



CHEMICAL BIOLOGY

- Moodle: <https://go.epfl.ch/CH-313>
 - Lecture slides (evening before the lecture)
 - Distributed presentation topics (assignments)
 - Forum (for questions and announcements)
- Examination (written, graded, detailed information will follow)
- Contact:
 - Moodle forum (for questions)
 - markus.jeschek@epfl.ch
- **“Concepts over details!”**
- **Interact! Ask! Discuss! Anytime!**

Group Presentations

- Critical discussion of primary literature
- Illustrative examples for topics from the lecture

- Why?
 - Repetition of core concepts, techniques etc.
 - Presentation skills and critical discussion of research
 - Insight into current research topics

- How?
 - Two students per group
 - Assignments distributed one week before delivery of presentation (via Moodle)
 - **Send slides: markus.jeschek@epfl.ch (Mon evening before presentation)**
 - **15 min presentation (both group members should present!) + Q&A**

EPFL Tipps for Group Presentations

- Rough structure
 - Short intro on general topic
 - Main presentation according to assignment
 - Brief outlook incl. points of criticism/open questions/personal opinion as kick-starter for the discussion
- Everybody should participate in the discussion, incl. constructive(!) feedback on presentation style
- Questionnaires with different points, feedback by peers
- Typical assignment:
 - You will receive a certain topic including a related publication
 - Introduce the topic using the publication
 - present the motivation behind the research, methodology, key results (not every graph!)
 - Additional questions will be provided hinting towards central points
 - Be encouraged to look/present beyond the questions and the provided paper

Group Presentations – Schedule

#	Name1	Name2	Presentation on...	Assignment on...
1	Winger Quentin	Jeremy	Sep 23, 2025	Sep 16, 2025
2	Ema	Ariane	Sep 30, 2025	Sep 23, 2025
3	Benjamin	Matthieu	Oct 7, 2025	Sep 30, 2025
4	Ivana	Ipek	Oct 14, 2025	Oct 7, 2025
5	Abigail	Robin	Oct 28, 2025	Oct 14, 2025
6	Mridhula	Elodie	Nov 4, 2025	Oct 28, 2025
7	Andrea	Florian	Nov 11, 2025	Nov 4, 2025
8	Melodie	?	Nov 18, 2025	Nov 11, 2025
9	Bastien	?	Nov 25, 2025	Nov 18, 2025
10	Nicole	Maria	Dec 2, 2025	Nov 25, 2025
11	Eva	?	Dec 9, 2025	Dec 2, 2025

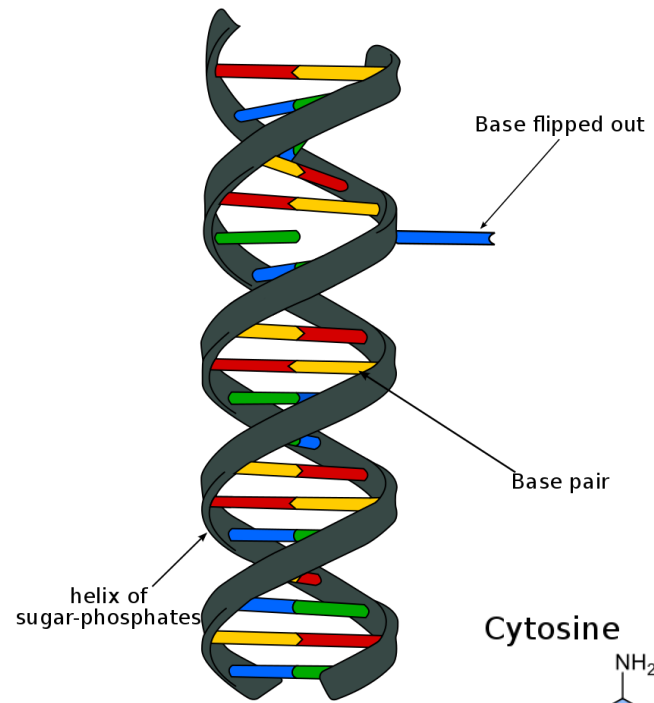
Course Topics – Overview

- Week 1 | Introduction + DNA
- **Week 2 | DNA**
- Week 3 | DNA
- Week 4 | RNA
- Week 5 | Protein/Enzymes
- Week 6 | Enzymes
- Week 7 | Enzymes
- Week 8 | Membranes
- Week 9 | Metabolism
- Week 10 | Metabolism
- Week 11 | Engineering
- Week 12 | Engineering
- Week 13 | Engineering
- Week 14 | LSAM Intro + Exam Preparation

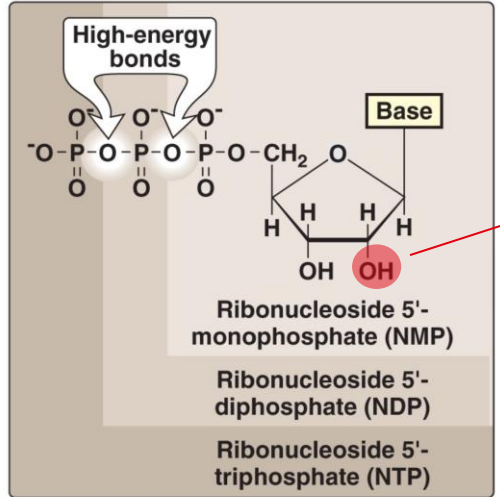
[tentative schedule]

DNA

Deoxyribonucleic acid (DNA)



DNA
Deoxyribonucleic acid

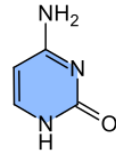


-OH → ribose (RNA)
-H → deoxyribose (DNA)

nucleotides
= nucleoside
+ phosphate(s)

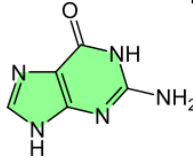
Copyright © 2011 Wolters Kluwer Health | Lippincott Williams & Wilkins

Cytosine



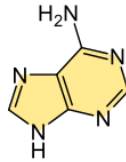
C

Guanine



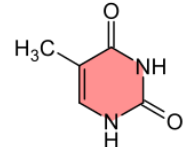
G

Adenine



A

Thymine

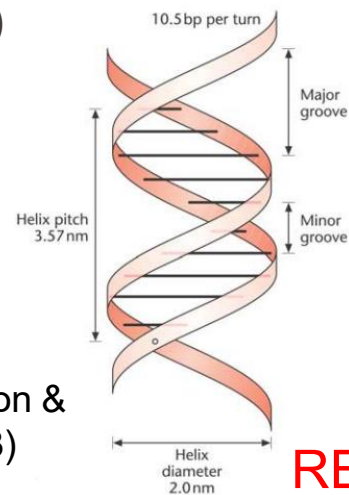
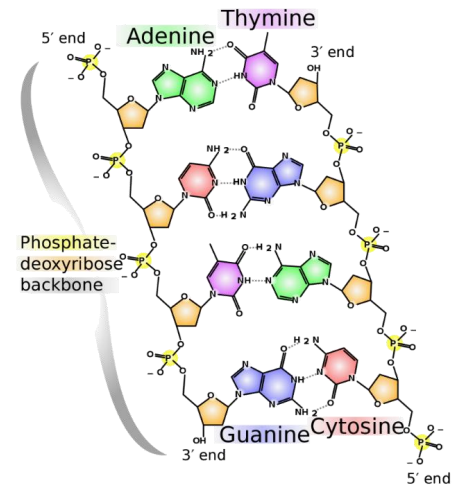


T

RECAP

Key Properties of DNA

- Long polymer
 - backbone: alternating phosphate and deoxyribose groups
 - negatively charged
 - double-helix structure
- Extremely stable (high T, solvents, hydrolysis etc.)
- Highly dense information storage (e.g. 215 petabytes/g)
- Codes for structure and function
 - non-coding RNAs
 - mRNAs → proteins
- **Directionality!** (3'-end, 5'-end)
- **Complementarity!**

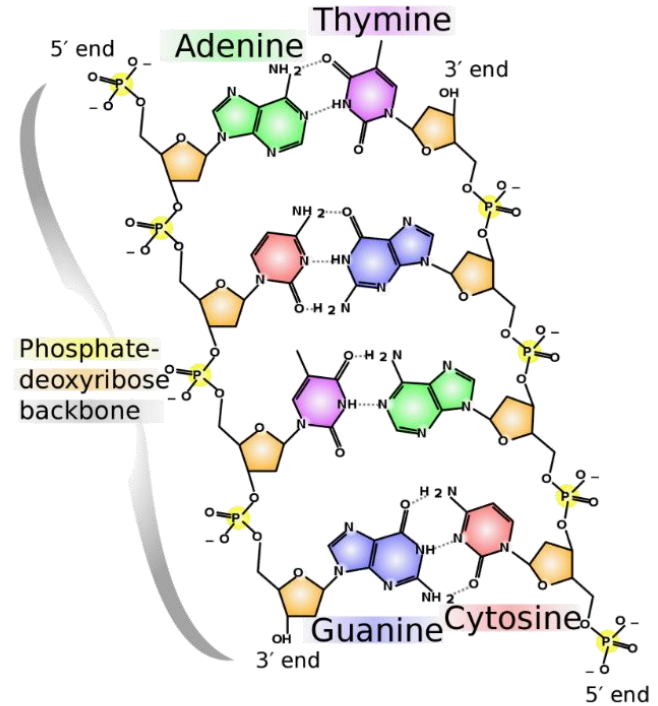


B-DNA (Watson & Crick, 1953)

