# GENOME INFORMATICS

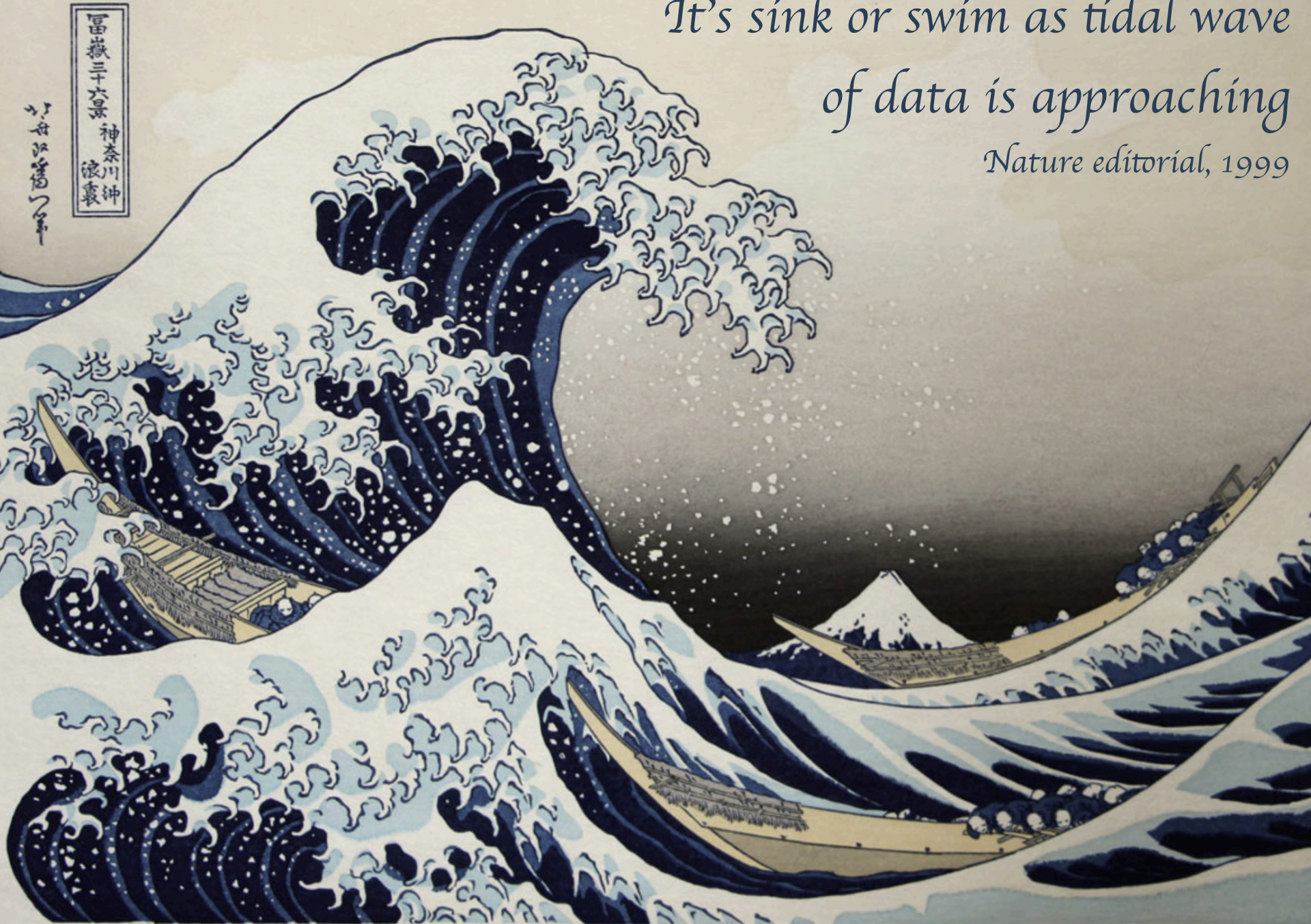http://bioinformatics.uni-muenster.de/teaching/Current/Genome_informatics/index.hbi

Prof. Dr. Wojciech Makałowski
Institute of Bioinformatics
University of Münster, Germany

