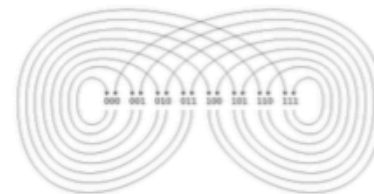
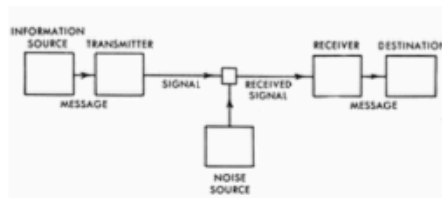


Part 1: Basics of Genomics

Olgica Milenkovic
University of Illinois, Urbana-Champaign

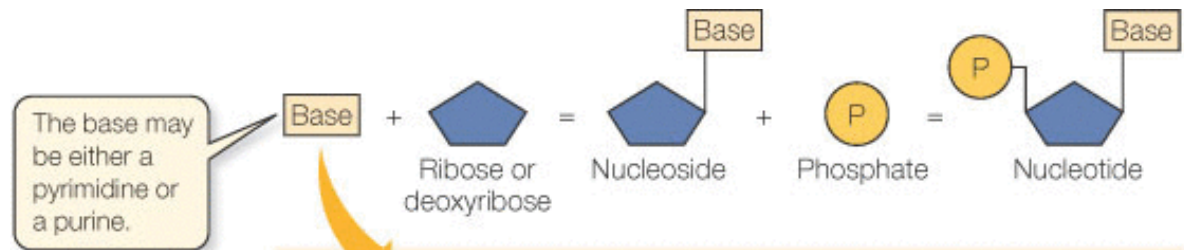
North American School of Information Theory, Texas, 2018

May 2018

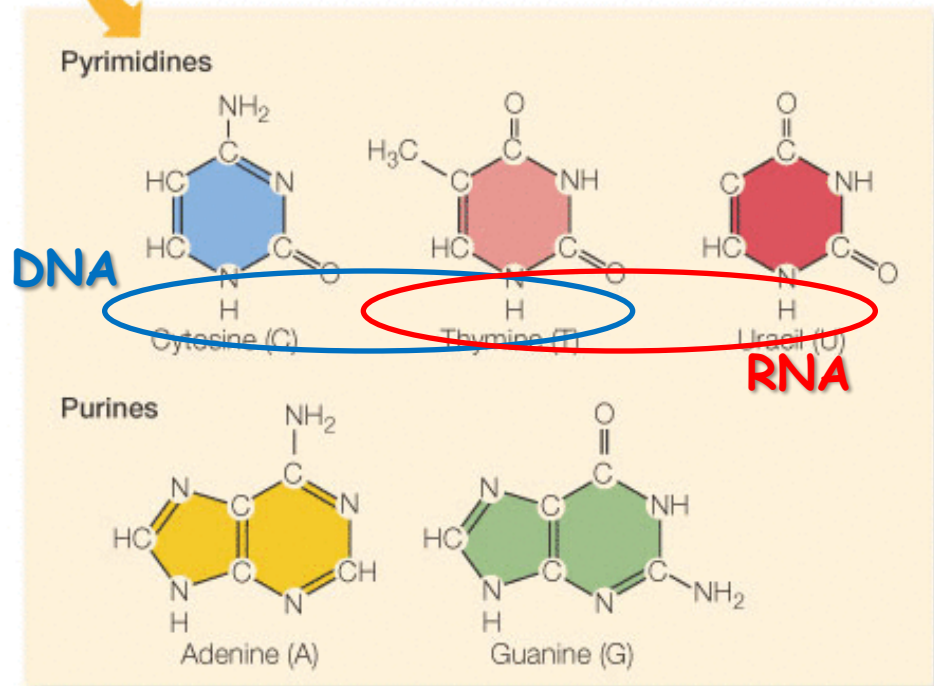


NUCLEIC ACIDS

- **DNA** (Deoxyribonucleic Acid) and **RNA** (Ribonucleic Acid): information storage molecules made up of “nucleotides”.



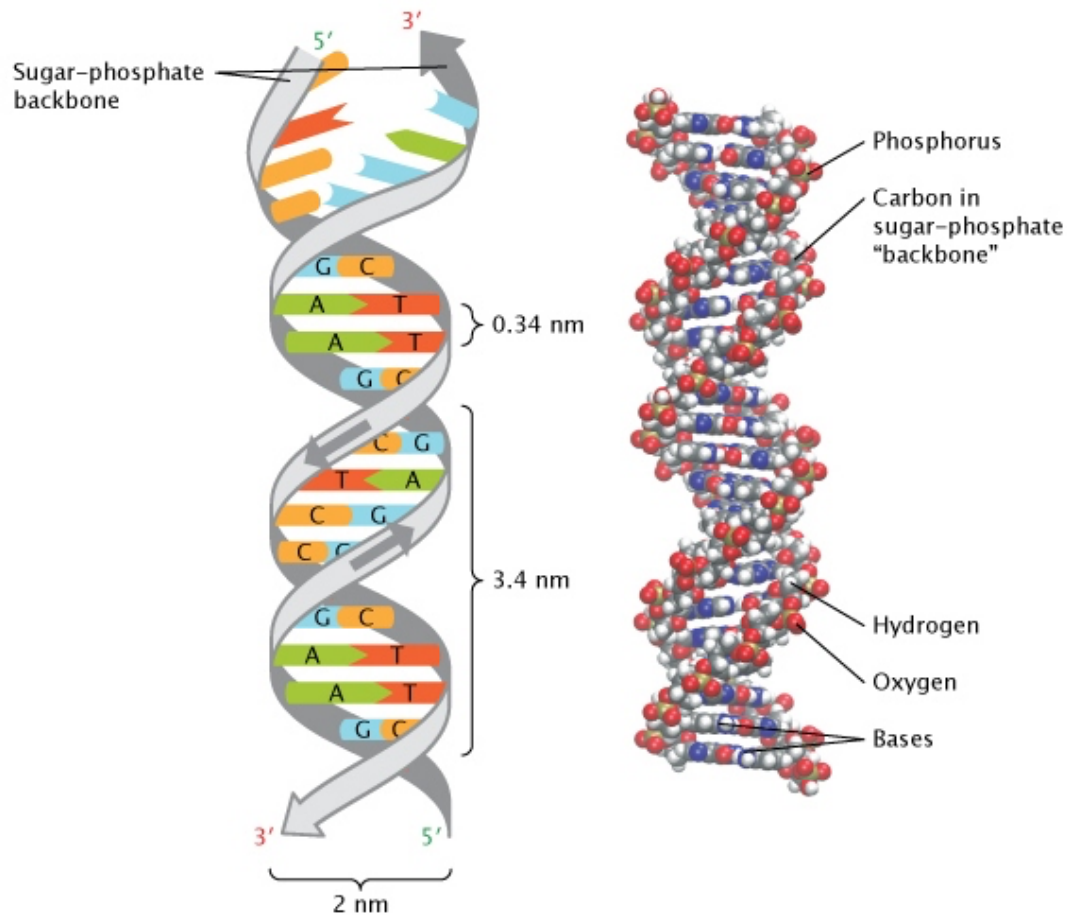
Base	Abbreviation
Adenine	A
Guanine	G
Cytosine	C
Uracil	U
Thymine	T



Source: Nature

NUCLEIC ACIDS

- **DNA** (Deoxyribonucleic Acid) and **RNA** (Ribonucleic Acid): information storage molecules made up of “nucleotides”.



Source: Nature

DISCOVERING THE STRUCTURE OF DNA

○ Chargaff's Rules (1949)

- Amount of each dNTP varies between organisms, but $[dA]=[dT]$ and $[dC]=[dG]$ in ALL organisms

Table 3-2 Data Leading to the Formulation of Chargaff's Rules

Source	Adenine to Guanine	Thymine to Cytosine	Adenine to Thymine	Guanine to Cytosine	Purines to Pyrimidines
Ox	1.29	1.43	1.04	1.00	1.1
Human	1.56	1.75	1.00	1.00	1.0
Hen	1.45	1.29	1.06	0.91	0.99
Salmon	1.43	1.43	1.02	1.02	1.02
Wheat	1.22	1.18	1.00	0.97	0.99
Yeast	1.67	1.92	1.03	1.20	1.0
<i>Hemophilus influenzae</i>	1.74	1.54	1.07	0.91	1.0
<i>E-coli</i> K2	1.05	0.95	1.09	0.99	1.0
Avian tubercle bacillus	0.4	0.4	1.09	1.08	1.1
<i>Serratia marcescens</i>	0.7	0.7	0.95	0.86	0.9
<i>Bacillus schatz</i>	0.7	0.6	1.12	0.89	1.0

SOURCE: After E. Chargaff et al., *J. Biol. Chem.* 177 (1949).

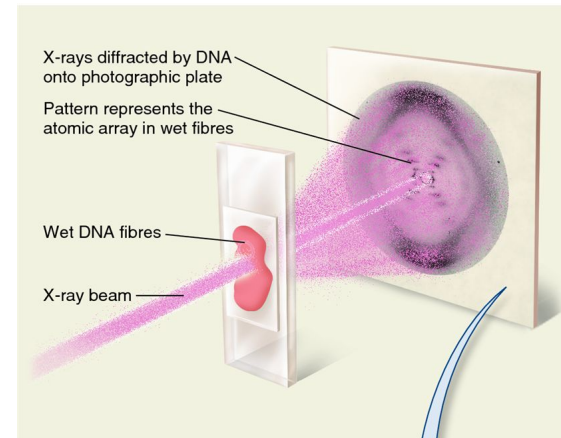
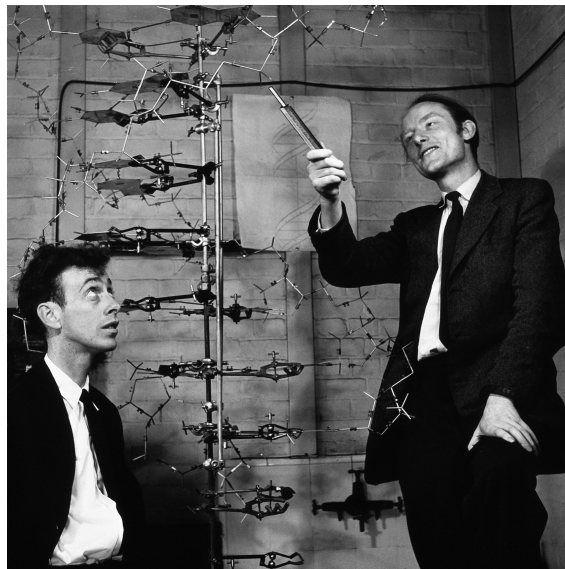
DISCOVERING THE STRUCTURE OF DNA

○ Rosalind Franklin & Maurice Wilkins

- X-ray diffraction suggested helix of uniform width with stacked bases, with sugar-phosphate on outside.

○ James Watson & Francis Crick

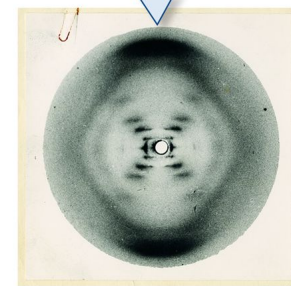
- Postulated double-helix model.



