



**The Abdus Salam
International Centre for Theoretical Physics**



2255-1

**2nd Conference on Systems Biology and New Sequencing
Techniques" (2-4 November), preceded by Introductory Lectures on
"Quantitative Approaches to Biological Problems" (31 October - 1 November)**

31 October - 4 November, 2011

DNA Sequencing Methods Past, Present, and Future

Arjang Hassibi
*University of Texas at Austin
USA*

DNA Sequencing Methods

Past, Present, and Future

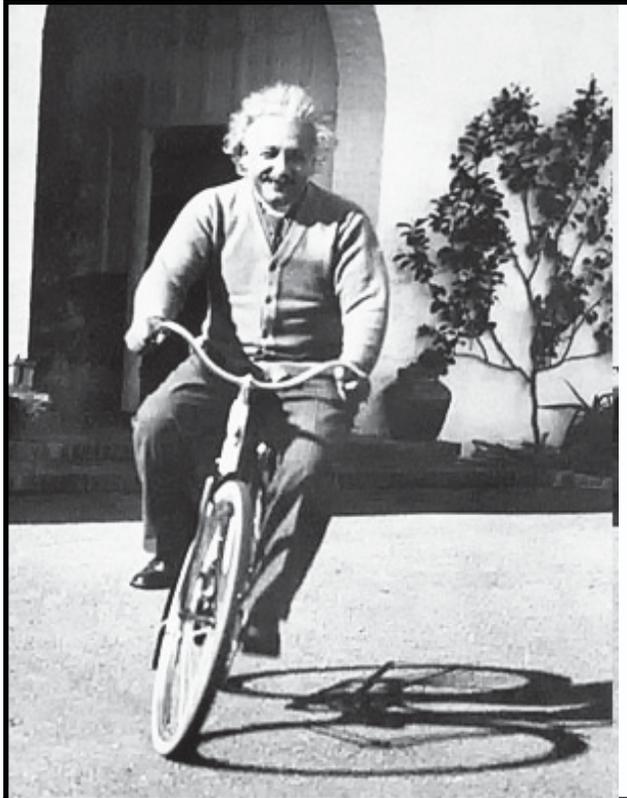
Arjang Hassibi

**Electrical and Computer Engineering Department
Institute for Cellular and Molecular Biology
University of Texas at Austin, TX**

November 1st, 2011

The Abdus Salam International Centre for Theoretical Physics

Why Sequence DNA?



Albert Einstein (1879-1955)
Greatest Achievement: Theory of Relativity

—



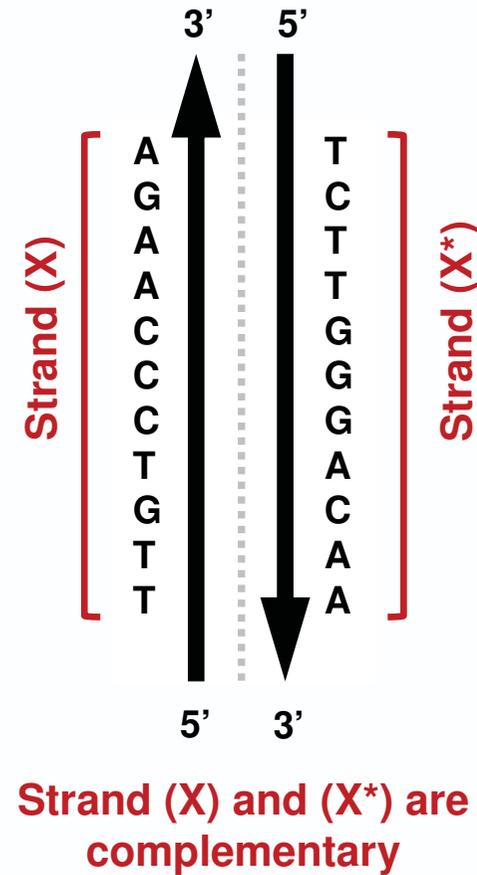
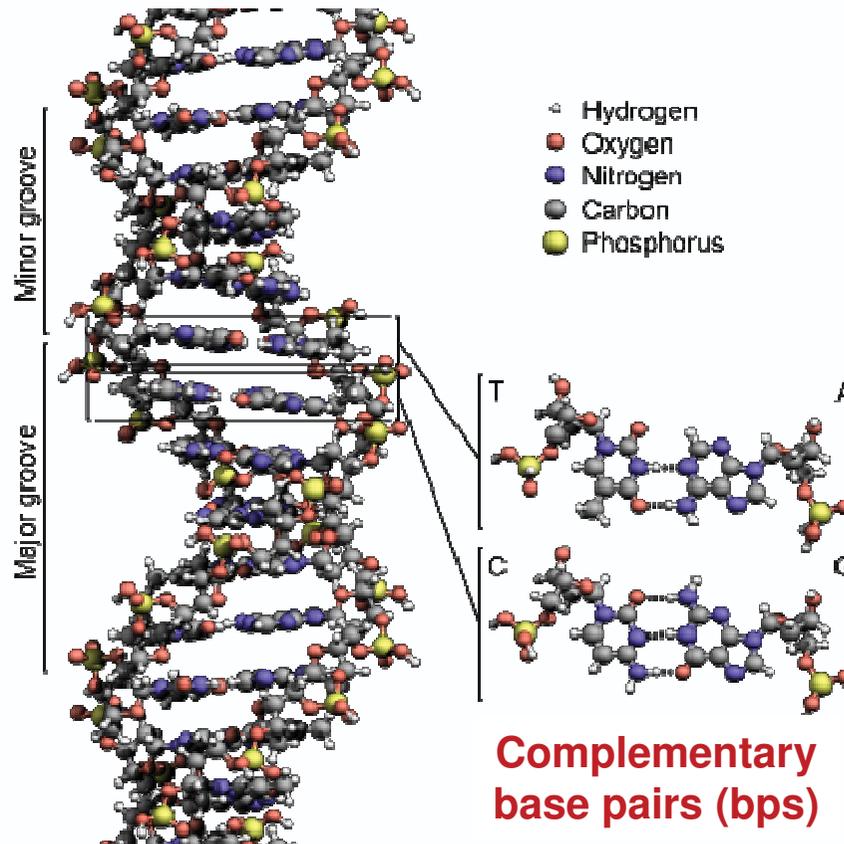
Bobo the Chimp (1995-Now)
Greatest Achievement: Shown Above

= 1.5%
DNA
Difference

Small details in the DNA can make a huge difference!

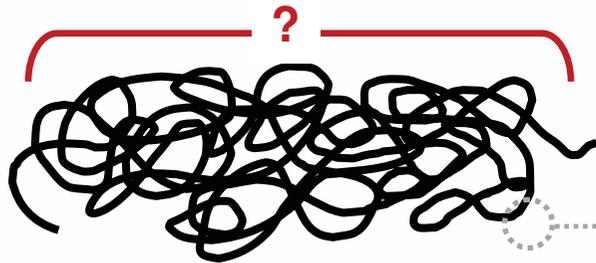
Problem Definition

Identify the structure of a long polymer molecule that consists of only four building blocks (A, C, G, and T)



Scope of the Problem

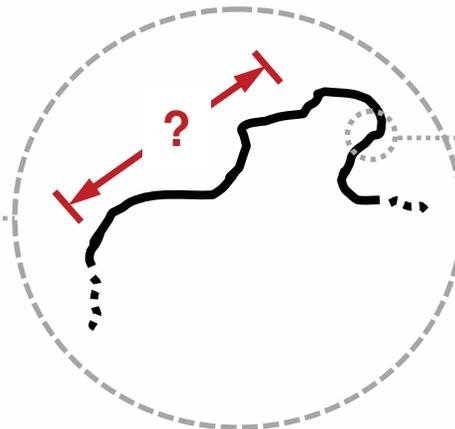
Whole Genome
Sequencing



$>10^9$ bps

Gigabytes (GB)
of Information

Targeted
Sequencing



10^6 bps

Megabytes (MB)
of Information

Genotyping /
SNP Detection



100's bps

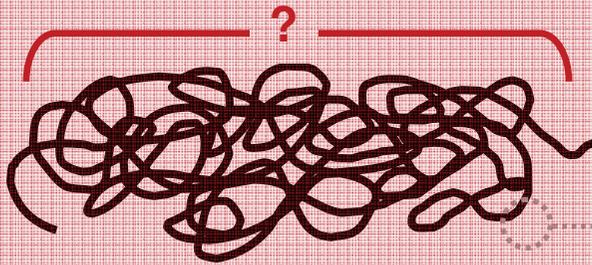
Kilobytes (kB)
of Information

Research / Discovery

Diagnostics

Scope of the Problem

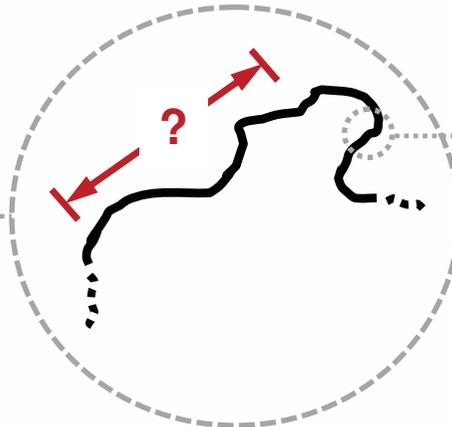
Whole Genome Sequencing



$>10^9$ bps

Gigabytes (GB)
of Information

Targeted Sequencing



10^6 bps

Megabytes (MB)
of Information

Genotyping / SNP Detection

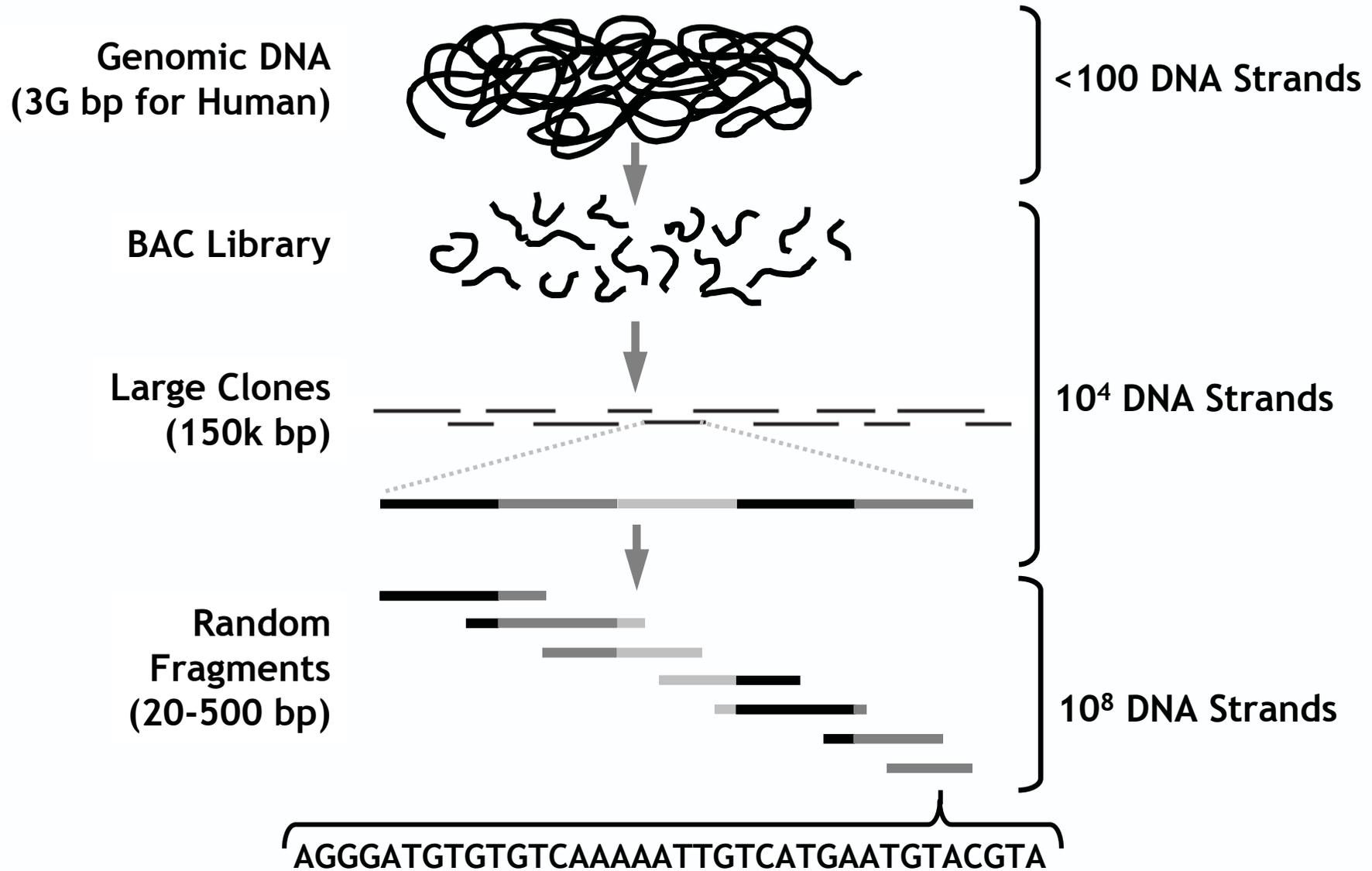
... ACTT?GGTCG...

100's bps

Kilobytes (kB)
of Information

The focus of this talk, i.e., how we can sequence long DNA strands

General Approach: Divide and Conquer

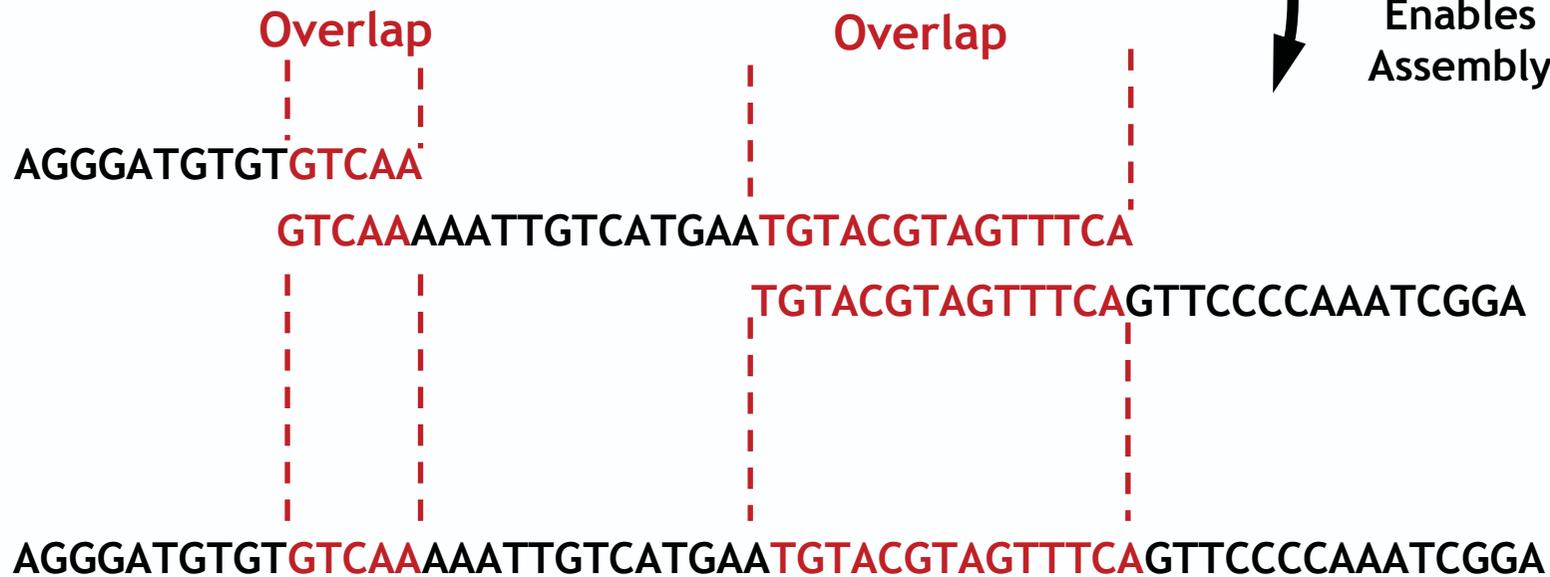


Redundancy vs. Accuracy (1)

①
Sequencing
Result

```
AGGGATGTGTGTCAA  
GTCAAAAATTGTCATGAATGTACGTAGTTTCA  
TGTACGTAGTTTCAGTTCCCAAATCGGA
```

②
Overlapping
Enables
Assembly



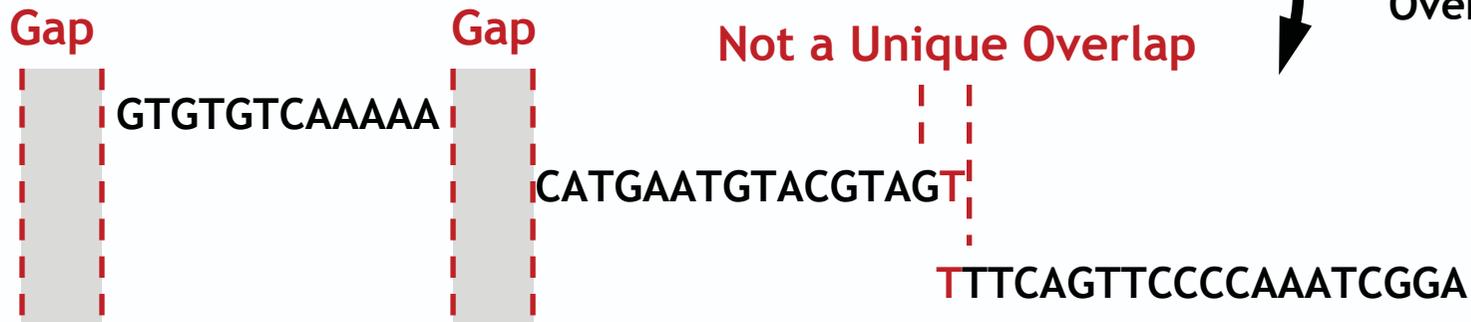
③ Sequence of the original DNA strand is identified

Redundancy vs. Accuracy (2)

① Sequencing Result

```
GTGTGTCAAAA  
CATGAATGTACGTAGT  
TTTCAGTTCCCAAATCGGA
```