RECAP

Complementarity of DNA

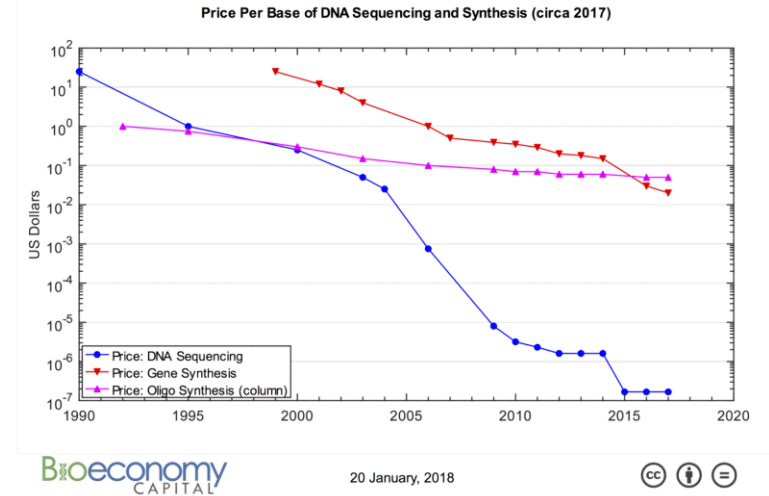
- Through non-covalent interaction, “base pairs” (bp)
 - A – T (two hydrogen bonds, ~4 kcal/mol),
 - C – G (three hydrogen bonds, ~6 kcal/mol)
- Essential biological property
 - DNA replication
 - transcription, Translation (genetic code)
 - gene regulation
 - etc.
- Basis for numerous biotechnological techniques and applications
 - PCR
 - gene synthesis, DNA assembly
 - sequencing
 - nanotechnology (e.g. DNA origami)
 - etc.



- DNA synthesis and sequencing
 - essential every-day tools across the entire life sciences
 - chemical tools play key roles in both
 - **knowing how it works is critical for correct use and applications!**

- DNA sequencing
 - “classical” (Sanger method)
 - “next-generation” techniques

- DNA synthesis
 - chemical synthesis of oligonucleotides
 - assembly of larger DNA molecules



Q: Why would you want to synthesize and sequence DNA? Name examples!

RECAP

DNA (“write”)

- Oligonucleotides
 - first step of DNA synthesis
 - single stranded!
 - 50 – 100 nt (rarely up to 300 nt)

DNA oligonucleotide (“oligo”)

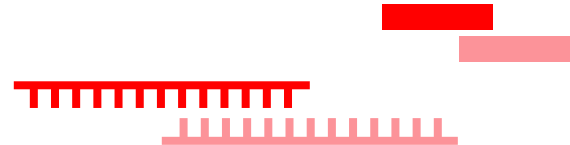
5' -ACGTACGTTTACTAG-3'



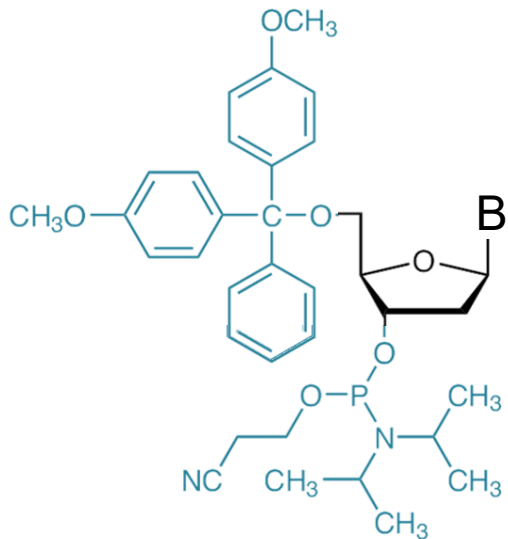
- Chemical oligo synthesis
 - Sequential coupling of single nucleotides
 - On solid support
 - 3' → 5' direction
 - Chemical protecting groups to avoid multiple couplings

5' -ACGTACGTTTACTAG-3'

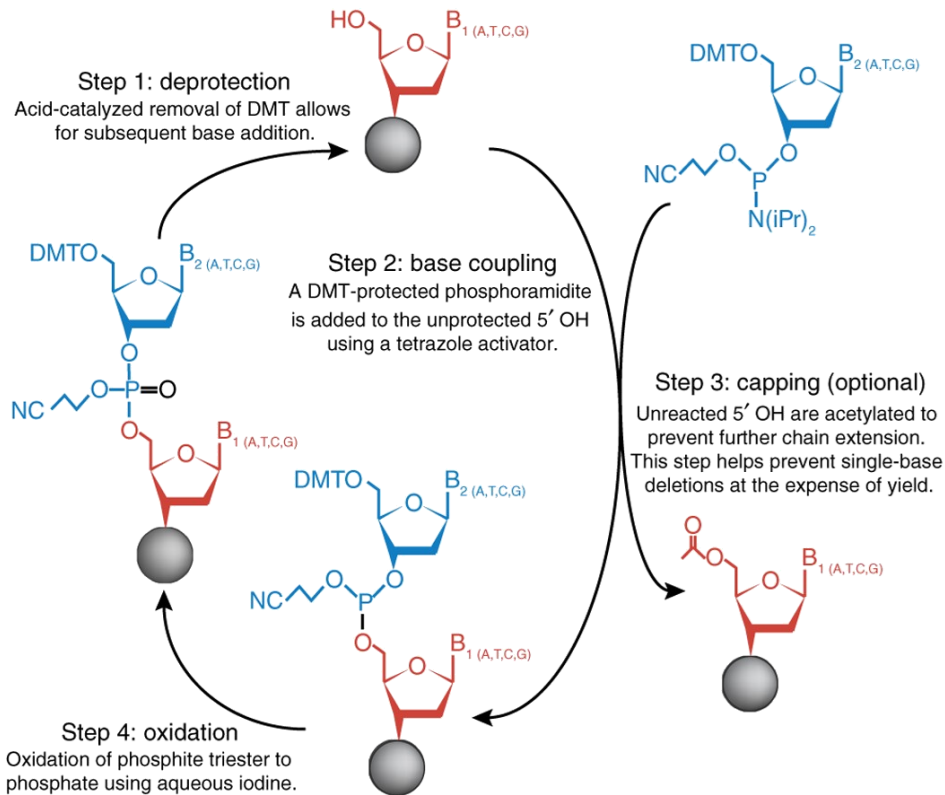
|||||
3' -AAATGATCTTACTAG-5'

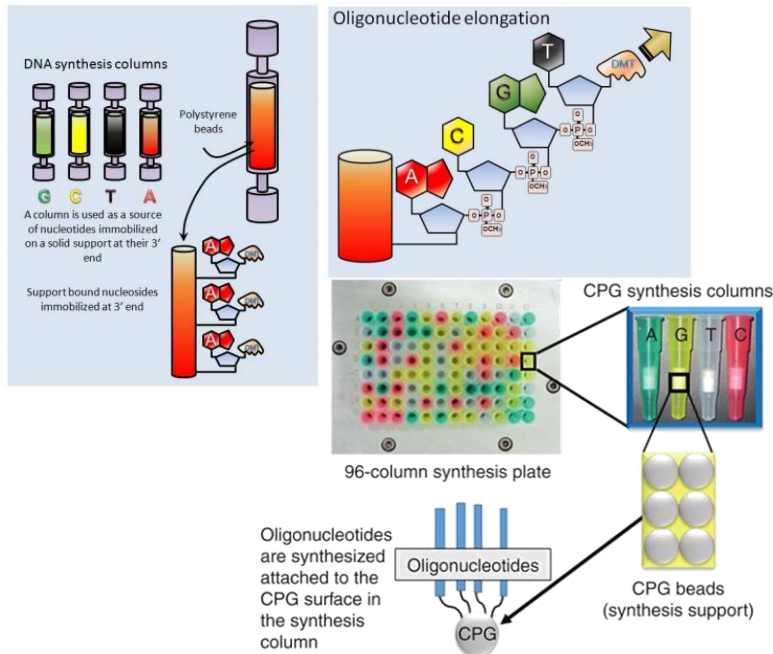


dimethoxytrityl (DMT)