(a) The method of X-ray diffraction

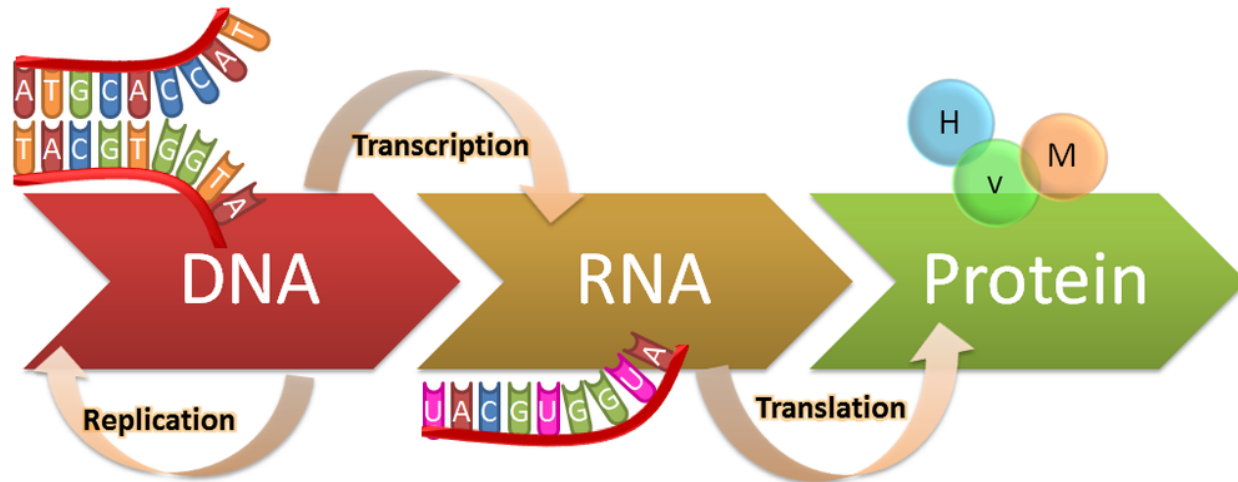


(b) Rosalind Franklin



(c) Franklin's X-ray diffraction pattern of wet DNA fibres

THE CENTRAL DOGMA

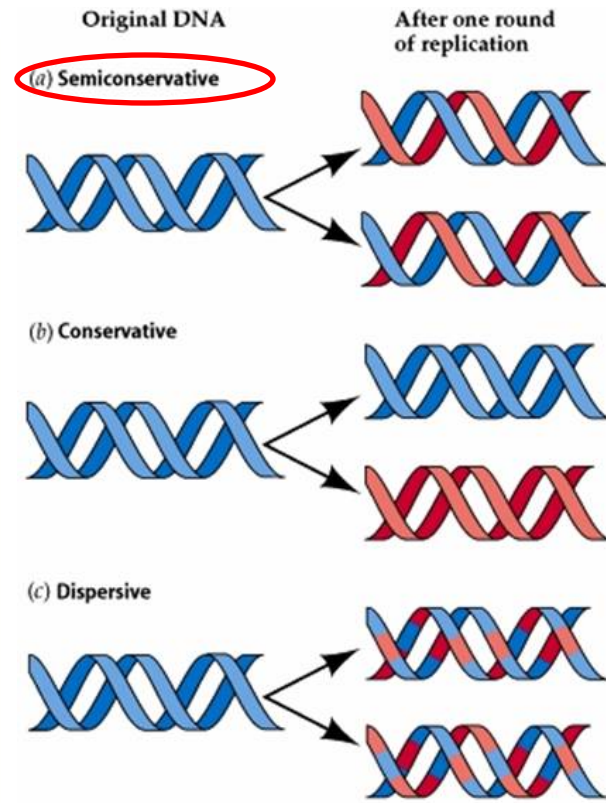
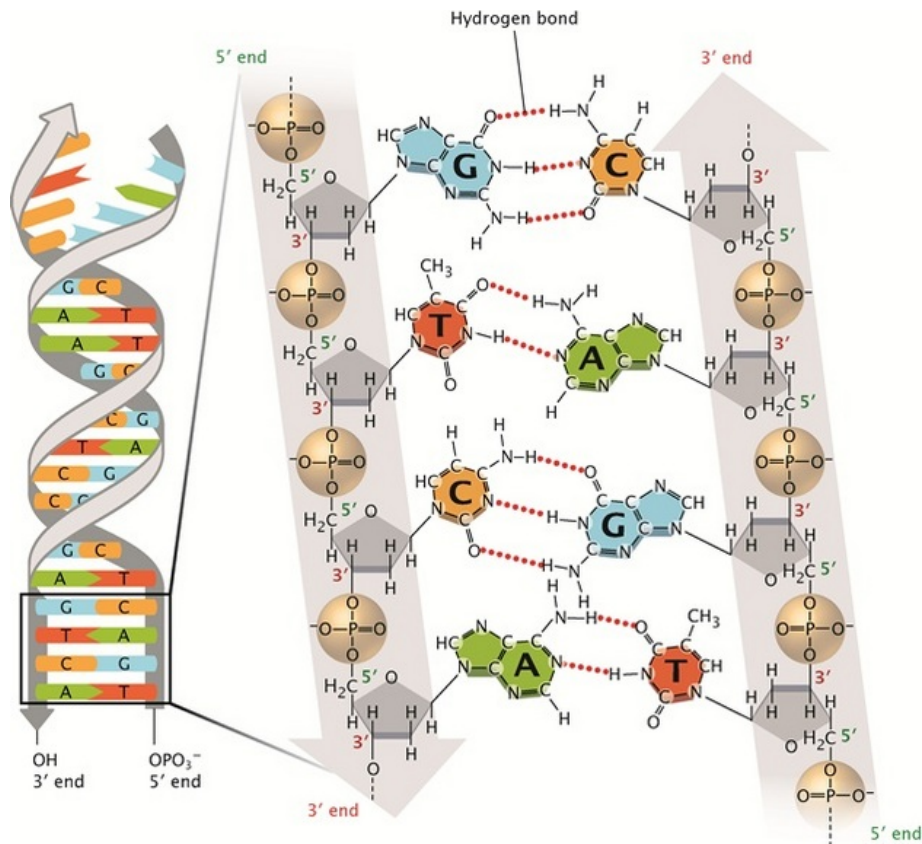


Source: Genius Media Group

- DNA encodes genetic information that “directs” the cell how to make proteins and RNAs.
- Information carried in the nucleotide sequence is copied into an RNA (TRANSCRIPTION).
- Information in RNA is used to build proteins (TRANSLATION).

DNA REPLICATION

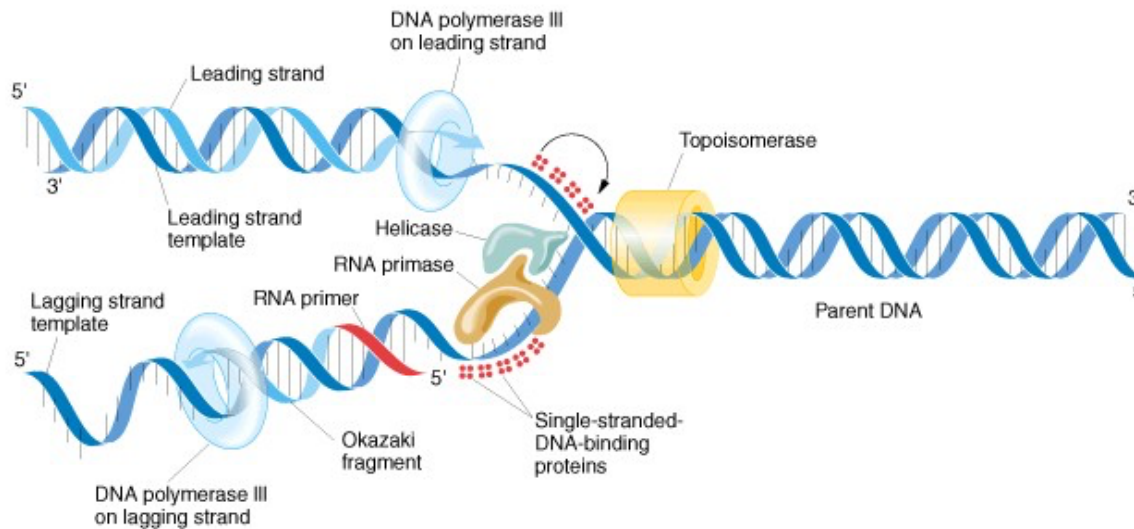
- Complementary base pairing (A-T, C-G).
- Semiconservative model of DNA replication.



Source: Nature

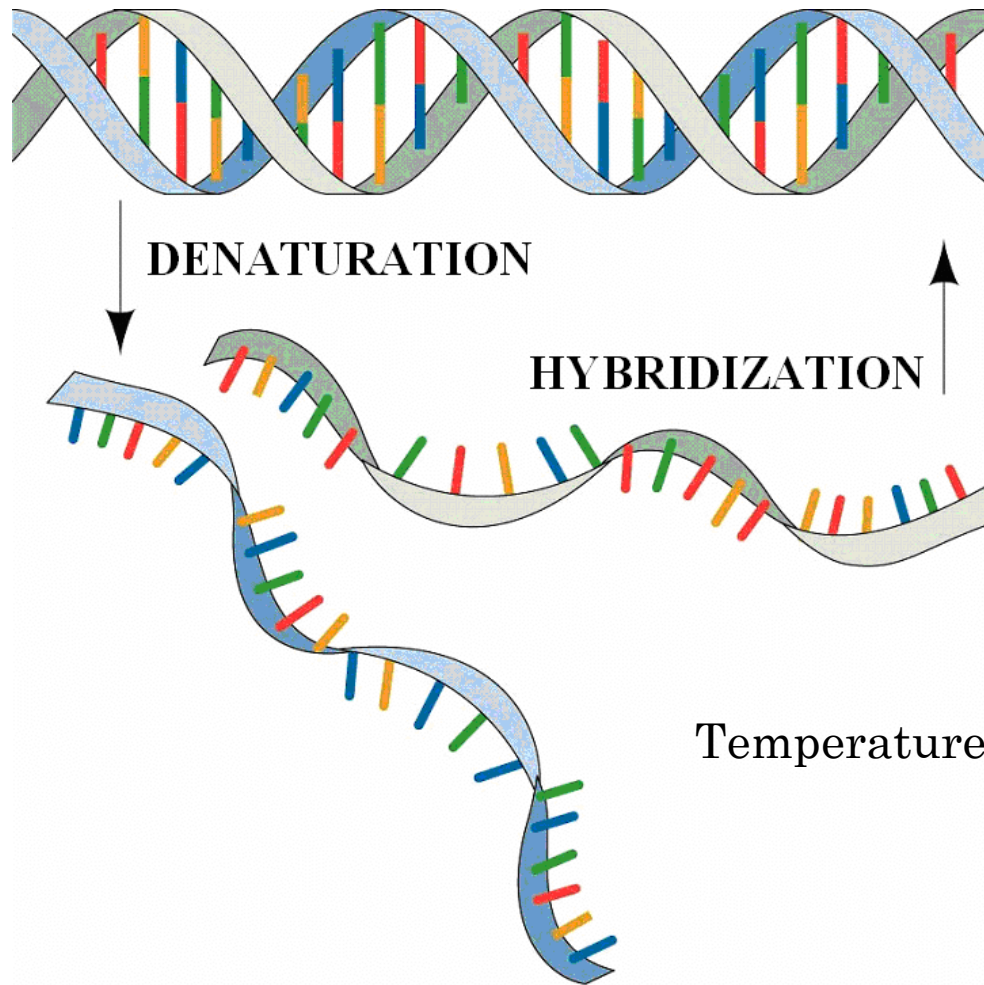
DNA REPLICATION REQUIRES PRIMERS!

- “A **primer** is a short strand of RNA or **DNA** (generally about 18-22 bases) that serves as a starting point for **DNA** synthesis. It is required for **DNA replication** because the enzymes that catalyze this process, **DNA polymerases**, can only add new nucleotides to an existing strand of **DNA**.”
- Source: [Wikipedia](#) 😊.



