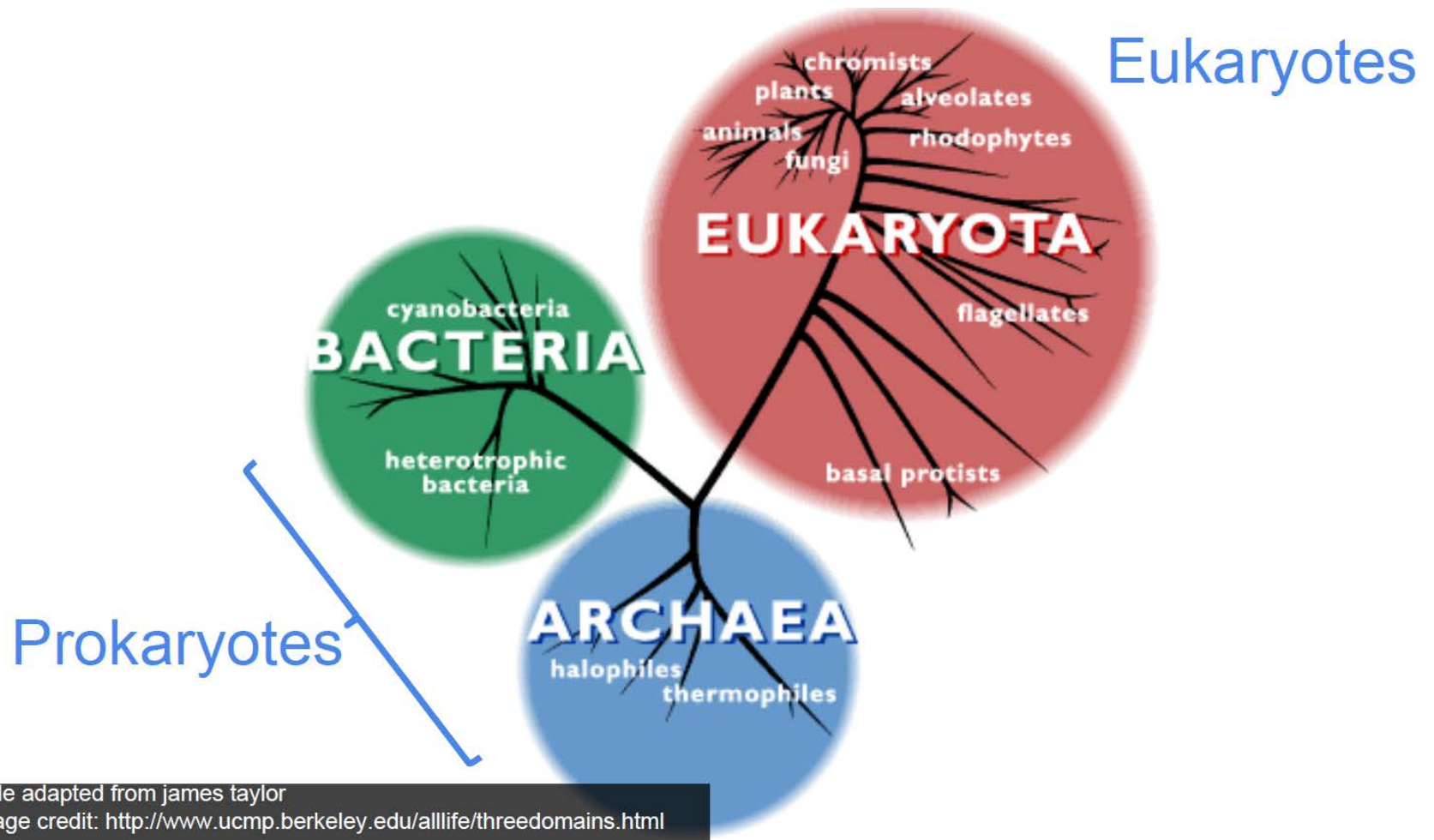


DNA/RNA sequencing

Computational Biology II

Credit to Coursera, Genomic data science specialization

Basic domains of life



Basic types of cells

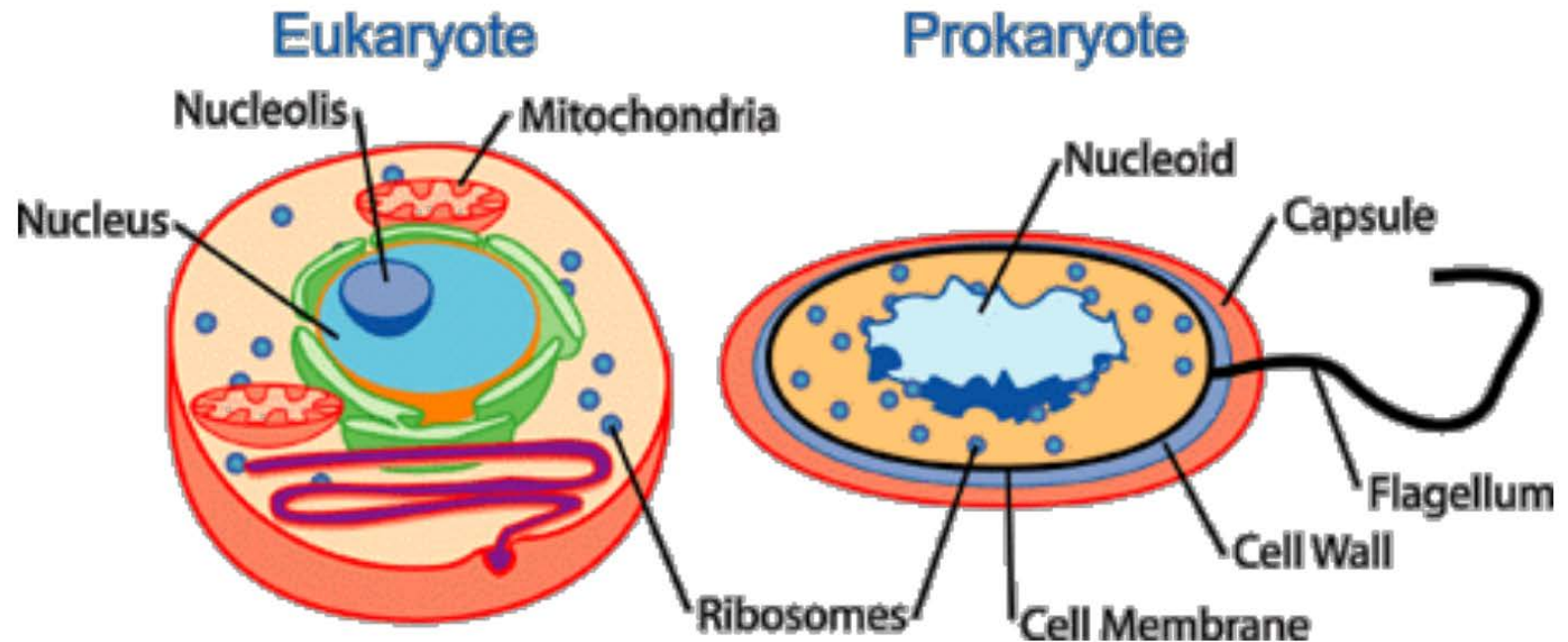


image credit: <http://commons.wikimedia.org/wiki/File:Celltypes.png>
slide adapted from james taylor