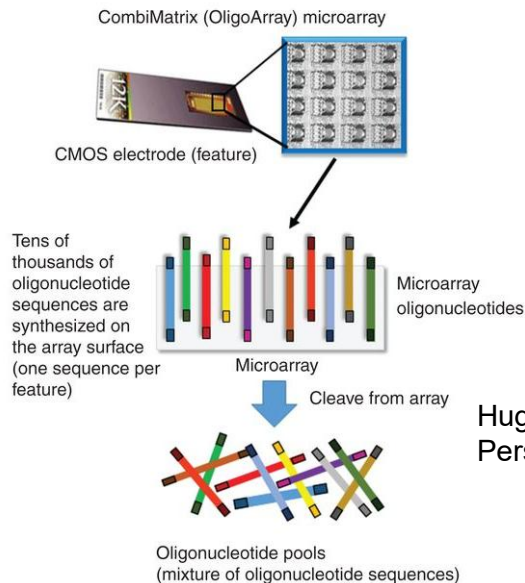


phosphoramidite





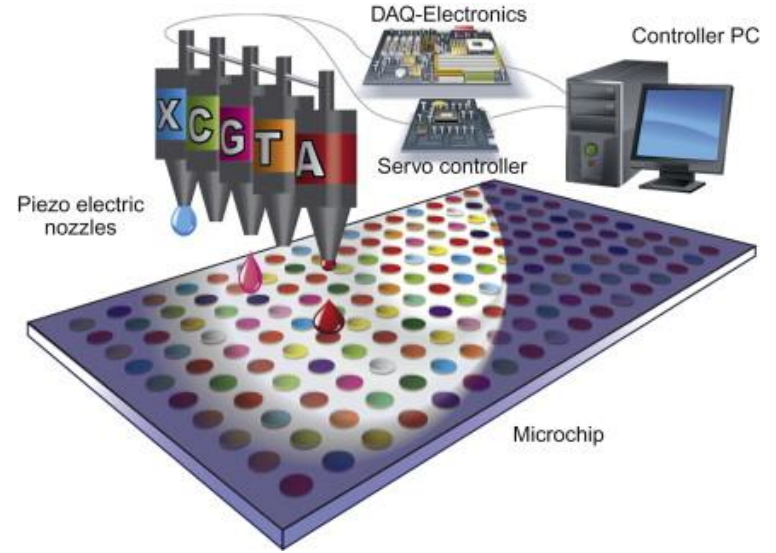
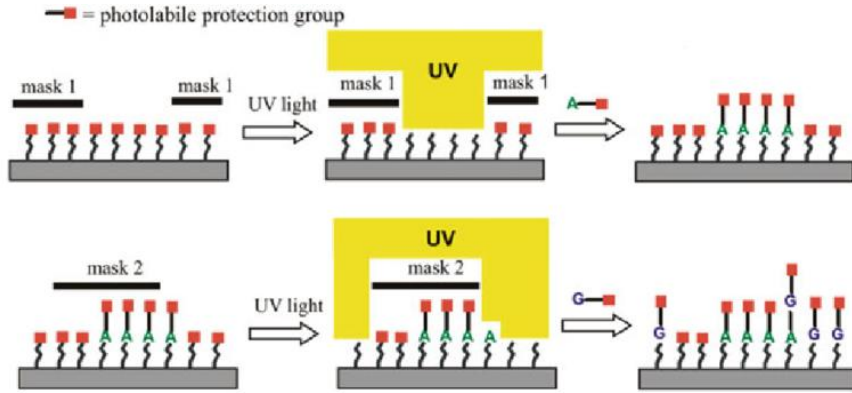
Microarray-based oligonucleotide synthesis



- “classical” , column-based
 - approx. 1-20 US-\$ per oligo
 - one sequence per column

- high-throughput , microarray-based
 - < 0.05 US-\$ per oligo
 - sequence pools/“libraries“ (max. 250,000 per chip)
 - lower quantities
 - e.g. photolithography or ink-jet technology

Microarrays - Photolithography & Ink-Jet Technology

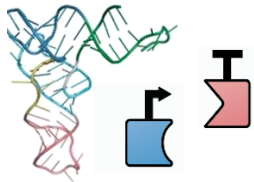


Synthesis Errors

- Deletions, depurination (A/G) → shortened product, cleavage
- Error rates: approx. 0.5-2% per nucleotide position
- Higher error rates for microarray synthesis
- Purification possible (electrophoresis, chromatography; expensive!)

Q: How high is the yield of full-length product for a 100 nt oligo at 1.5% per-position error?

Typical Size of Genetic “Parts”



Genetic «Part»

Approximate size range

promoters, terminators, tRNAs

50 – 200 bp *de novo synthesis*

genes, gene fragments



0.2 – 3.0 kbp

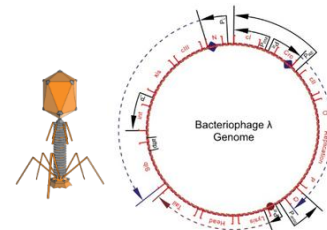


operons, genetic circuits

5 – 20 kbp

gene clusters, phage genomes, BACs

20 – 500 kbp

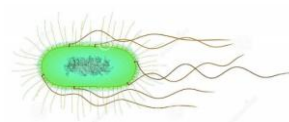


yeast chromosomes, small bacterial genomes

500 – 2000 kbp

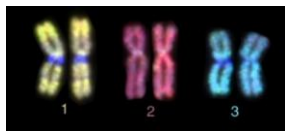
common bacterial genomes (*E. coli* !)

2.0 – 5.0 Mbp



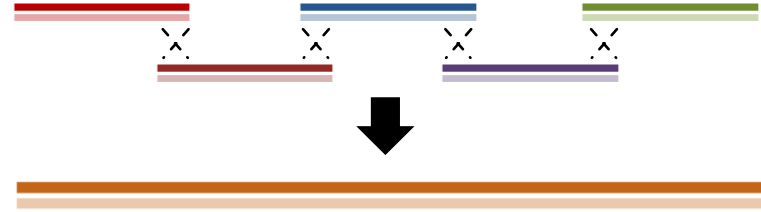
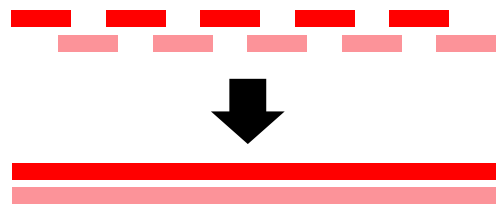
human chromosomes

50 – 250 Mbp



Recommendation: Refresh your memory on what these parts are/do!

DNA Synthesis – Overview



chemical synthesis
(de novo)

enzymatic DNA assembly

PCR-like → *in vitro* assembly → *in vivo* assembly

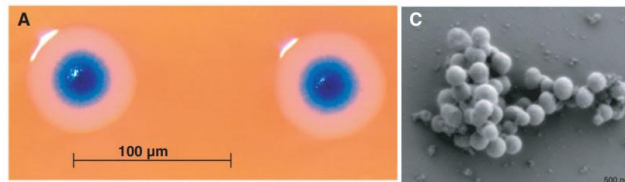
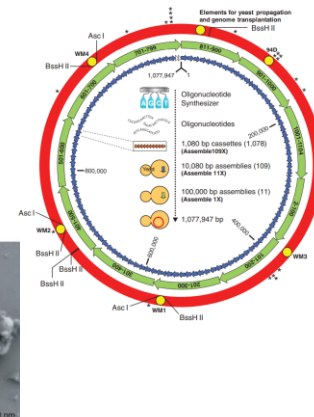
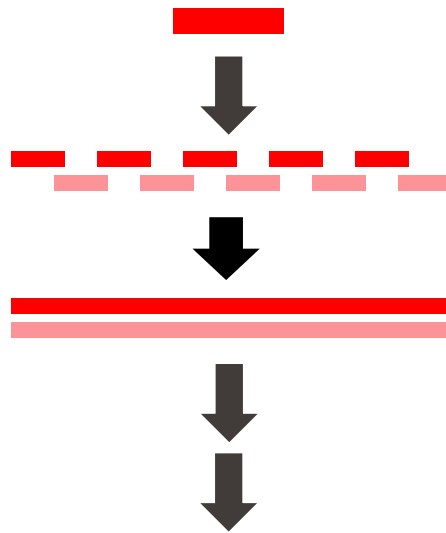
50 – 100 nt

500 – 5000 bp

10 – 200 kb

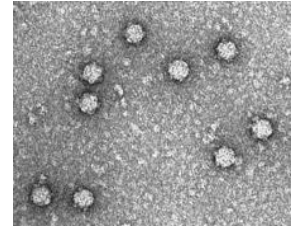
up to 1 Mbp
(few Mbp with tricks)

- Synthetic procedures to obtain pieces of DNA > 200 bp
- Step-wise joining of smaller pieces
- Verification by sequencing!
 - errors from *de novo* synthesis
 - assembly errors
- Modular process (construction of DNA libraries possible)
- General trend towards high throughput
 - parallelization, miniaturisation
 - total length, number of pieces \uparrow
 - price per assembly \downarrow

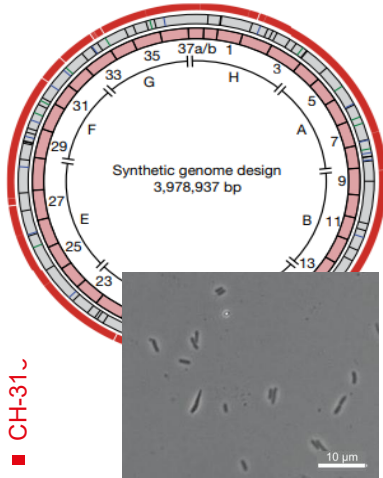


Generating a synthetic genome by whole genome assembly: ϕ X174 bacteriophage from synthetic oligonucleotides

Hamilton O. Smith, Clyde A. Hutchison III[†], Cynthia Pfannkoch, and J. Craig Venter[‡]



15440–15445 | PNAS | December 23, 2003 | vol. 100 | no. 26



nature

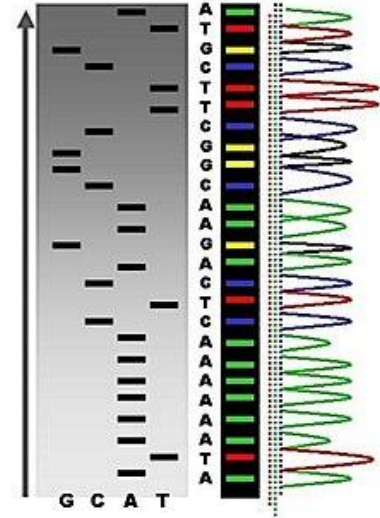
Total synthesis of *Escherichia coli* with a recoded genome

Julius Fredens^{1,4}, Kaihang Wang^{1,2,4}, Daniel de la Torre^{1,4}, Louise F. H. Funke^{1,4}, Wesley E. Robertson^{1,4}, Yonka Christova¹, Tiongsun Chia¹, Wolfgang H. Schmied¹, Daniel L. Dunkelmann¹, Václav Beránek¹, Chayasith Uttamapinant^{1,3}, Andres Gonzalez Llamazares¹, Thomas S. Elliott¹ & Jason W. Chin^{1*}

514 | NATURE | VOL 569 | 23 MAY 2019

DNA (“read”)

- Analytical process to determine the sequence of nucleotides (nucleobases) in a DNA molecule
- Key technology in molecular biology/biotech (“era of genomics”)
 - early methods: 1970s
 - wider availability: 1980s/1990s
 - “next-generation” sequencing: late 1990s/early 2000s, ongoing!
- Selected applications
 - molecular biology
 - evolutionary biology
 - epidemiology, virology
 - quality control (DNA synthesis, cloning etc.)
 - diagnostics, forensics
 - etc. etc.