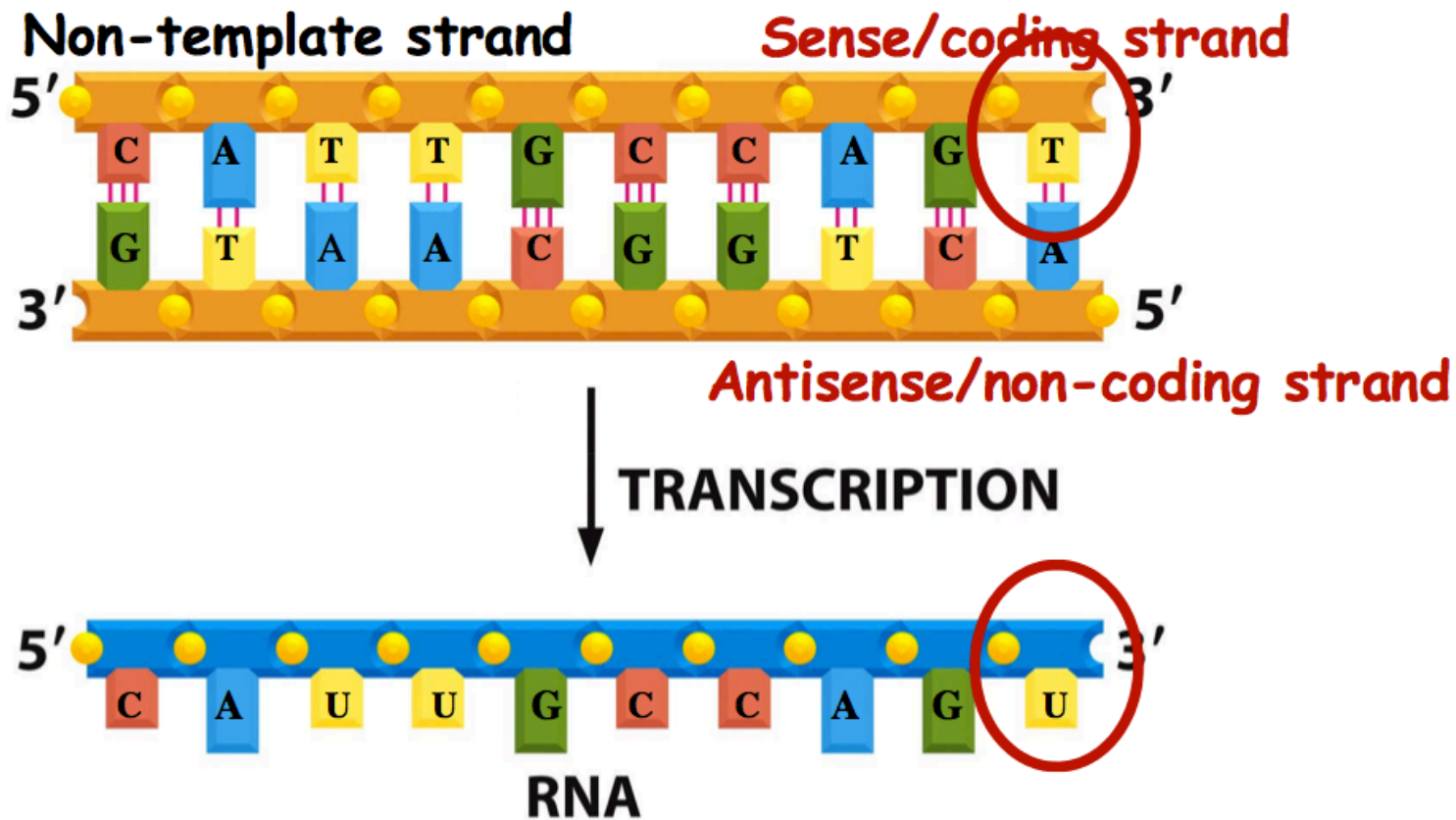
Source: Nature

DNA HYBRIDIZES AND DENATURES



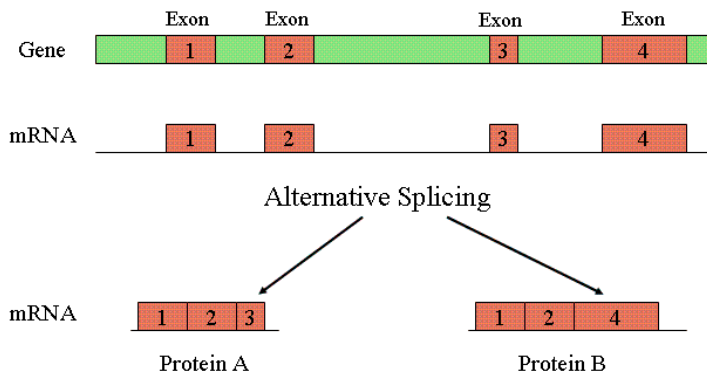
TRANSCRIPTION (DNA \rightarrow RNA)

- RNA (usually) single-stranded.
- T \rightarrow U.



TRANSLATION (RNA → PROTEIN)

- Alternative Splicing in eukaryotes.



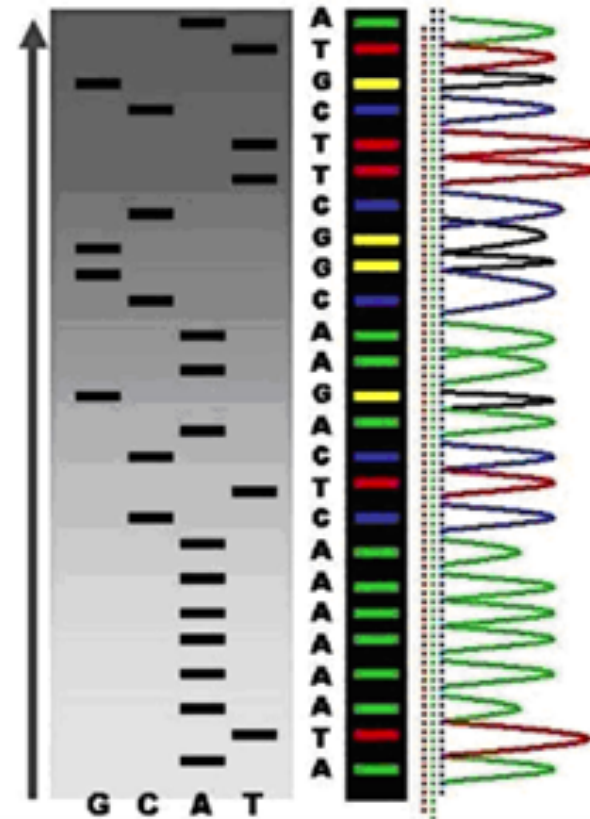
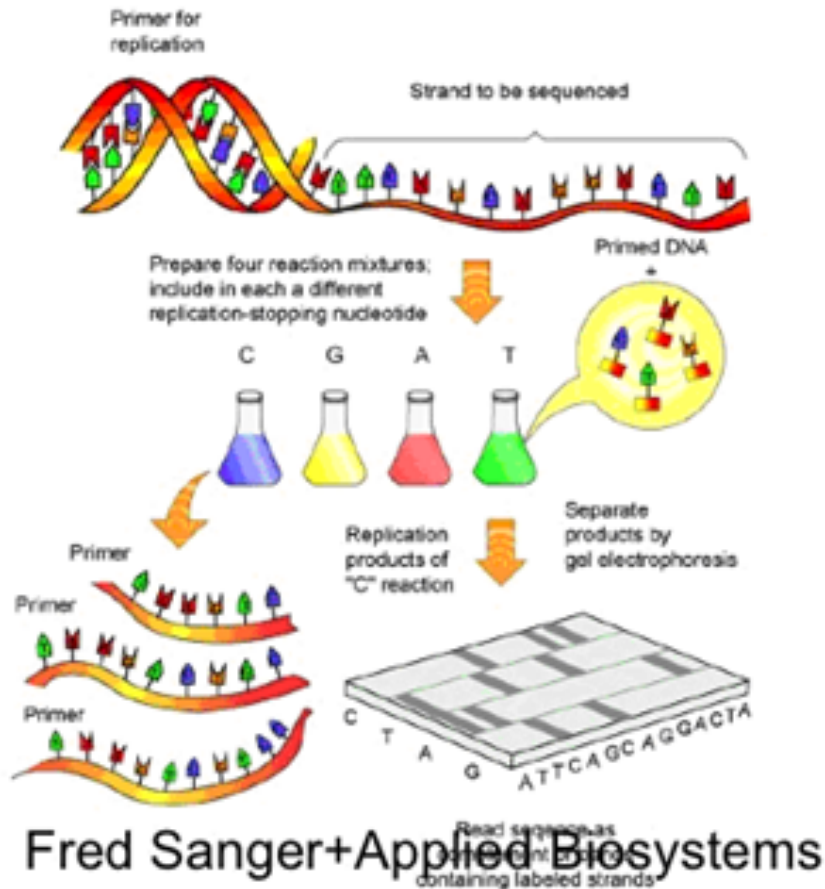
Source: NCBI

- Codon table of Amino Acids (protein)
 - Degenerate (redundant).

		2nd base in codon				
		U	C	A	G	
1st base in codon	U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys STOP Trp	U C A G
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G
	A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
	G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G
						3rd base in codon























READING DNA: SANGER SEQUENCING

Sequencing by capillary electrophoresis



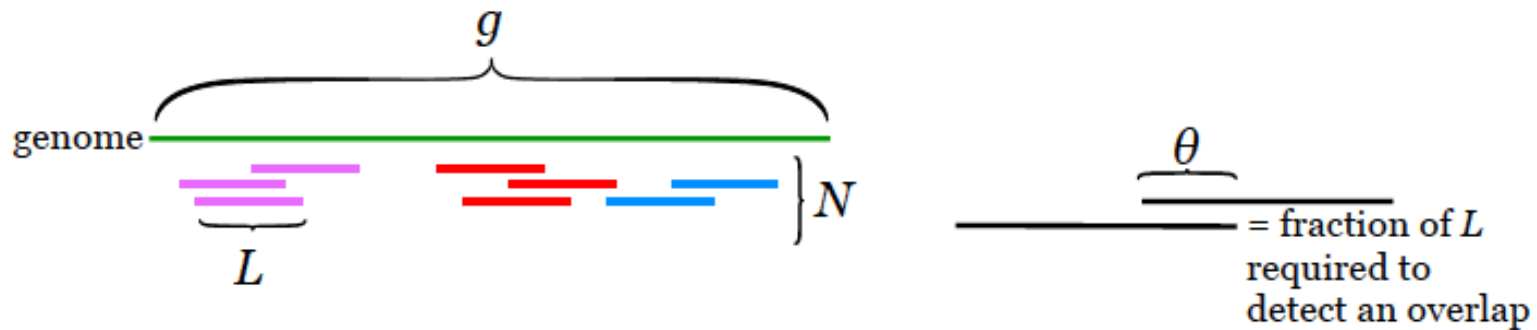
READING DNA: SANGER SEQUENCING

Gel:

	G	GCGAATGCGTCCACACGCTACAGGT G
	T	GCGAATGCGTCCACACGCTACAGGT
	G	GCGAATGCGTCCACACGCTACAG G
	G	GCGAATGCGTCCACACGCTACAG
	A	GCGAATGCGTCCACACGCTAC A
	C	GCGAATGCGTCCACACGCTAC
	A	GCGAATGCGTCCACACGCT A
	T	GCGAATGCGTCCACACGCT
	C	GCGAATGCGTCCACACG C
	G	GCGAATGCGTCCACACG
	C	GCGAATGCGTCCACAC
	A	GCGAATGCGTCCAC A
	A	GCGAATGCGTCCAC A
	C	GCGAATGCGTCCAC
	A	GCGAATGCGTCC A
	C	GCGAATGCGTCC
	C	GCGAATGCGT C
	T	GCGAATGCG T
	G	GCGAATGCG
	C	GCGAATG C
	G	GCGAAT G
	T	GCGAAT

READING DNA: COVERAGE AND READ LENGTH

How many reads do we need to be sure we cover the whole genome?



An **island** is a contiguous group of reads that are connected by overlaps of length $\geq \theta L$.

(Various colors above)

Want: Expression for expected # of islands given N, g, L, θ .