It's sink or swim as tidal wave
of data is approaching
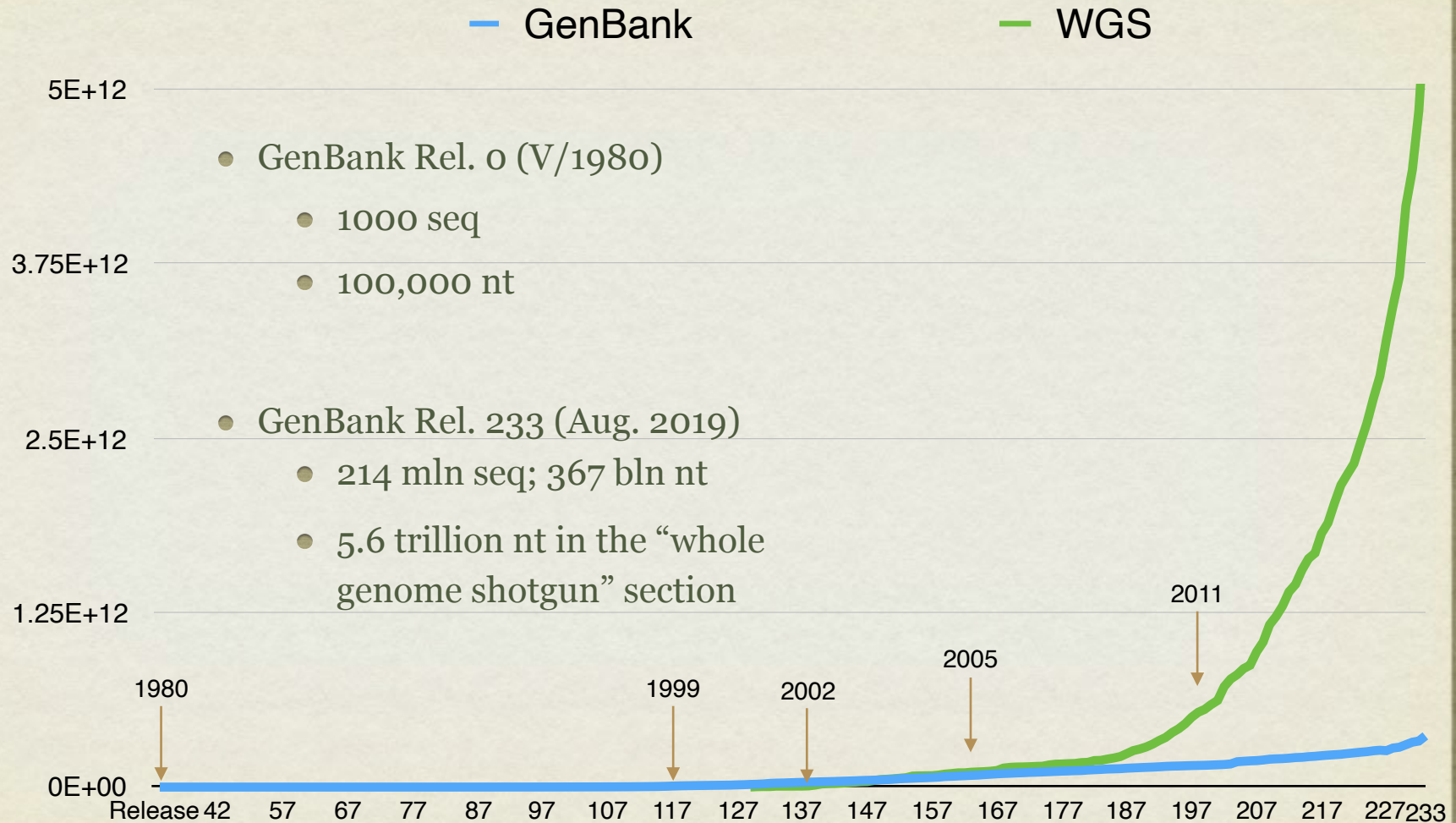
*Nature editorial, 1999*

Unfortunately, it's not a tidal wave, it's a tsunami!