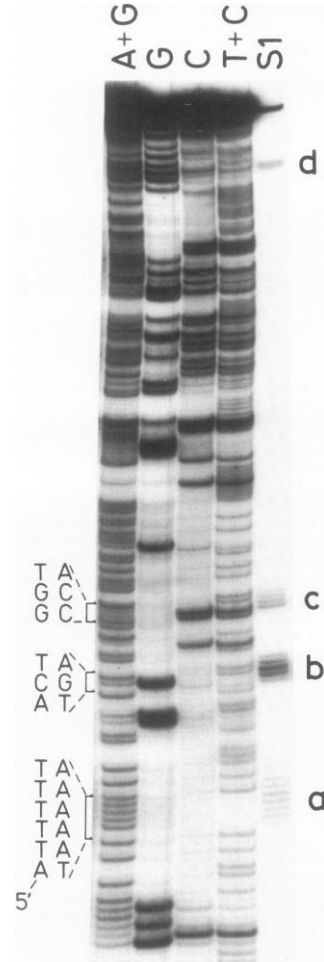


**Second generation
sequencer: 454 GS20**



- Allan Maxam & Walter Gilbert (1976–1977)
- Chemical modification of nucleobases → selective cleavage
 - (1) labelling of target DNA on 5'- or 3'-end (radioactive, dyes etc.)
 - (2) splitting into four samples and base-specific modification
 - A+G: depurination with formic acid
 - G: methylation with dimethyl sulfate
 - C+T: hydrolysis with hydrazine
 - C: hydrolysis with hydrazine in presence of NaCl
 - (3) cleavage at the modified base with hot piperidine
 - (4) fragments are separated/resolved in acrylamide gel electrophoresis
- Base modifications (step 2) occur stochastically
 - lower efficiency at greater lengths! (max. 200-300 bp)
- No longer in use!

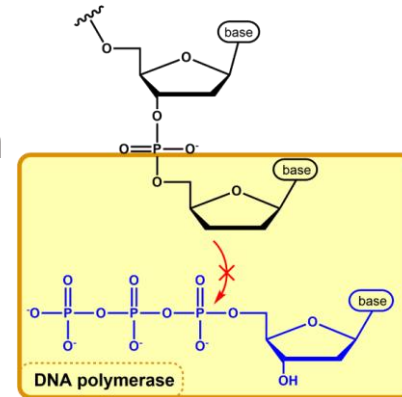
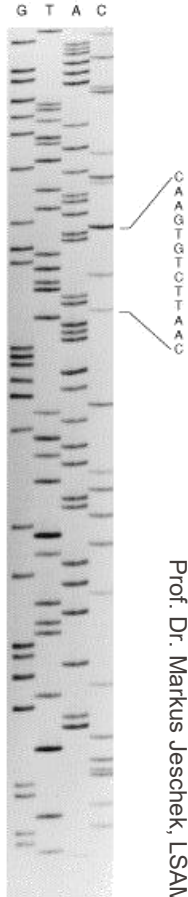
Maxam-Gilbert Sequencing



The Gold Standard – Sanger Sequencing

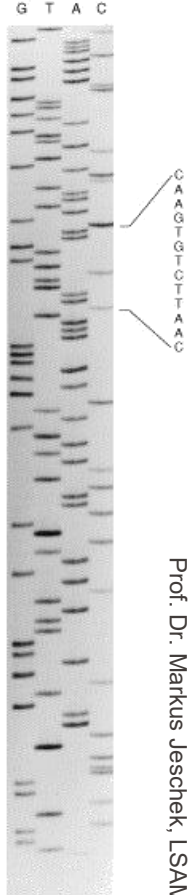
- Chain-termination/dideoxynucleotide method
- Simpler, faster, less hazardous, automatable
- History:
 - developed in 1977 by Frederick Sanger
 - first commercialized in 1986
 - still in use today!
 - standard low-throughput sequencing in “every lab”

- PCR-based method
- Premature termination of DNA polymerase reaction
- So-called “stop nucleotides”



Q: Why does the polymerization stop?

- Classical procedure:
 - (1) splitting of DNA into four samples
 - (2) addition of one di-deoxynucleotide triphosphates (ddNTP) to each sample (together with mix of all four dNTPs)
 - (3) addition of primer and DNA polymerase → elongation and stochastic (!) incorporation of ddNTP → termination
 - (4) separation of fragments in acrylamide gel electrophoresis (visualization through radioactively labelled primer or nucleotides)



Nature Vol. 265 February 24 1977

687

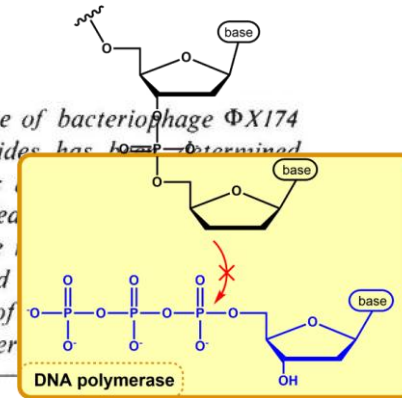
articles

Nucleotide sequence of bacteriophage Φ X174 DNA

F. Sanger, G. M. Air¹, B. G. Barrell, N. L. Brown¹, A. R. Coulson, J. C. Fiddes, C. A. Hutchison III², P. M. Slocombe³ & M. Smith⁴

MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK

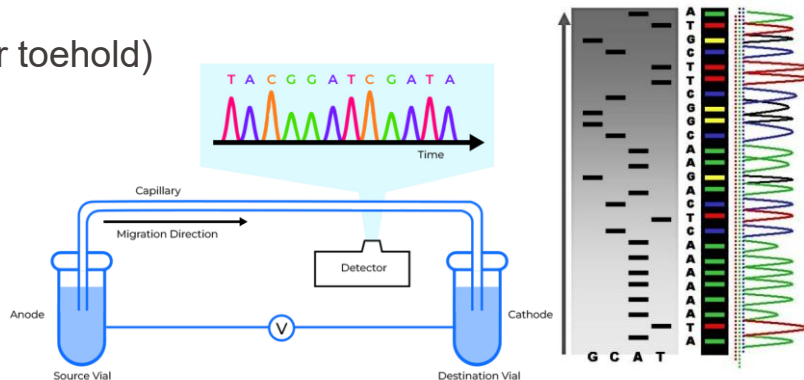
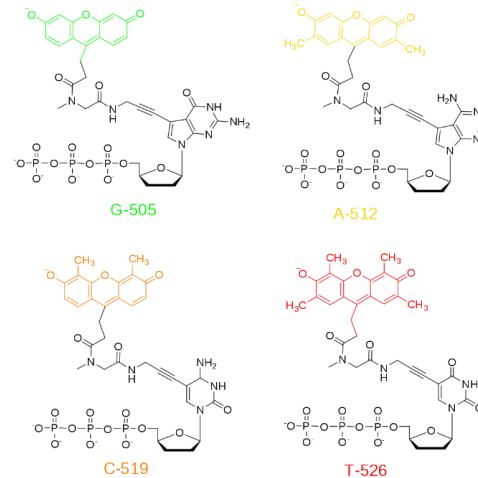
A DNA sequence for the genome of bacteriophage Φ X174 of approximately 5,375 nucleotides has been determined using the rapid and simple 'plus-minus' method. The sequence identifies many of the features of the production of the proteins of the phage, including initiation and termination of the genes for the proteins and RNAs. Two pairs of genes in the same region of DNA using different



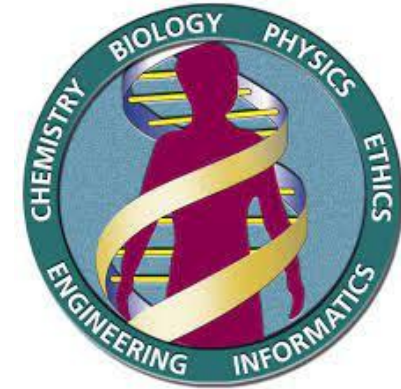
- Fluorescently labelled ddNTPs (“dye-terminator sequencing”)
 - no radioactive labels/autoradiography
 - one instead of four reactions
- Capillary electrophoresis
- Automation, miniaturization etc.
- 600 – 1200 bp, < 4 CHF per sample, <24 h result delivery**
- high accuracy (> 99.9% per positions)
- Limitations:
 - weak signal in the beginning (~50 bp from primer toehold)
 - length limit
 - repetitive sequences

Q: What could be reasons for the length limitation?

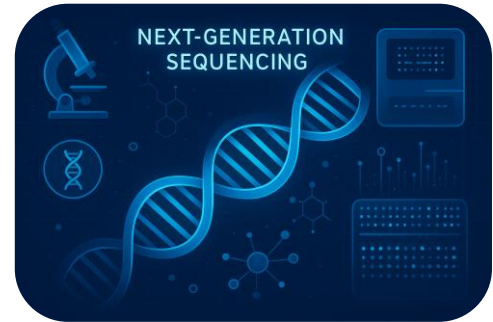
How can one sequence fragments larger than ~ 1kbp?