From C. Kingsford lecture notes

READING DNA: COVERAGE AND READ LENGTH

$\lambda := N/g =$ probability a read starts at a given position
(assuming random sampling)

Pr(k reads start in an interval of length x)

x trials, want k “successes,” small probability λ of success

Expected # of successes = λx

Poisson approximation to binomial distribution:

$$\text{Pr}(k \text{ reads in length } x) = e^{-\lambda x} \frac{(\lambda x)^k}{k!}$$

Expected # of islands = $N \times \text{Pr}(\text{read is at rightmost end of island})$

$$\begin{aligned} \underline{\quad (1-\theta)L \quad \theta L \quad} &= N \times \text{Pr}(0 \text{ reads start in } (1-\theta)L) \\ &= N e^{-\lambda(1-\theta)L} \frac{\lambda^0}{0!} \quad (\text{from above}) \\ &= N e^{-\lambda(1-\theta)L} \\ &= N e^{-(1-\theta)LN/g} \quad \leftarrow LN/g \text{ is called the } \mathbf{coverage} \mathbf{c}. \end{aligned}$$

READING DNA: COVERAGE AND READ LENGTH

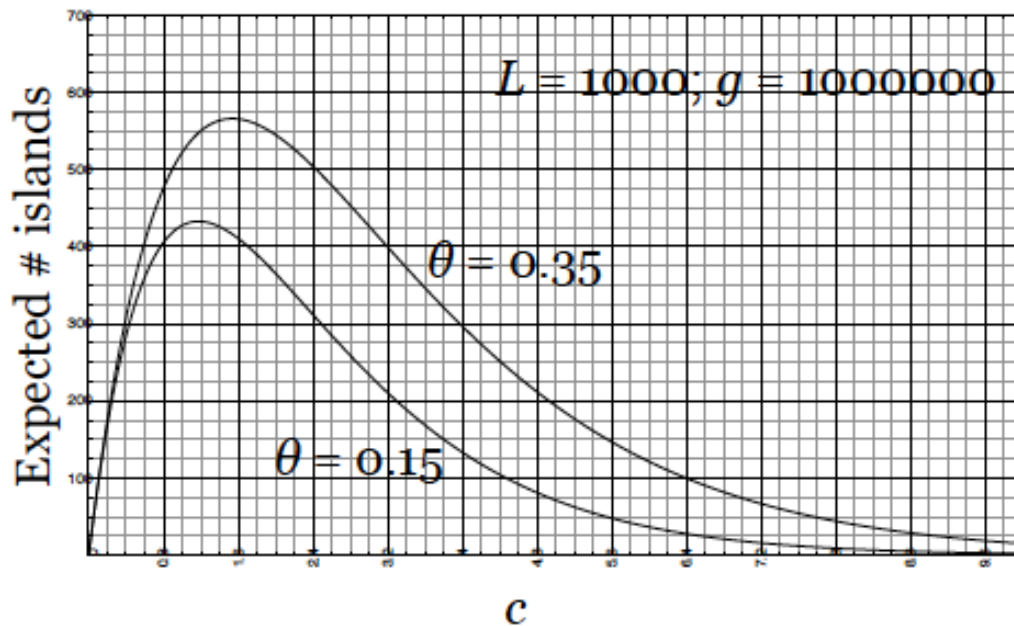
Rewrite to depend more directly on the things we can control: c and θ

$$\text{Expected \# of islands} = Ne^{-(1-\theta)LN/g}$$

$$= Ne^{-(1-\theta)c}$$

$$= \frac{L/g}{L/g} Ne^{-(1-\theta)c}$$

$$= \frac{g}{L} ce^{-(1-\theta)c}$$



From C. Kingsford lecture notes

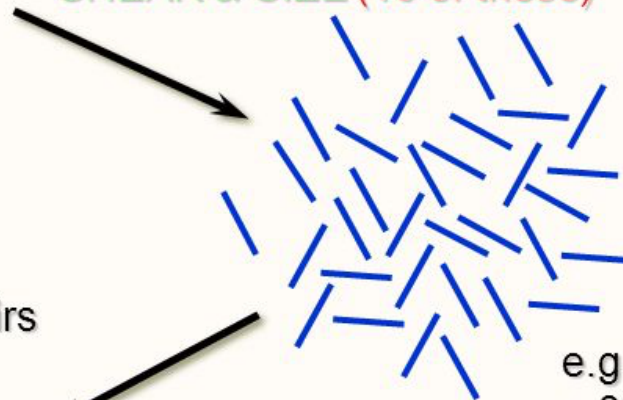
READING DNA: LESSONS

Mate-Pair Shotgun DNA Sequencing

DNA target sample

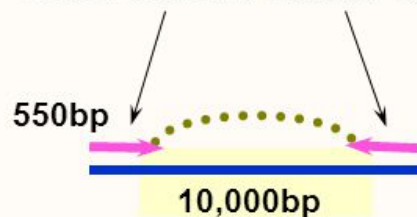


SHEAR & SIZE (16 of these)



e.g., 10Kbp
 $\pm 8\%$ std.dev.

End Reads / Mate Pairs



CLONE (16 of these)
& END SEQUENCE (automated)

DNA SEQUENCING: ILLUMINA PLATFORMS



50 – 600Gb
2 – 11 days
2 x 100bp max

10 – 180Gb
7 – 40 hours
2 x 150bp max

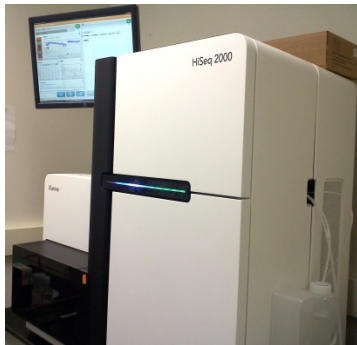


Larger projects, fewer runs

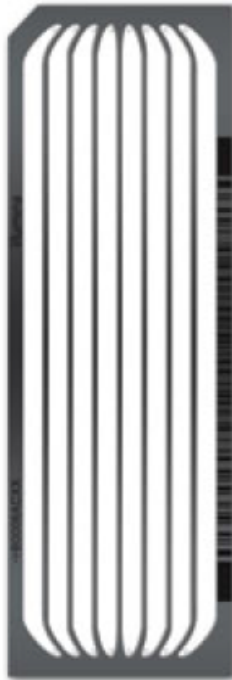
Smaller projects, quick results

Courtesy of Alvaro Hernandez, UIUC

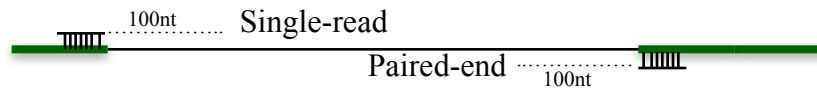
ILLUMINA HISeq 2000



Flowcell
(8 lanes)



ILLUMINA LIBRARY FRAGMENT:



* Reads (sequences) = 100nt in length

* **Number of reads per lane:**

* 150 to 200 million single-reads (15 to 20 GB)

* 300 to 400 million paired-reads (30 to 40 GB)

- Run length: 7 days for single-reads,
13 days for paired-reads

Courtesy of Alvaro Hernandez, UIUC

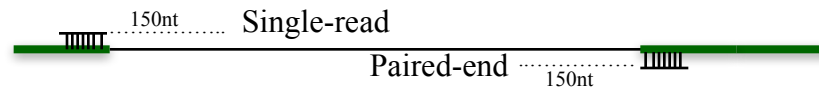
ILLUMINA HiSEQ 2500



Flowcell
(2 lanes)



ILLUMINA LIBRARY FRAGMENT:



*Reads (sequences) = 150-300nt in length

Number of reads per lane:

* 100 to 150 million single-reads (15 to 22 GB)

* 200 to 300 million paired-reads (30 to 45 GB)

- Run length: <1 day for single-reads,
2 days for paired-reads

DNA SEQUENCING: ILLUMINA NOVASEQ



NovaSeq 5000 and 6000
cost **\$850,000** and **\$985,000 (2017)**

In "Rapid Run Mode," the **Illumina HiSeq 2500** instrument is capable of generating approximately 150 millions **reads** passing filter **per** lane, or up to 300 million **reads** passing filter **per** lane for paired-end sequencing.

NovaSeq: Output size up to 6Tb, up to 6 billion reads per run, length of reads 2x150

Courtesy of Alvaro Hernandez, UIUC

FILE FORMAT: FASTA

- Text-based format for nucleotide or peptide sequences
- Line 1: description
 - “>” symbol
 - Sequence identifier

Database	Format
GenBank	gb <i>accession</i> <i>locus</i>
EMBL Data Library	emb <i>accession</i> <i>locus</i>
DDBJ, DNA Database of Japan	dbj <i>accession</i> <i>locus</i>
NBRF PIR	pir <i>entry</i>
Protein Research Foundation	prf <i>name</i>
SWISS-PROT	sp <i>accession</i> <i>entry name</i>
Brookhaven Protein Data Bank	pdb <i>entry</i> <i>chain</i>
Patents	pat <i>country</i> <i>number</i>
GenInfo Backbone Id	bbs <i>number</i>
General database identifier	gnl <i>database</i> <i>identifier</i>
NCBI Reference Sequence	ref <i>accession</i> <i>locus</i>
Local Sequence Identifier	lcl <i>identifier</i>

- Description [optional]

FILE FORMAT: FASTA

- Line 2: sequence data (1 or more lines)
 - Protein or nucleic acids sequences
 - Amino acids: A-Z, * (translation stop), - (gap)
 - Nucleic acids: A, C, G, T, U, R, Y, K, M, S, W, B, D, H, V, N (any A C G T U), X (masked), - (gap)
 - Each line less than 80 characters
- File extensions

Extension	Meaning	Notes
fasta (.fas)	Generic fasta	Can be .fa, .seq, .fsa
fna	Fasta nucleic acid	Generic nucleic acids
ffn	FASTA nucleotide coding regions	Coding regions for a genome
faa	Fasta amino acid	.mpfa: multiple protein fasta
frn	FASTA non-coding RNA	Non-coding RNA regions

FILE FORMAT: FASTA

- Example: “random.fna”

```
>SEQUENCE_1
TGGCAATCTTGCTTCTGTTTACGGCTGGCATAGTTACGACA
GGTCTTTTTTCT
>SEQUENCE_2
CCGGTTTCTTCAACCTTAGTTCTGGTAGCAGAATCAAGATA
CATGTTTTTCGT
>SEQUENCE_3
GACGGCGTCAGCTGCAACAACACTGTGCGCGCCATTGCCCTG
CCGGGGCGATC
```

- Example: “NP_852610.1”

```
>gi|31563518|ref|NP_852610.1| microtubule-associated
proteins 1A/1B light chain 3A isoform b [Homo sapiens]
MKMRFFSSPCGKAAVDPADRCKEVQQIRDQHPSKIPVIIERY
KGEKQLPVLDKTKFLVPDHVNMSELVKI
```


FILE FORMAT: FASTQ

- Text-based format for sequences and quality scores
- Line 1
 - “@” symbol
 - Sequence identifier
 - Description [optional]
- Line 2
 - Raw sequence letters
- Line 3
 - “+” symbol
 - Sequence identifier
 - Description [optional]
- Line 4
 - Quality values for sequences in Line 2

FILE FORMAT: FASTQ

- Quality value Q is an integer-valued function of p , the probability that the corresponding base call is incorrect
- Phred quality score:
 - $Q_{sanger} = -10 \log_{10} p$
- Quality values in increasing order of quality (ASCII):

```
!"#$%&'()*+,-  
./0123456789:;<=>?@ABCDEFGHIJKLMNPOQRSTUVWXYZ[\]  
^_`abcdefghijklmnopqrstuvwxyz{|}~
```

FILE FORMAT: FASTQ

- Illumina sequence identifiers

- Example: @HWUSI-EAS100R:6:73:941:1973#0/1

HWUSI-EAS100R	Unique instrument name
6	Flowcell lane
73	Tile number within flowcell lane
941	'x'-coordinate of the cluster within tile
1973	'y'-coordinate of the cluster within tile
#0	Index number for multiplexed sample
/1	Member of a pair