Eukaryotic cells

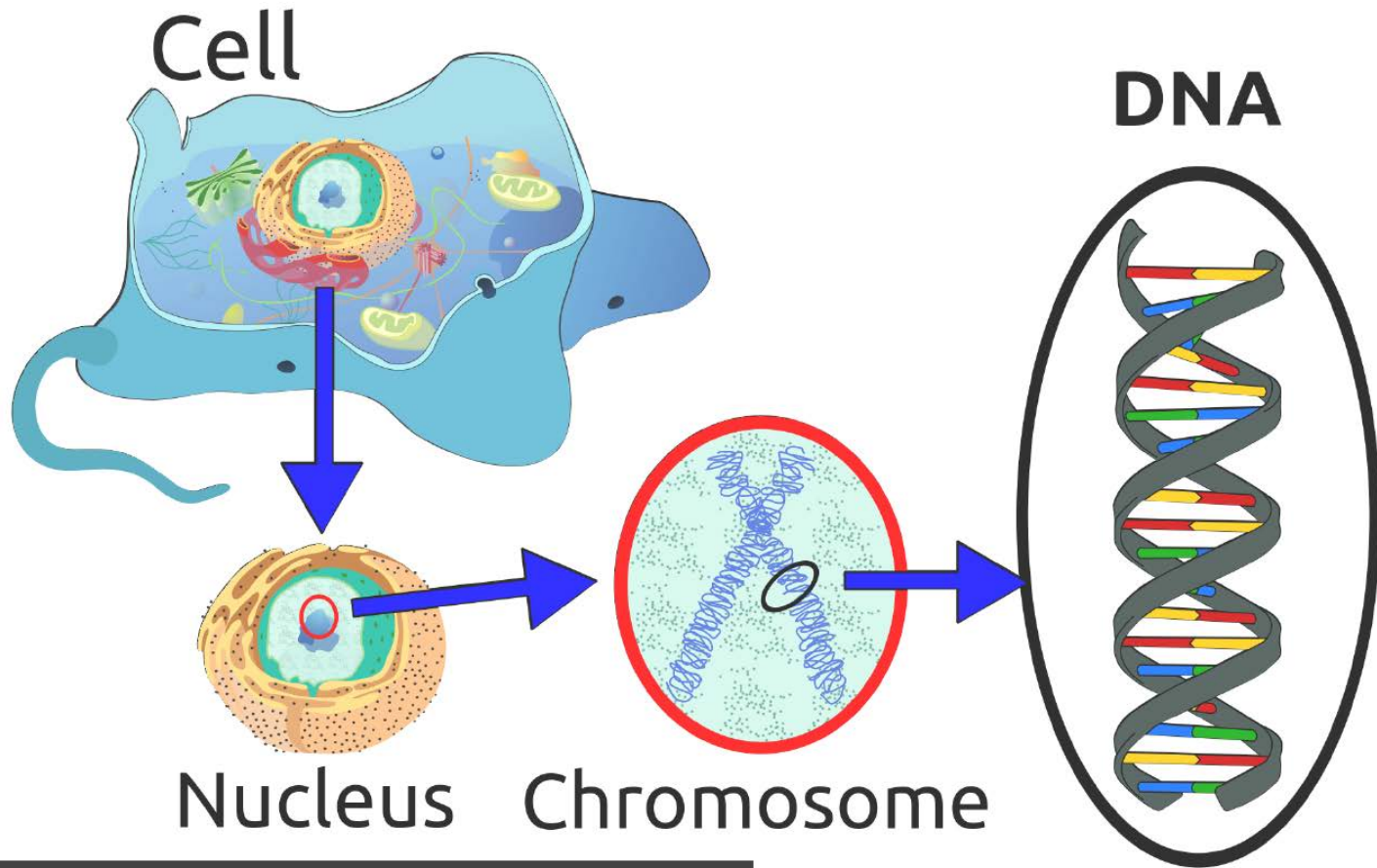