- 1990-2003
- consortium involving USA (coordinating), UK, Japan, France, Germany, China
- **3 billion US-\$ and 20 universities** (NIH funded)
- Based on Sanger sequencing and primer walking
- “End point” in 2003: ~85% of the genome (i.e. ~2.7 of 3.2 Mbp of the haploid genome)
- In parallel: privately funded project, Craig Venter (shotgun sequencing!)
- Jan 2022: final gapless assembly
- **Today: <1000 US-\$ per human genome (NGS)**

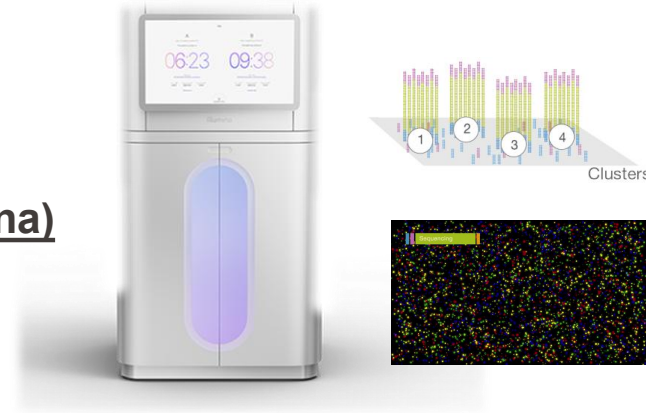


- “Massively parallel sequencing“, “2nd/3rd generation sequencing“, “deep sequencing“, “high-throughput sequencing”...
- Different “recent” technologies (from late 1990s/early 2000s)
- application scope (selection)
 - whole-genome sequencing (“shotgun” approach)
 - diagnostics
 - transcriptomics, amplicon sequencing (libraries!)
 - single-cell sequencing approaches
 - etc.
- **key technology for modern day biotech and medicine**



- “Second generation“ methods

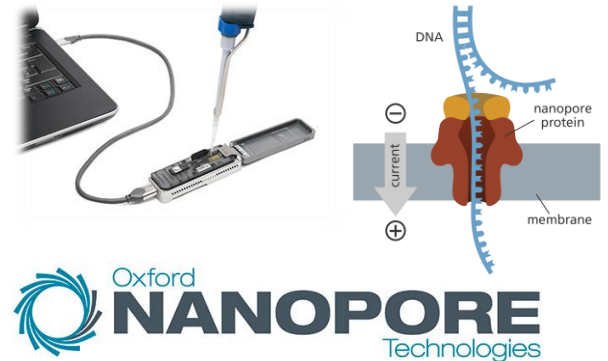
- (1) Pyrosequencing (Roche 454)
- (2) Sequencing by synthesis (Illumina)
- (3) Sequencing by ligation
- (4) Ion semiconductor sequencing



illumina®

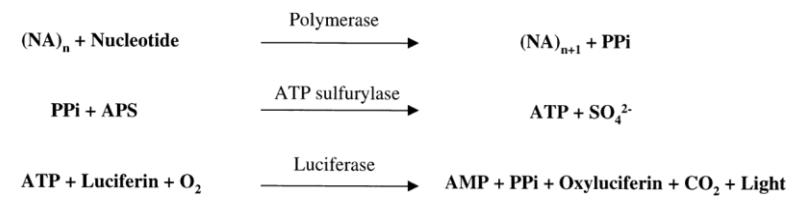
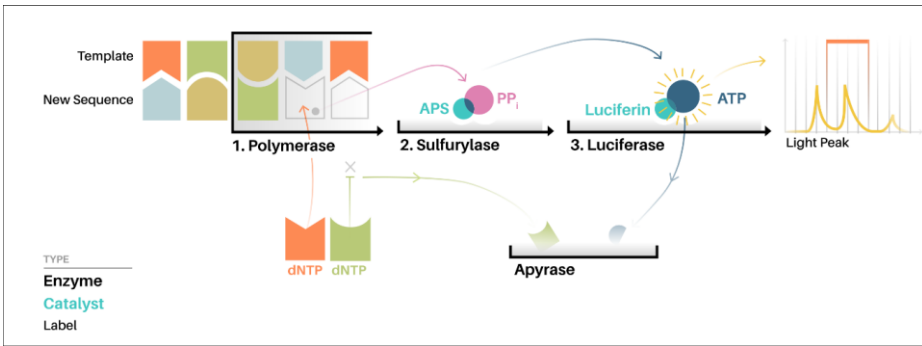
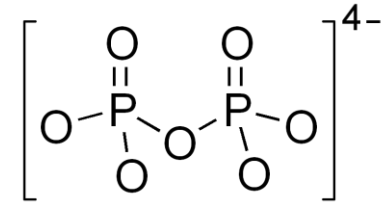
- “Third generation“ methods (single molecule)

- (1) Single-molecule real-time sequencing (SMRT; Pacific Biosciences)
- (2) Nanopore sequencing (MinION etc.; Oxford Nanopore Technologies)



Oxford
NANOPORE
Technologies

- **Sequencing by synthesis** (following Sanger approach)
- dNTP incorporation detected via enzyme-coupled assay
 - (1) **DNA polymerase** adds one nt releasing pyrophosphate (PPi)
 - (2) **sulfurylase** converts adenosine 5'-phosphosulfate (APS) and PPi to ATP
 - (3) ATP used by **luciferase** to convert luciferin under light emission
 - (4) unincorporated dNTPs/ATP degraded by **apyrase**



- 1st benchtop HTP sequencer
- Target DNA ligated to “**adaptors**”
- Single-stranded fragments captured on beads (**one per bead!**)
- Emulsion PCR to amplify fragments **monoclonally**
- Decorated beads loaded into picotiter plate (**Poisson distribution!**)
- $\sim 10^6$ Wells = reaction vessels ($\sim 29 \mu\text{m}$ diameter)
- dNTPs added one by one

- Limitations:

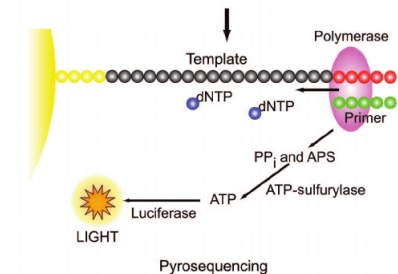
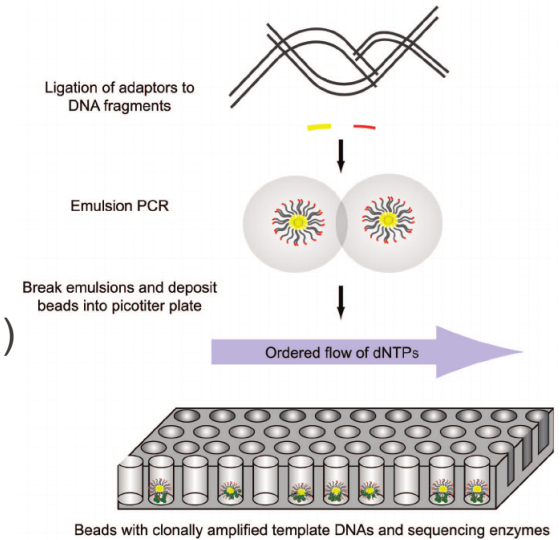
- no termination \rightarrow multiple incorporations, homopolymer read errors
- expensive

- discontinued 2013

Q: Why is the amplification by PCR needed?

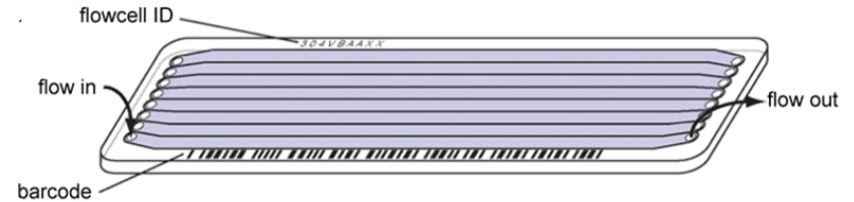
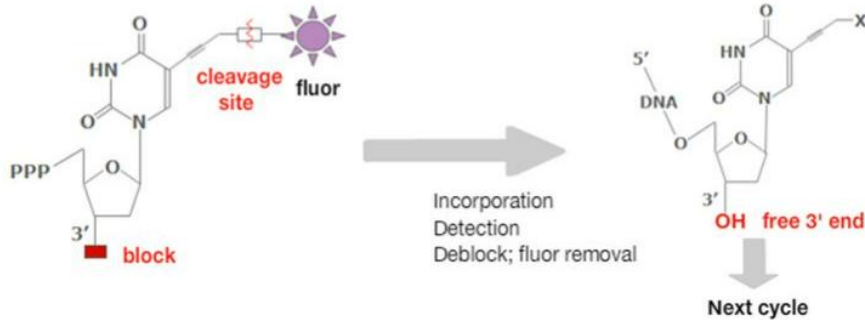
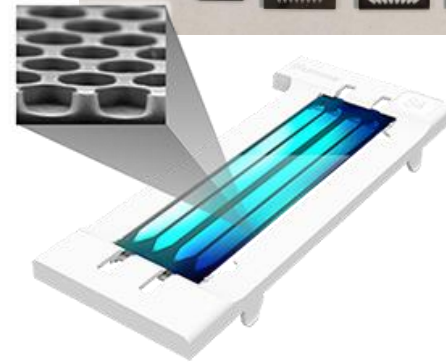
What is “monoclonality” and why is it needed?

What is a Poisson distribution?