- File extensions

- .fq, .fastq

FILE FORMAT: FASTQ

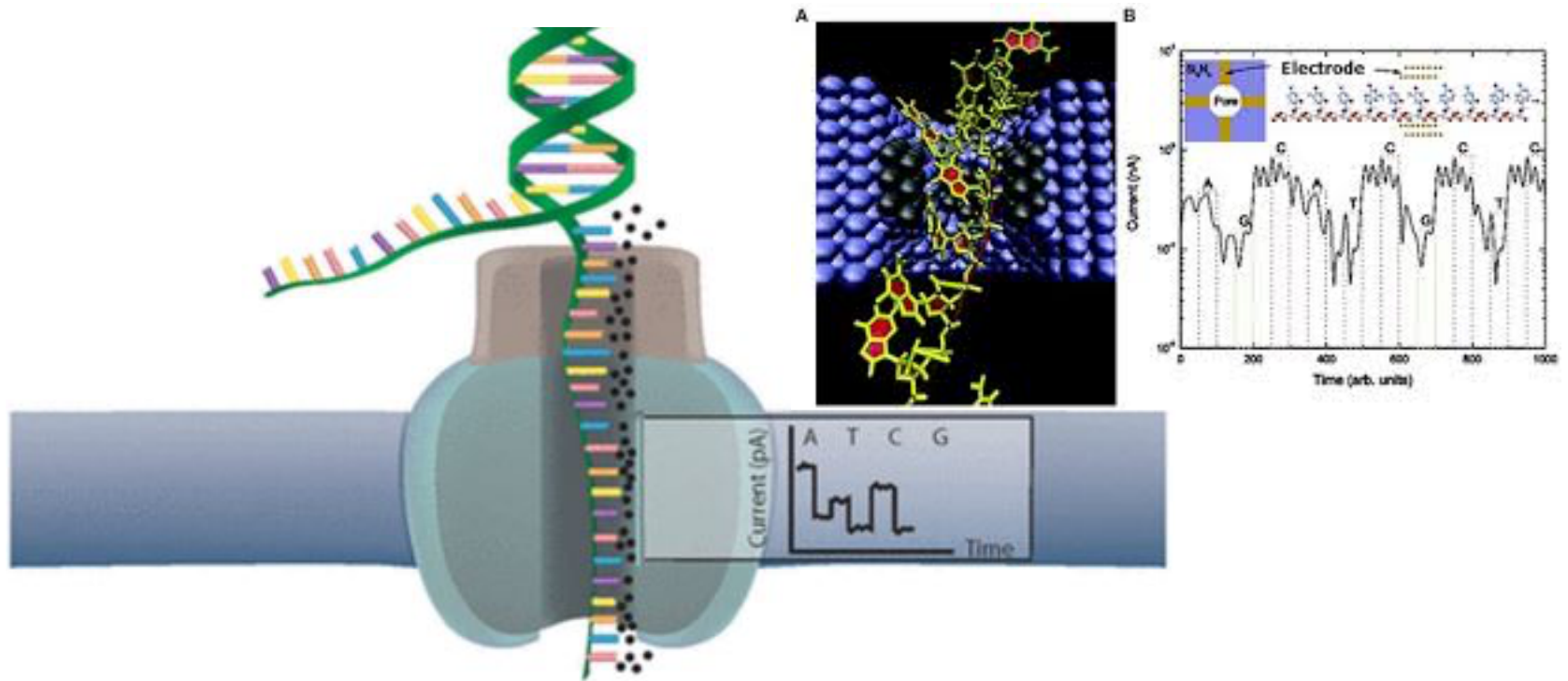
- Example (Illumina): “sample1.fq”

```
@HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNNTTACCTTNNNNNNNNNNN
TAGTTTCTTGAGATTTGTTGGGGGAGACATTTTGTGATTGCCTTGAT
+HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
efcffffcfeeffcffffdf`feed]`_Ba_^__[YBBBBBBBBBBRTT\]][]ddd`ddd^dddadd
^BBBBBBBBBBBBBBBBBBBBBBBBBBBB
```

- Example (NCBI read archive)

```
@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGGCCGCTGCCGATGGCGTCAAATCCCACC
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC
```

THIRD GENERATION SEQUENCERS



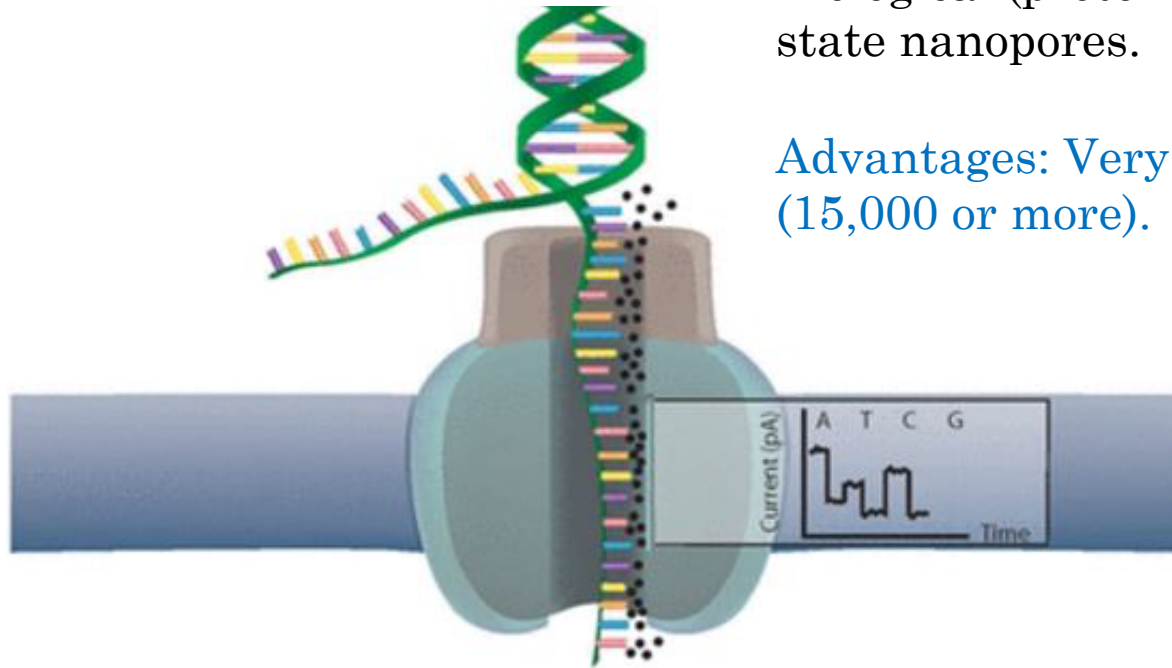
Credit: LabBiotech

Biological (protein) versus solid state nanopores.

THIRD GENERATION SEQUENCERS

Biological (protein) versus solid state nanopores.

Advantages: Very long reads (15,000 or more).



Credit: LabBiotech

Advantages: Very long reads (15,000 or more).

Disadvantages: Many substitution, insertion and deletion errors.

BASE CALLING

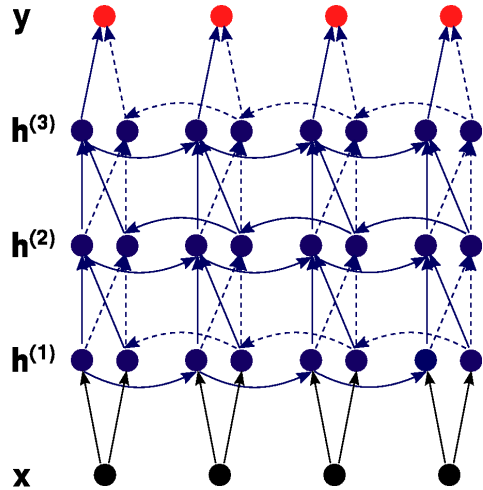
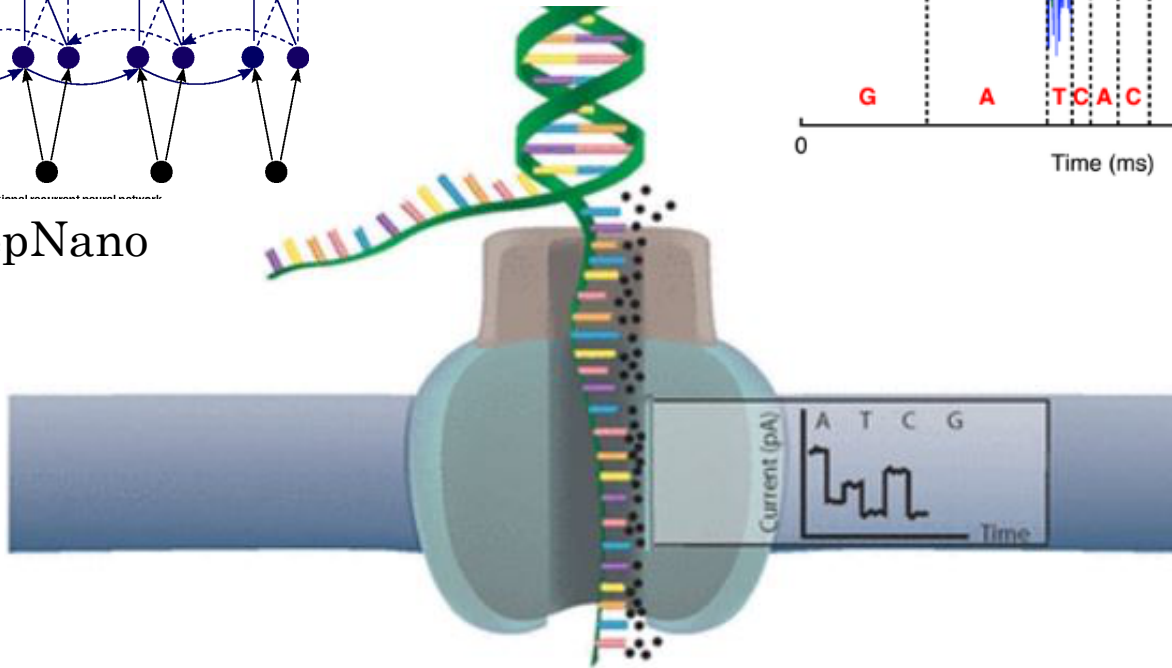
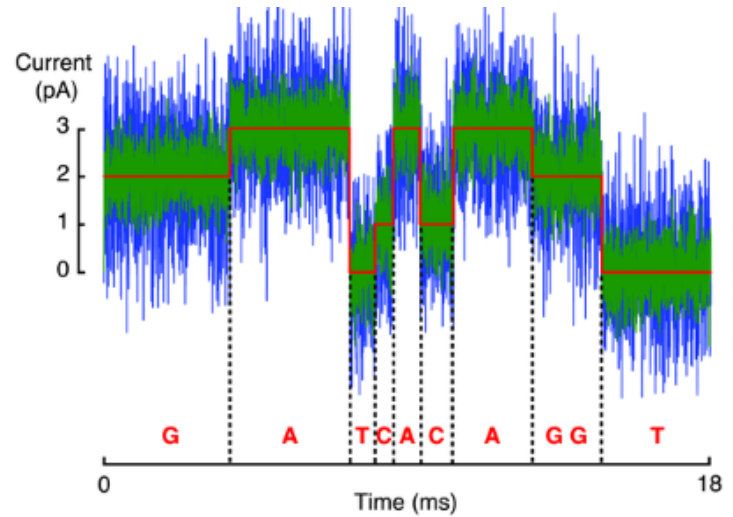


Fig. 2. Schematic of a bidirectional recurrent neural network.

DeepNano



Credit: LabBiotech

MINION AND GRIDION



Cheap(er) and Portable



RECONSTRUCTING SEQUENCES FROM TRACES ALIGNMENT

Position Consensus	20	22	23	24	25	26	32	42	45	50	53	56	69	81	93	103	121	122	123	124	147	148	149	153	155	156	170	183	Total
AXE	A	-	-	-	-	-	A	C	C	A	C	C	T	C	A	A	-	-	-	-	G	T	T	T	C	G	A	G	13
DW	A	-	-	-	-	-	A	C	C	A	C	C	T	C	A	A	-	-	-	-	G	T	T	G	C	G	A	G	12
FNZ	A	-	-	-	-	-	A	C	C	A	C	C	T	C	A	A	-	-	-	-	G	T	T	G	C	G	A	G	12
NQ	A	-	-	-	-	-	A	C	C	G	C	C	T	C	A	A	-	-	-	-	G	T	T	G	C	G	C	G	12
ZT	A	-	-	-	-	-	A	C	C	G	C	C	T	C	A	A	-	-	-	-	G	T	T	G	C	G	C	G	12
MEM	A	-	-	-	-	-	A	C	C	G	T	C	T	C	A	A	-	-	-	-	G	T	T	G	C	G	C	G	13
GD	A	A	A	A	G	C	A	C	T	A	C	C	C	G	A	A	-	-	-	-	G	C	T	G	T	G	C	G	7
IM	A	A	A	A	G	C	A	C	T	A	C	C	C	G	A	A	-	-	-	-	G	C	T	G	T	G	C	G	7
NRS	A	A	A	A	G	C	A	C	T	A	C	C	C	G	A	A	-	-	-	-	G	C	T	G	T	G	C	G	7
ZPJ	A	A	A	A	G	C	A	C	T	A	C	C	C	G	A	A	-	-	-	-	G	C	T	G	T	G	C	G	7
AB	A	A	A	A	G	C	A	C	T	A	C	C	C	C	A	A	A	T	T	A	G	T	T	G	C	G	C	G	0
XC	A	A	A	A	G	C	A	C	T	A	C	C	C	C	A	A	A	T	T	A	G	T	T	G	C	G	C	G	0
MR	T	A	A	A	G	C	C	C	T	A	C	C	C	C	A	A	A	T	T	A	G	T	T	G	C	G	C	G	2
UBG	T	A	A	A	G	C	C	C	T	A	C	C	C	C	A	A	A	T	T	A	G	T	T	G	C	G	C	G	2
FYC	A	A	A	A	G	C	A	A	T	A	C	T	C	C	T	G	A	T	T	A	G	T	T	G	C	T	C	-	6
JRR	A	A	A	A	G	C	A	A	T	A	C	T	C	C	T	G	A	T	T	A	G	T	T	G	C	T	C	-	6
UM	A	A	A	A	G	C	A	A	T	A	C	T	C	C	T	G	A	T	T	A	G	T	T	G	C	T	C	-	6
WZ	A	A	A	A	G	C	A	A	T	A	C	T	C	C	A	G	A	T	T	A	G	T	T	G	C	T	C	-	5
DGO	A	A	A	A	G	C	A	A	T	A	C	C	C	C	A	A	A	T	T	A	-	-	-	G	C	T	C	-	6
Summary:	17A	13A	13A	13A	13G	13C	17A	14C	13T	16A	18C	15C	13C	15C	16A	15A	9A	9T	9T	9A	18G	14T	18T	18G	15C	14G	16C	14G	135
	2T	6-	6-	6-	6-	6-	2C	5A	6C	3G	1T	4T	6T	4G	3T	4G	10-	10-	10-	10-	1-	4C	1-	1T	4T	5T	3A	5-	