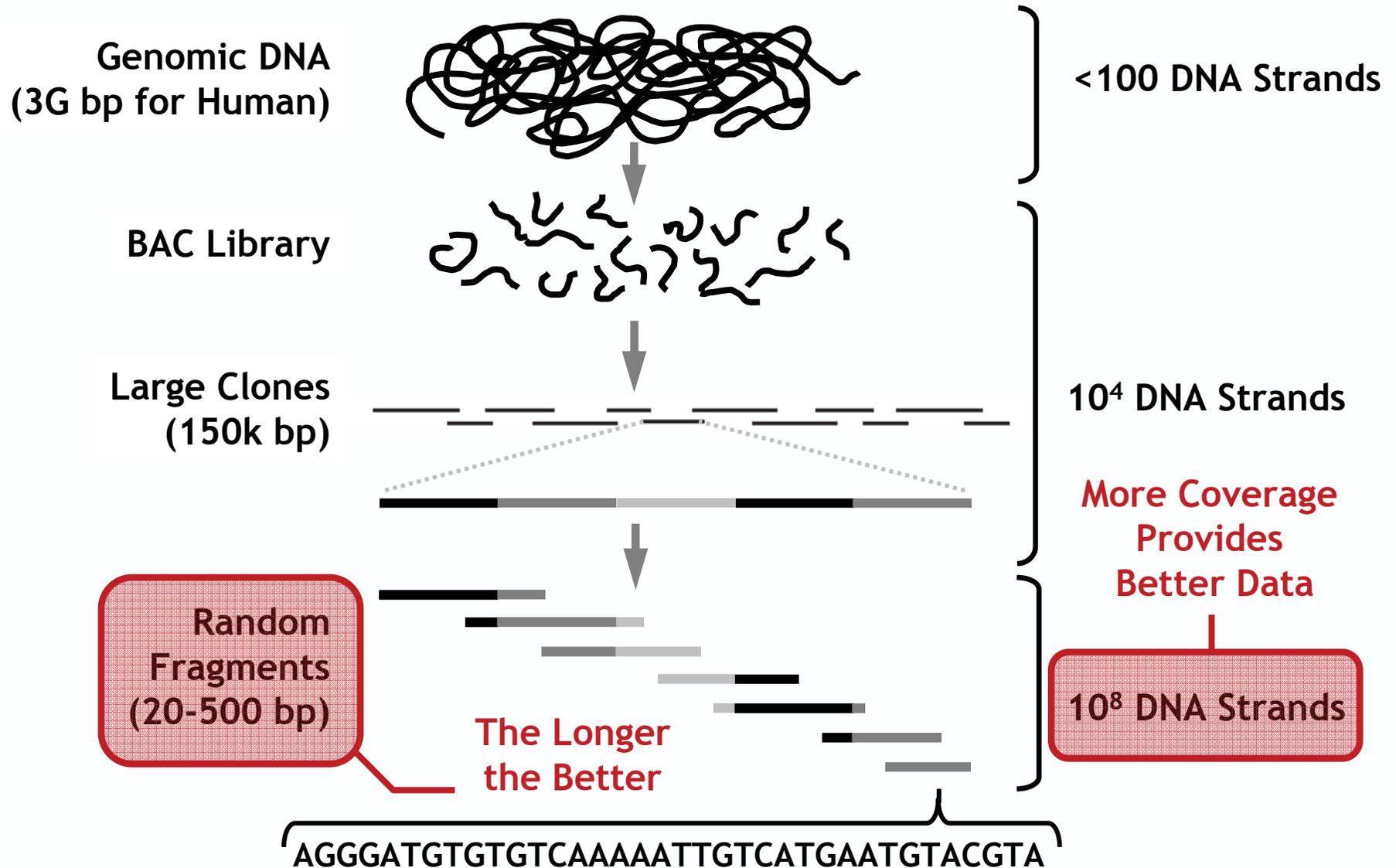
② Cannot Find Overlap



AGGGATGTGTGTCAAAAATTGTCATGAATGTACGTAGTTTCAGTTCCCAAATCGGA

FAILED! Not enough information to identify the sequence

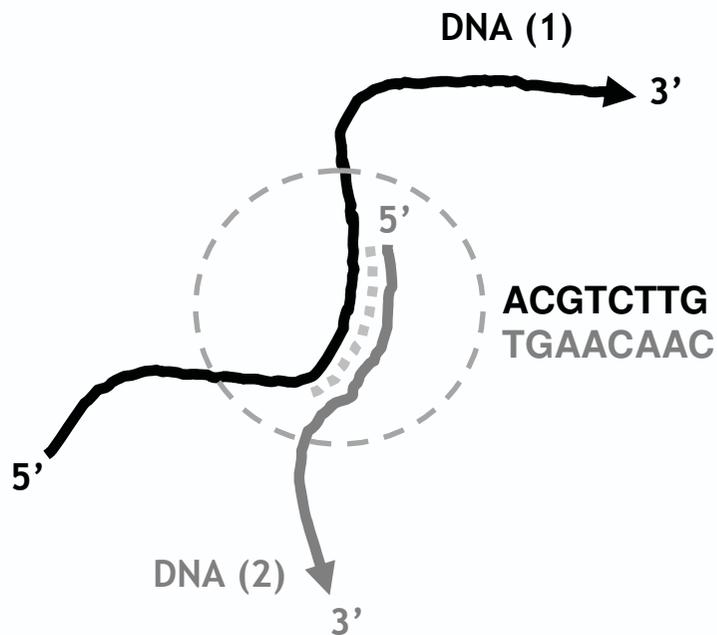
General Approach: Divide and Conquer



Our DNA Toolbox (1)

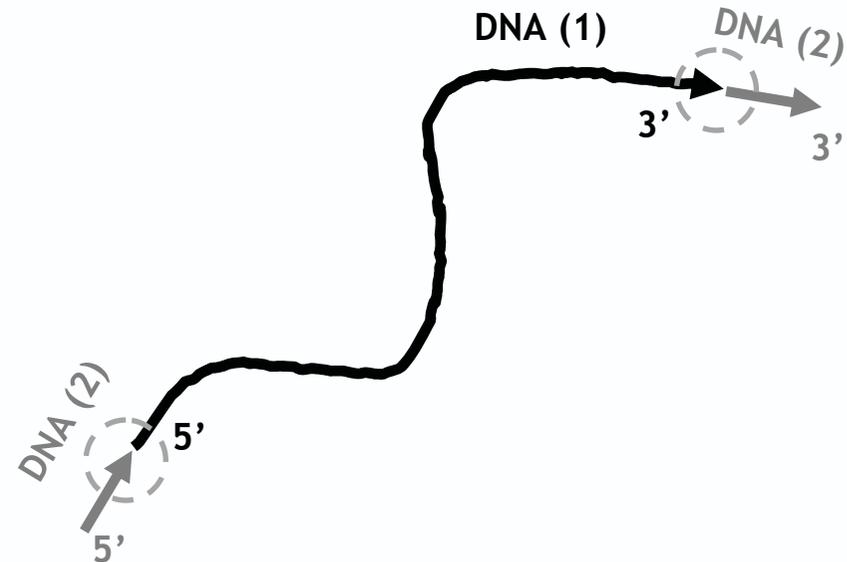
“Sticking” DNA fragments to one another

DNA Hybridization



Sequence-specific
(Energy $\sim 0.2kT/bp$ at 300°K)

Covalent Attachment

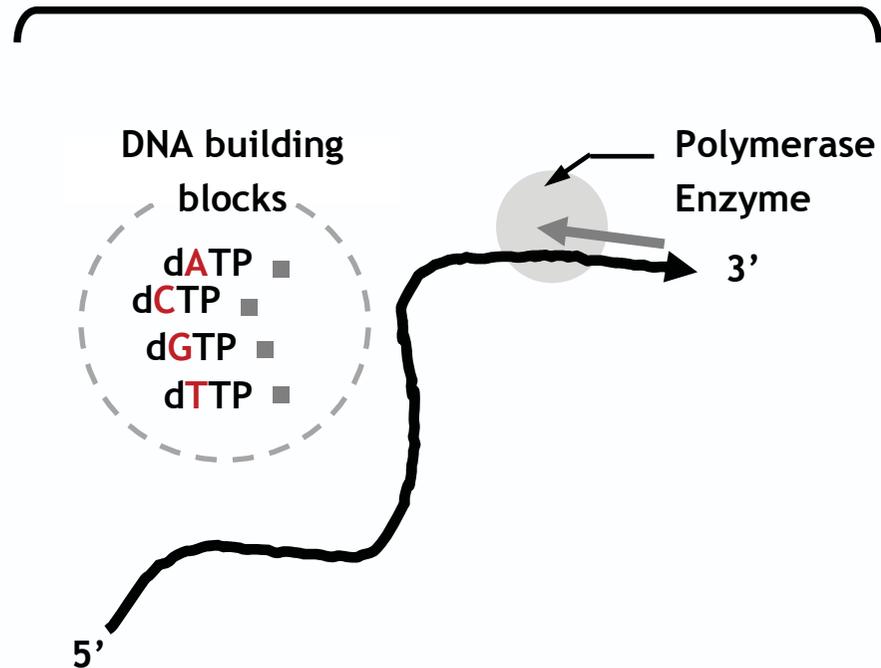


Non-specific and requires activation chemistries

Our DNA Toolbox (2)

Replicating/Copying DNA

DNA Polymerization

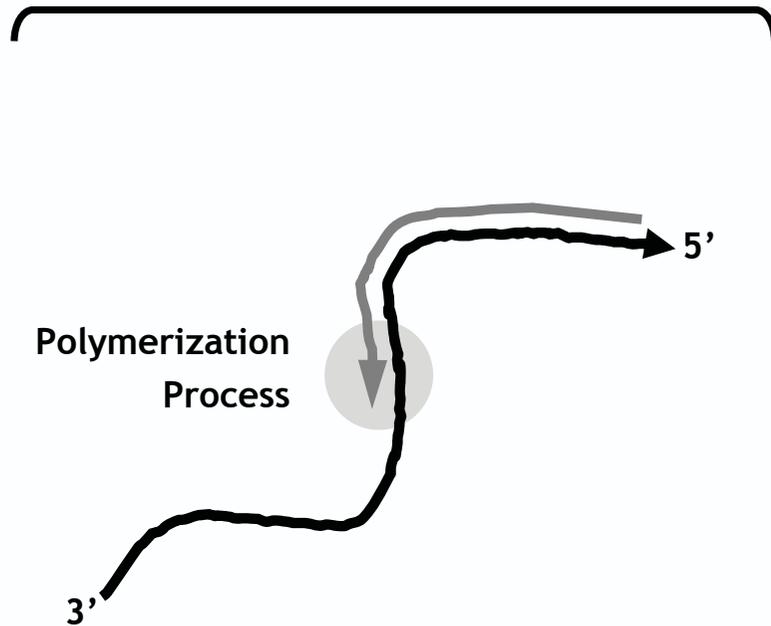


DNA can copy itself
Polymerase and nucleotides
(A, C, G, and T) should be available)

Our DNA Toolbox (2)

Replicating/Copying DNA

DNA Polymerization

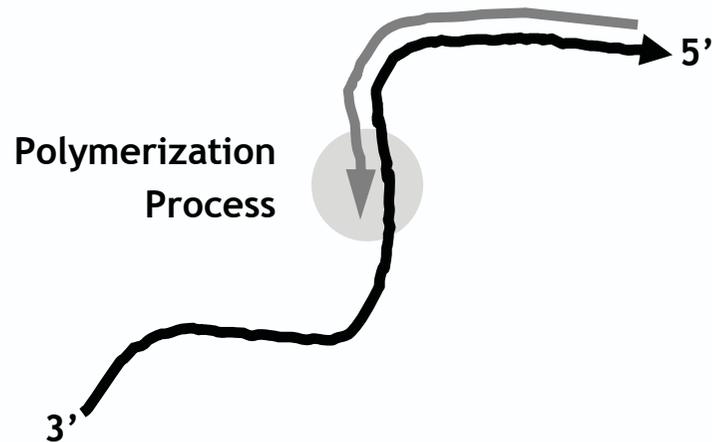


DNA can copy itself
Polymerase and nucleotides
(A, C, G, and T) should be available)

Our DNA Toolbox (2)

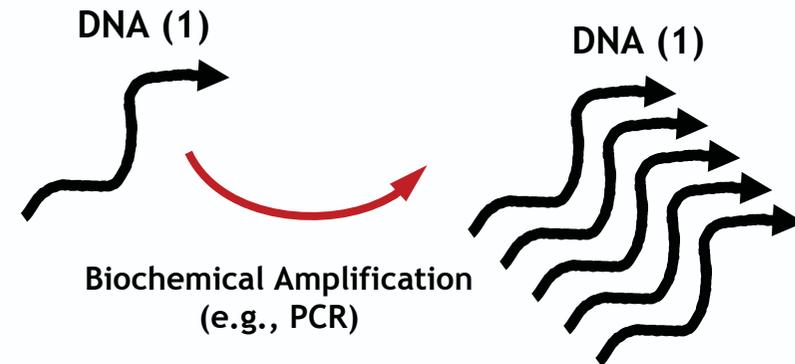
Replicating/Copying DNA

DNA Polymerization



DNA can copy itself
Polymerase and nucleotides
(A, C, G, and T) should be available)

DNA Amplification / Copying



Amplification gain can be $>10^6$

Sanger Sequencing (1977¹)



1958



1980

Frederick Sanger (1918-Present)

Nobel Prize in chemistry in 1958
(Structure of Proteins)

Nobel Prize in chemistry in 1980
(DNA Sequencing)

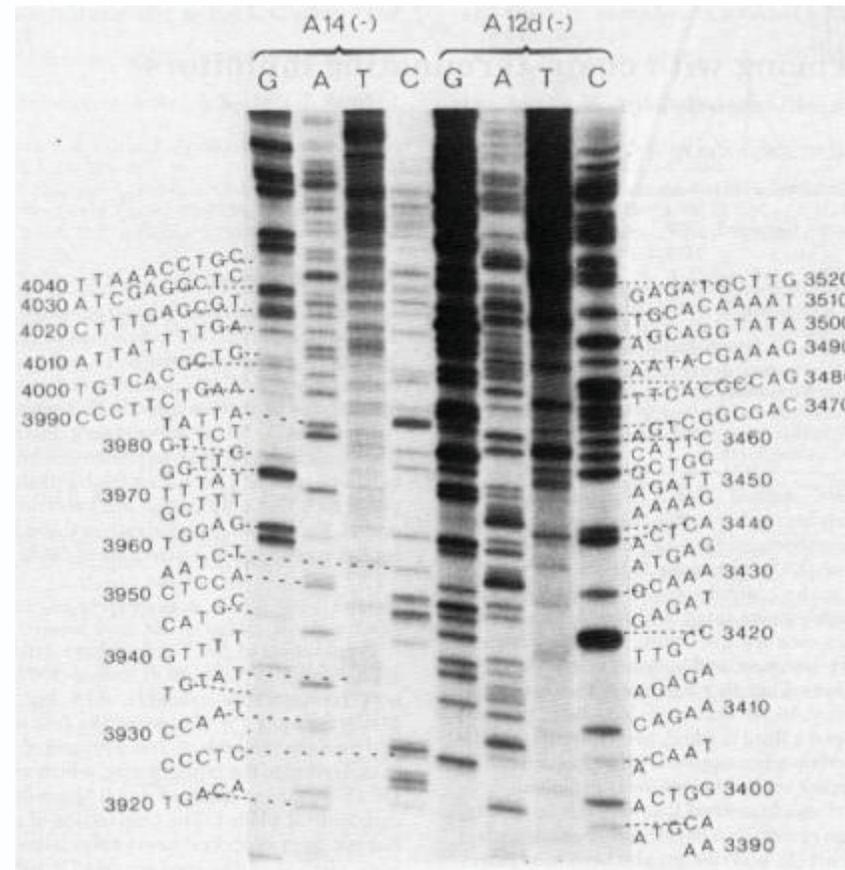
DNA sequencing with chain-terminating inhibitors

(DNA polymerase/nucleotide sequences/bacteriophage ϕ X174)

F. SANGER, S. NICKLEN, AND A. R. COULSON

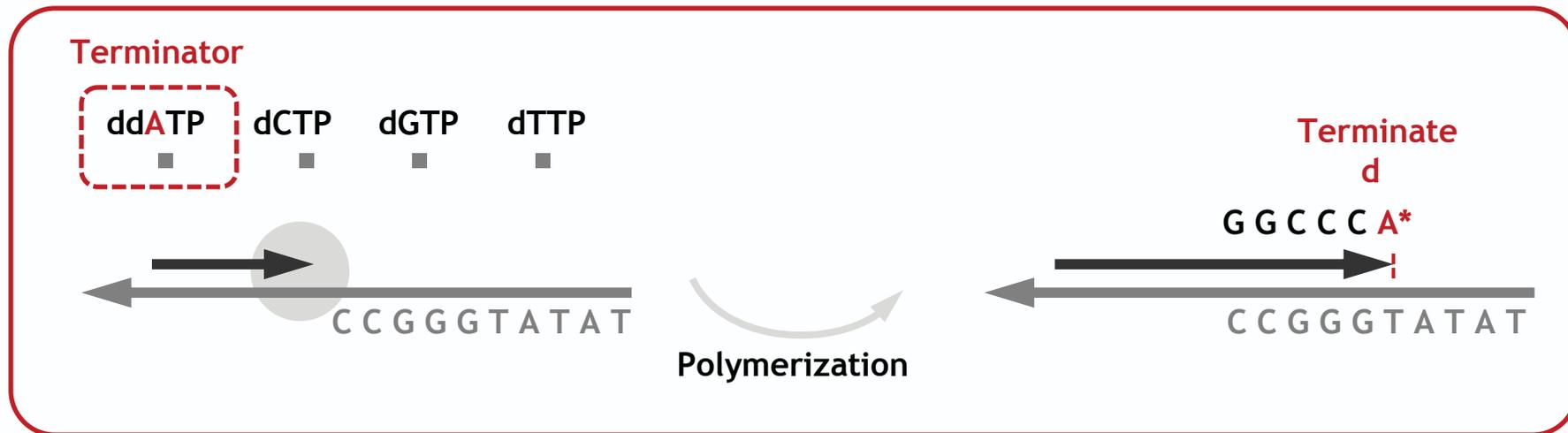
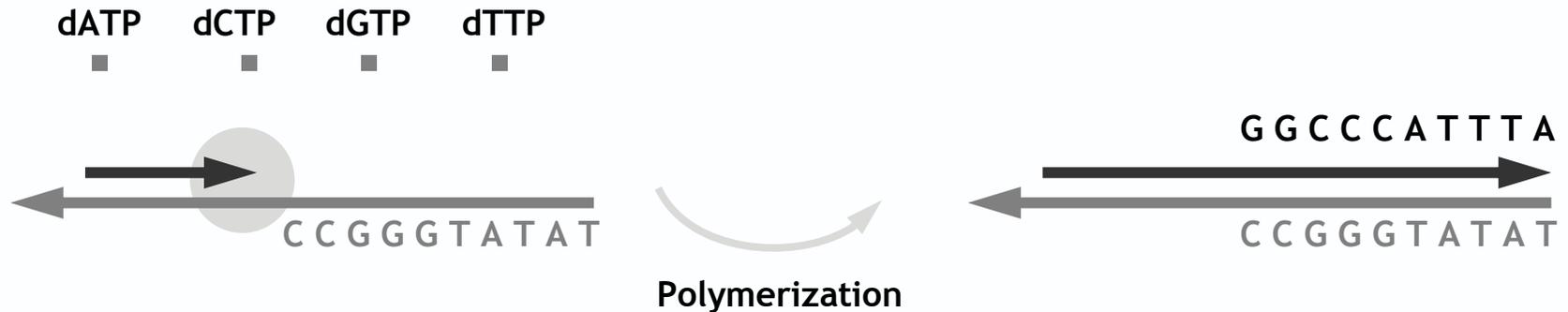
Medical Research Council Laboratory of Molecular Biology, Cambridge CB2 2QH, England

Contributed by F. Sanger, October 3, 1977



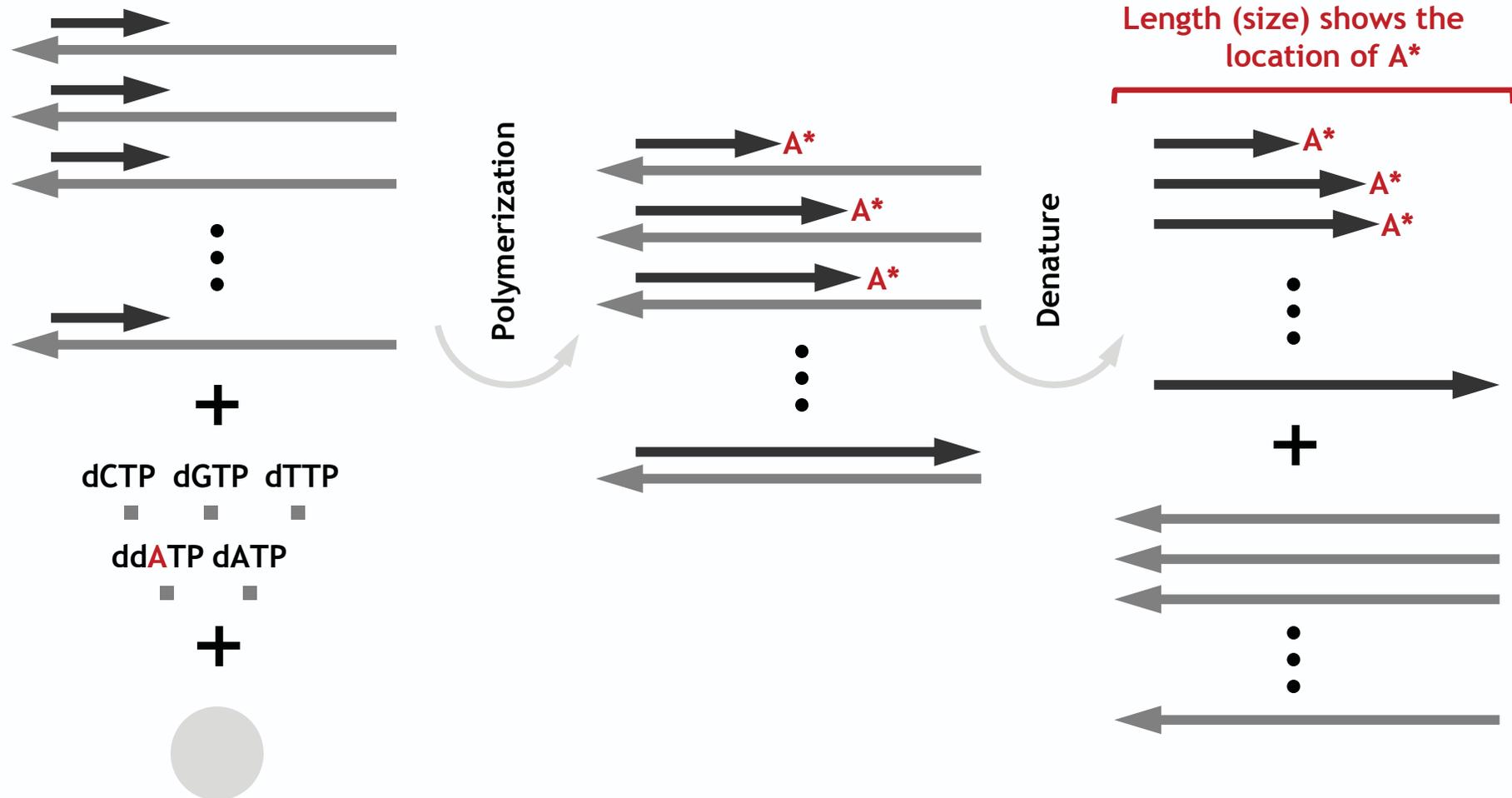
Chain Termination

By using **dideoxynucleotides (ddNTPs)** instead of deoxynucleotides (dNTPs), we can stop (terminate) polymerization