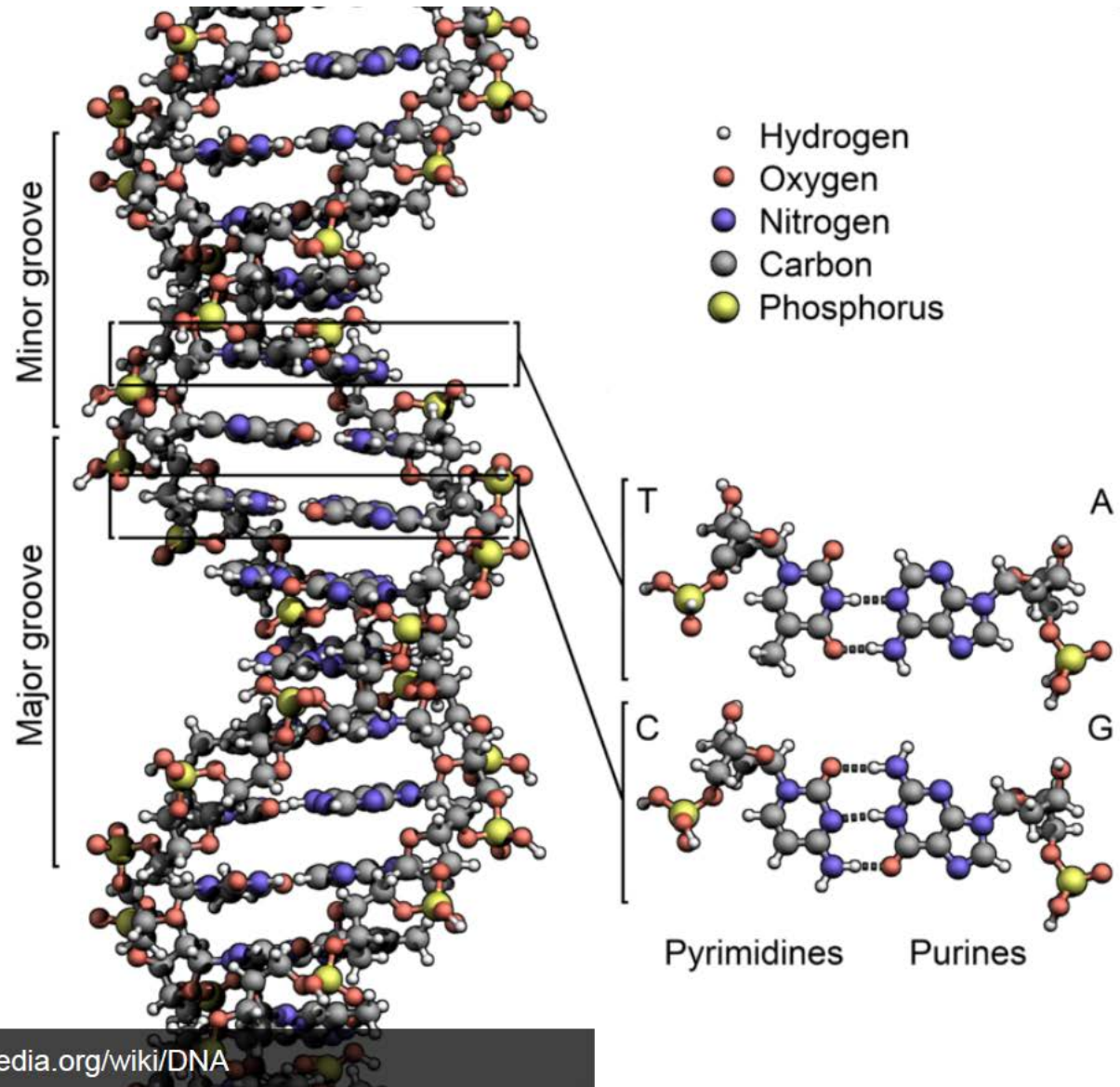
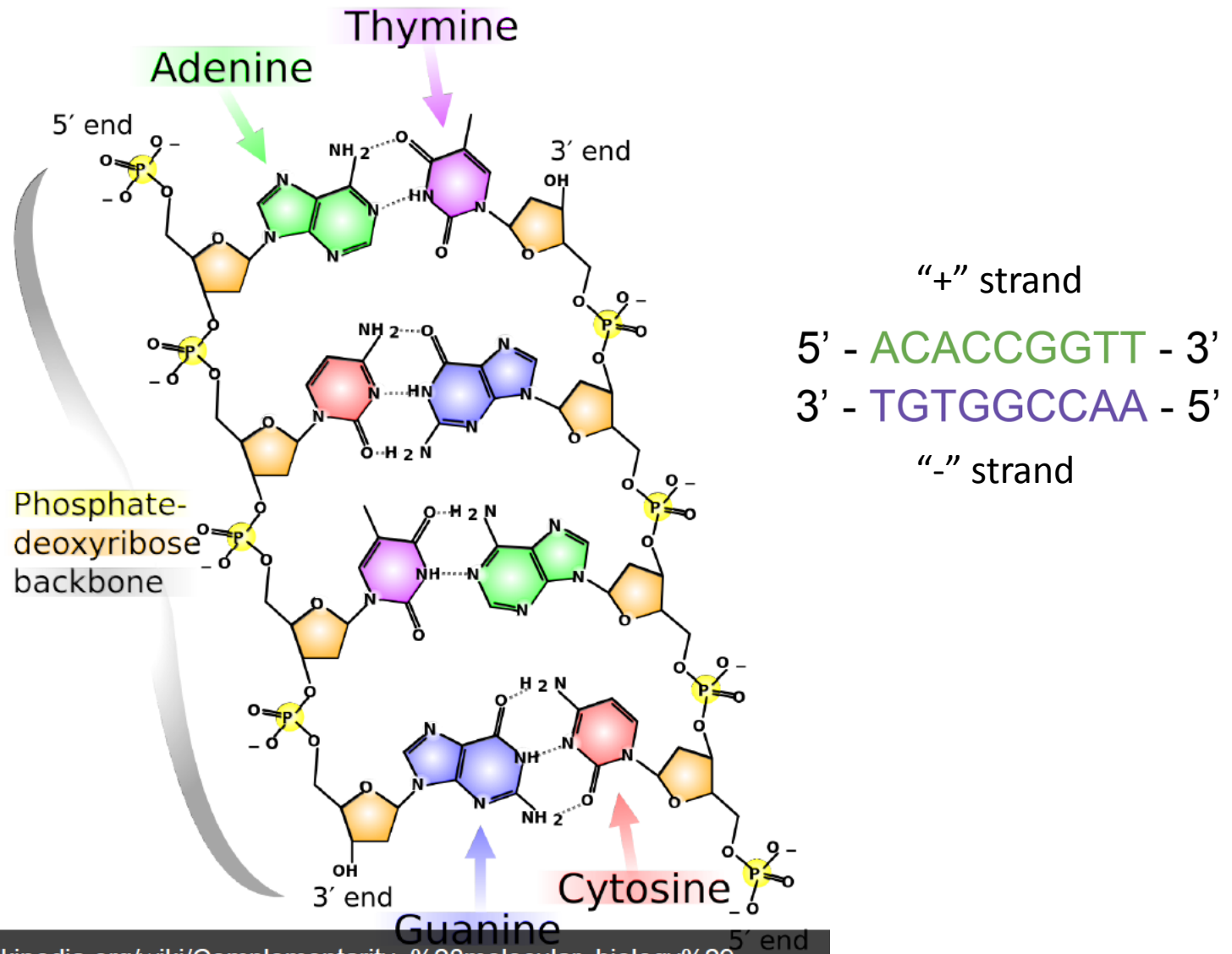