Dynamic programming, greedy algorithms etc
 Basic computation of weighted Levenshtein distance
 CLUSTAL OMEGA, MUSCLE, etc

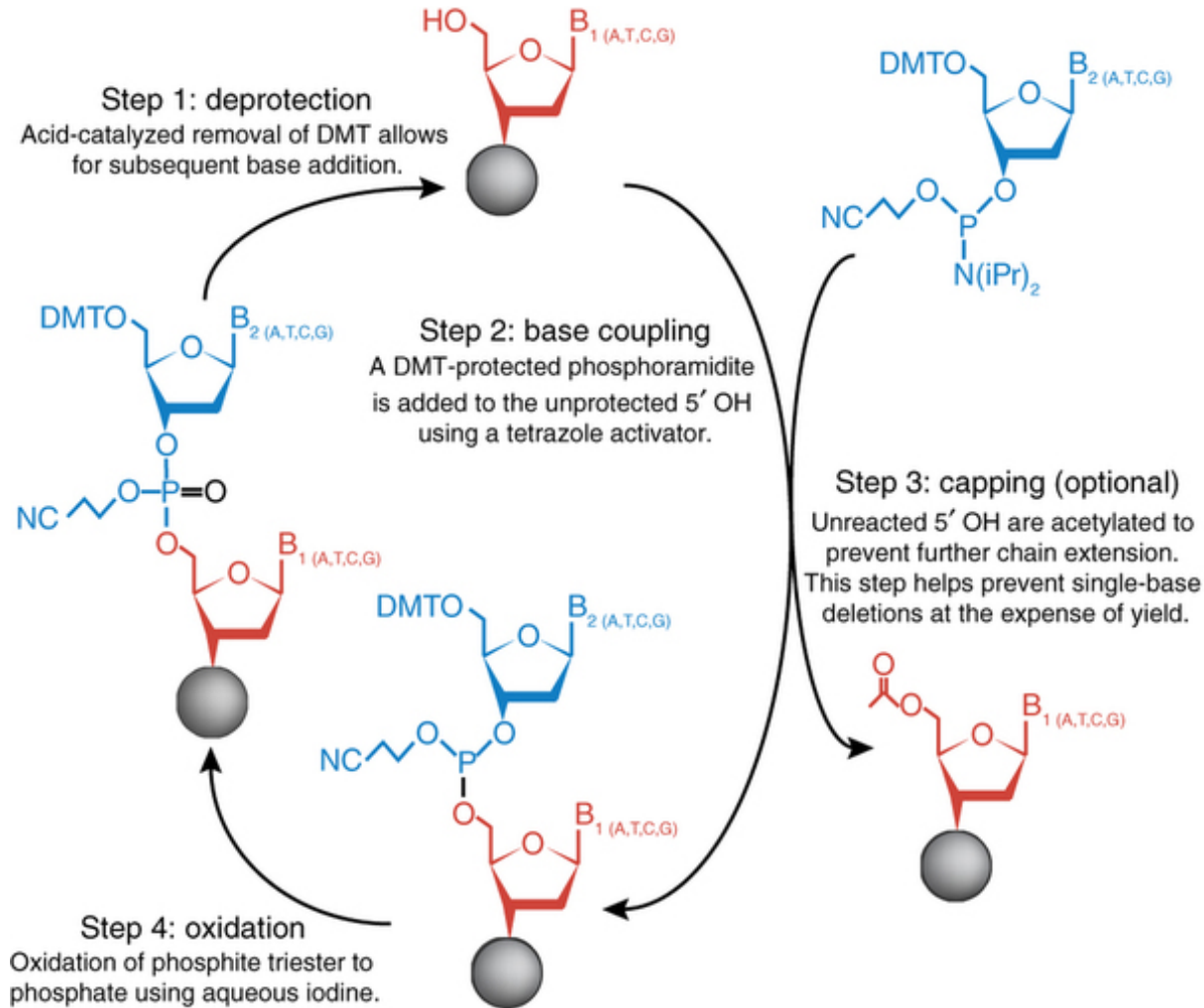
RECONSTRUCTING SEQUENCES FROM TRACES ALIGNMENT

The screenshot shows a software window titled "DNA Sequence Alignment". On the left, there are input fields for "Sequence 1" (GAATTCAGTTA), "Sequence 2" (GGATCGA), "Similarity Score" (1), "Non Similarity Score" (0), and "Gap Penalty" (-1). An "Align" button is located below these fields. On the right, a similarity matrix is displayed. The matrix has a blue cell at the top-left corner (row 1, column 1). Below the matrix, the aligned sequences are shown in blue text: GAATTCAGTTA and GGA-TC-G--A.

▶	-	G	A	A	T	T	C	A	G	T	T	A	
	-	0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-10	-11
	G	-1	1	0	-1	-2	-3	-4	-5	-6	-7	-8	-9
	G	-2	0	1	0	-1	-2	-3	-4	-4	-5	-6	-7
	A	-3	-1	1	2	1	0	-1	-2	-3	-4	-5	-5
	T	-4	-2	0	1	3	2	1	0	-1	-2	-3	-4
	C	-5	-3	-1	0	2	3	3	2	1	0	-1	-2
	G	-6	-4	-2	-1	1	2	3	3	3	2	1	0
*	A	-7	-5	-3	-1	0	1	2	4	3	3	2	2