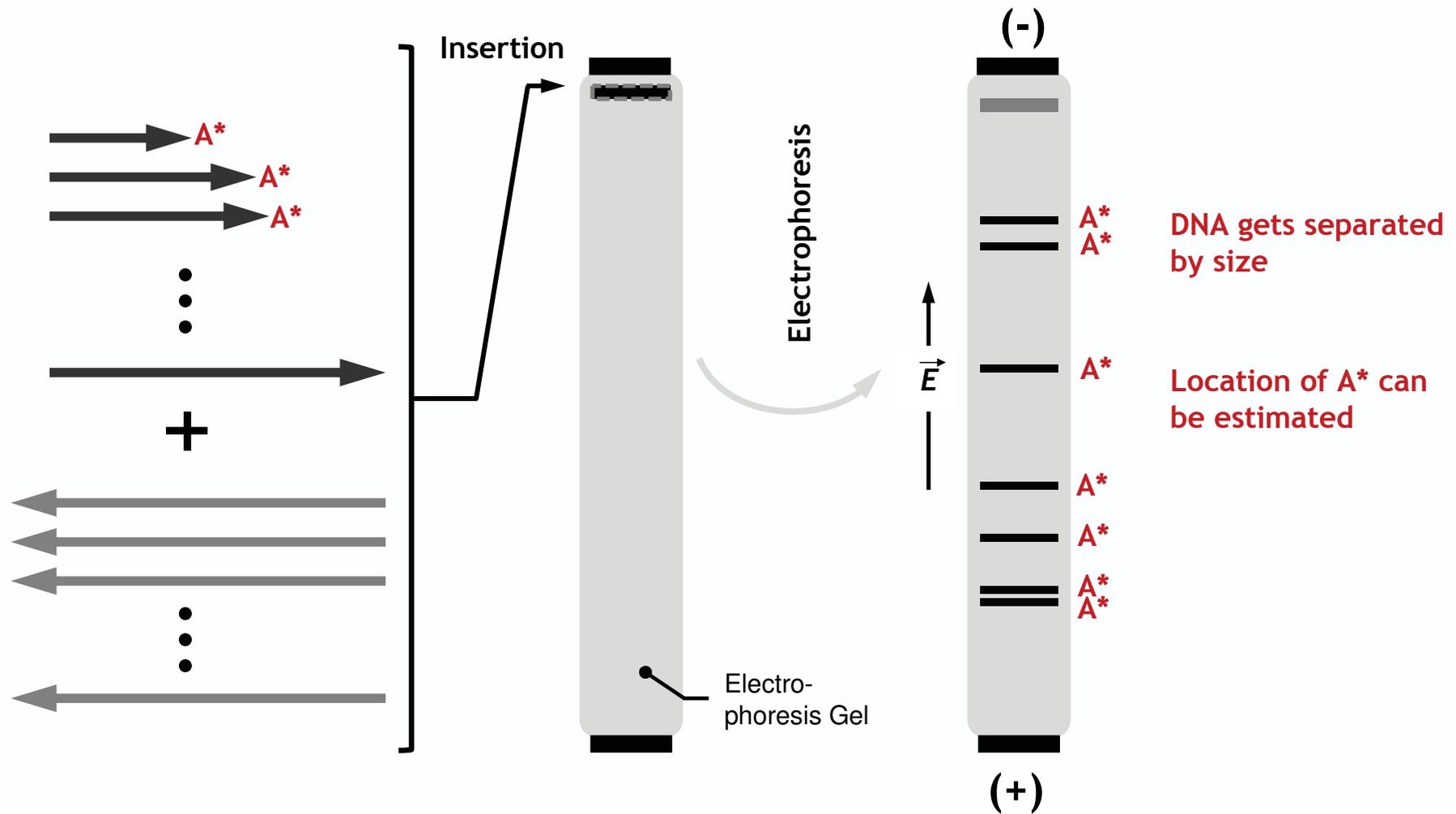


Sequence-Dependant Fragments

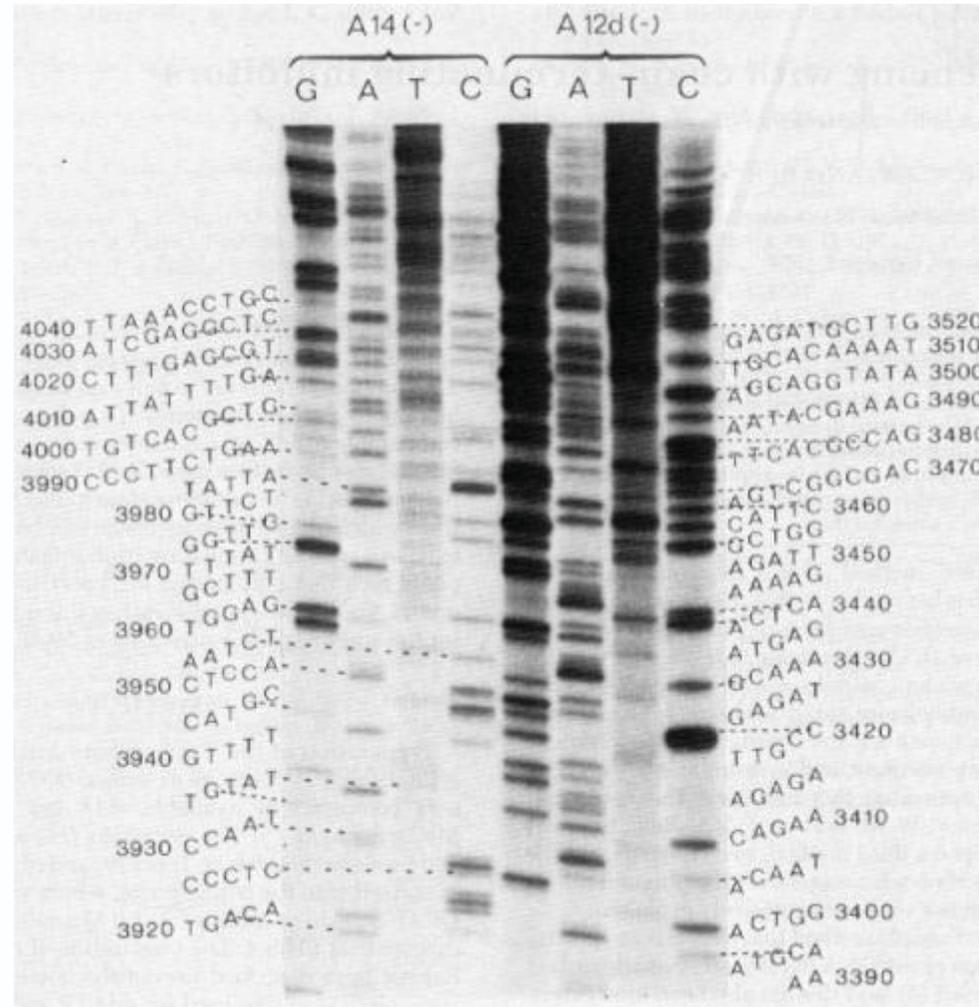
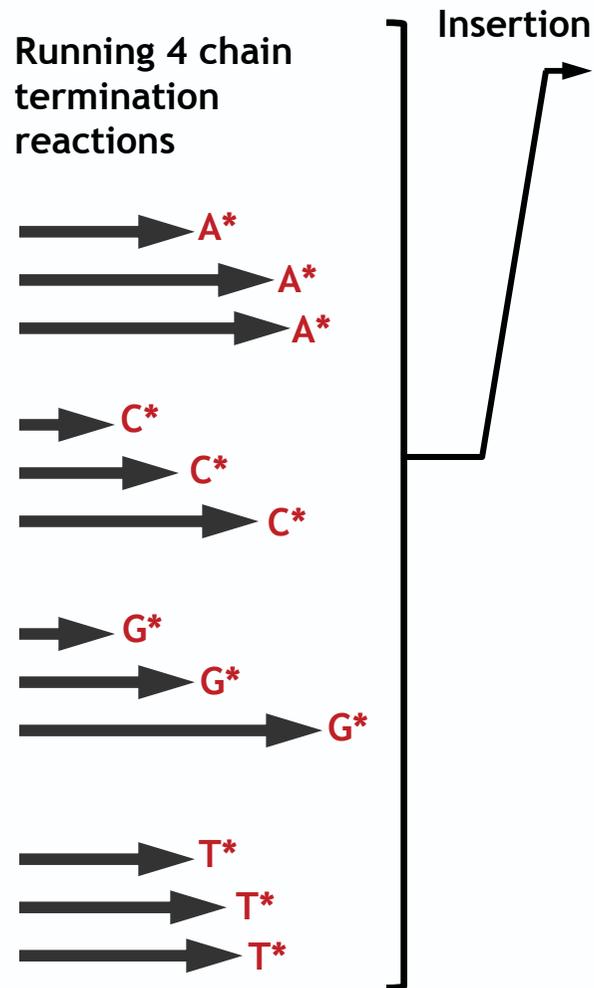
Polymerization with a mixture of dNTPs and ddTTPs



Finding the Length by Electrophoresis



Imaging Using Radioactive Isotopes

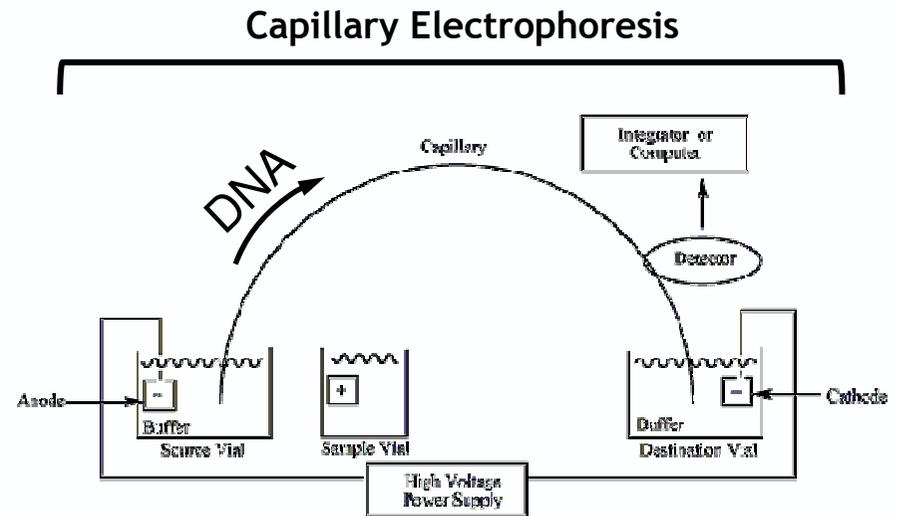


Sanger Sequencers in Human Genome Project (HGP¹)



Sanger DNA sequencer in 2010
ABI (Life Technologies) 3730

96 capillary electrophoresis
~ 1000 bps reads per channel
~ 100 kbps per day



It will take years to sequence a human genome!

¹ IHGSC "Finishing the euchromatic sequence of the human genome". *Nature* 431 (7011): 931-94, (2004).

Sequence-by-Synthesis (SBS)¹

Edward D. Hyman independently came up with this method in 1988.
He made no money out of it and seldom gets any credit

A New Method of Sequencing DNA

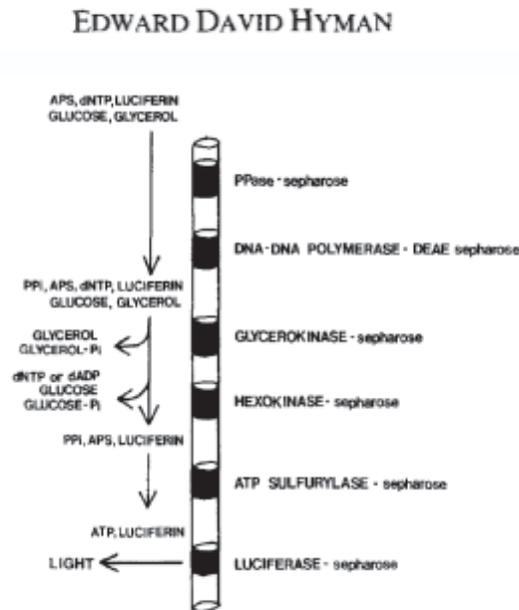
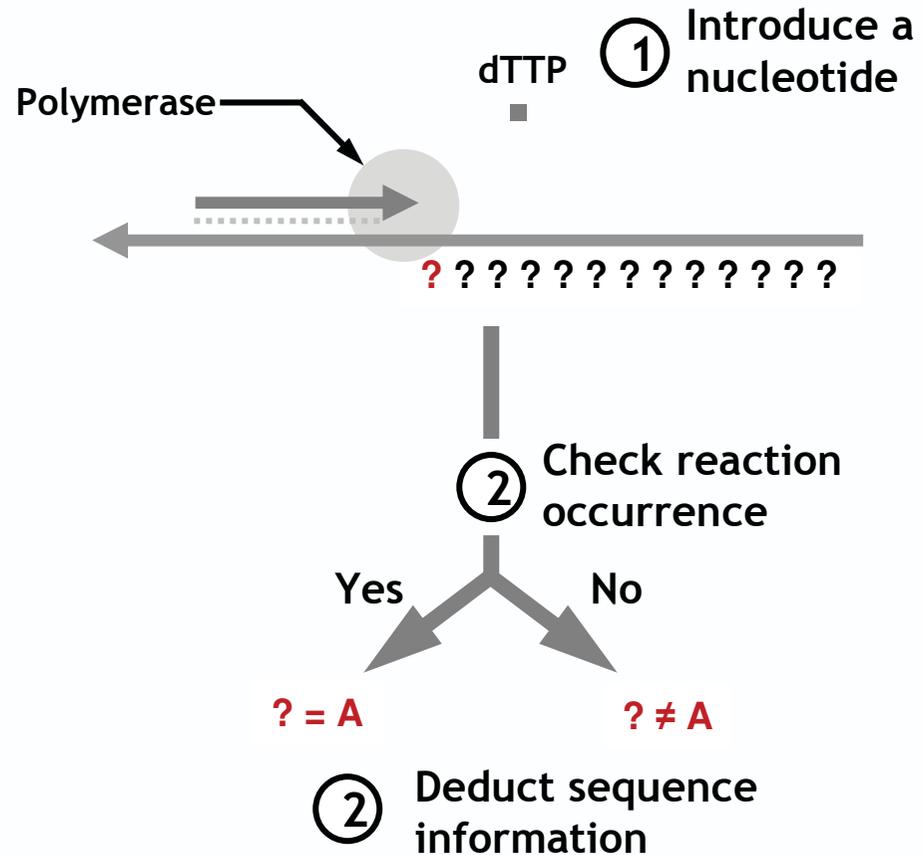
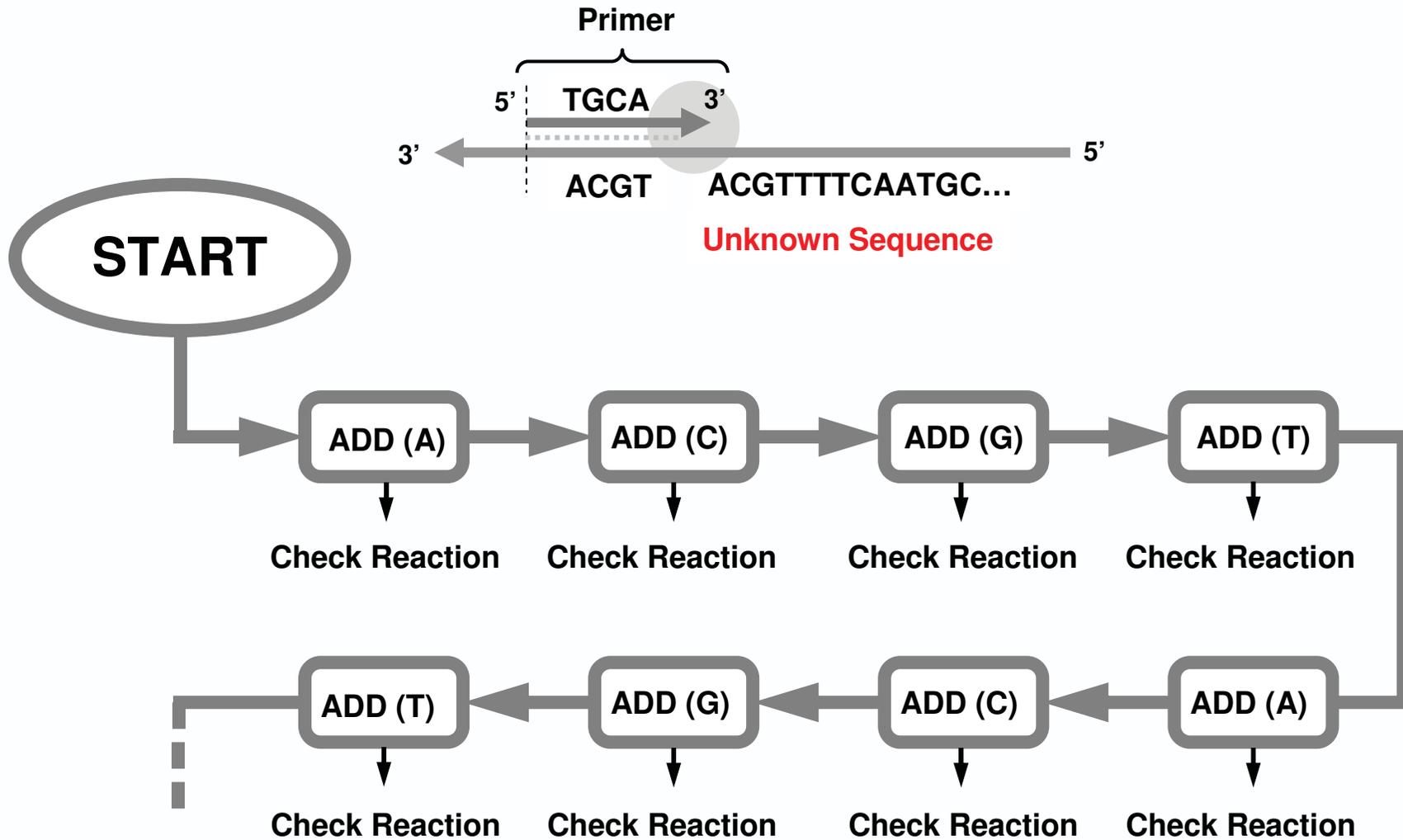


FIG. 1. Schematic diagram of DNA sequencer.