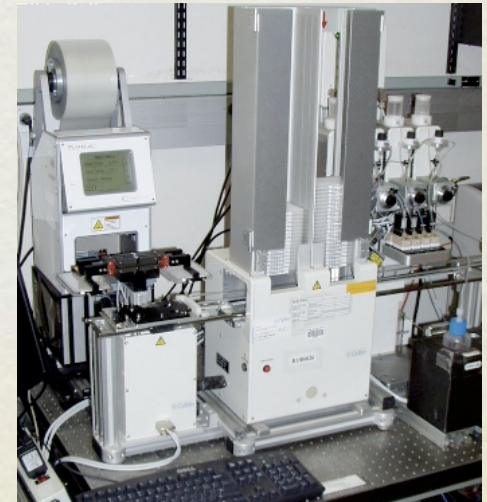# GROWTH OF BIOMEDICAL INFORMATION - GENBANK

**— GenBank**  **— WGS**

- GenBank Rel. 0 (V/1980)
  - 1000 seq
  - 100,000 nt

- GenBank Rel. 233 (Aug. 2019)
  - 214 mln seq; 367 bln nt
  - 5.6 trillion nt in the "whole genome shotgun" section

# TECHNOLOGY MEETS BIOLOGY

# IMPROVING TECHNOLOGY



Number of Humans Genomes Sequenced Over the Next 5 and 10 Years

# GETTING SEQUENCES

```
TGCATCGATCGTAGCTAGCTAGCGCATGCTAGCTAGCTAGCTAGCTACGATGCATCG
TGCATCGATCGATGCATGCTAGCTAGCTAGCTAGCATGCTAGCTAGCTAGCTATTGG
CGCTAGCTAGCATGCATGCATGCATCGATGCATCGATTATAAGCGCGATGACGTCAG
CGCGCGCATTATGCCGCGGCATGCTGCGCACACAGTACTATAGCATTAGTAAAAA
GGCCGCGTATATTTTACACGATAGTGCGGCGCGGCGCGTAGCTAGTGCTAGCTAGTC
TCCGGTTACACAGGTAGCTAGCTAGCTGCTAGCTAGCTGCTGCATGCATGCATTAGT
AGCTAGTGTAGCTAGCTAGCATGCTGCTAGCATGCAGCATGCATCGGGCGCGATGCT
GCTAGCGCTGCTAGCTAGCTAGCTAGCTAGGCGCTAATTATTTATTTTGGGGGGTTA
AAAAAAAAATTTCGCTGCTTATACCCCCCCCCACATGATGATCGTTAGTAGCTACT
AGCTCTCATCGCGCGGGGGGATGCTTAGCGTGGTGTGTGTGTGTGGTGTGTGTGGTC
CTATAATTAGTGCATCGGCGCATCGATGGCTAGTCGATCGATCGATTTTATATATCT
AAAGACCCCATCTCTCTCTCTTTTCCCTTCTCTCGCTAGCGGGCGGTACGATTTACC
GGCCGCGTATATTTTACACGATAGTGCGGCGCGGCGCGTAGCTAGTGCTAGCTAGTC
AGCTCTCATCGCGCGGGGGGATGCTTAGCGTGGTGTGTGTGTGTGGTGTGTGTGGTC
TGCATCGATCGATGCATGCTAGCTAGCTAGCTAGCATGCTAGCTAGCTAGCTATTGG
CTATAATTAGTGCATCGGCGCATCGATGGCTAGTCGATCGATCGATTTTATATATCT