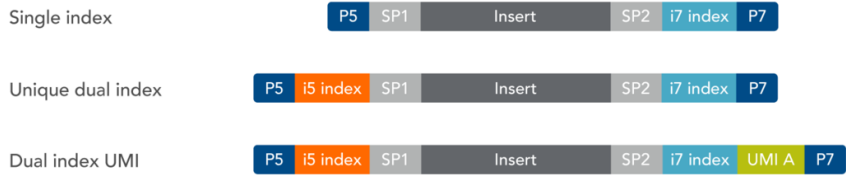
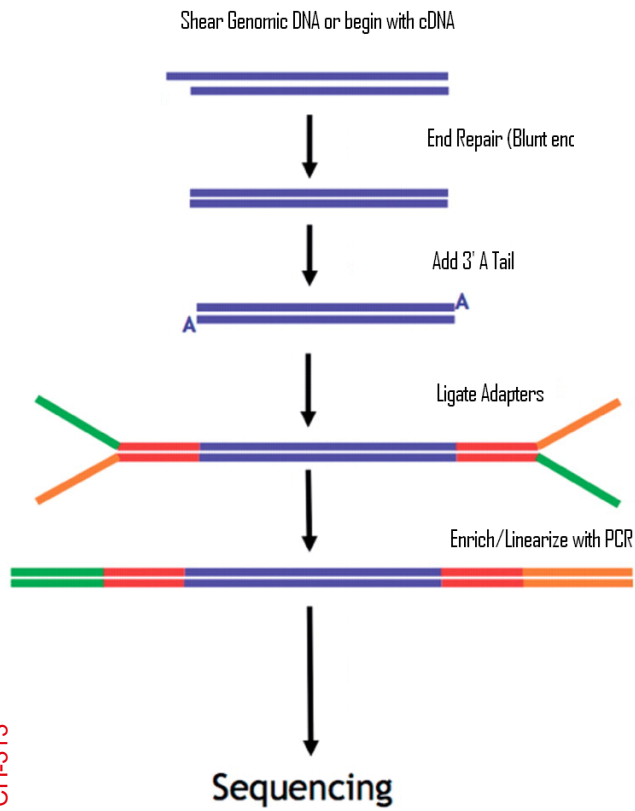


Sequencing by Synthesis (Illumina)

- “Extension of Sanger principle”
- Most widely used NGS technology today
- Sequencing on chip surface (“flow cells”)
- **Fluorescently labelled, “blocked” dNTPs**
- **Reversible termination:** blocked dNTPs converted into dNTPs to continue elongation after each nucleotide/“cycle”

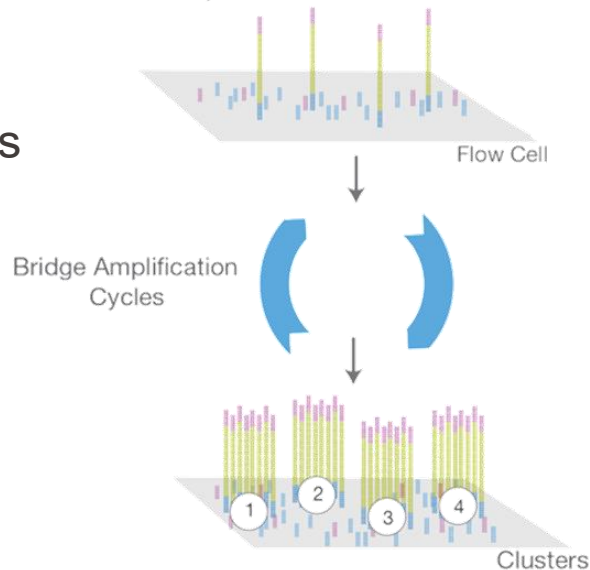
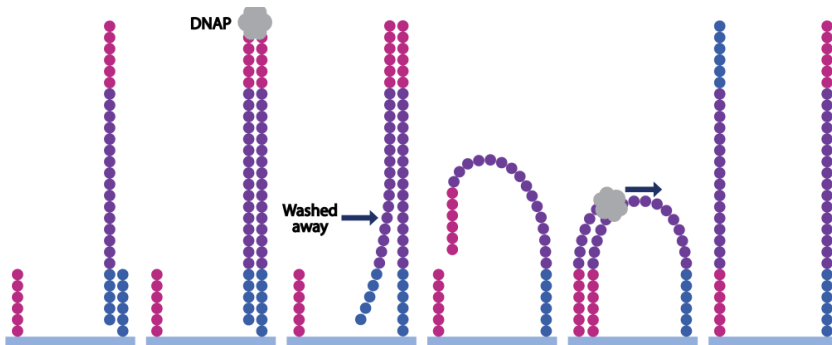


Illumina Seq. – Sample preparation (1)



- Flow cell binding sequence:** Platform-specific sequences for library binding to instrument
- Sequencing primer sites:** Binding sites for general sequencing primers
- Sample indexes:** Short sequences specific to a given sample library
- Molecular index/barcode:** Short sequence used to uniquely tag each molecule in a given sample library
- Insert:** Target DNA or RNA fragment from a given sample library

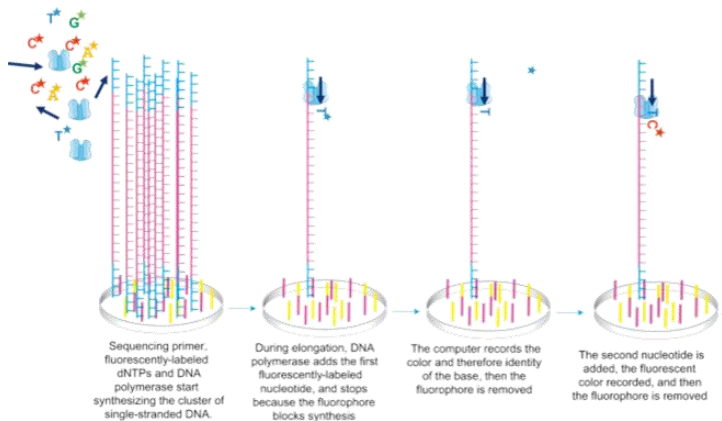
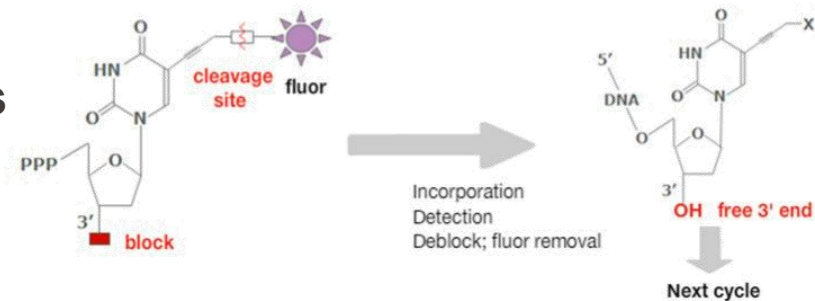
- Isothermal DNA amplification on flow cell surface
- Via “bridges” using two types of immobilized oligos
- Resulting clusters:
 - **monoclonal** (~1000 copies per cluster)
 - several hundred thousand per mm²
 - defined by 2D-coordinates on flow cell
 - reverse strands washed off before sequencing



Library is loaded into a flow cell and the fragments are hybridized to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

Illumina Seq. – Sequencing (3)

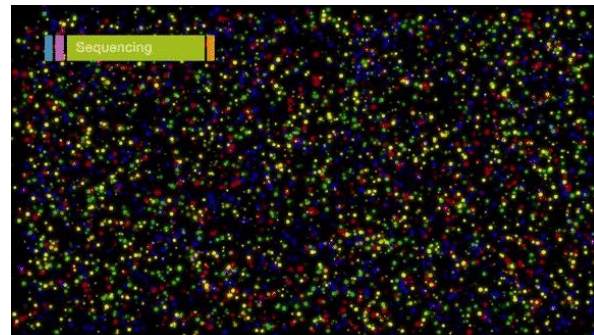
- DNA polymerase extends sequencing primer
- **Fluorescent, reversible termination dNTPs**
- Cycle (once per nt/position)
 - (1) elongation by one nucleotide
 - (2) imaging
 - (3) de-blocking+removal of fluorophore



Cycle 1

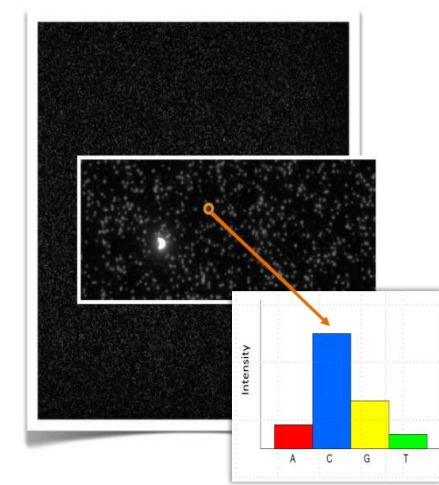
Cycle 2

Cycle 3



- Images analysed in real time (each cycle+channel)
 - “base calling” and quality scores obtained for each cluster
 - images deleted
 - FASTQ file produced
- FASTQ files
 - line 1: unique read ID (instrument/flow cell ID, lane, cluster coordinates, index, PE: forward or reverse read etc.)
 - line 2: sequence**
 - line 4: quality score (“Phred”, ASCII code)

```
@HISEQ:126:H14YJADXX:1:1101:1118:2101 1:N:0:ATCACG
CTCCATAGTCAGAAACTTCAGCATGACAGTACCTCATGCTGCATCAGGTGATCATGAAAAGATTAC
+
@@?ADDDD?ADHDIIIIIIIEIIIGEFHC<?FH4C9E9BGAFIGH<DG9BD?@DGEGHHG<DCBB
```