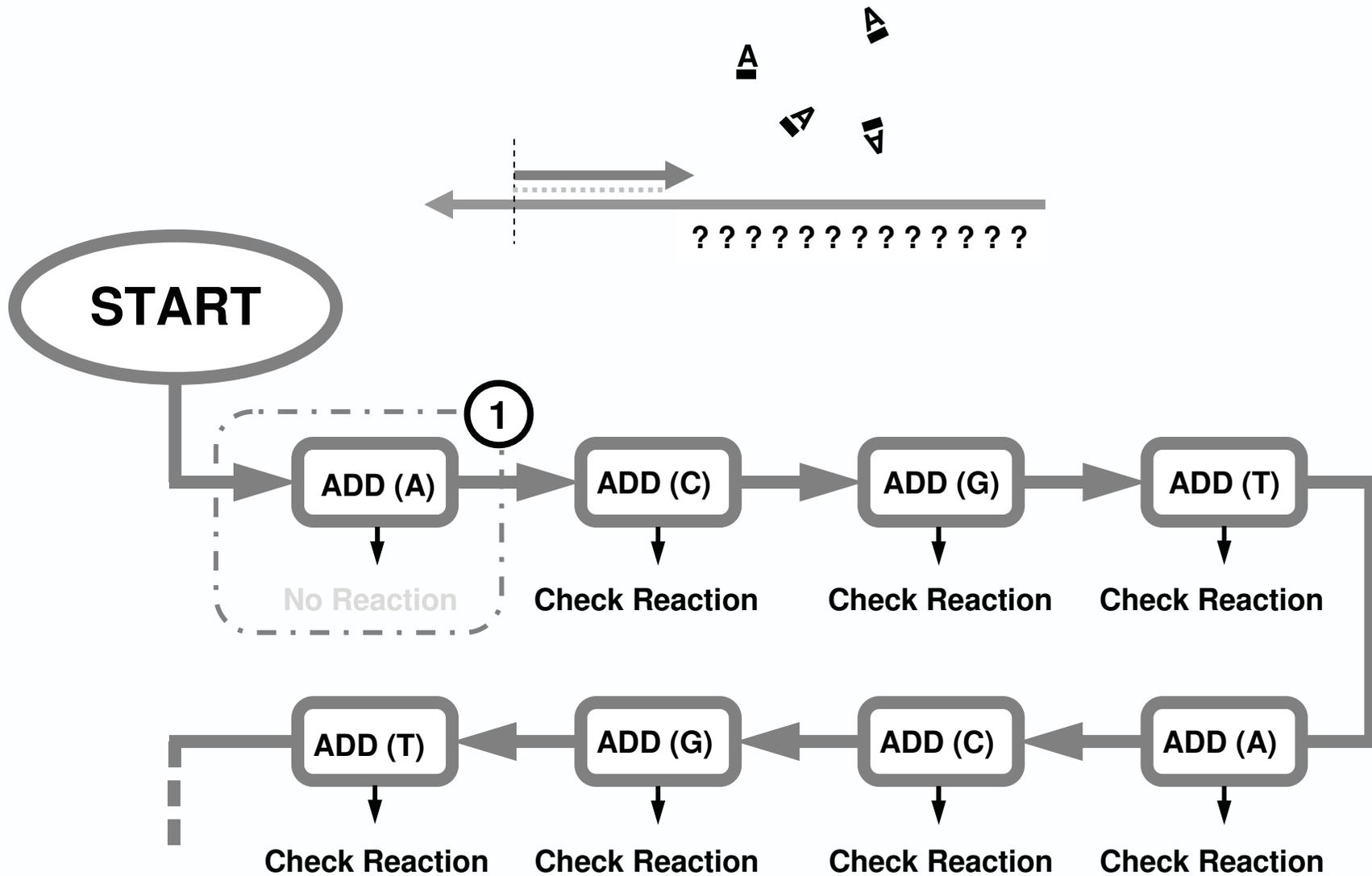


¹ ED Hyman "A new method of sequencing DNA," *Analytical Chemistry*, 174, 423-436, (1988).

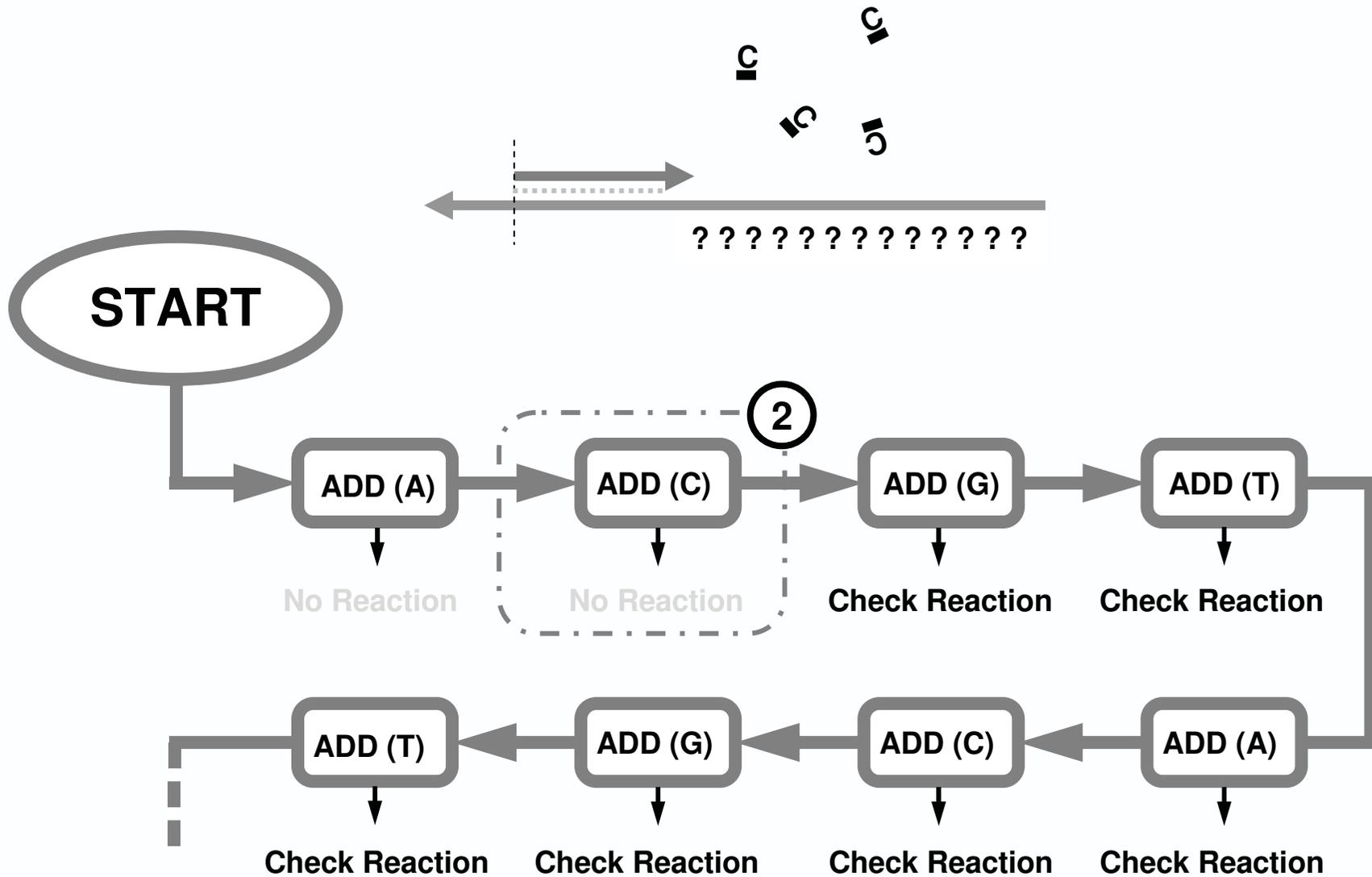
SBS Detection Algorithm



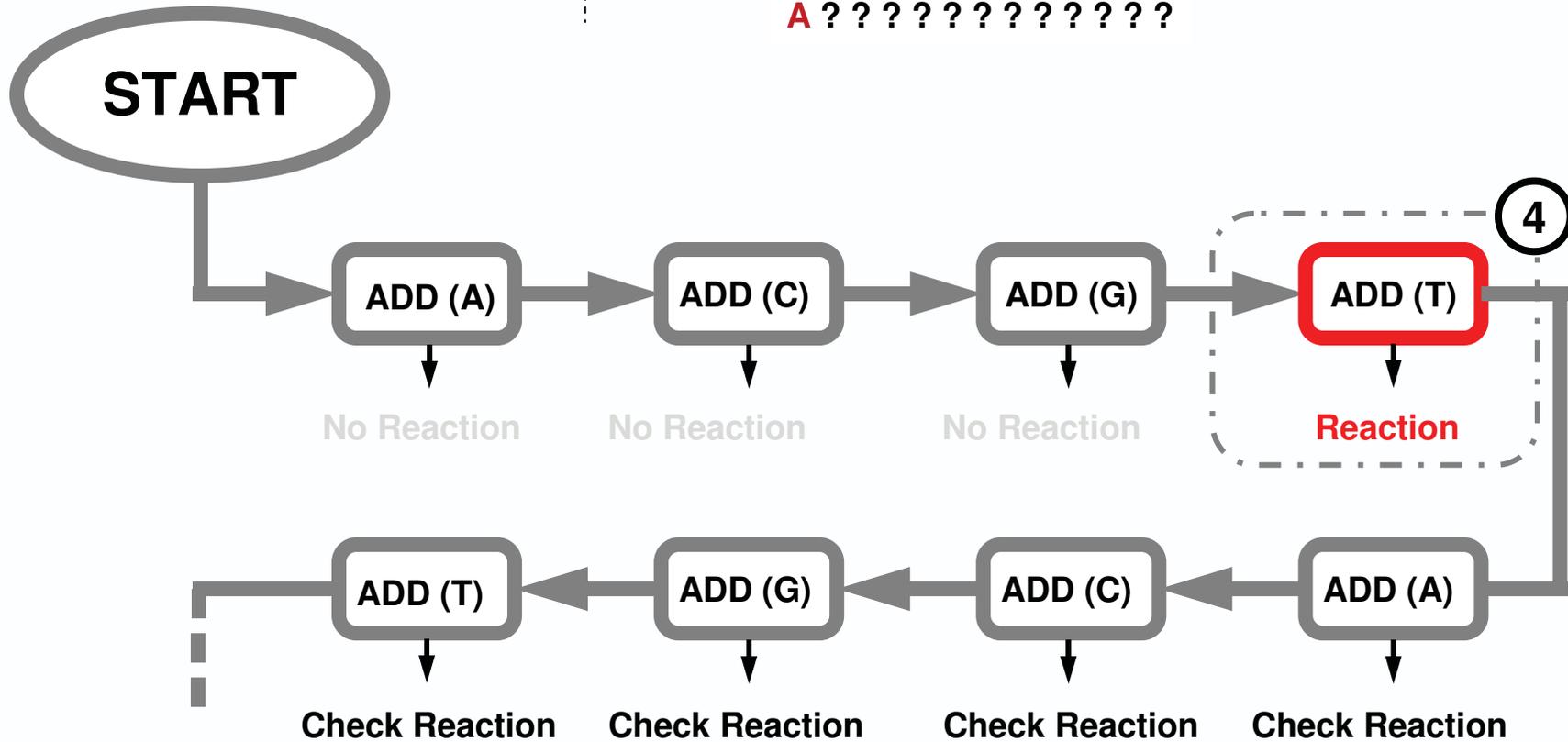
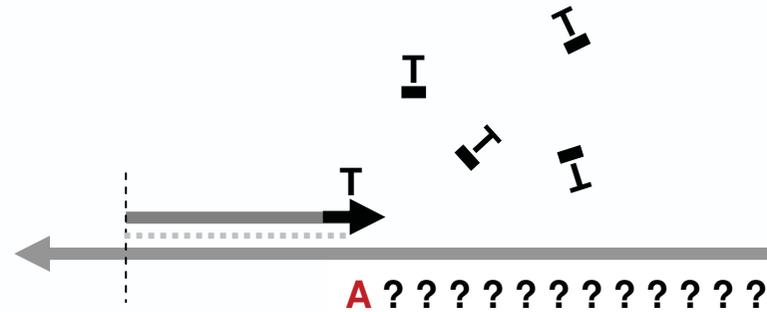
SBS Detection Algorithm



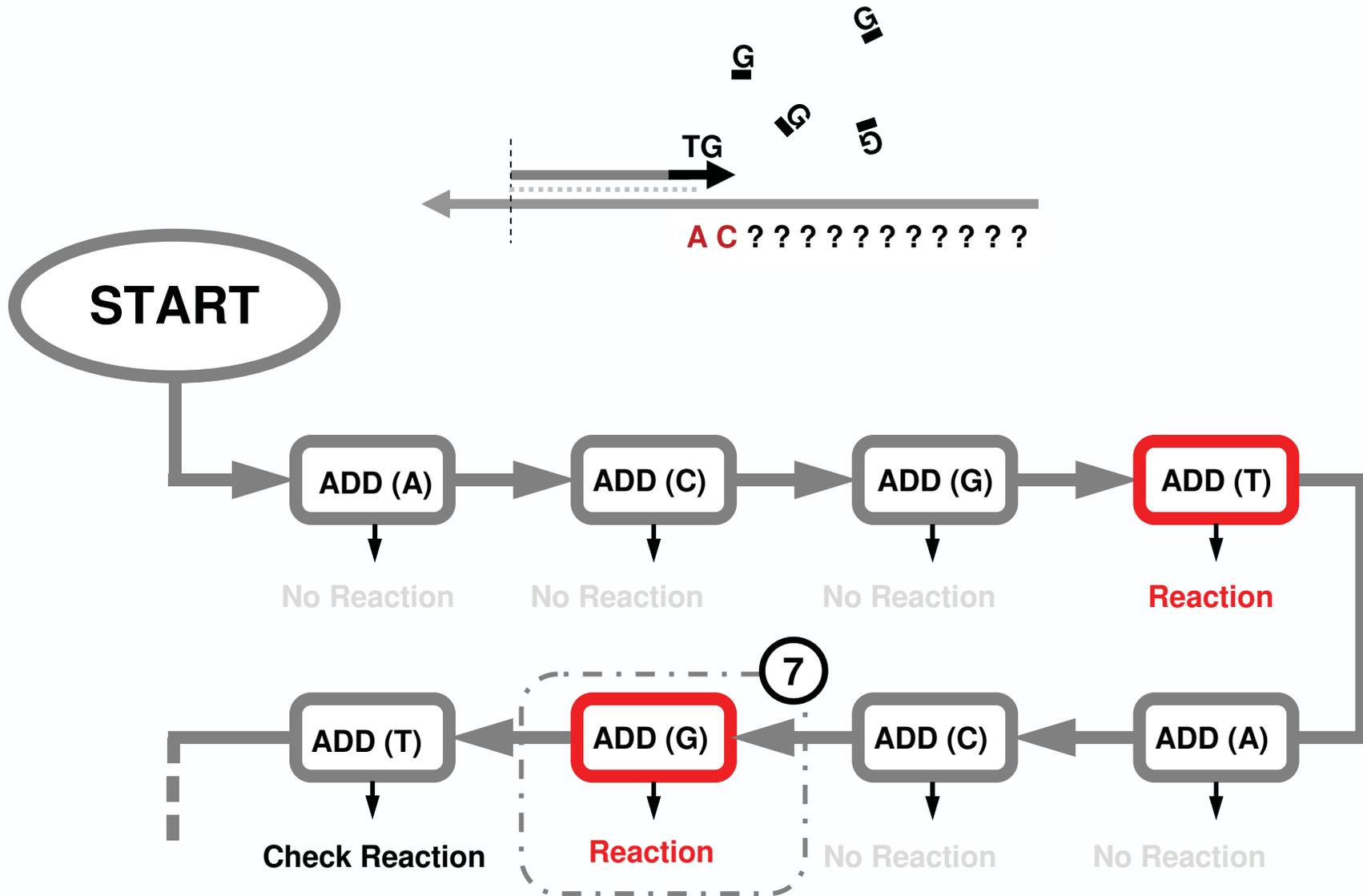
SBS Detection Algorithm



SBS Detection Algorithm

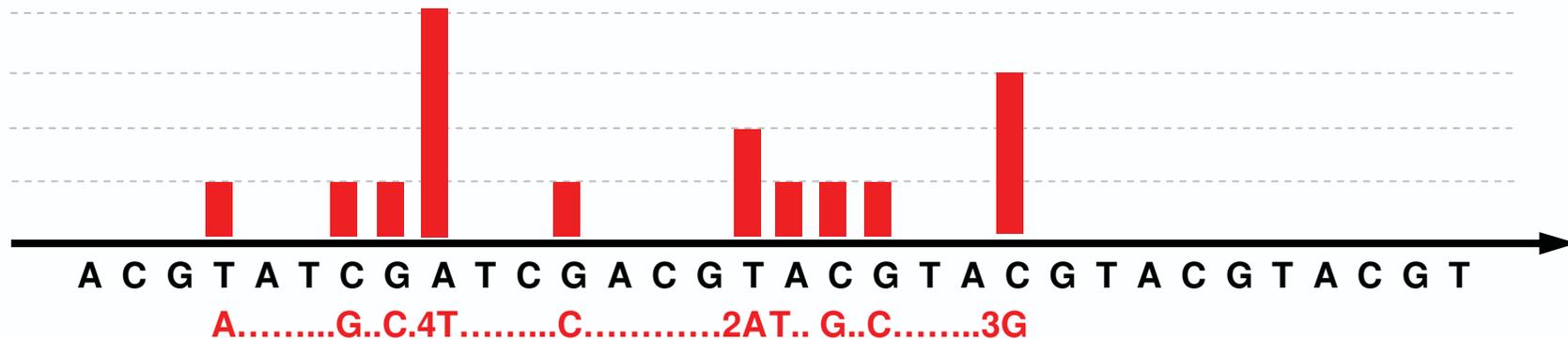


SBS Detection Algorithm



SBS Reaction Diagram

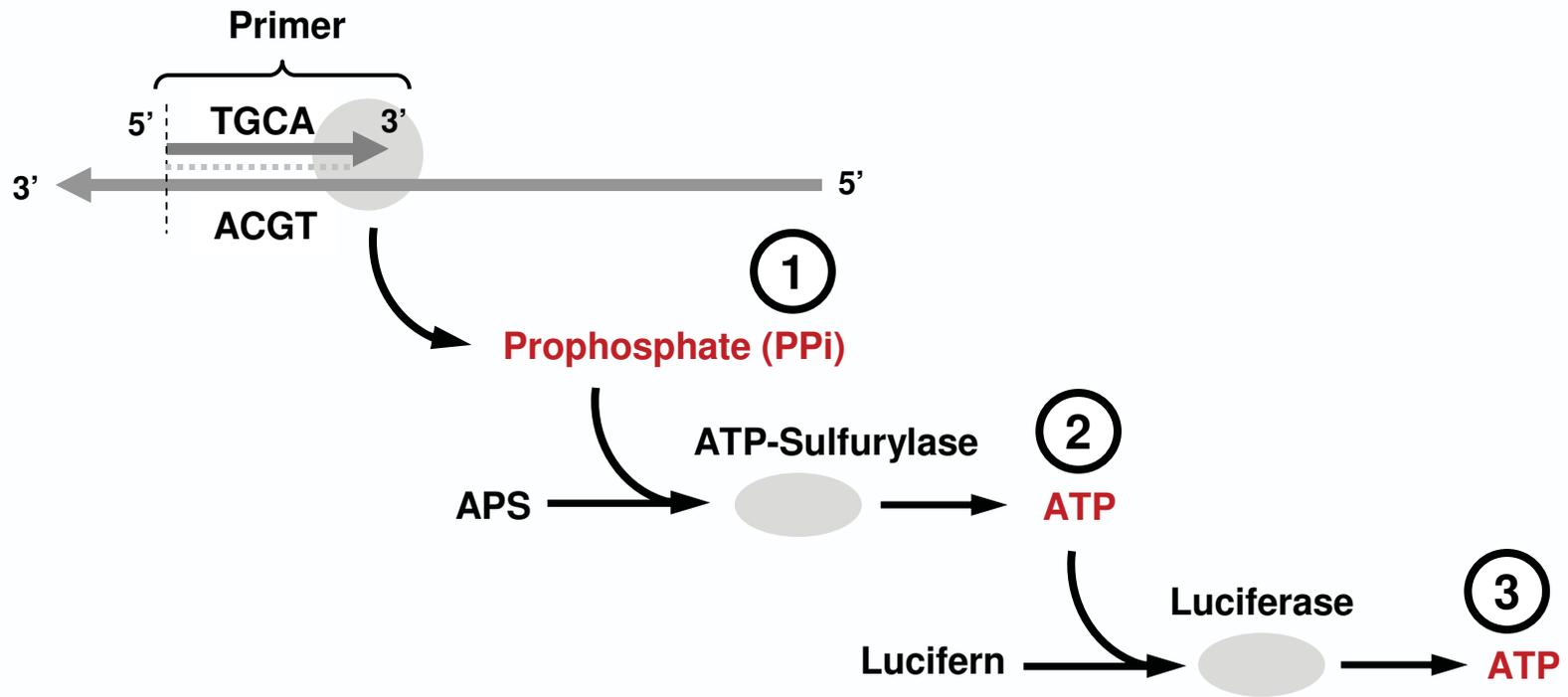
Presence of a reaction and its magnitude is used to sequence the unknown DNA fragment



