
Up to 15 Gb



NextSeq Series +  
Up to 120 Gb



NextSeq Series +  
Up to 120 Gb



HiSeq Series +  
Up to 750 Gb



HiSeq X Series†  
Up to 800 Gb



NovaSeq Series +  
Up to 2 Tb

# Sequencing

illumina®

Platform	Reads per run	Read length (mode or average)	Bases per run (gigabases)
ABI Sanger	96	800	0.0000768
454	1 millions	700	0.7
IonTorrent	75 millions	200	15
SOLiD	3 billions	75	320
<b>Illumina</b>	<b>600 millions to 6 milliards</b>	<b>100 to 300</b>	<b>7.5 to 2 000</b>



MiniSeq System

50,000 \$



MiSeq Series +

99,000 \$



NextSeq Series +

250,000 \$



NextSeq Series +

250,000 \$



HiSeq Series +

690,000 \$-  
900,000 \$



HiSeq X Series†

1,000,000 \$  
(X10= 10,000,000 \$)



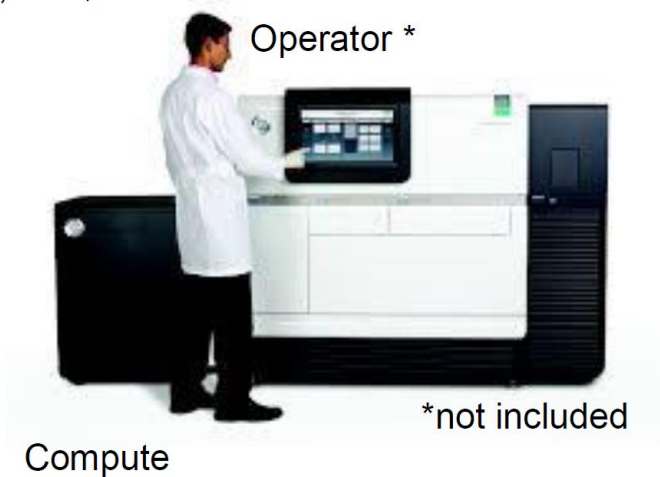
NovaSeq Series +

850,000 \$-  
985,000 \$

# Single molecule Long Read Sequencing



RSII 700,000 \$



Platform	Year	Reads per run	Read length (mode or average)	Bases per run (gigabases)
ABI Sanger	2002	96	800	0.0000768
454	2011	1 millions	700	0.7
SOLiD	2013	3 milliards	75	320
IonTorrent	2015	75000000	200	15
Illumina	2016	600 millions to 6 milliards	100 to 300	7.5 to 2 000
PacBio	2014	660000	13500	20

Sequel 350,000 \$



# Single molecule Long Read Sequencing



MinION Mk1: portable, real time biological analyses

Platform	Reads per run	Read length (mode or average)	Bases per run (gigabases)
ABI Sanger	96	800	0.0000768
454	1 millions	700	0.7
SOLiD	3 milliards	75	320
IonTorrent	75000000	200	15
Illumina	600 millions to 6 milliards	100 to 300	7.5 to 2 000
PacBio	660000	13500	12.000
<b>Oxford Nanopore</b>	<b>4400000</b>	<b>9545</b>	<b>42</b>

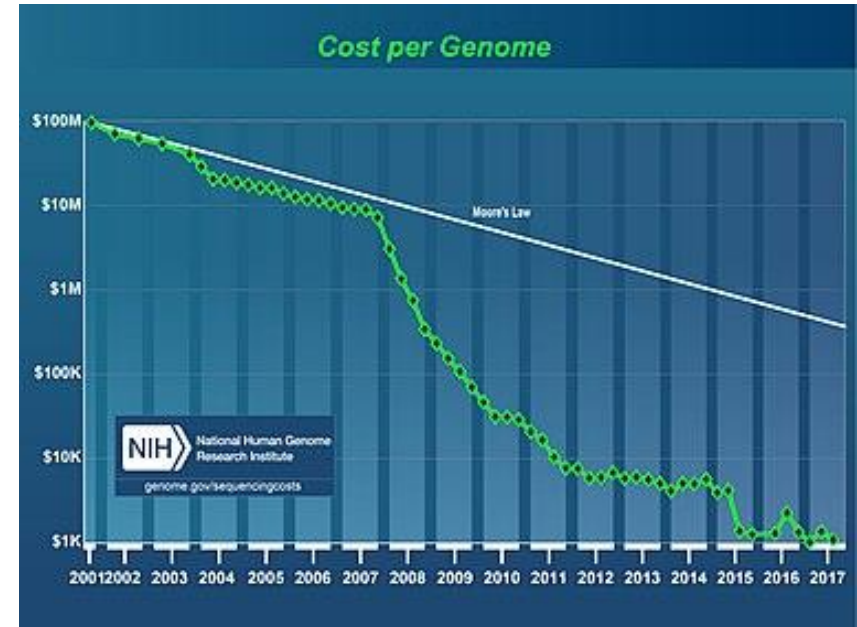
MinION 1,000 \$



# Sequencing Technologies are getting Faster and Cheaper

The use of nucleic acid sequencing has increased exponentially as the ability to sequence has become accessible to research and clinical labs all over the world

Several Sequencing Technologies Applications “-Seq” : RNA-Seq, Chip-Seq, SingleCell-Seq, etc.



This demand has driven the development of High Throughput Sequencing (HTS), that are becoming exponentially cheaper.