image credit: <http://en.wikipedia.org/wiki/DNA>

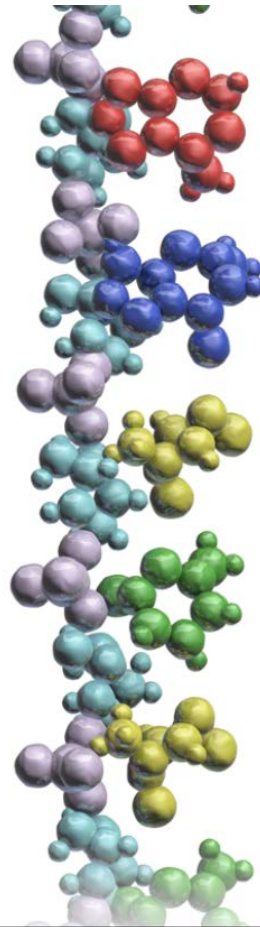
DNA nucleotides



DNA structure



RNA structure



Adenine

Guanine

Cytosine

Uracil

Cytosine

Single stranded nucleic acid

RNA sequence

5' - ACACCGGTT 3'
3' - TGTGGCAA - 5'

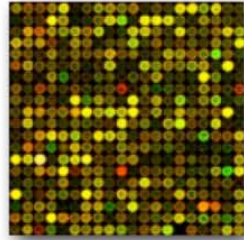


ACACCGGUU

DNA sequencing



Sanger DNA
sequencing
1977-1990s



DNA Microarrays
Since mid-1990s

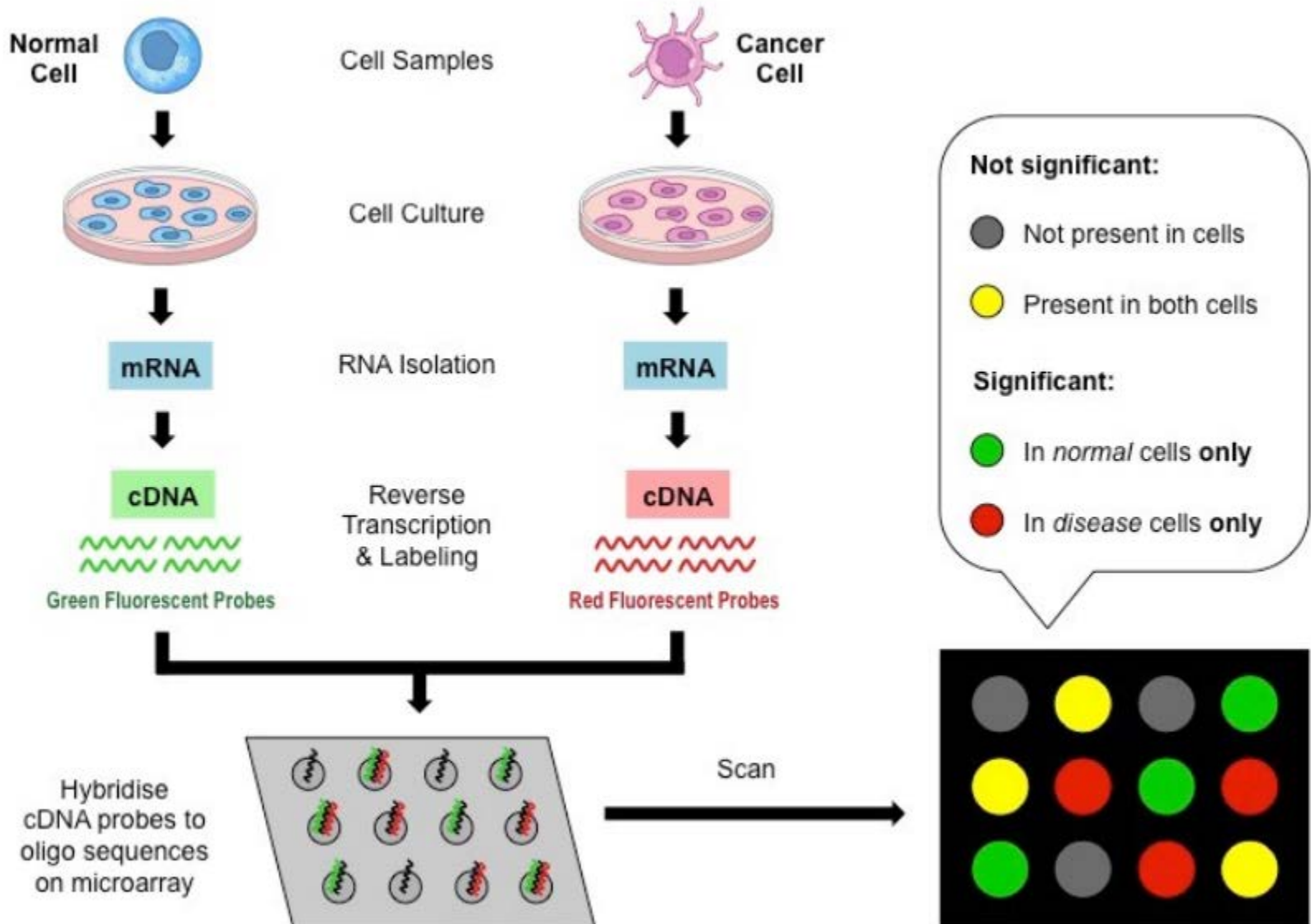


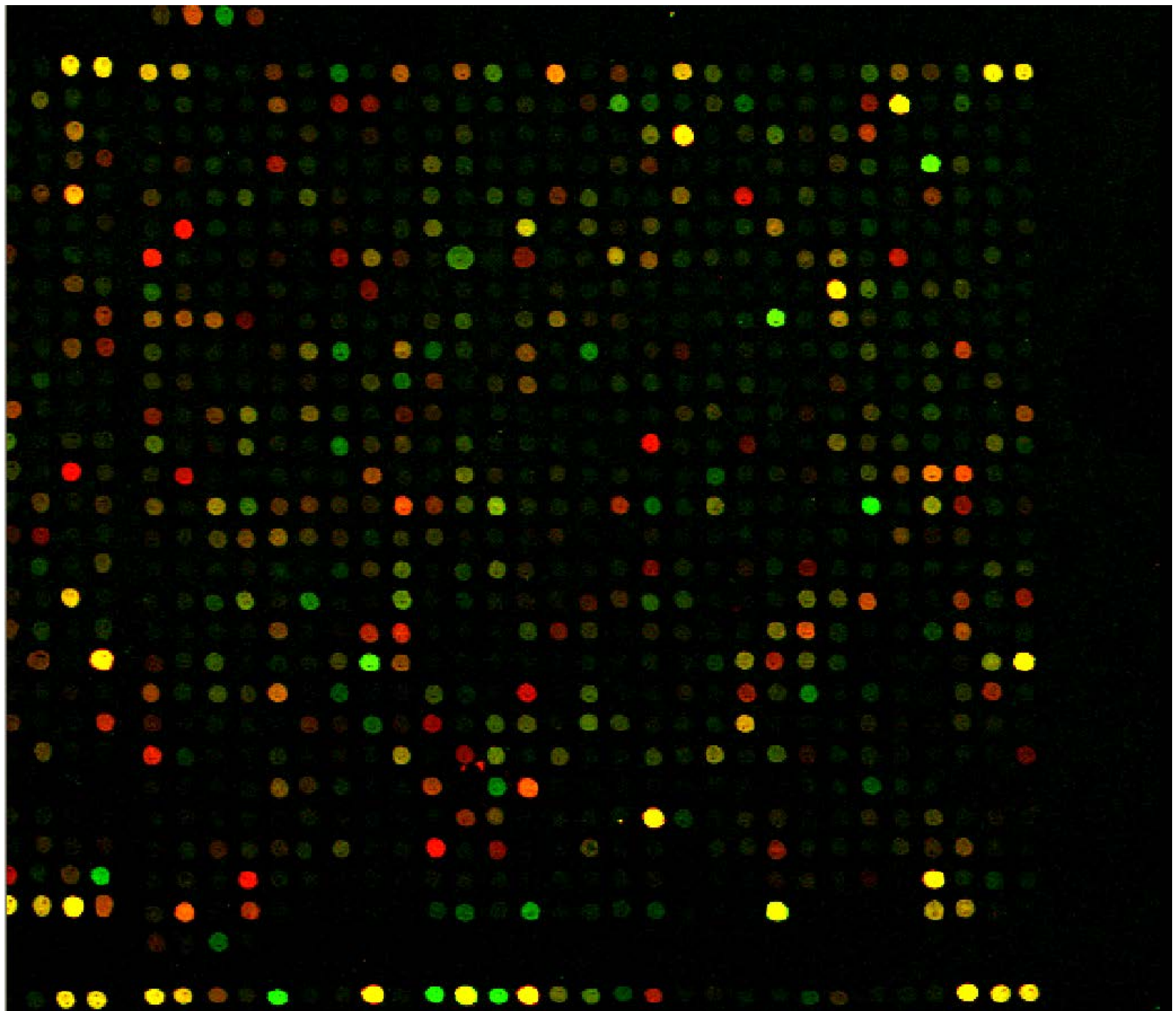
2nd Generation
DNA Sequencing
Since ~2007