Sequence: **ACGTTTTCAATGCGGG**

Pyrosequencing¹

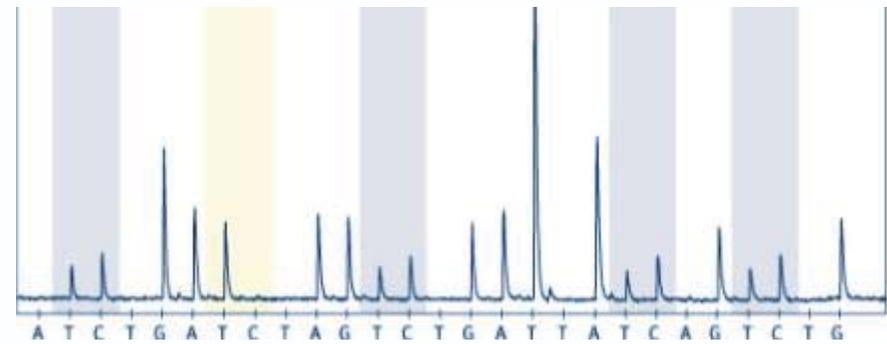
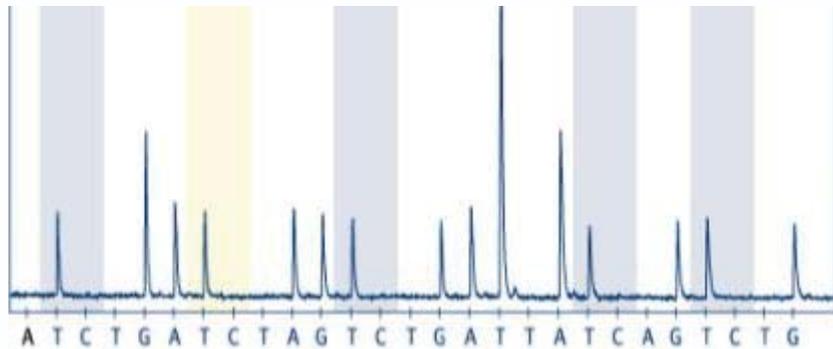
In 1998, Mostafa Ronaghi and Pål Nyrén improved Hyman's SBS method, created a robust assay and coined the term Pyrosequencing



In Pyrosequencing successful reaction trigger an enzymatic cascade to generate photons with $\lambda=562\text{nm}$

¹ M. Ronaghi, M. Uhlén, and P. Nyrén, "A sequencing method based on real-time pyrophosphate,". *Science*, 1998.

Pyrosequencing Data



Qiagen's PyroMark (2011)

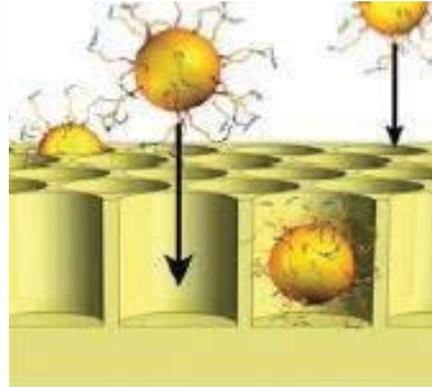
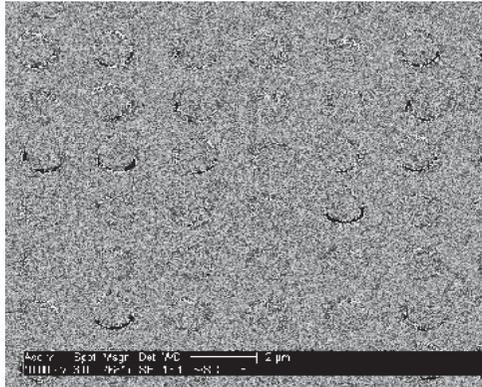
24 microtitre well format
~100 bps read length
~50 kbps per day

**It will take years to sequence a
human genome!**



Random Bead Arrays + Pyrosequencing

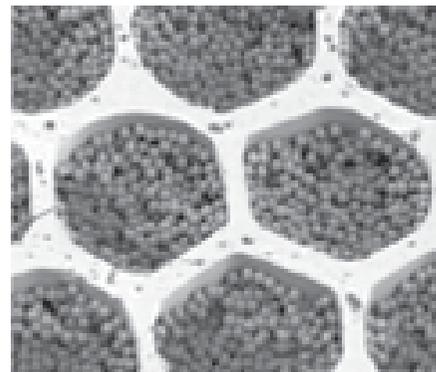
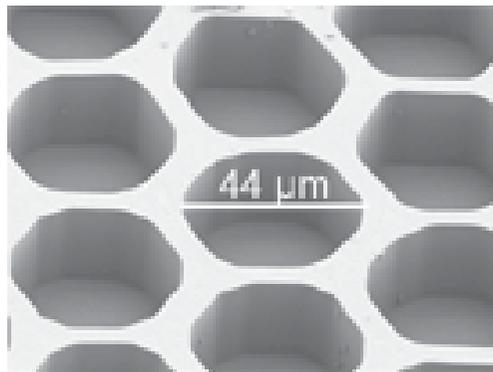
In 2005, the company 454, a subsidiary of Curagen created a high-throughput Pyrosequencing system enabled by random bead arrays



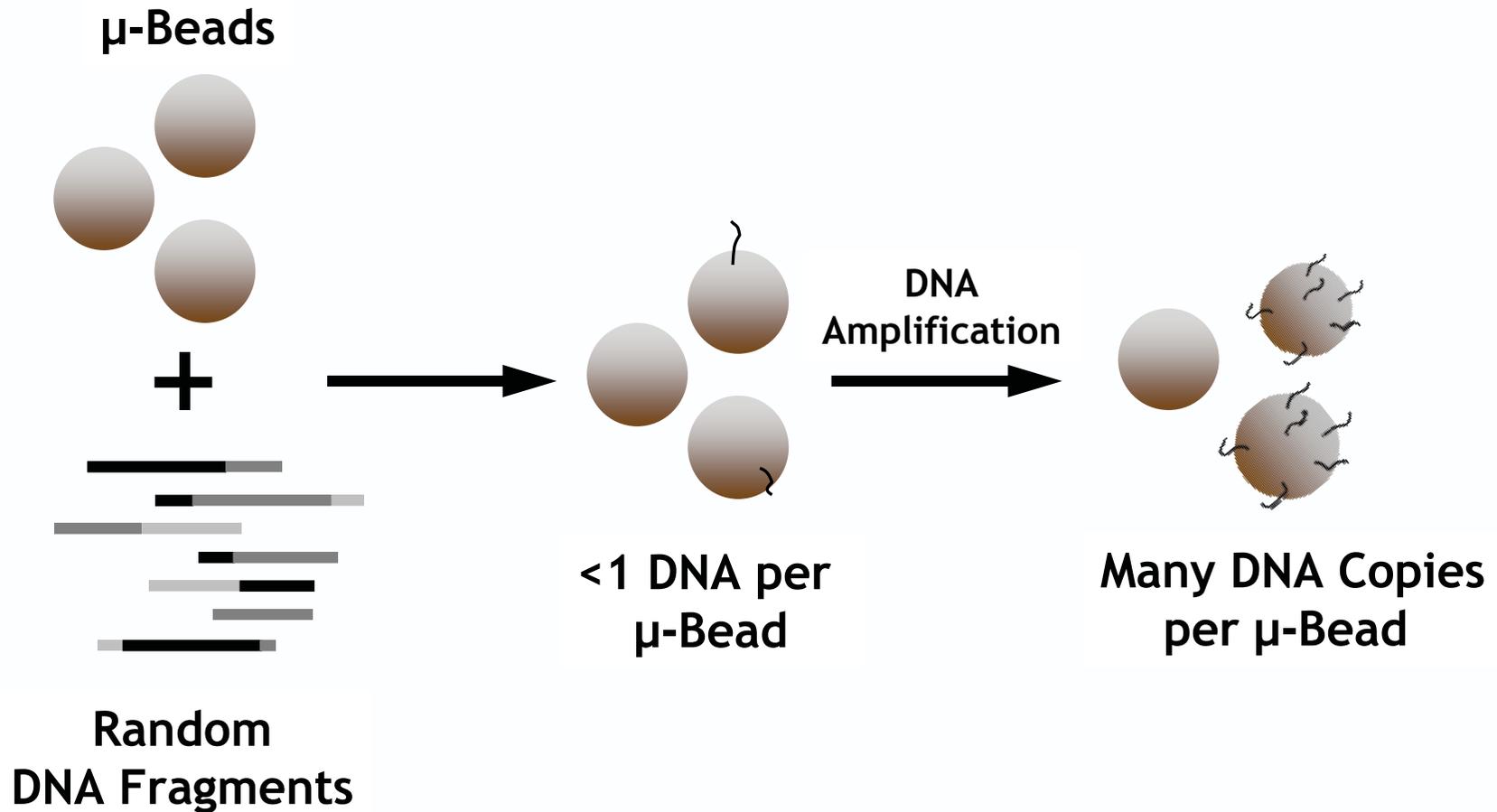
**454 (Roche) GS FLX+
(2011)**

**>1 M micro-wells
~700 bps read length
~700 Mbs per day**

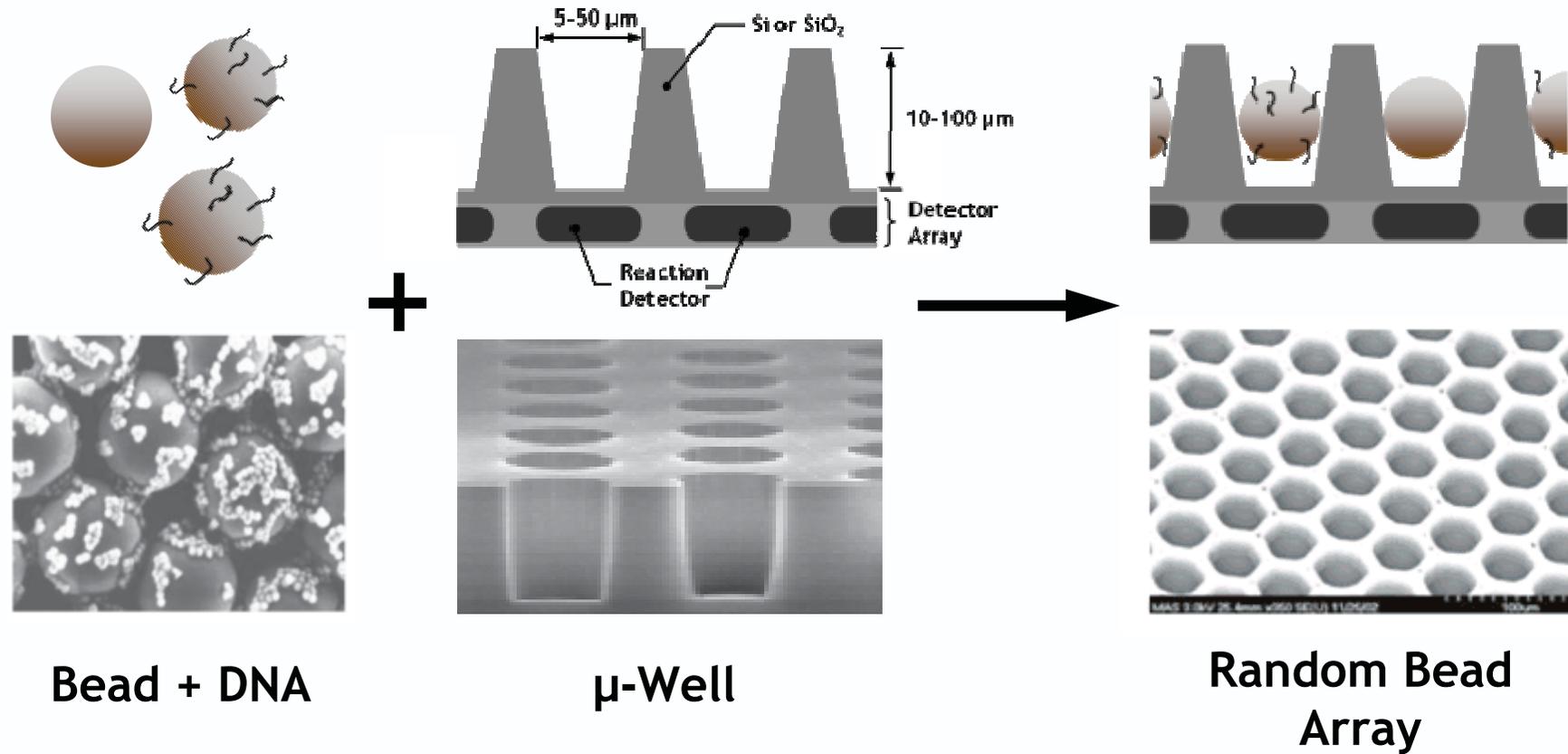
**It will take 2 weeks to sequence
a human genome**



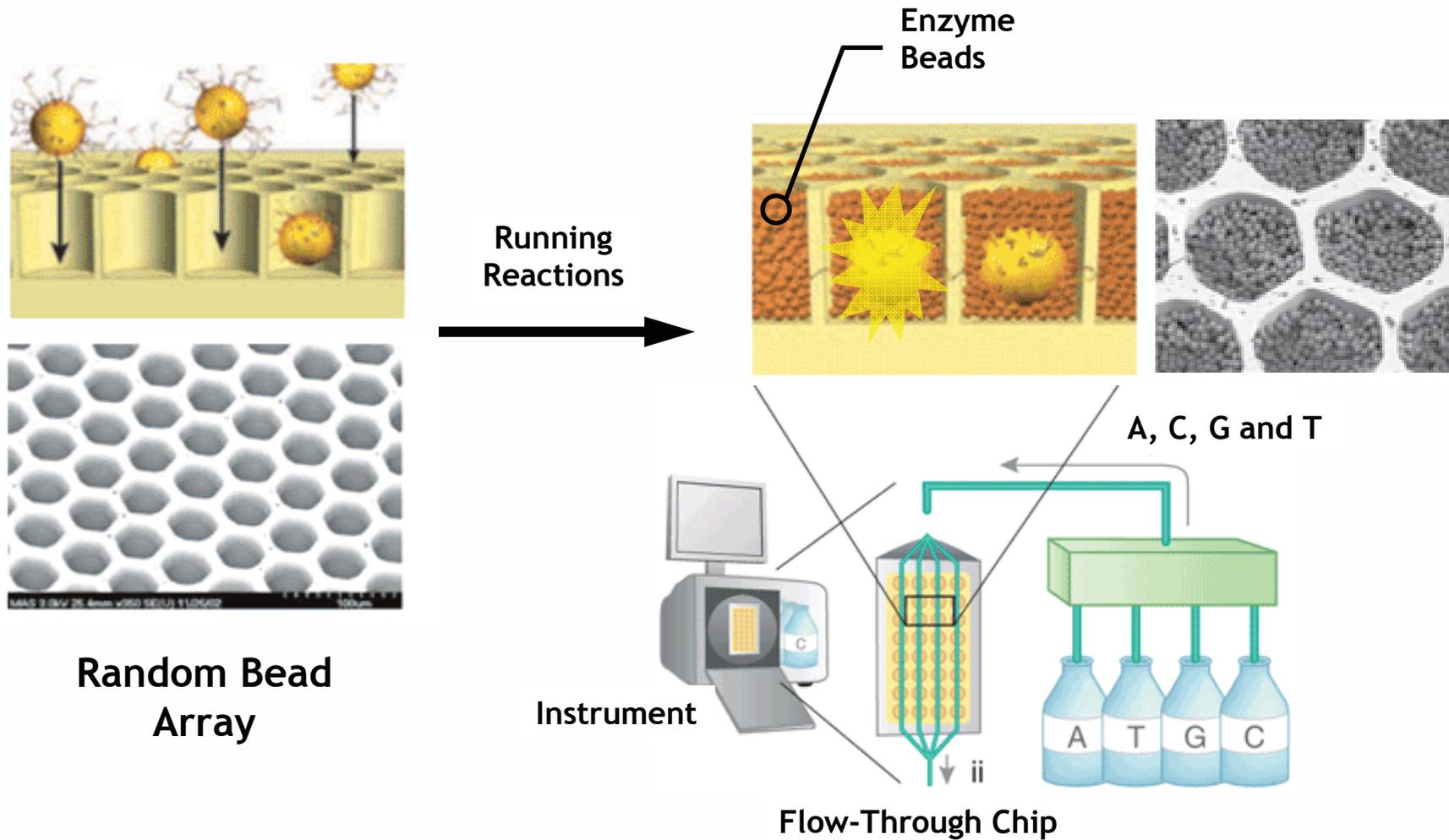
Step 1: Create DNA Clusters on μ -Beads



Step 2: Loading the Array

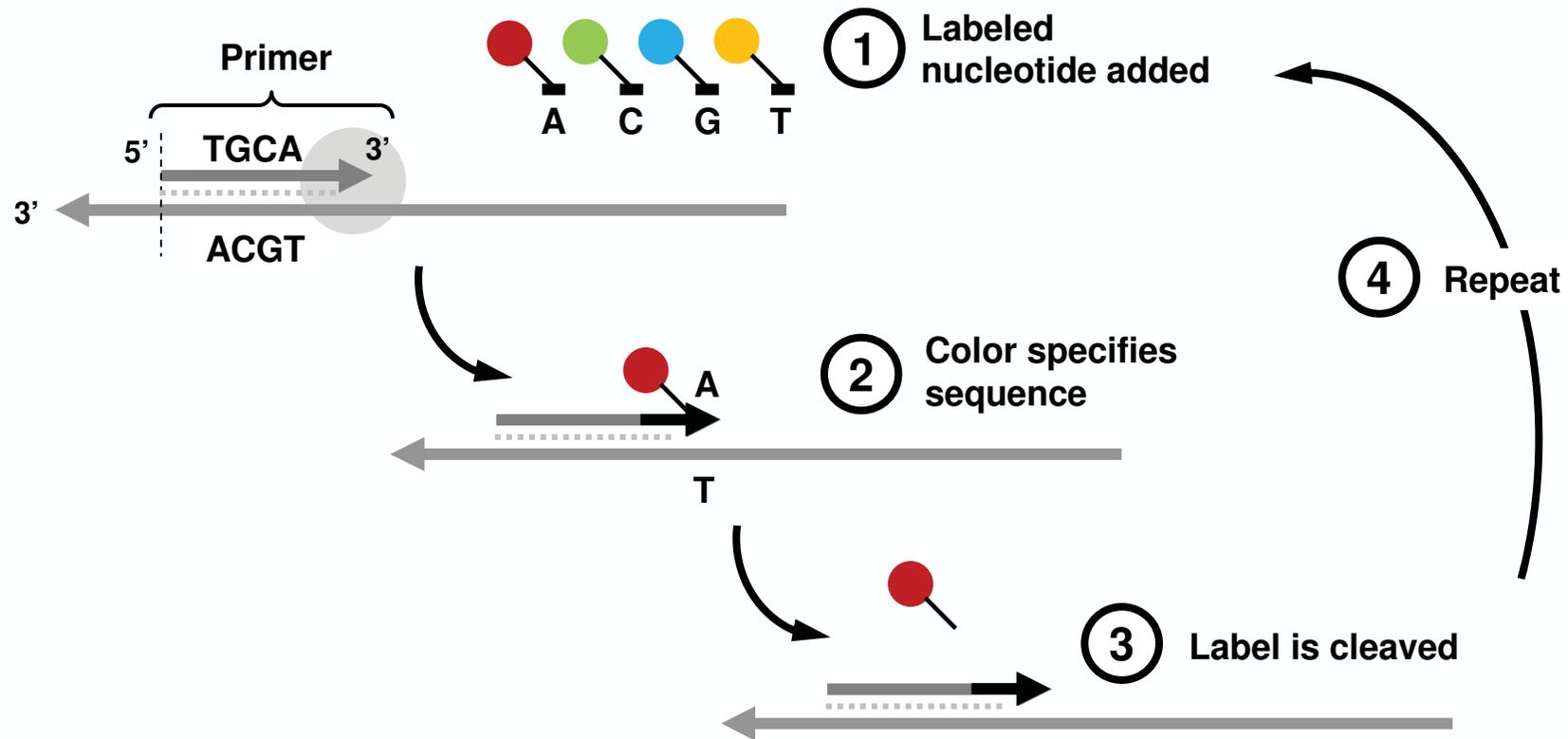


Step 3: Pyrosequencing



Illumina (Solexa) DNA Sequencing

Merges the technologies from Solexa¹ (UK), Manteia² (Switzerland) to create a high-throughput SBS system