The exponential growth of genomic data unfortunately is not followed by an exponential growth of storage

# A new Human Genome every 6 min

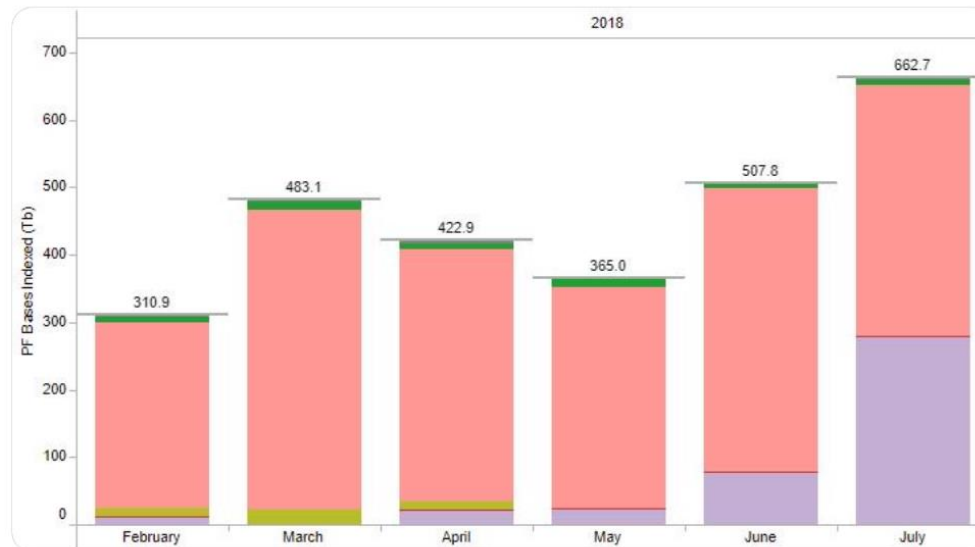


**Broad Genomics**  
@BroadGenomics

Segui

The sequencing lab just rang in its biggest month of sequence data generation EVER! 663 Terabases in the month of July. Equivalent to a human genome every 6 minutes.

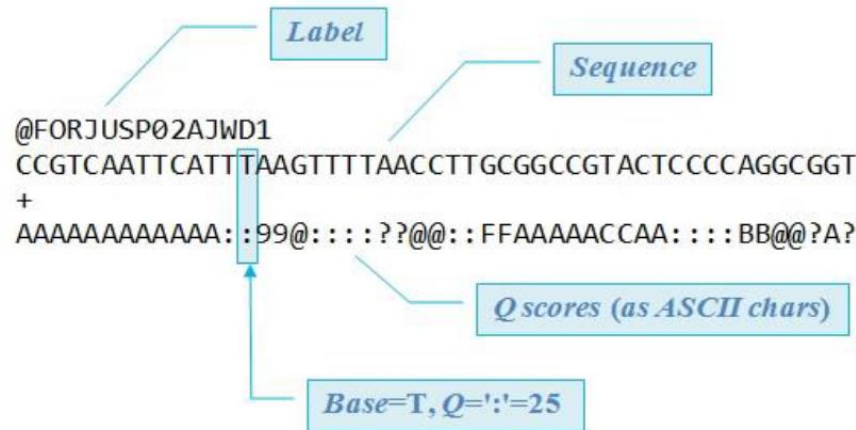
#NewBenchmark #QualityandQuantity  
#CantStopWontStop





# What is the output of HTS

## FastQ File



- In many applications **compression** is required to reduce space and improve throughput. e.g. The most popular reads mapper (BWA) can operate directly on the compressed FastQ.
- The **DNA sequence** exposes a **high redundancy**, especially on large reads collections with high coverage, and thus it is **highly compressible**
- **Quality values** have **much higher entropy** than a genomic sequence because their alphabeth usually span a much larger range of values (e.g. [1 - 40]), and they are not repetitive.
- When FASTQ files are compressed **quality values account for 70%** of the total space.



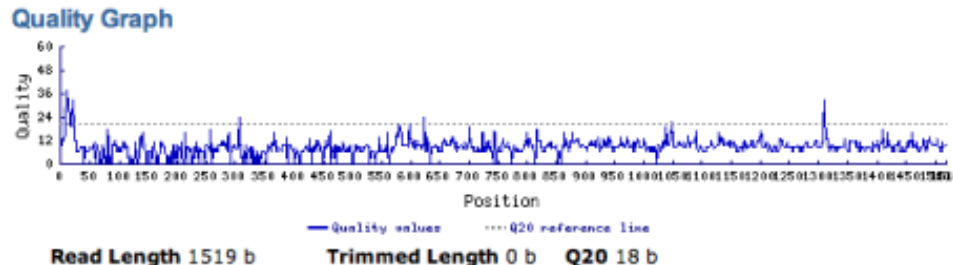
# Background on Quality Values

Quality Value	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%

- $Q(i)$ : Phred-scaled probability that the  $i$ -th base of the read being wrong
- $Q(i) = -10 \log_{10} \text{Prob}\{\text{the base } i \text{ of read is wrong}\}$

## Applications:

- SNP/Mutation Detection
- Removal of low-quality reads
- Reads Mapping
- Detection of Overlapping reads
- Error Correction
- Compression
- etc.



Quality graphs obtained from Geospiza's iFinch ([www.geospiza.com](http://www.geospiza.com)).