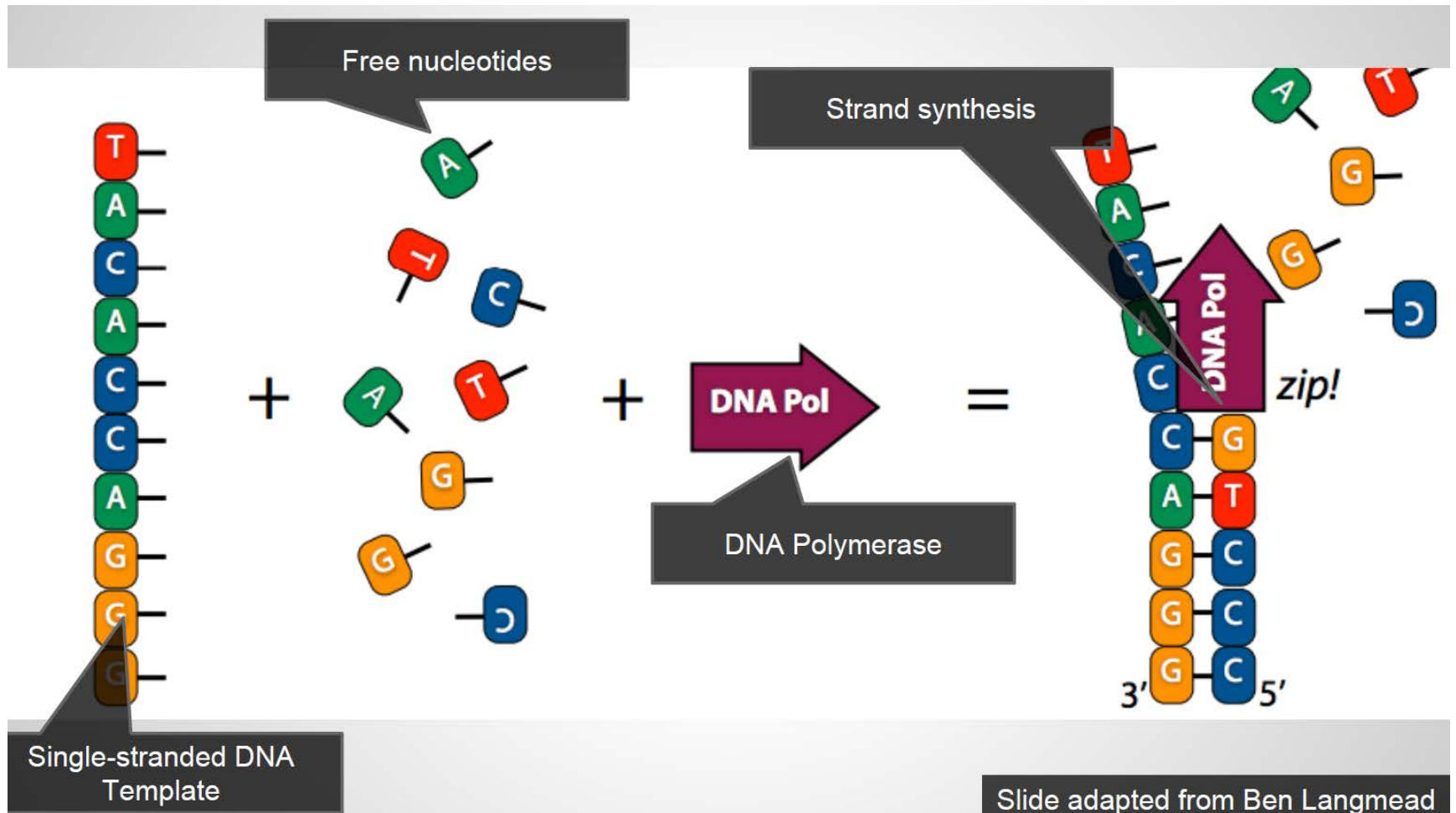
3rd Generation
& single molecule
Sequencing
Since ~ 2010