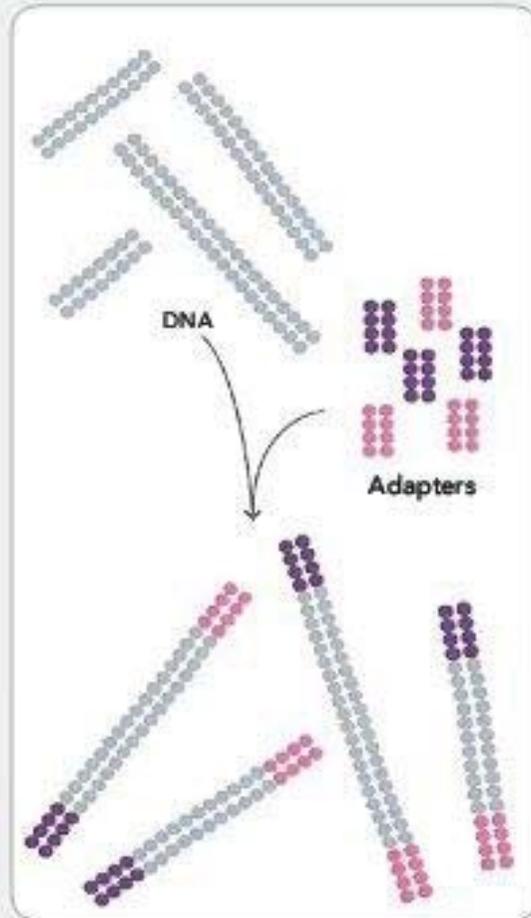


¹ Turcatti *et al.*, "A new class of cleavable fluorescent nucleotides: synthesis and optimization as reversible terminators for DNA sequencing by synthesis," *Nucleic Acids Res.*, (2008).

² Adessi *et al.*, "Solid phase DNA amplification: characterization of primer attachment and amplification mechanisms," *Nucleic Acids Res.*, (2000).

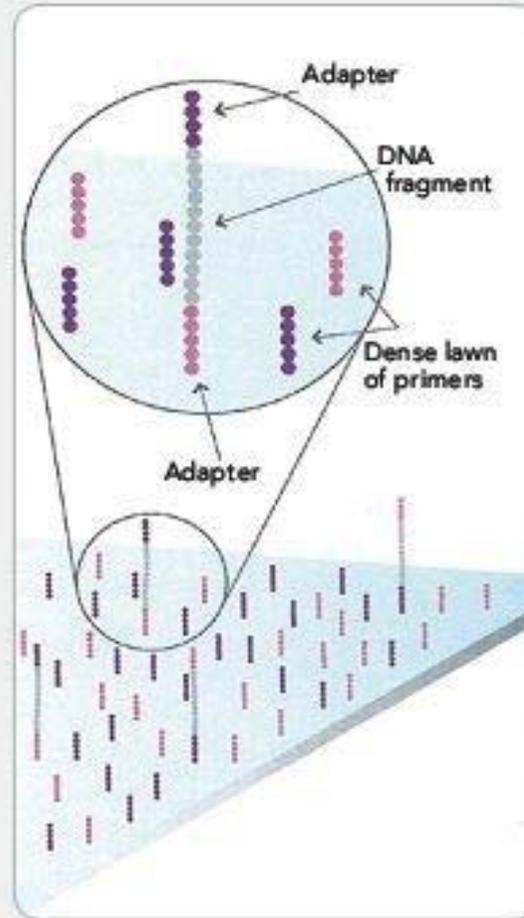
Step 1: Creating Fragments

1. PREPARE GENOMIC DNA SAMPLE



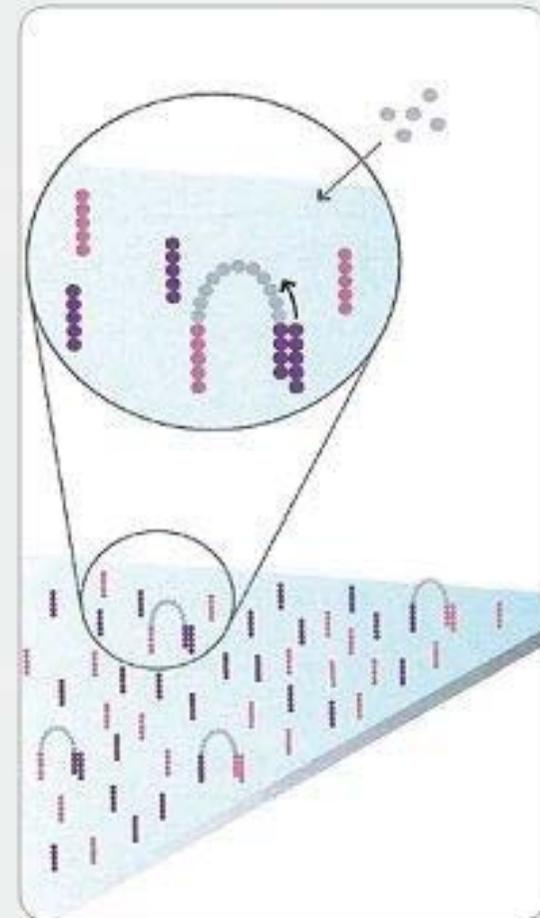
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



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

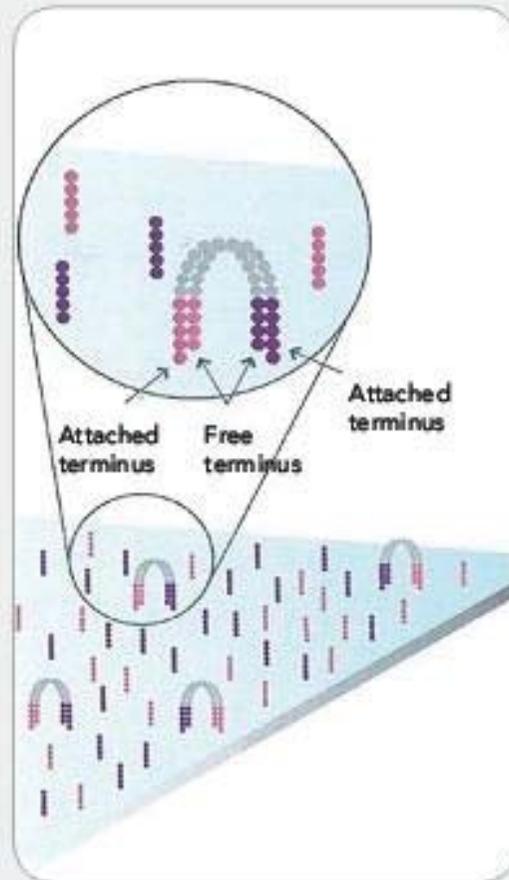
3. BRIDGE AMPLIFICATION



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

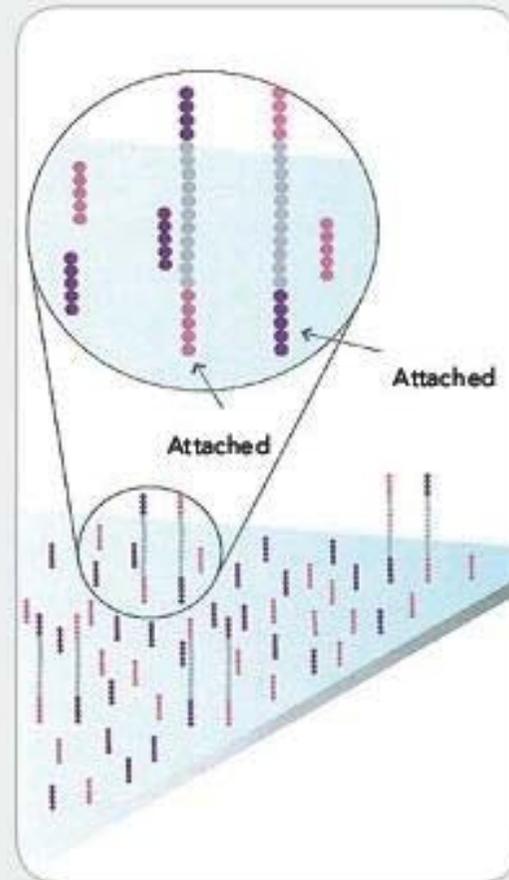
Step 2: Bridge Amplification

4. FRAGMENTS BECOME DOUBLE STRANDED



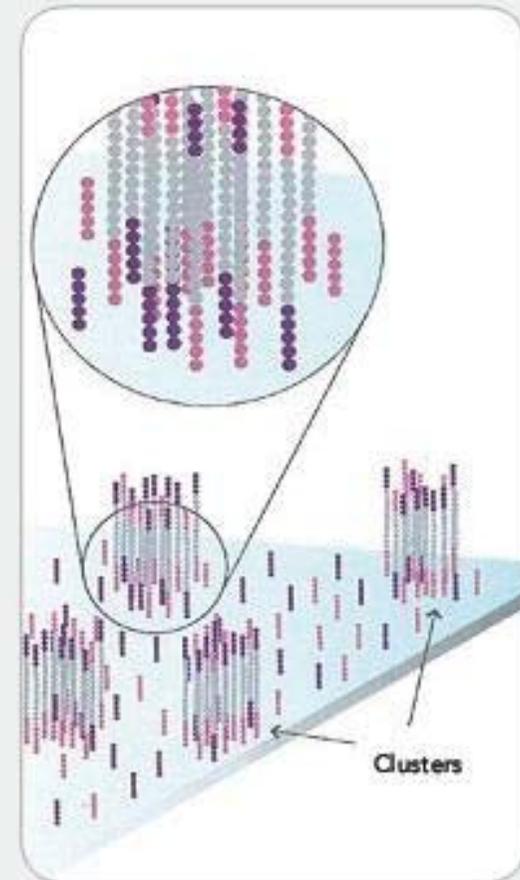
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



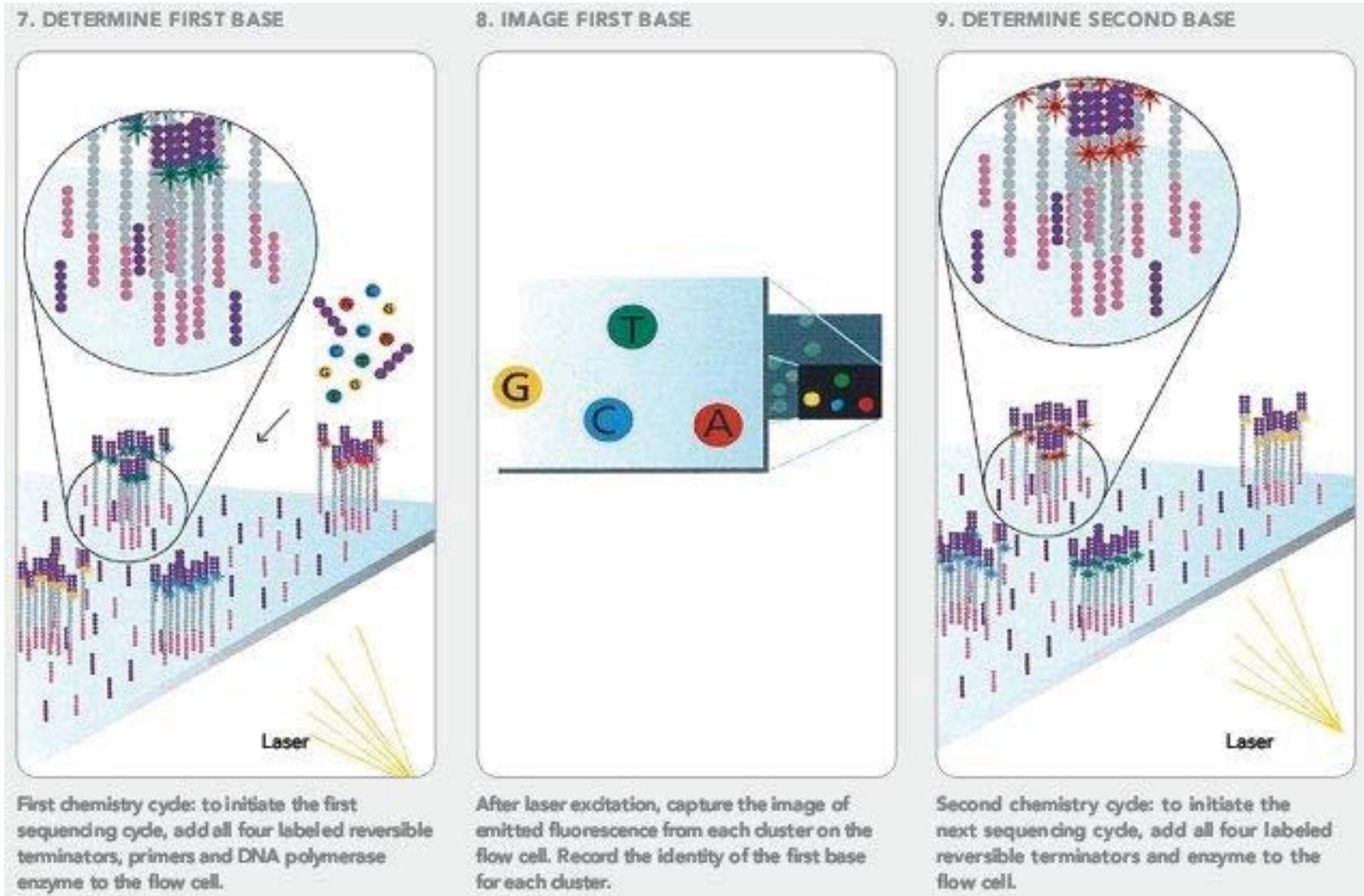
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



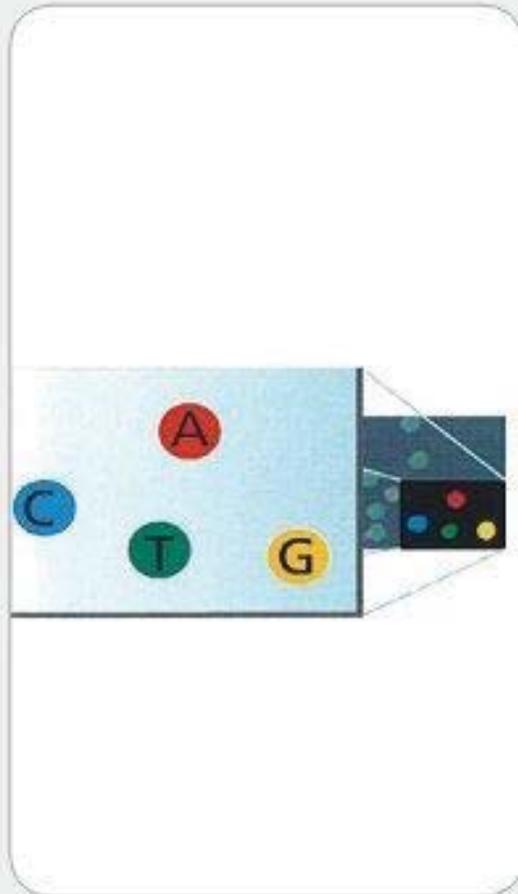
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

Step 3: SBS Using Labels



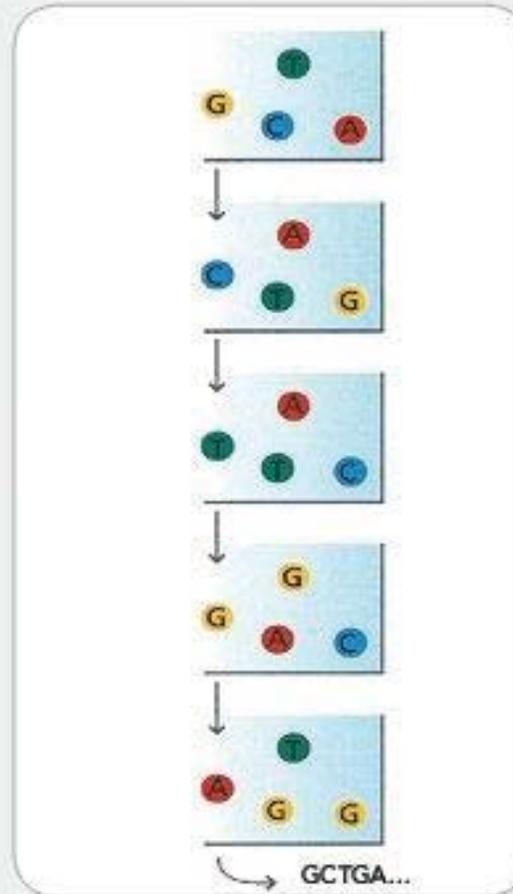
Step 4: Data Analysis

10. IMAGE SECOND CHEMISTRY CYCLE



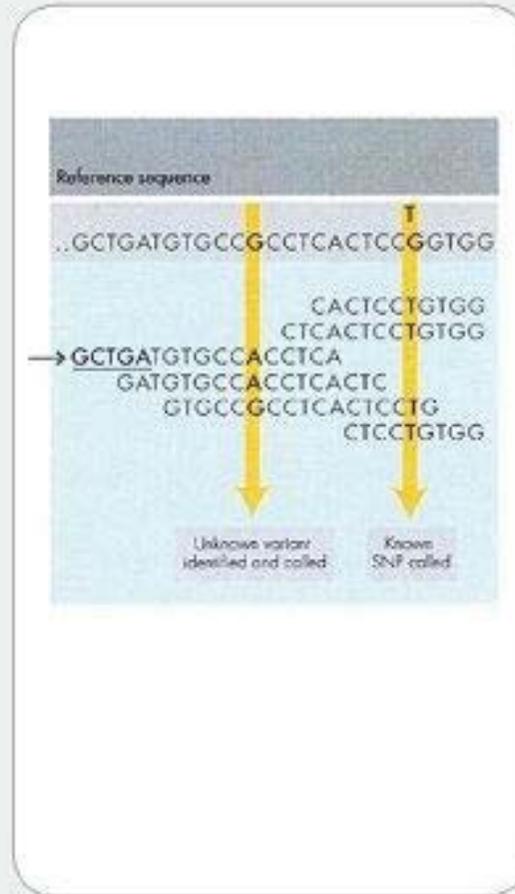
After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



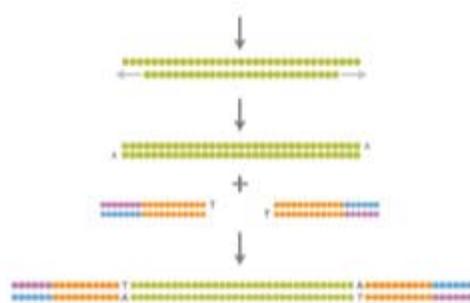
Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

12. ALIGN DATA



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

ILLUMINA (SOLEXA) DNA SEQUENCING



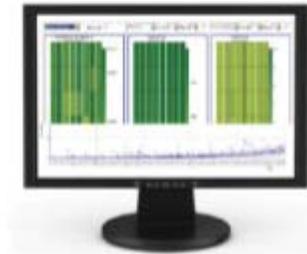
Library Preparation
~2 h [15 min hands-on (Nextera)]
< 6 h [< 3 h hands-on (TruSeq)]



Cluster Generation
~5 h (<10 min hands-on)



Sequencing by Synthesis
~1.5 to 11 days



CASAVA
2 days (30 min hands-on)

HiSeq 2000 (2011)

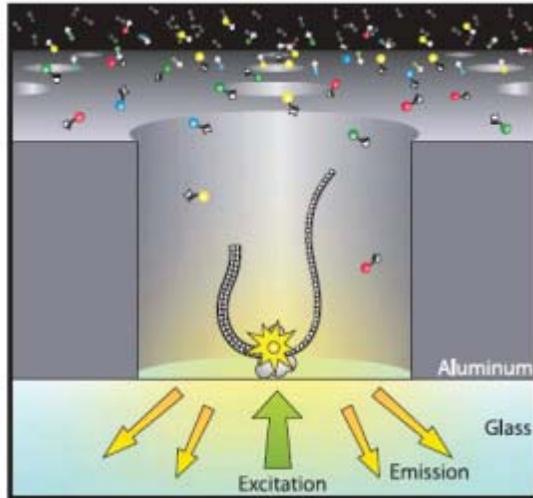
~ 100 bps read length
~ 100 Gbps per 2 days

It will take 2 days to sequence a human genome!