GAATTCAGTTA
GGA-TC-G--A

FROM READING TO WRITING

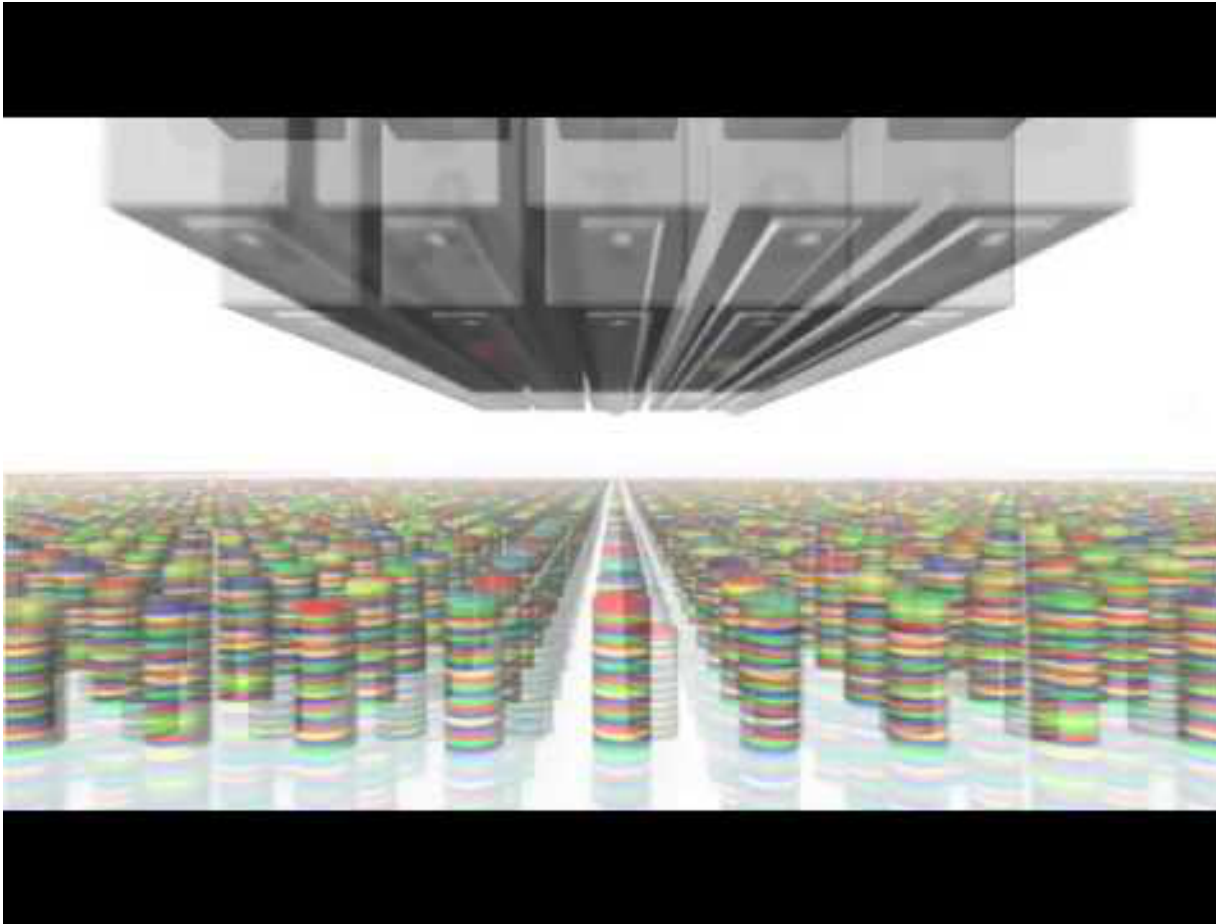


Expensive

Relatively slow but parallelizable

Commercially available from Agilent Twist IDT

FROM READING TO WRITING



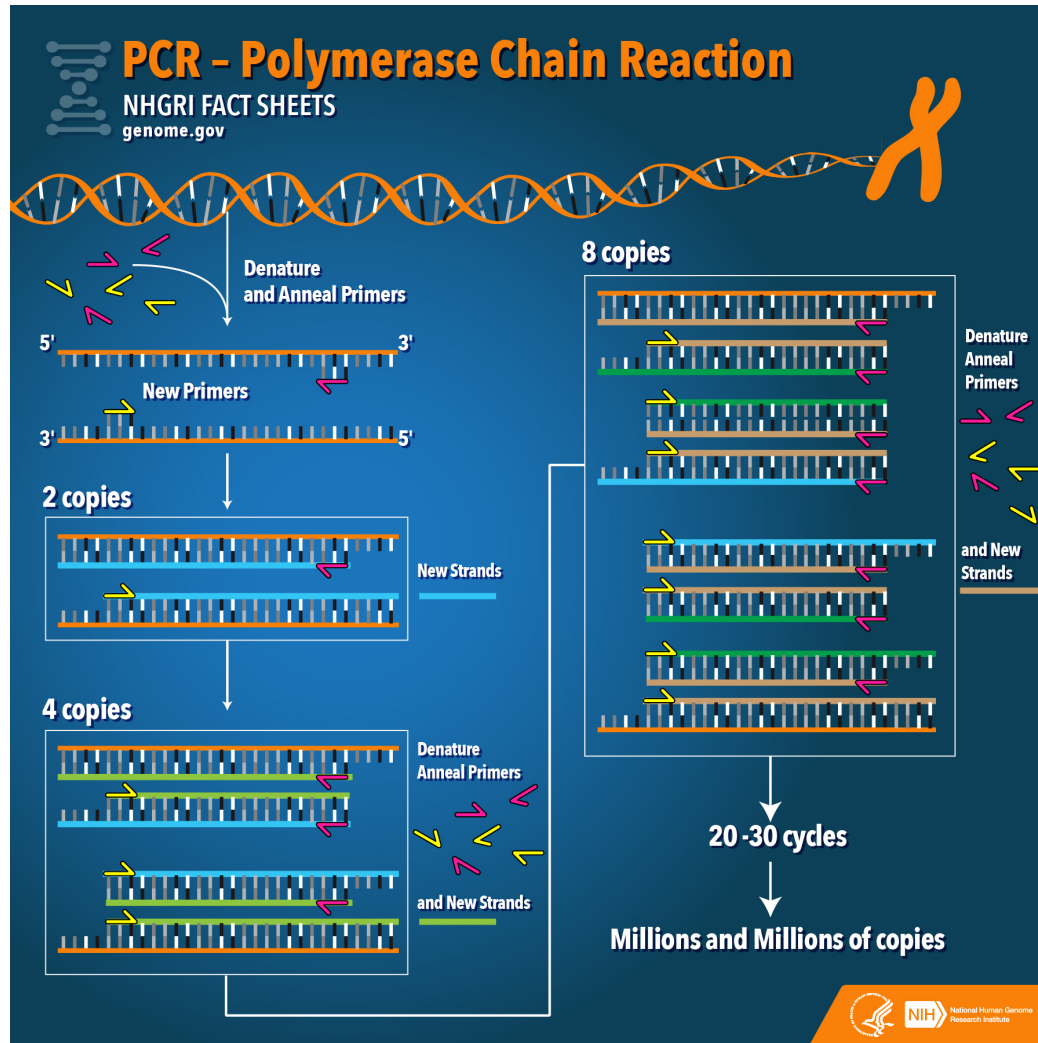
Expensive

Relatively slow but
parallelizable

Commercially
available from
Agilent
Twist
IDT

Credit: Agilent (INKJet Technology)

FROM WRITING TO COPYING



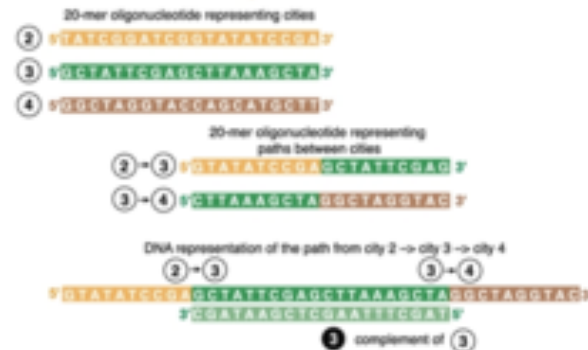
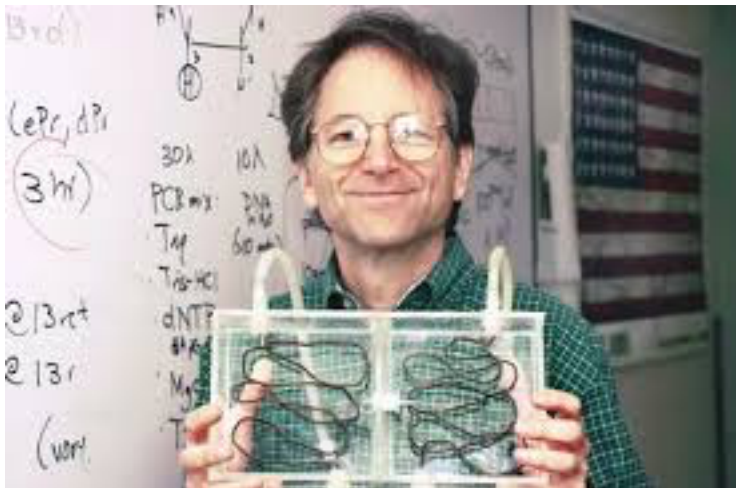
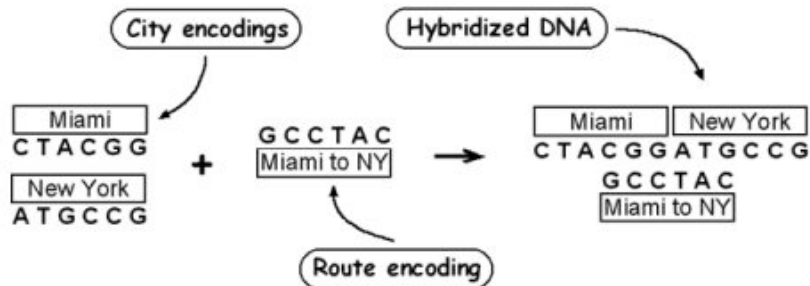
PCR machine

Thermofisher

(\$200)

Credit: NIH

AND COMPUTING: ADLEMAN'S EXPERIMENT



AND COMPUTING: STRAND DISPLACEMENT

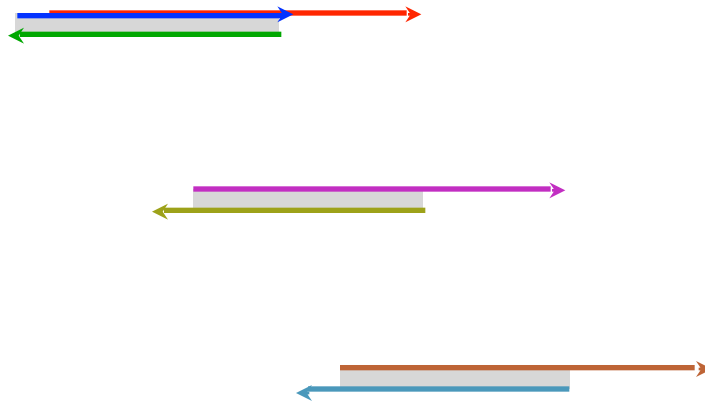
Input



Curtesy of D. Soloveichik



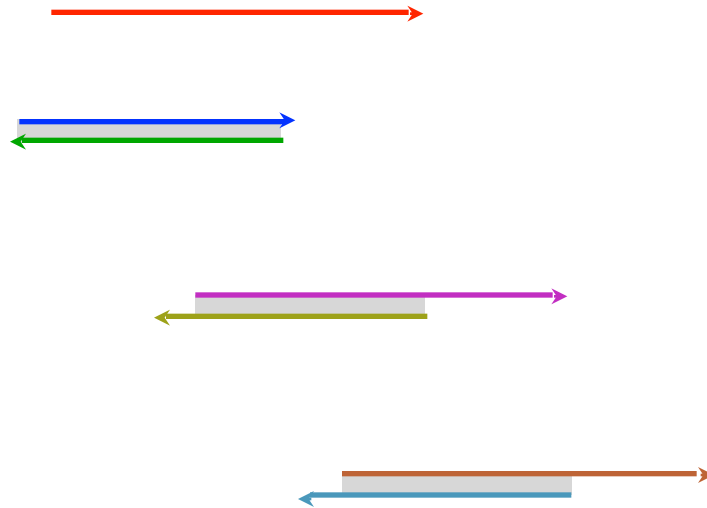
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



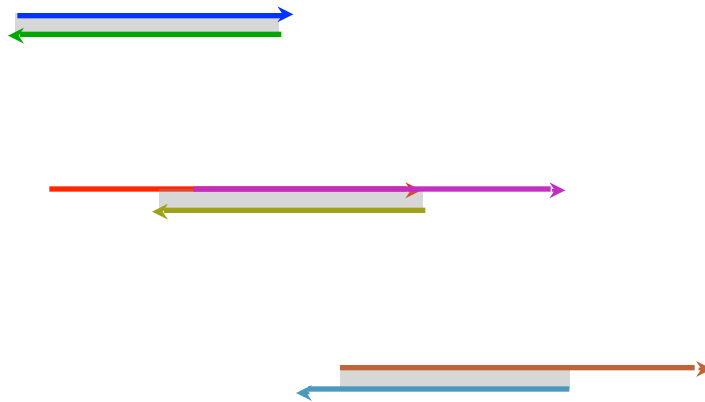
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



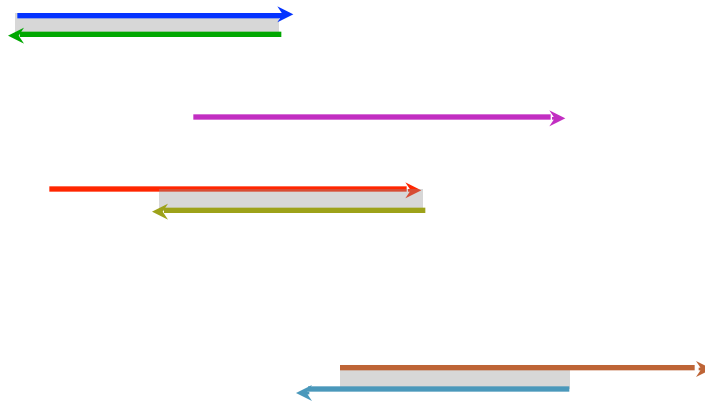
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



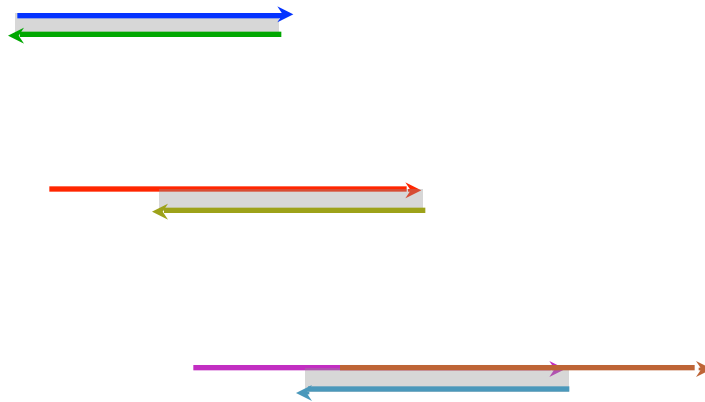
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



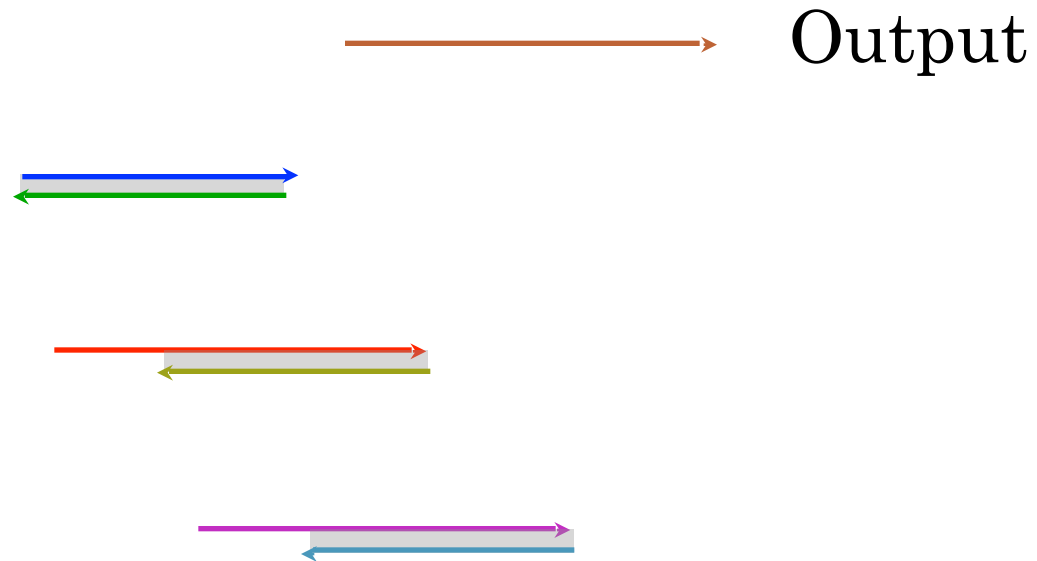
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