Microarrays

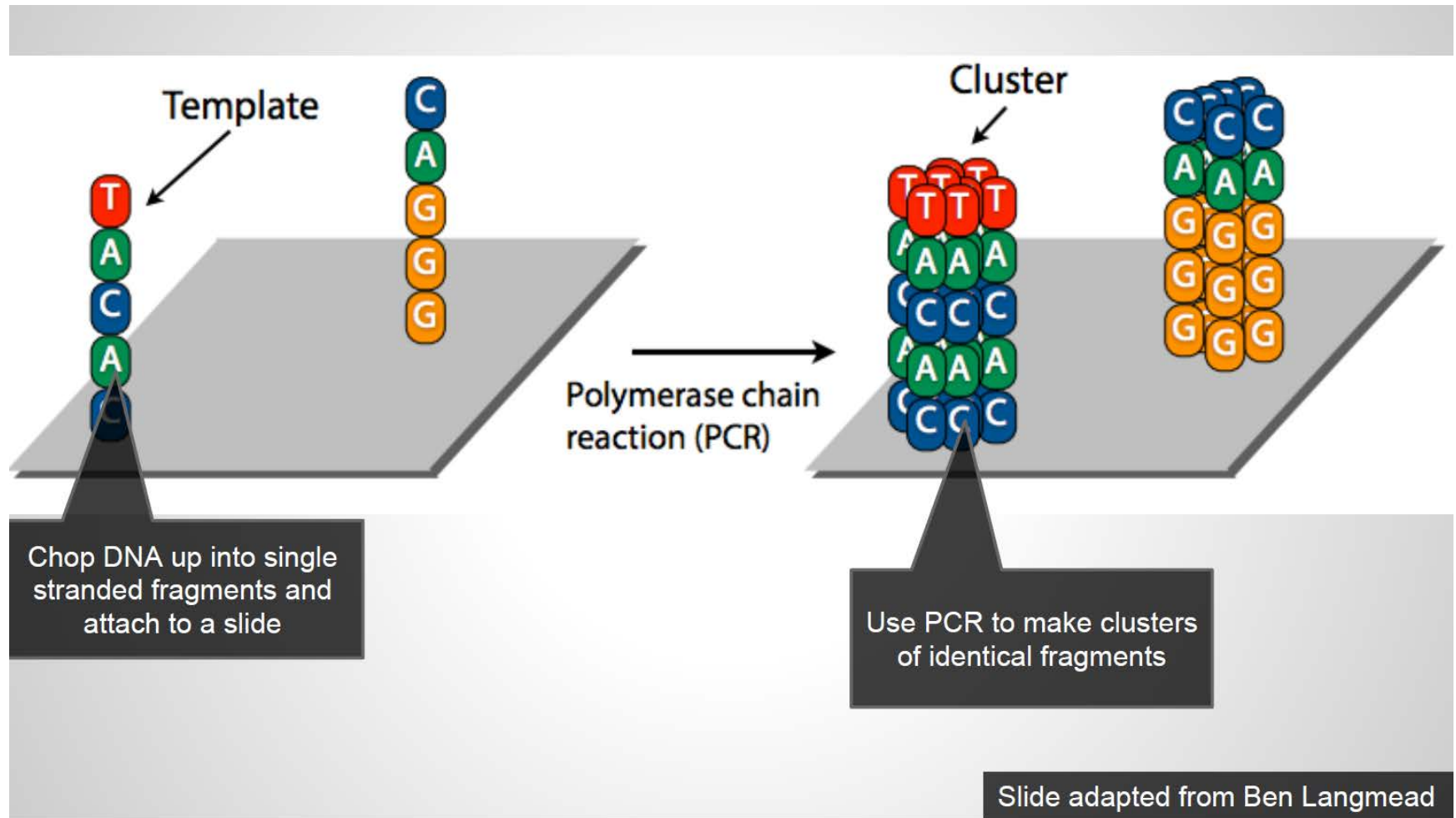




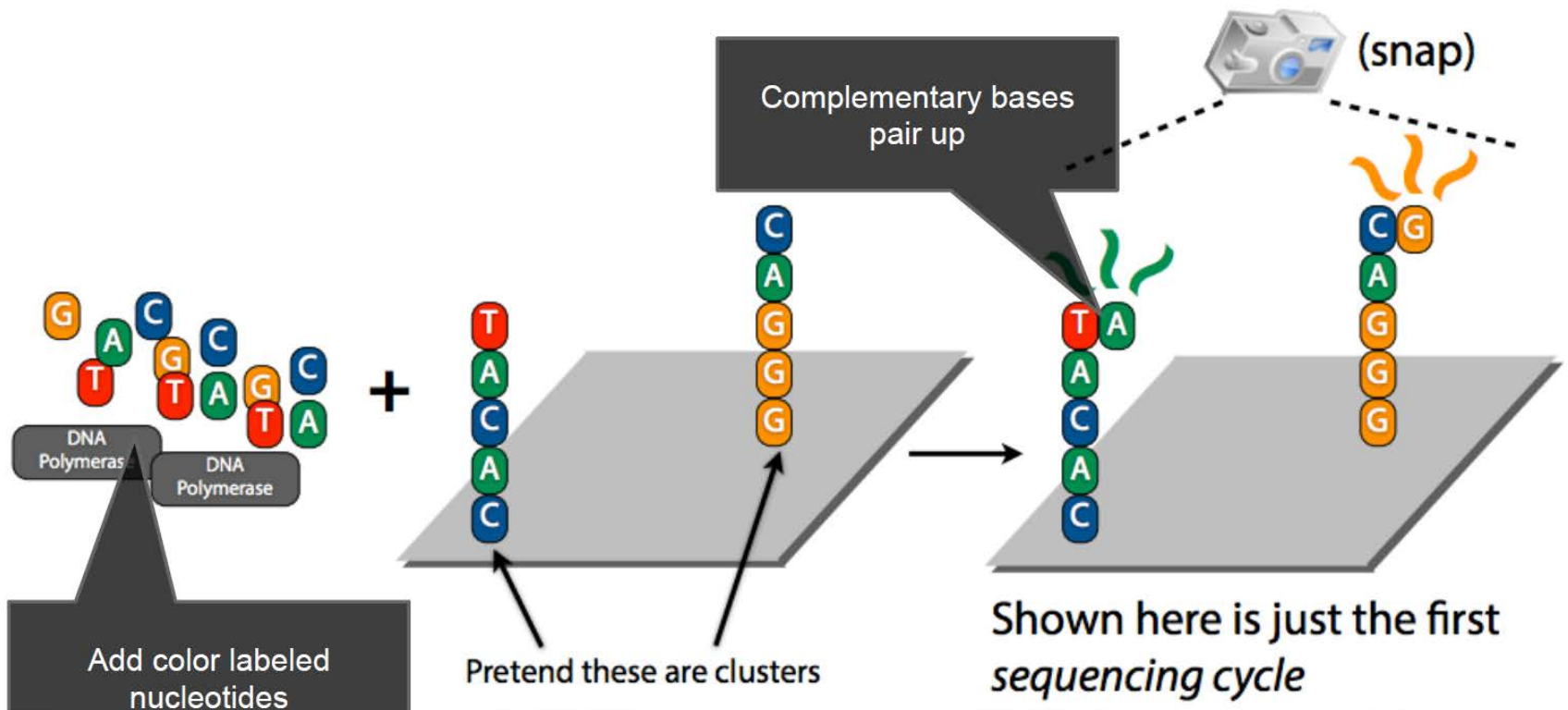
Polymerase chain reaction (PCR)