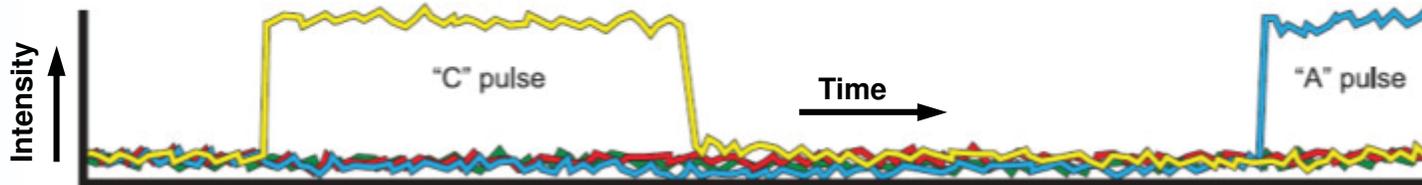
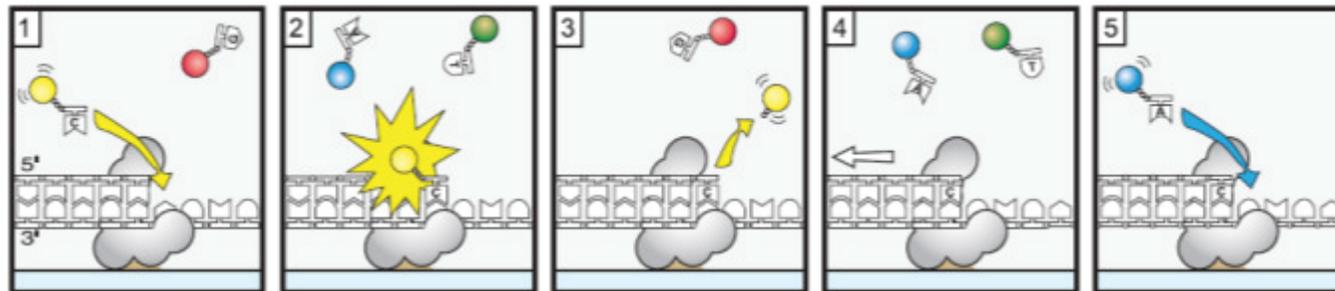
Upcoming Sequencers

Disclaimer: The presenter cannot verify the performance of these systems in October 2011

Pacific Biosciences¹ (PacBio)

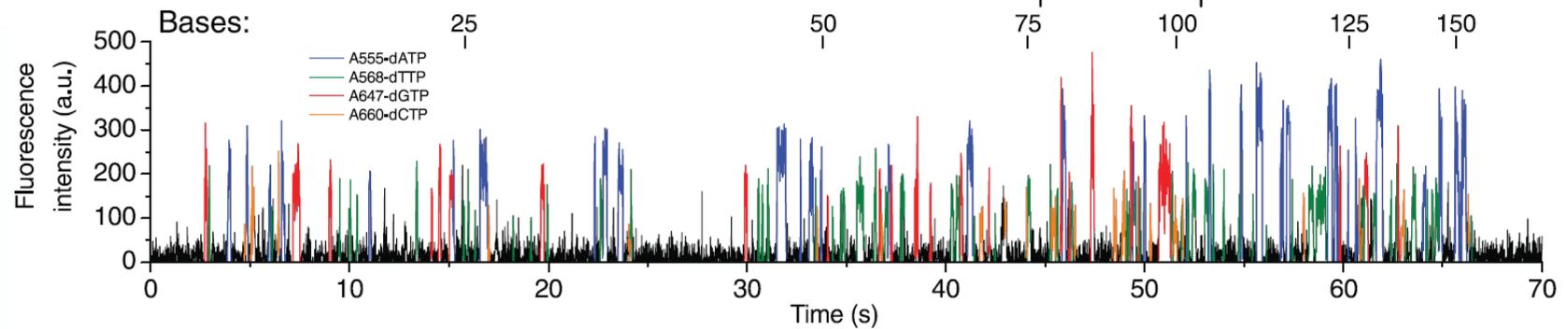
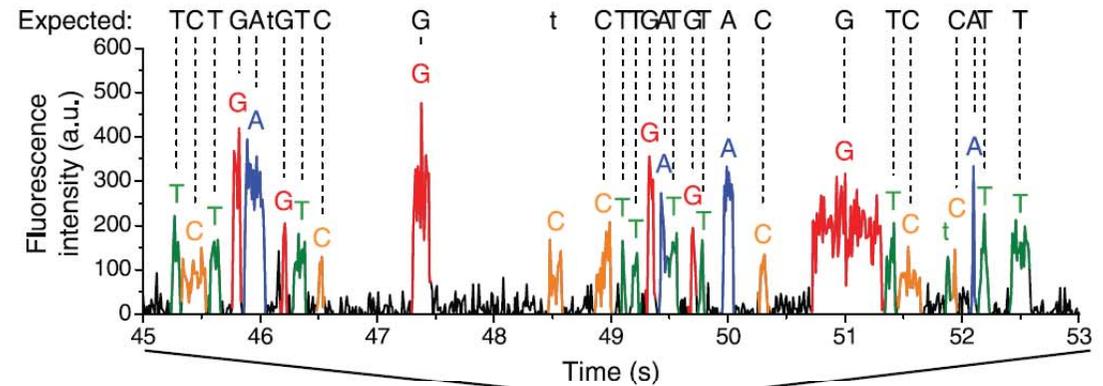


Polymerization of a single molecule is detected in real-time using a labeled nucleotides



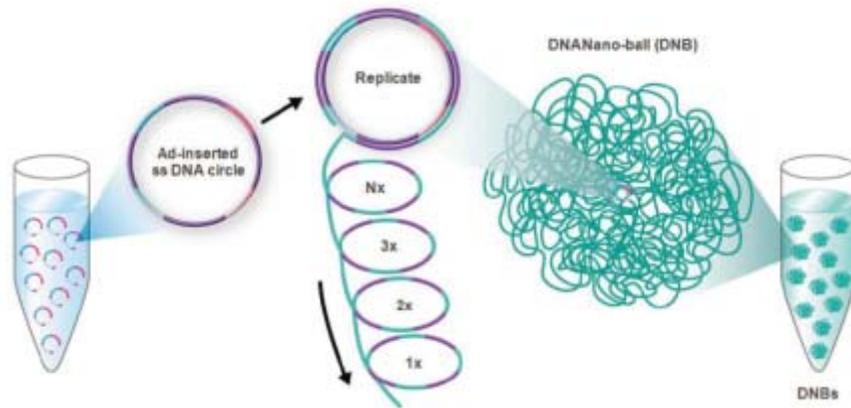
¹ J. Eid *et al.*, "Real-time DNA sequencing from single polymerase molecules," *Science*, 2009.

Pacific Biosciences (PacBio)



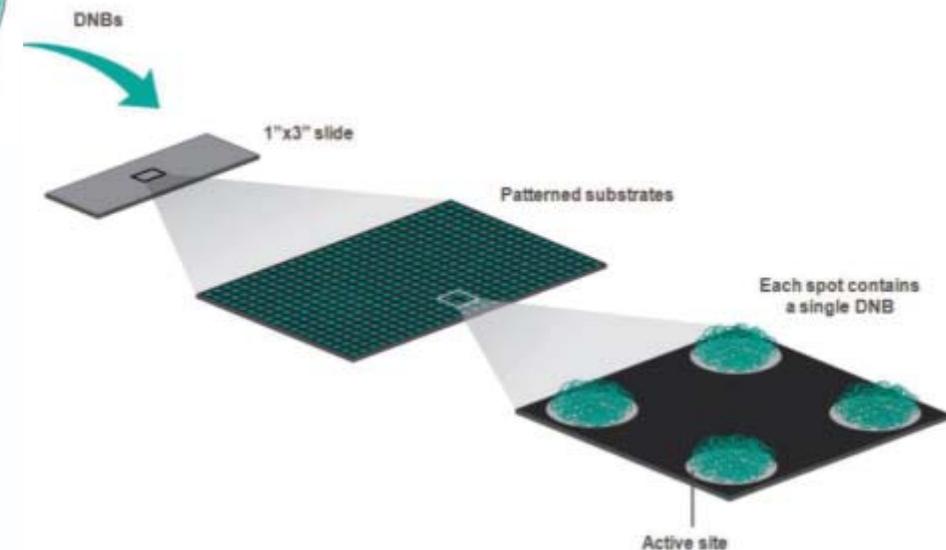
¹ J. Eid *et al.*, "Real-time DNA sequencing from single polymerase molecules," *Science*, 2009.

DNA Nanoball Sequencing¹ (Complete Genomics)



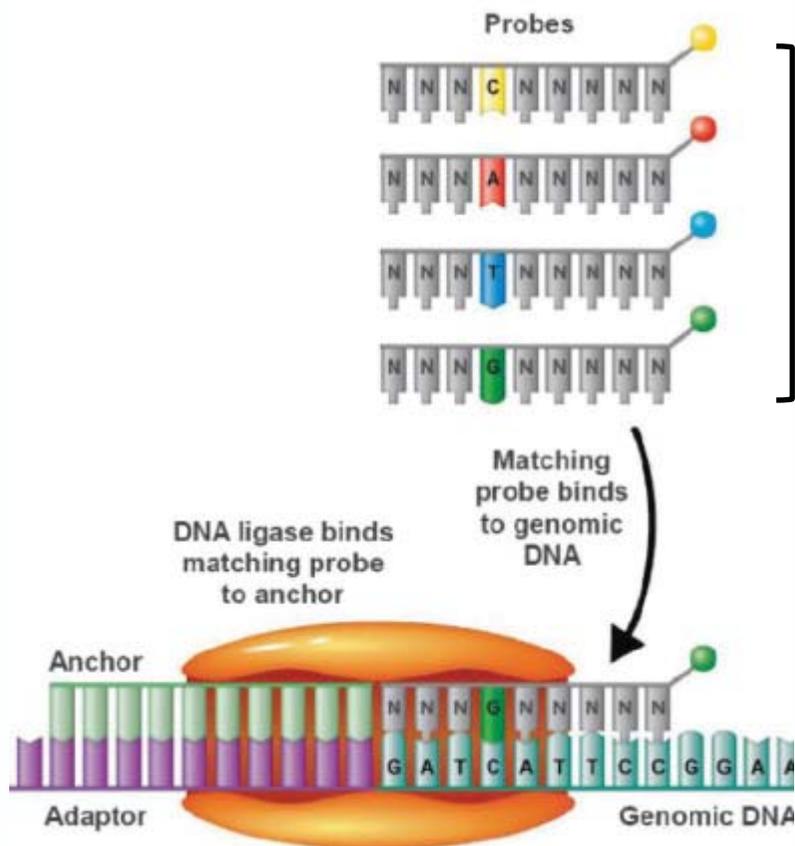
① Unknown DNA is amplified into repetitive long sequence to form a “nanoball”

② The “nanoballs” get randomly attached on a substrate with ~2.8B spots



¹ R. Drmanac *et al.*, “Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays,” *Science*, 2010.

DNA Nanoball Sequencing (Complete Genomics)

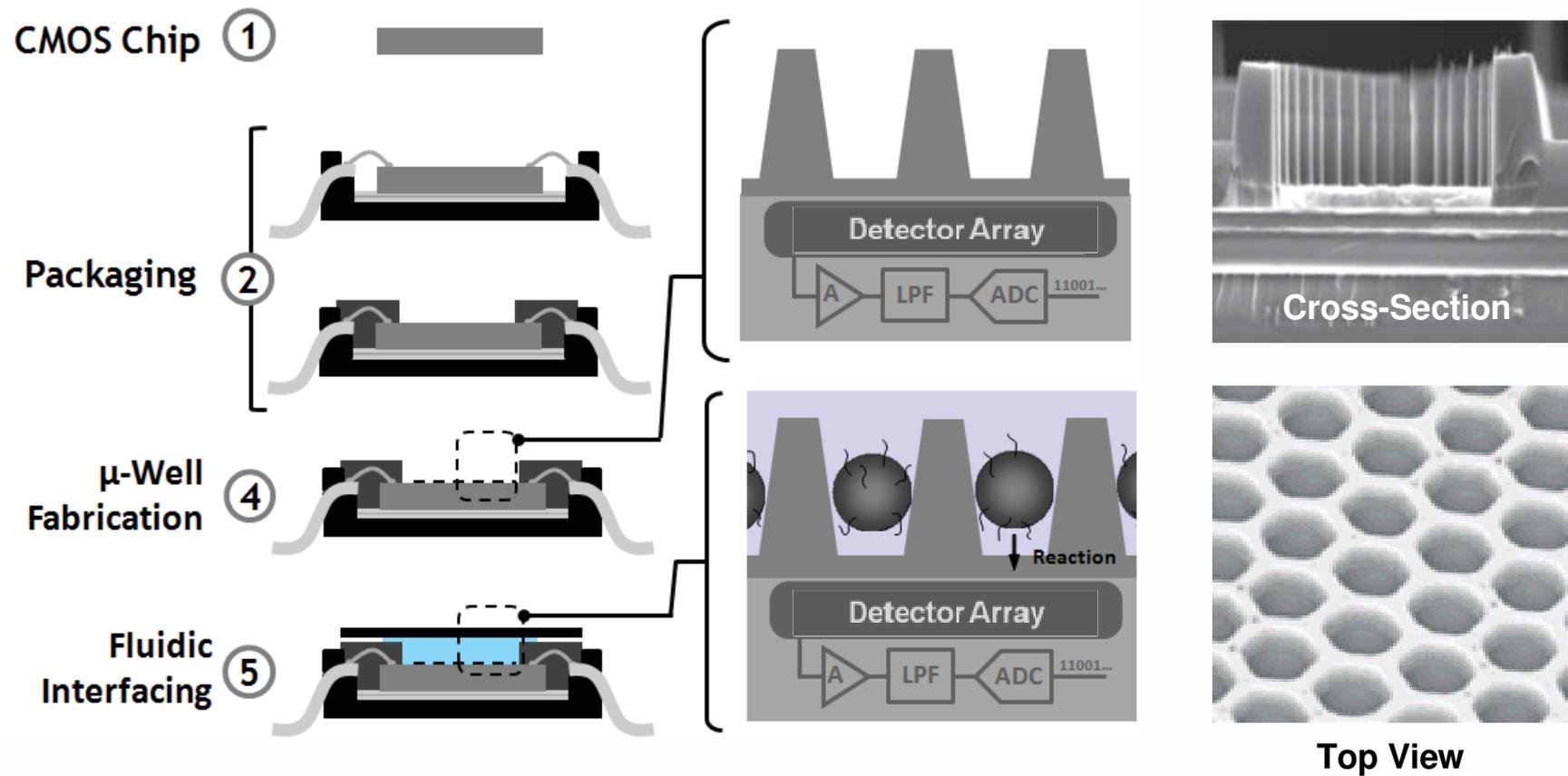


“N” is degenerate nucleotide, i.e., A, C, G, or T

Hybridization and ligation are used to find the sequence of specific locations on the unknown strand

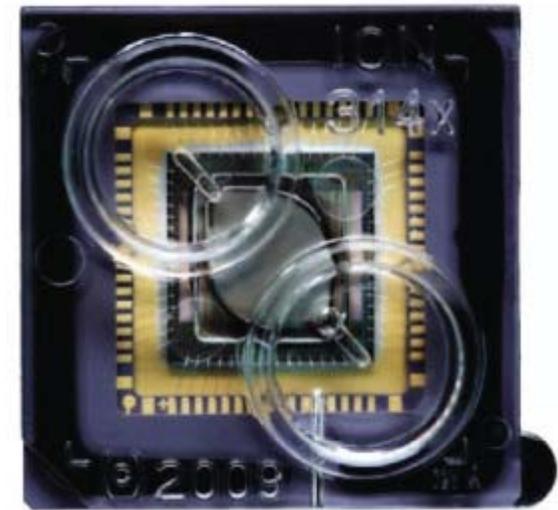
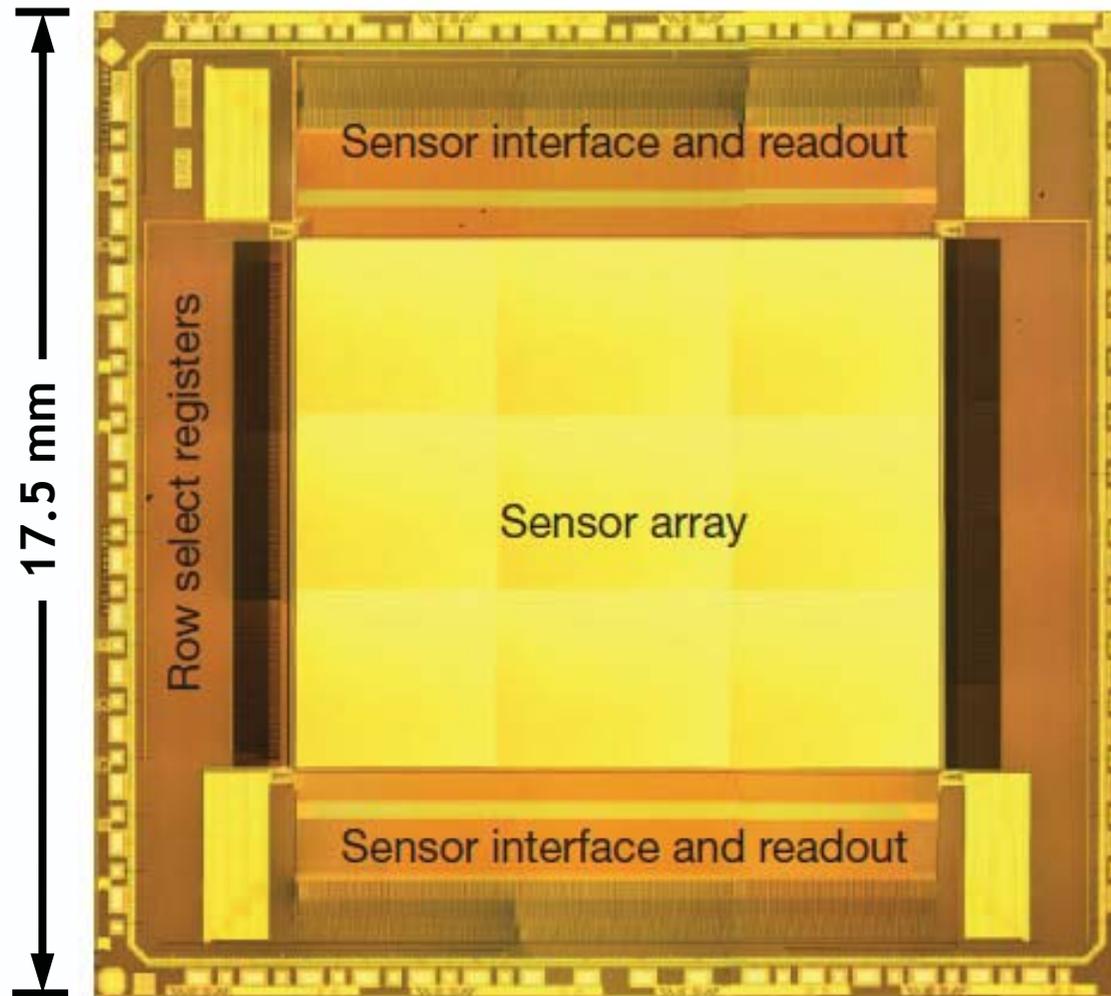
Semiconductor Sequencing¹ (Ion Torrent)