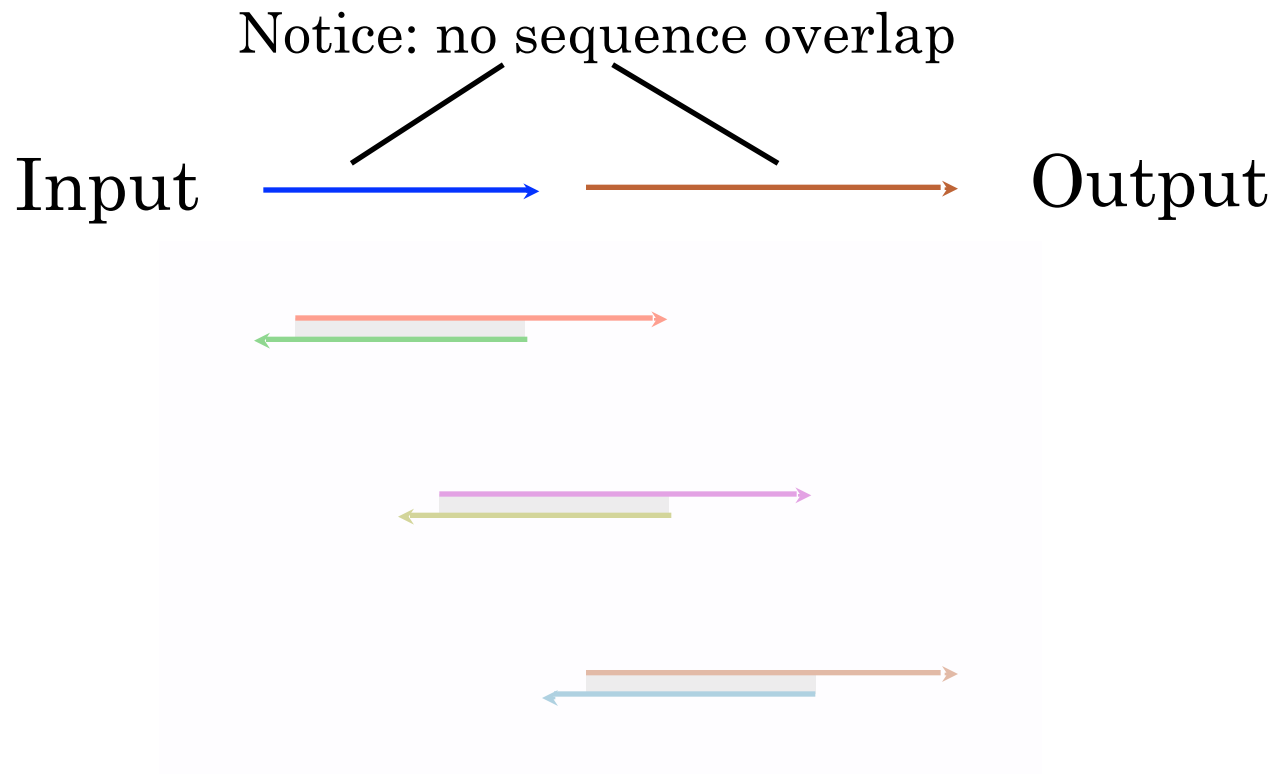
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



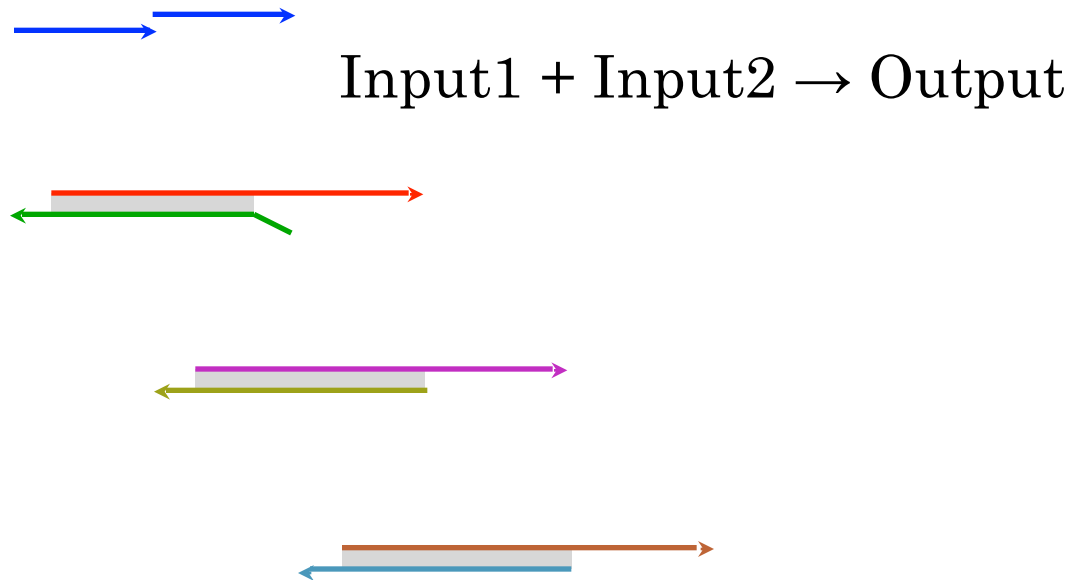
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik



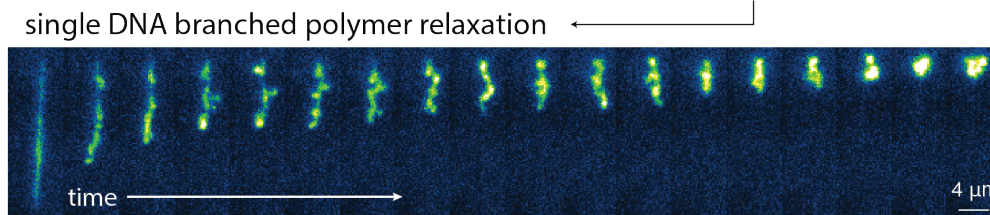
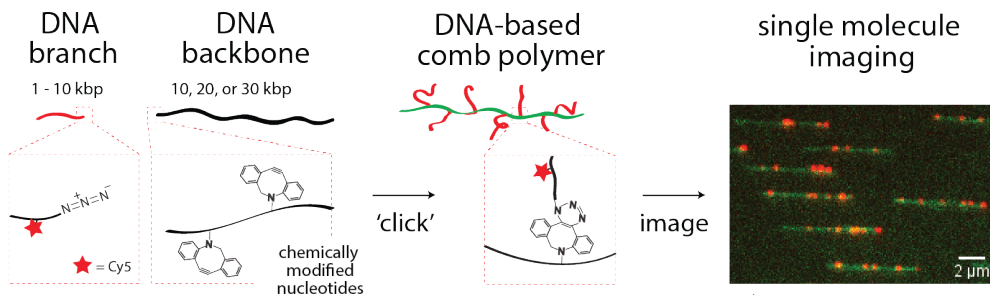
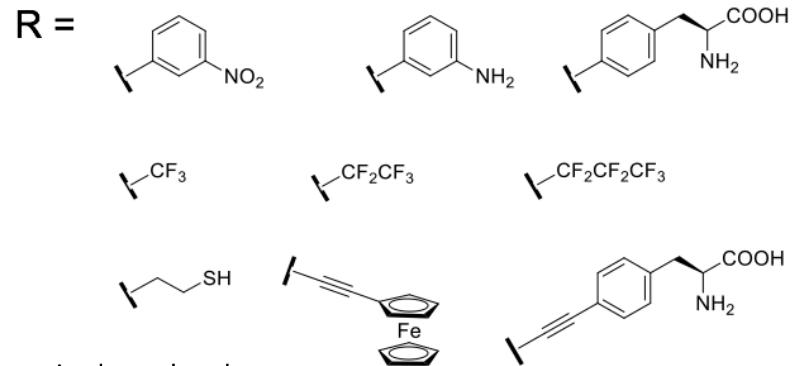
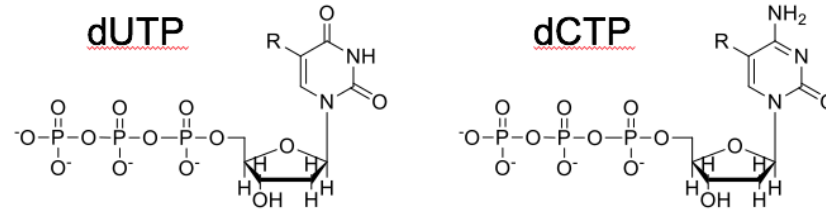
AND COMPUTING: STRAND DISPLACEMENT



Curtesy of D. Soloveichik

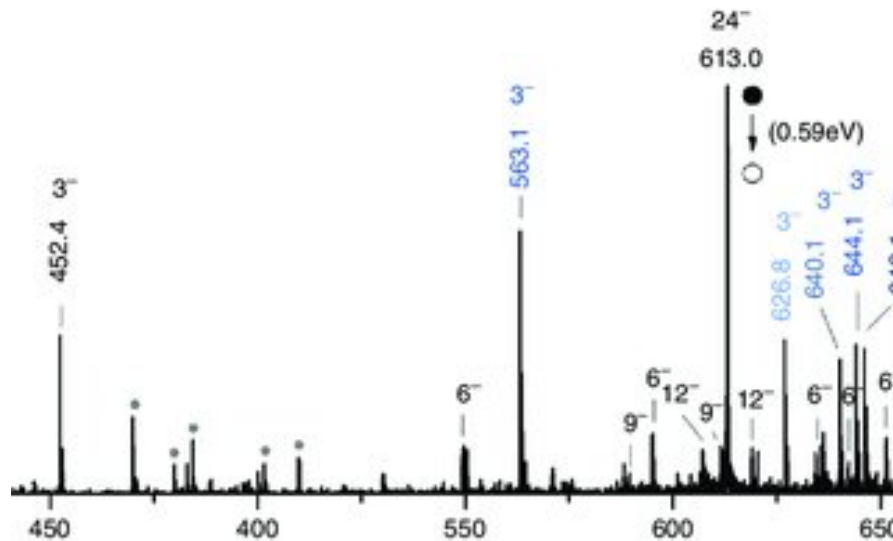
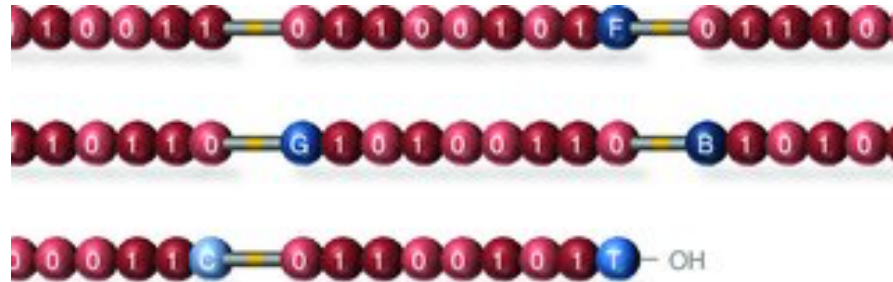
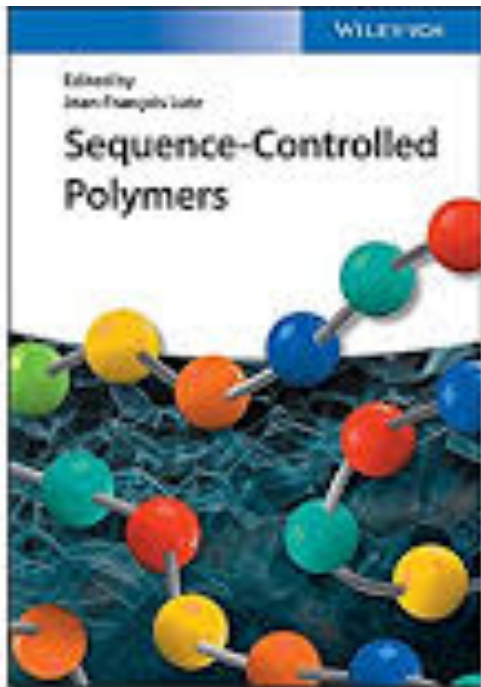


BEYOND DNA: CHEMICALLY MODIFIED DNA



Credit: C. Schroeder

BEYOND DNA: SYNTHETIC POLYMERS



Credit: J. F. Lutz