Polymerase chain reaction (PCR) II

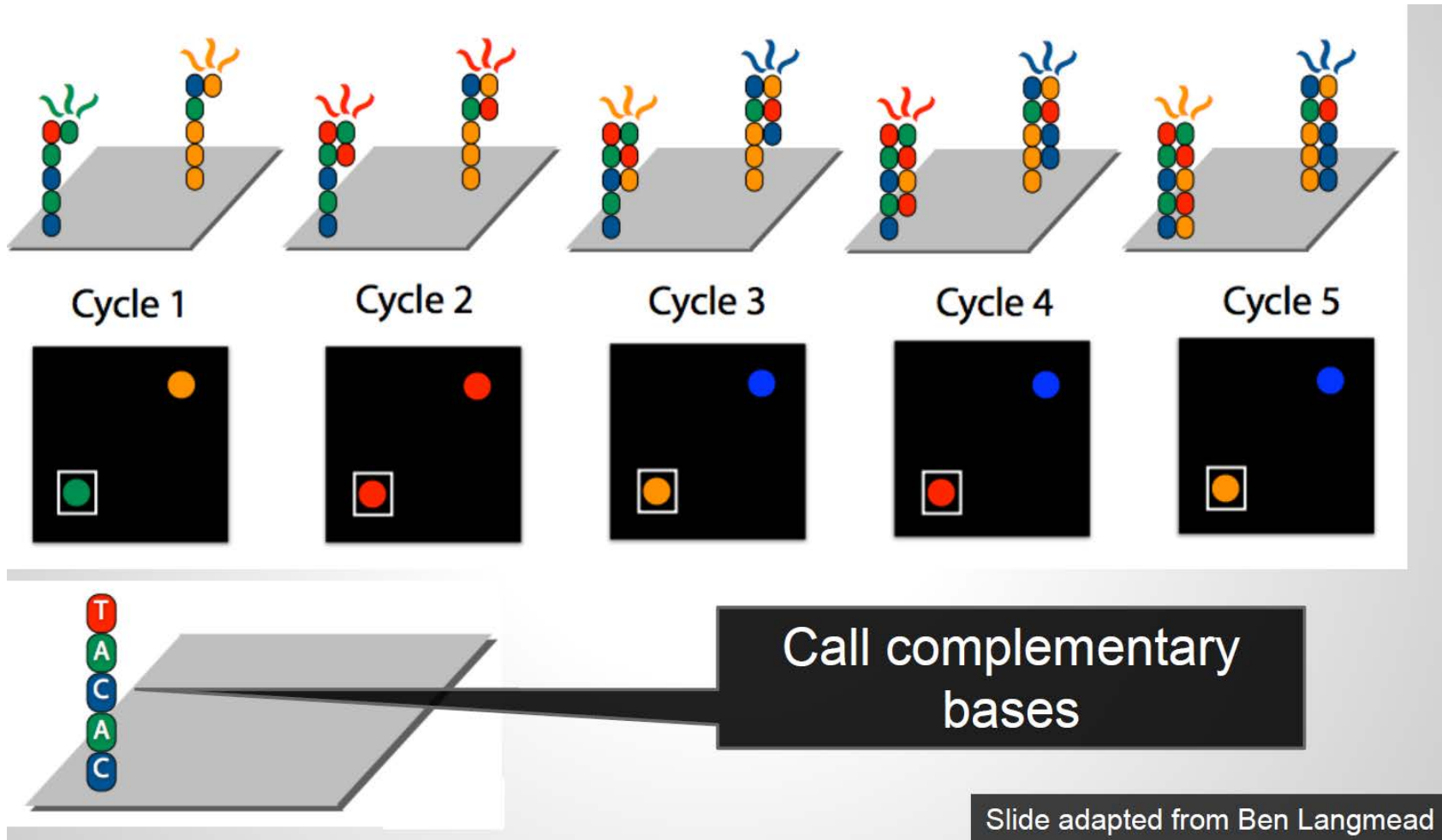


Next generation sequencing

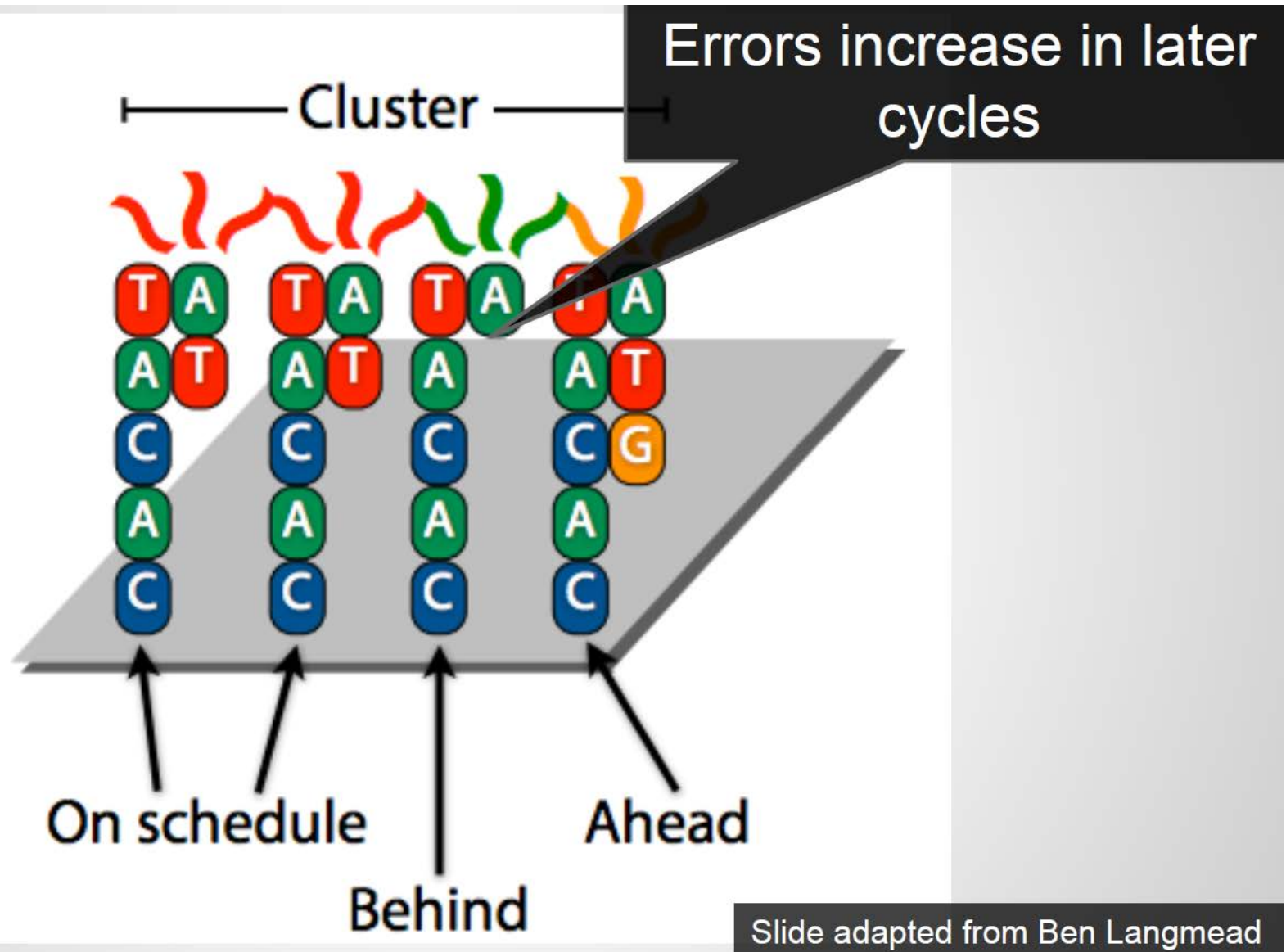


Slide adapted from Ben Langmead

Sequencing cycles



Sequencing errors



Output FASTQ file

[illegible][illegible]

Output FASTQ file II

Usually these identifiers give information about the precise location and geometry of the cluster within the sequencing instrument.