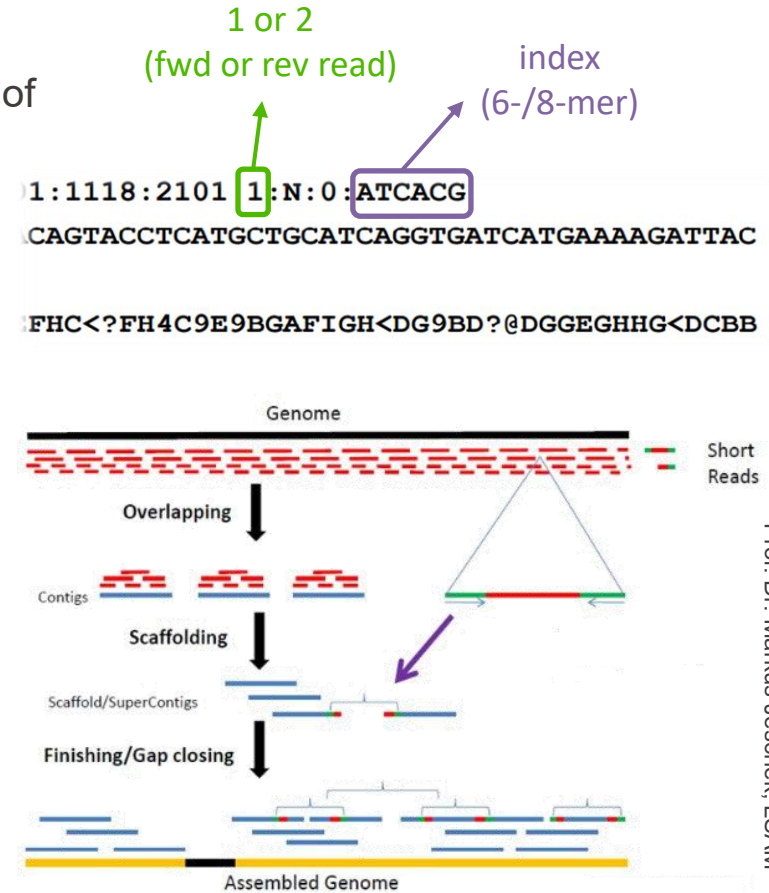
Phred quality scores are logarithmically linked to error probabilities

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%

$$\text{Phred} = -10 \log_{10} p$$

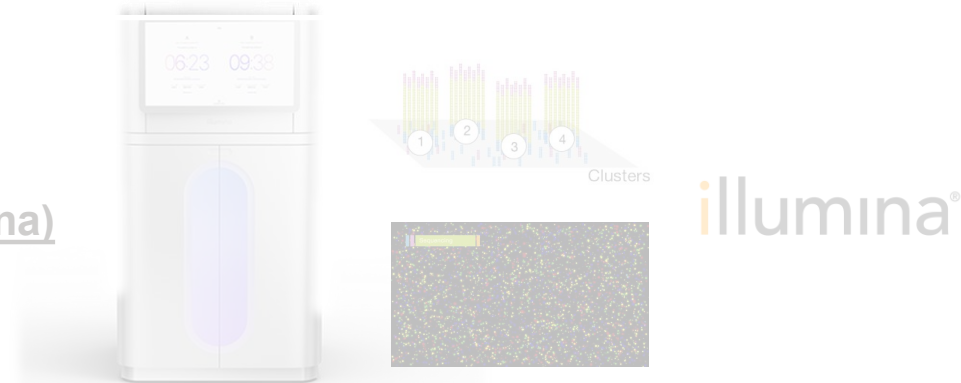
p = Probability call is incorrect

- From FASTQ files
 - optional: de-multiplexing of samples (deconvolution of index combinations)
 - optional: pairing of reads (paired-end reading)
- E.g. for whole-genome assembly...
 - short reads “stitched” together in silico via overlaps
 - contigs → scaffolds → entire assembly
 - “oversampling” for statistical coverage (30-200x)
 - mapping to reference genomes
- Different procedures for each application (amplicon sequencing etc.)



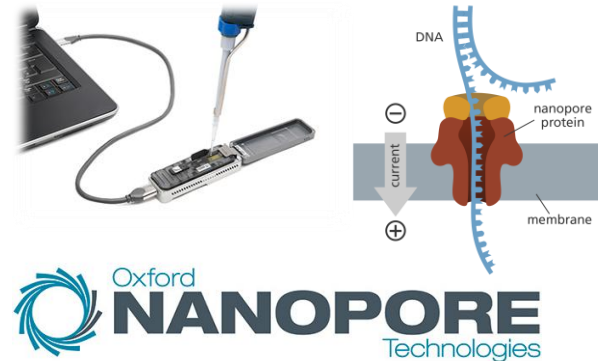
- “Second generation“ methods

- (1) Pyrosequencing (Roche 454)
- (2) Sequencing by synthesis (Illumina)
- (3) Sequencing by ligation
- (4) Ion semiconductor sequencing



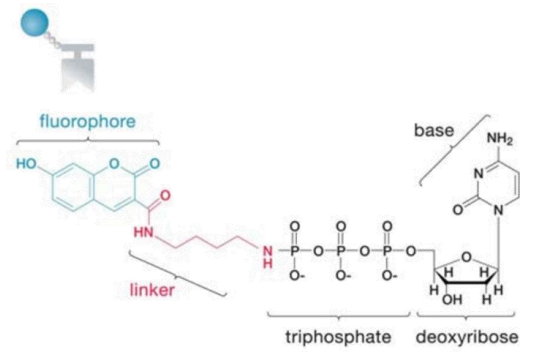
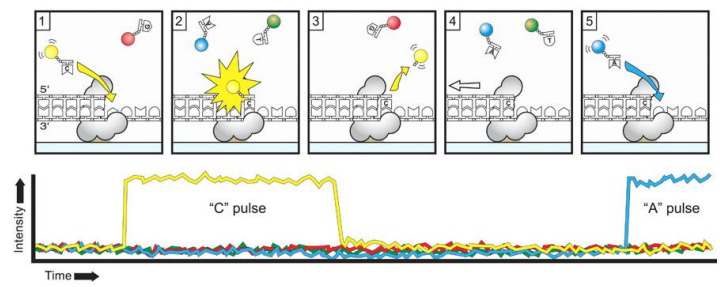
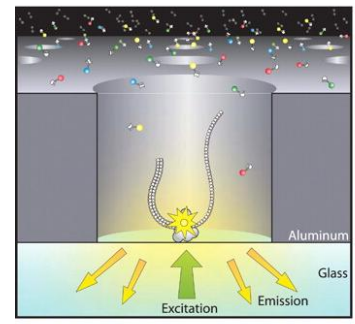
- “Third generation“ methods (single molecule)

- (1) **Single-molecule real-time sequencing (SMRT; Pacific Biosciences)**
- (2) Nanopore sequencing (MinION etc.; Oxford Nanopore Technologies)



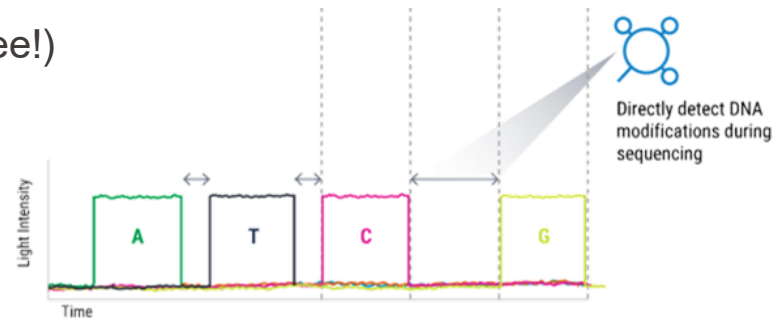
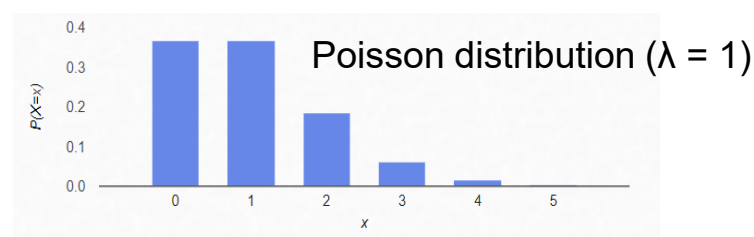
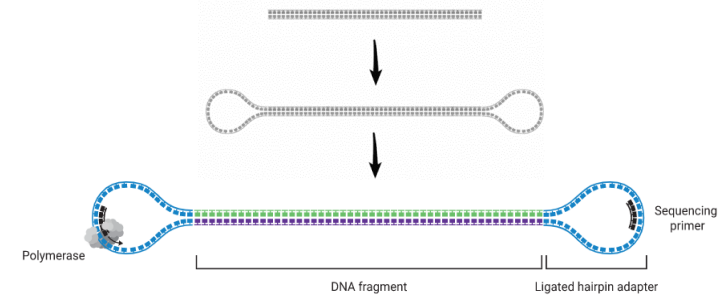
EPFL Single-molecule Real-time (SMRT) Sequencing 41

- Pacific Biosciences (PacBio)
- **Zero-mode waveguides (ZMWs)** contain...
 - single DNA polymerase molecule, immobilized
 - single molecule of DNA template
 - “phospholinked” nucleotides
 - reaction vessel: (70 nm diameter, 100 nm depth)
- Principle: labelled nucleotides incorporated by DNA pol. → labels are held a little longer in the ZMW than average diffusion → detected as a “flash”