Detecting polymerization electronically



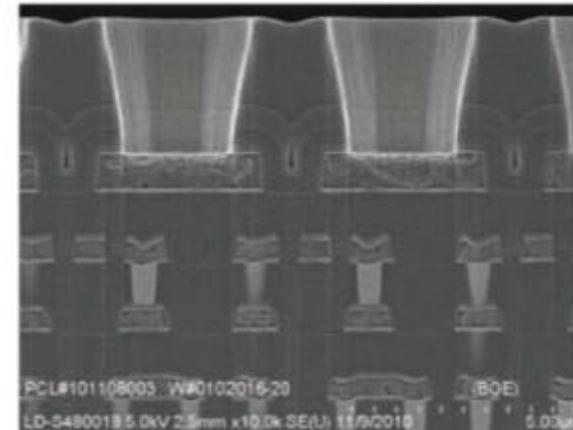
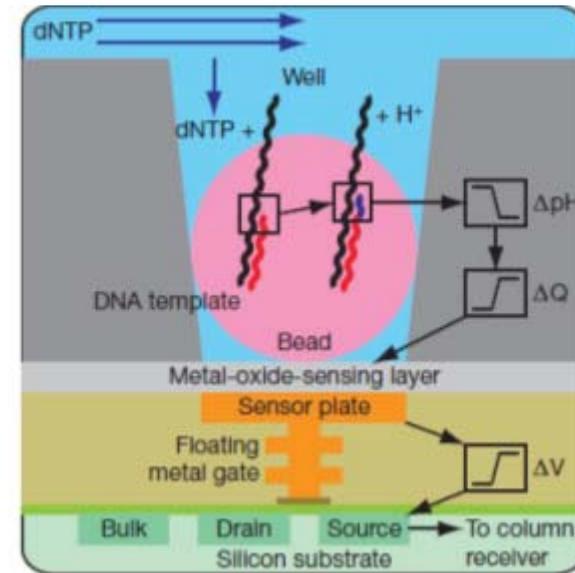
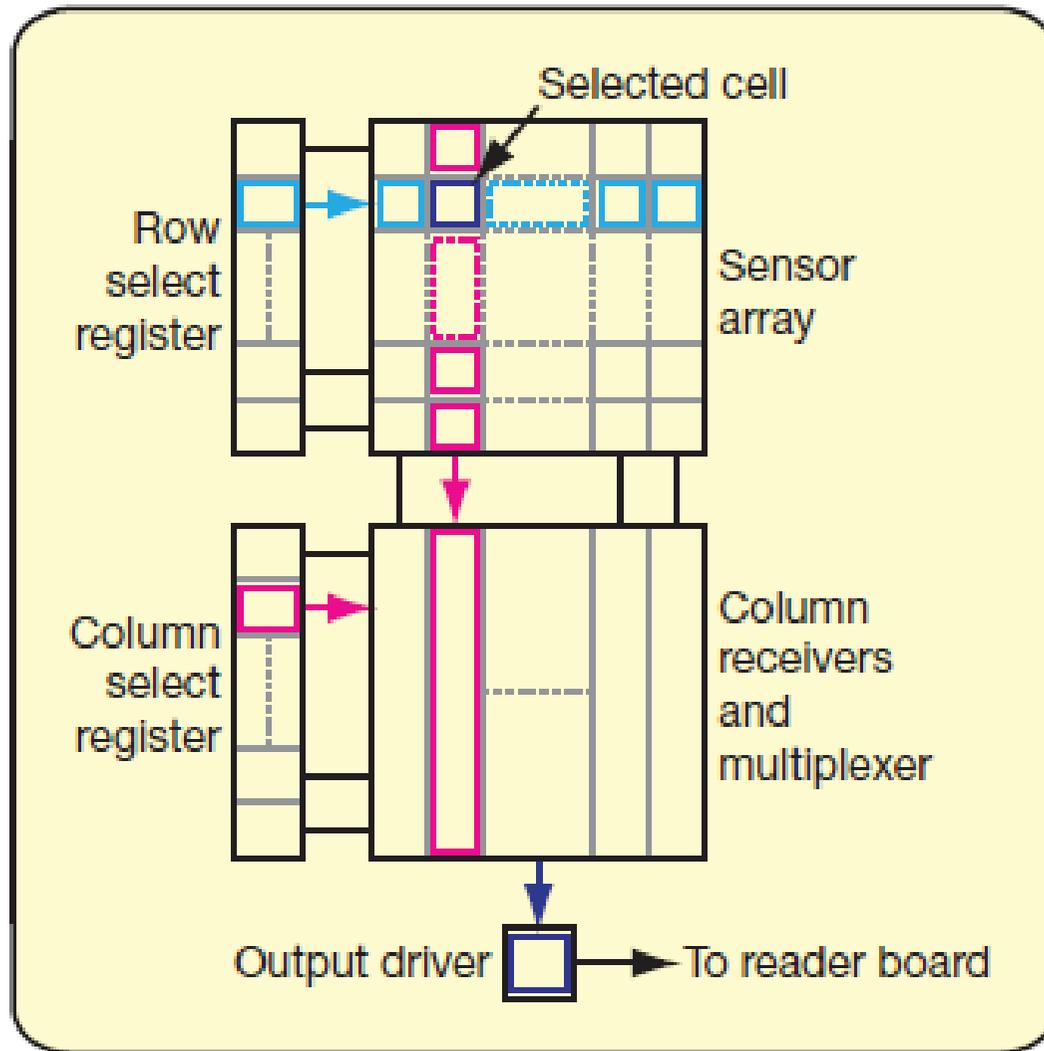
¹ Rothberg *et al.*, *Nature* (2011)

Ion Torrent Sequencing Chip¹



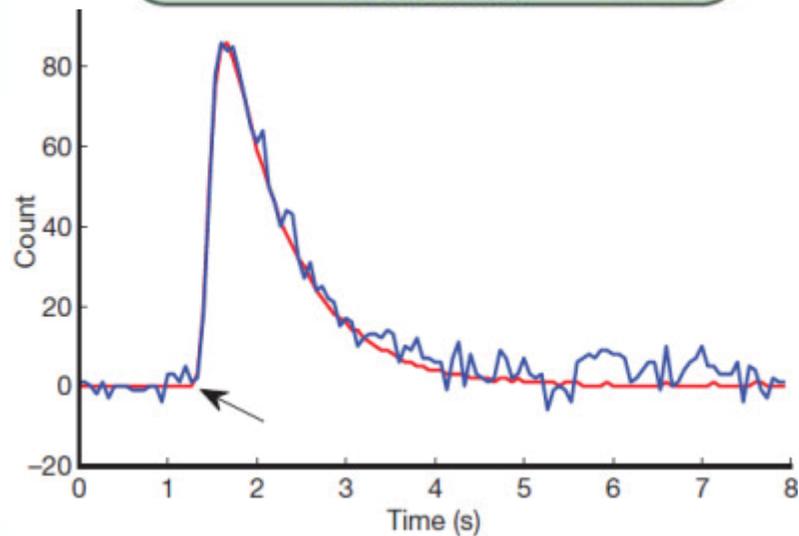
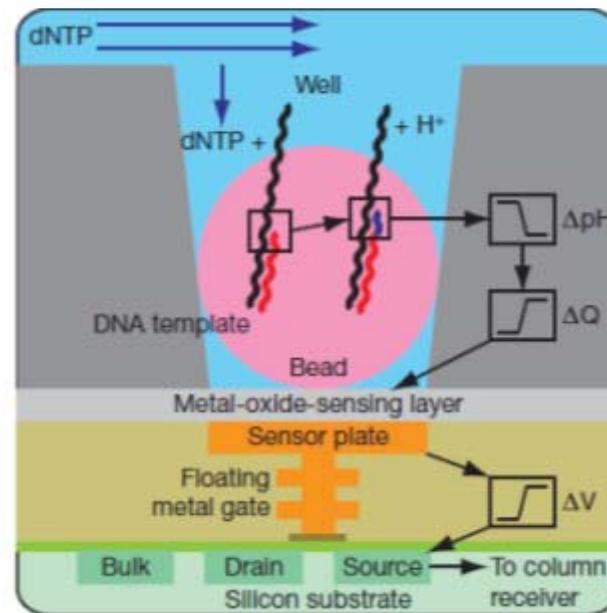
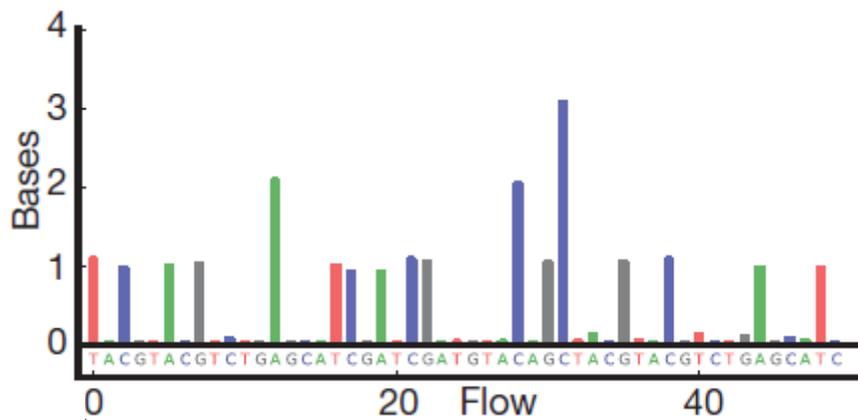
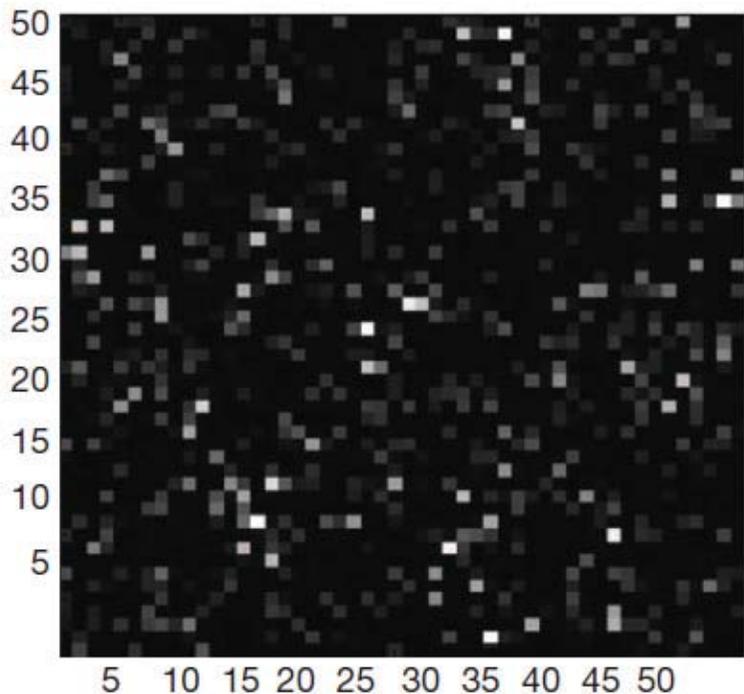
¹ Rothberg *et al.*, *Nature* (2011)

Architecture¹



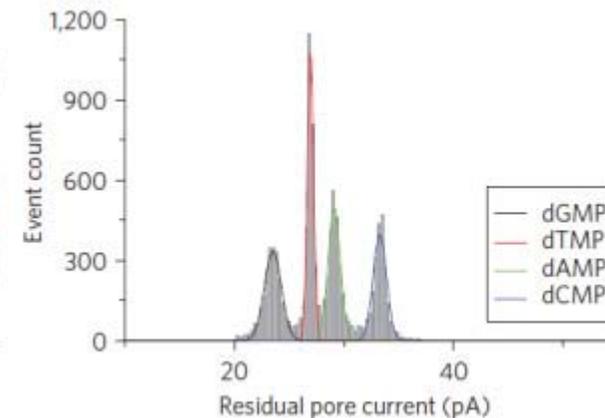
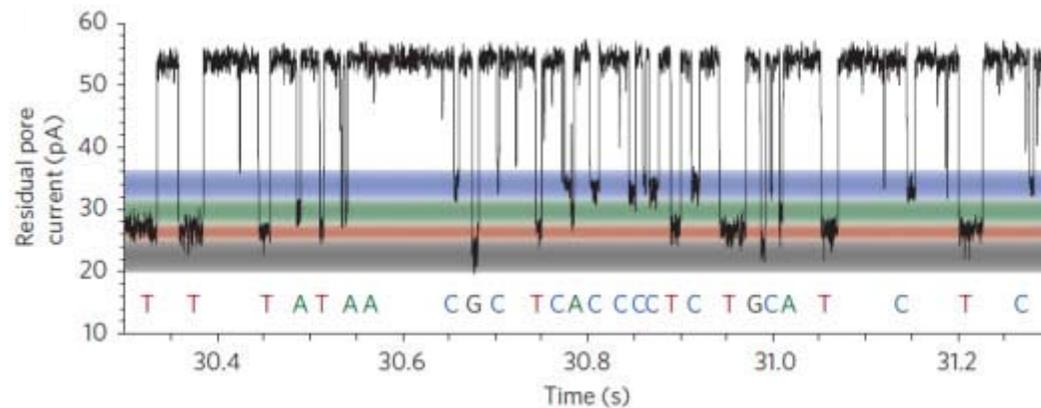
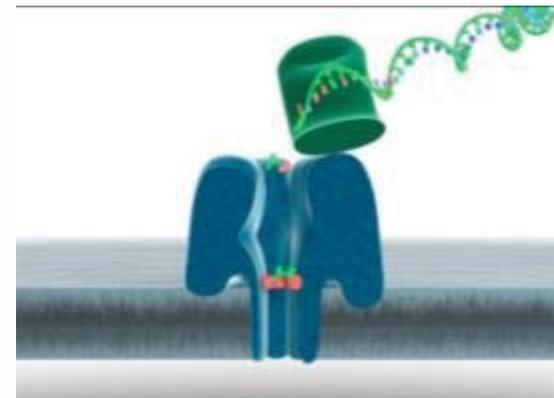
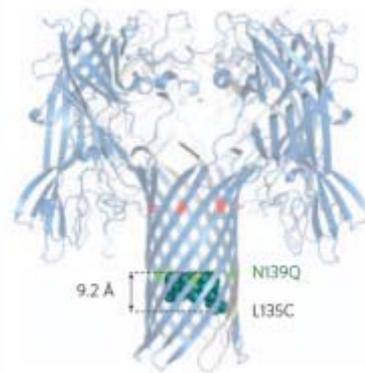
¹ Rothberg *et al.*, *Nature* (2011)

Electronic SBS Signals¹



Nanopore Sequencing¹ (Oxford Nanopore)

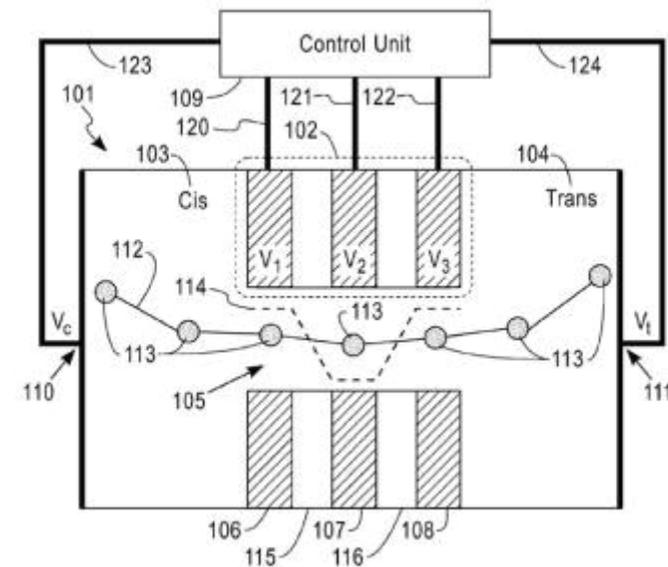
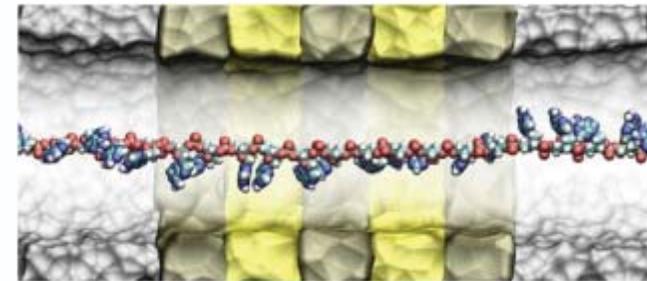
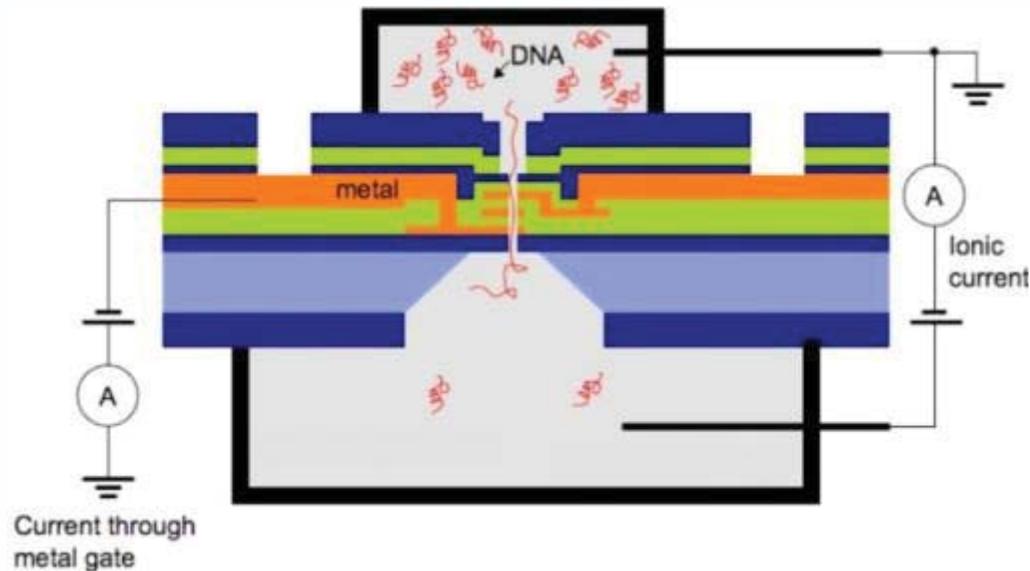
Digested DNA bases can create dissimilar blockage currents in a organic nanopore structure



¹ Clarke *et al.*, "Continuous base identification for single-molecule nanopore DNA sequencing," *Nature Nano*, 2009.

Nanopore Sequencing¹ (IBM/Roche)

Use solid-state probes with embedded electrodes to simultaneously ratchet and sequence DNA



¹ Polonsky *et al.*, "Nanopore in metal-dielectric sandwich for DNA position control," *Applied Physics Letters*, 2007.

Conclusion

- Sequencing matters
- It is a complicated problem
- It requires much more than biology and perhaps deep understanding of physics and engineering
- There are many crazy methods (and nebulous ideas) out there
- You should try to come up with your own method ... seriously!