Header (preceded by '@') **Sequence** **+Header or '+'** **Base qualities** **Read in pair**

@UNC14-SN744:186:D078MACXX:1:1101:1203:1868/1
NGAGAAAAGAGGATTATTGCTGAGTGGCAGCACCAGCCCCAAAGGGAA
+
#1-DDDDFFHHHHHHJJJJJJJJJJJGIIJJJJJIIJJGIGIGIJJGIIJJJJBC
@UNC14-SN744:186:D078MACXX:1:1101:1340:1972
CAGGATTTGGCCTTAGCTTCTGGGCCTATCGGCTGCCTTCCCTCTACT
+
CCCFFFFFHHHHHHJIIJJJJJJJJJJIIJJJJJIIJJJJFHGHBBHIGDG<FG

The diagram illustrates the structure of a FASTQ file. It shows two reads. The first read is a single-end read with a header line starting with '@', followed by the sequence line, a plus sign, and the base quality line. The second read is a paired-end read, indicated by the 'Read in pair' label, and is enclosed in a red box. It follows the same format but includes a mate identifier in the header line. Arrows point from the labels to the corresponding parts of the text: 'Header' points to the '@' symbol, 'Sequence' points to the sequence string, '+Header or '+' points to the plus sign, 'Base qualities' points to the quality string, and 'Read in pair' points to the second read's header line.

Base quality score

- Let p_b = probability that the call at base b is incorrect
- Quality value: $Q_{\text{sanger}} = -10 \log_{10} p_b$ (integer)
- Sanger (Phred quality scores): 0..93 (ASCII characters 33..126)

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

- In practice, the maximal quality value is ~40.
- Quality values below 20 are typically considered low.

Representation: Single sequences, Fasta/Pearson format

- **Fasta/Pearson** – the most common format for Sanger sequences, and for genomic and gene sequences

Header (1 line, starts with '>')

>gi|23238195|ref|NM_000878.2| Homo sapiens interleukin 2 receptor, beta (IL2RB), mRNA
GCAGCCAGAGCTCAGCAGGGCCCTGGAGAGATGGCCACGGTCCCAGCACCGGGGAGGACTGGAGAGCGCG
CGCTGCCACCGCCCCATGTCTCAGCCAGGGCTTCCTTCCTCGGCTCCACCCTGTGGATGTAATGGCGGCC
CCTGCTCTGTCTGGCGTCTGCCCCCTCCTCATCCTCCTGCCCCCTGGCTACCTCTTGGGCATCTGCAG
CGGTGAATGGCACTTCCCAGTTCACATGCTTCTACAACCTCGAGAGCCAACATCTCCTGTGTCTGGAGCCA
AGATGGGGCTCTGCAGGACACTTCCTGCCAAGTCCATGCCTGGCCGGACAGACGGCGGTGGAACCAAACC
TGTGAGCTGCTCCCCGTGAGTCAAGCATCCTGGGCCTGCAACCTGATCCTCGGAGCCCCAGATTCTCAGA
AACTGACCACAGTTGACATCGTCACCCTGAGGGTGCTGTGCCGTGAGGGGGTGCGATGGAGGGTGATGGC
CATCCAGGACTTCAAGCCCTTTGAGAACCCTTCGCCTGATGGCCCCCATCTCCCTCCAAGTTGTCCACGTG
GAGACCCACAGATGCAACATAAGCTGGGAAATCTCCCAAGCCTCCCACTACTTTGAAAGACACCTGGAGT
TCGAGGCCCGGACGCTGTCCCAGGCCACACCTGGGAGGAGGCCCCCCTGCTGACTCTCAAGCAGAAGCA
GGAATGGATCTGCCTGGAGACGCTCACCCCAGACACCCAGTATGAGTTTCAGGTGCGGGTCAAGCCTCTG
CAAGGCGAGTTCACGACCTGGAGCCCCCTGGAGCCACCCCCCTGGCCTTCAGGACAAAGC

Sequence (nucleotide or protein) (1 or several lines)

Representation: Multiple sequences, Multi-Fasta format

➔ **>gi|28178860|ref|NM_000586.2| Homo sapiens interleukin 2 (IL2), mRNA**
CGAATTCCCCTATCACCTAAGTGTGGGCTAATGTAACAAAGAGGGATTTACCTACATCCATTTCAGTCAG
TCTTTGGGGGTTTAAAGAAATTCCAAAGAGTCATCAGAAGAGGAAAAATGAAGGTAATGTTTTTTCAGAC
AGGTAAAGTC. ATAAAAAAAAAAAAA

➔ **>gi|31982837|ref|NM_008366.2| Mus musculus interleukin 2 (Il2), mRNA**
ATCACCCCTTGCTAATCACTCCTCACAGTGACCTCAAGTCCTGCAGGCATGTACAGCATGCAGCTCGCATC
CTGTGTCACATTGACACTTGTGCTCCTTGTC AACAGCGCACCCACTTCAAGCTCCACTTCAAGCTCTACA
GCGGAAGCACAGCAGCAGCAGCAGCAGCAGCAG AGCTCTCCTCT

➔ **>gi|16758691|ref|NM_053836.1| Rattus norvegicus interleukin 2 (Il2), mRNA**
GAAGTCCTGCAAGCATGTACAGCATGCAGCTCGCATCCTGTGTTGCACTGACGCTTGTCTCCTTGTCAG
CAGCGCACCCACTTCAAGCCCTGCAAAGGAAACACAGCAGCACCTGGAGCAGCTGTTGCTGGACTTACAG
GTGCTCCTGAGAGGGATC AGTATTTAGAAGAGTCGATGAA

➔ **>gi|30794289|ref|NM_180997.1| Bos taurus interleukin 2 (IL2), mRNA**
GGGCTATCTGTTTCGGTCGTTTCATGTCAGCAATGTACAAGATACTCTTGTCTTGCATTGCACTAACTC
TTGCACTCGTTGCAAACGGTGCACCTACTTCAAGCTCTACGGGGAACACAATGAAAGAAGTGAAGTCATT
GCTGCTGGATTTACAGTTGCTTTTGGAGAAAGTTAAAAATCCTGAGAACCTCAAGCTCTCCAGGATGCAT
. ATTTAATAAAGTTGATGAATAAAAAAC

Summary