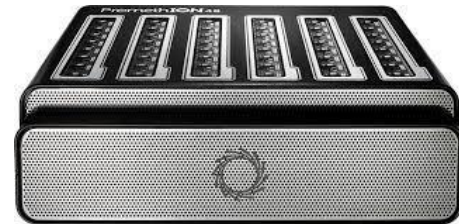
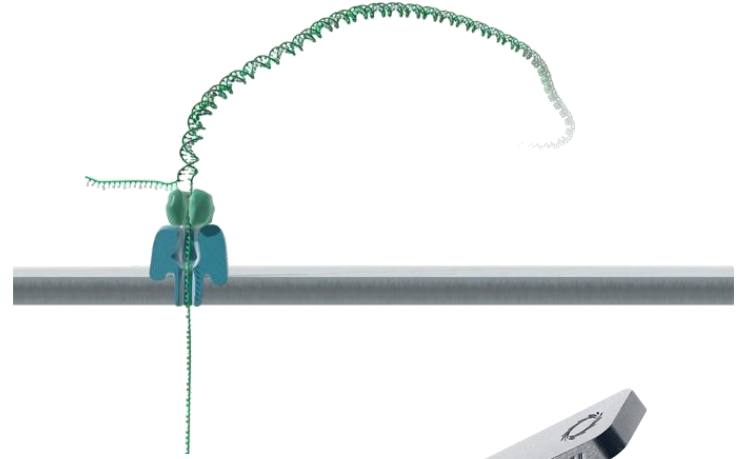
EPFL Single-molecule Real-time (SMRT) Sequencing 42

- Brief procedure:
 - linear DNA fragments ligated to bell-shaped adapters
 - DNA loaded on SMRT cell (Poisson distribution!)
 - Reading is performed in a circle (multiple times)
- Long-read technology (30 kb average, good for repetitive sequences!)
- Up to 4,000,000 reads
- High error rates (5-15% per base)
- Fast (\ll 24 h)
- Some modified bases can be directly detected (label-free!)
 - mainly methylated bases (e.g. N⁶-methyladenine, N⁴-methylcytosine)
 - epigenetics!



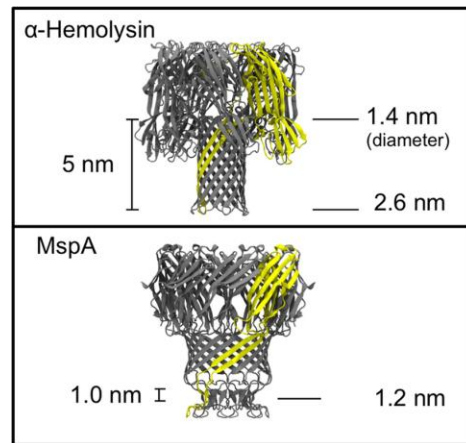
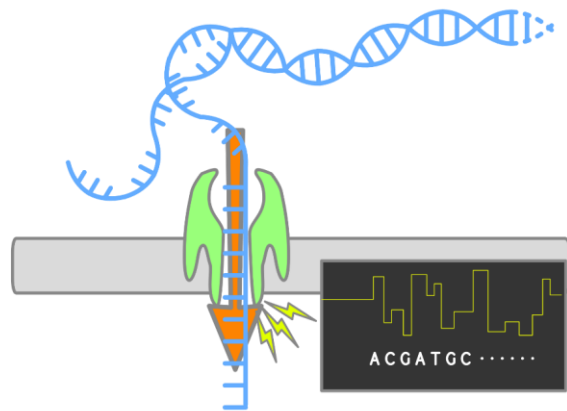
EPFL Nanopore Sequencing

- Oxford Nanopore Technologies (MinION, PromethION etc.)
- Label-free, single-molecule technique
- Portable equipment, real-time results
- Long-read technology (> 4Mbp successfully demonstrated)
- 512 – 2,675 pores per flow cell, repeated passage possible
- Error rates of 3-8% (lately improving drastically)
- Modified bases and other molecules (RNA, proteins) can be directly “sequenced“



EPFL Nanopore Sequencing - Principle

- single-stranded DNA/RNA molecules are “pushed” through nanopore via processive enzyme (e.g. DNA helicase)
- pore embedded in membrane and surrounded by electrolyte
- electric field across the membrane → electrophoretic motion of ions through pore
- if a larger molecule (e.g. DNA strand) occupies pore, ion flux is disrupted (detectable by voltage change in real time)
- voltage changes are specific for base/molecules
- pores: biological (α -hemolysin, MspA) or solid-state (metal, metal alloy)
- mostly synthetic membranes

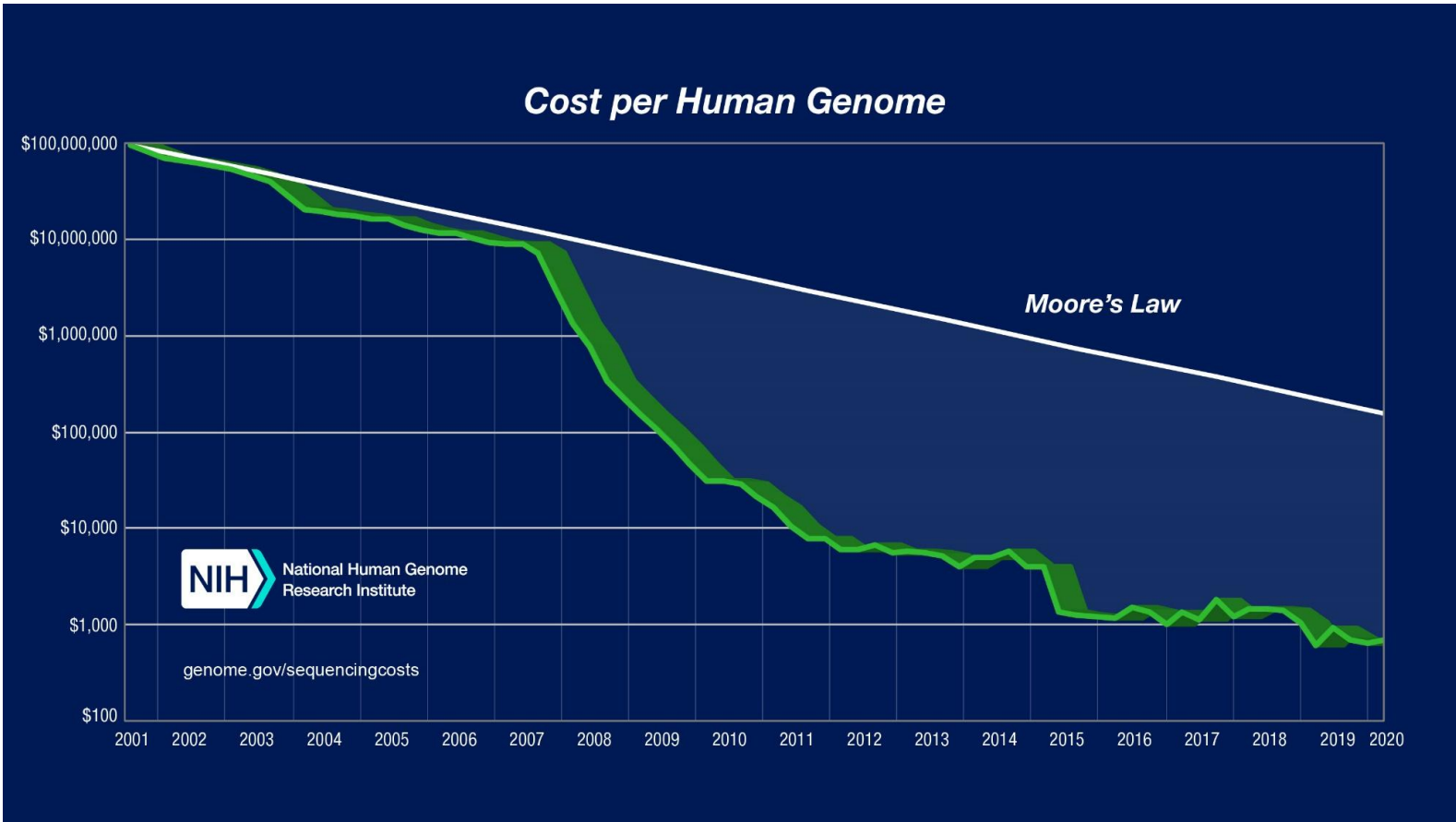


<https://www.youtube.com/watch?v=E9-Rm5AoZGw>

EPFL NGS – Comparison of Main Methods (approximate numbers)

Method	max. read length [bp]	error rate [%]	max. reads per run	time per run [h]	single molecule?	costs per Gb [US-\$]	Remarks
Pyrosequencing (Roche)	700	0.1-1	1M	24	no	10k	discontinued, expensive, homopolymer errors
Sequencing by synthesis (Illumina)	50-600	0.1-1	52B	4-48	no	2-150	expensive equipment, cheap Gb price, low error rates
Ion semiconductor (Ion Torrent)	600	~0.5	80M	2	no	50-1000	cheap equipment, very fast, homopolymer errors
SMRT sequencing (PacBio)	30k-100k	5-15	4M	0.5-20	yes	5-50	expensive equipment, long reads, fast, methylation
Nanopore Sequencing (Oxford)	> 4000k	3-8	dep. on length (~few 100k)	72	yes	5-100	handheld, cheap equipment, longest reads, other molecules
Sanger	1200	0.01	1	0.2-3	no	2-3M	gold standard, low throughput

EPFL DNA Sequencing – Costs



Questions?

Thank you!