CGCTAGCTAGCATGCATGCATGCATCGATGCATCGATTATAAGCGCGATGACGTCAG
TCCGGTTACACAGGTAGCTAGCTAGCTGCTAGCTAGCTGCTGCATGCATGCATTAGT
```
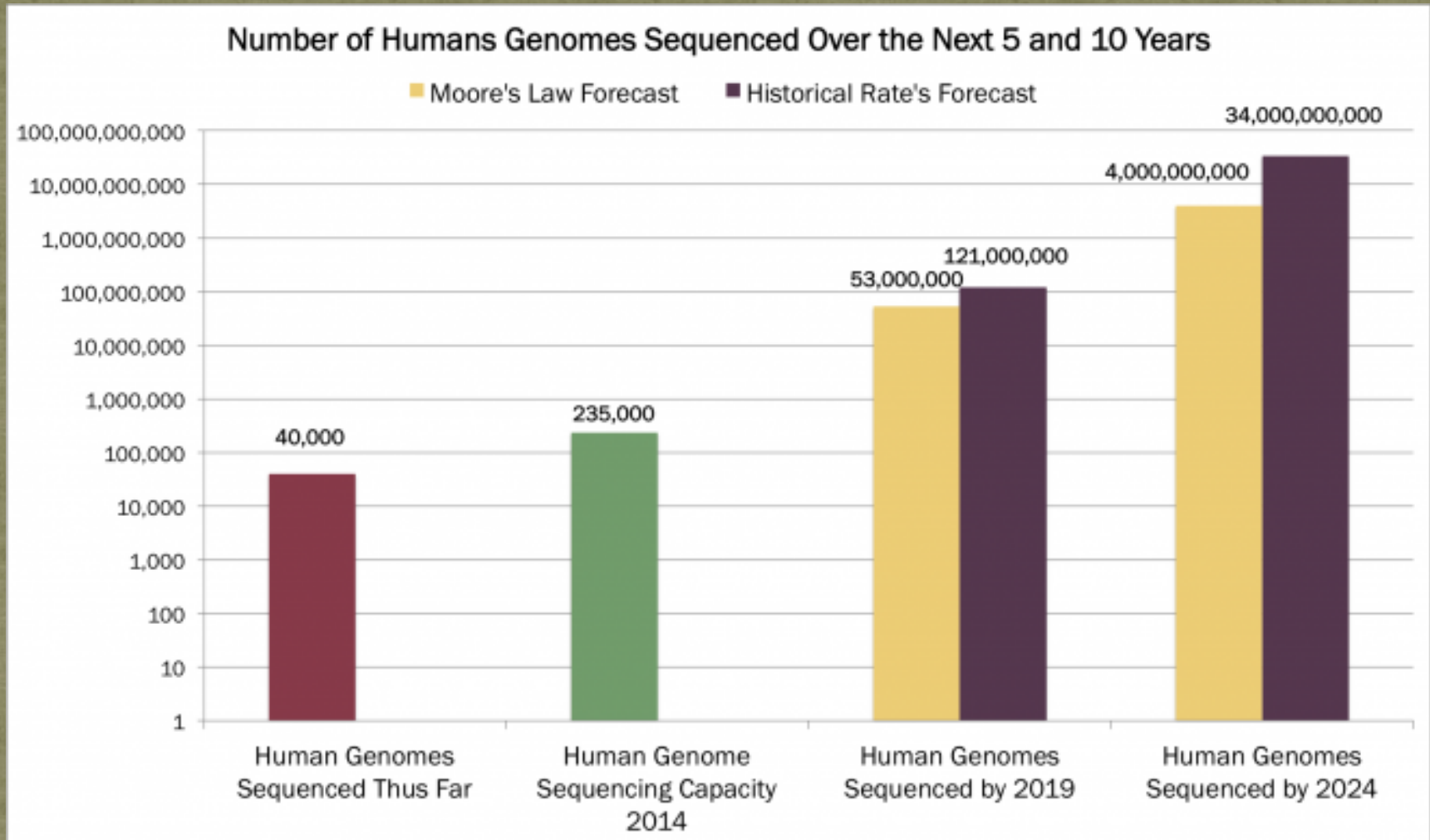
# READING ≠ UNDERSTANDING

Carmina qui quondam studio florente peregi, flebilis heu maestos cogor inire modos.

Ecce mihi lacerae dictant scribenda Camenae et ueris elegi fletibus orarigant.

**Boethius,** *Consolatio Philosophiae*

# READING ≠ UNDERSTANDING

We shall best understand the probable course of natural selection by taking the case of a country undergoing some physical change. If the country were open were open on its borders, new forms would certainly immigrate, and this also would bla, bla bla become extinct inhabitants.

Charles Darwin - *The Origin of Species*

# READING ≠ UNDERSTANDING

We shall best understand the probable course of natural selection by taking the case of a country undergoing some physical change. If the country were open were open on its borders, new forms would certainly immigrate, and this also would bla, bla bla become extinct inhabitants.

Charles Darwin - *The Origin of Species*

# CHALLENGE: HOW FROM THIS...

```
TGCATCGATCGTAGCTAGCTAGCGCATGCTAGCTAGCTAGCTAGCTACGATGCATCG
TGCATCGATCGATGCATGCTAGCTAGCTAGCTAGCATGCTAGCTAGCTAGCTATTGG
CGCTAGCTAGCATGCATGCATGCATCGATGCATCGATTATAAGCGCGATGACGTCAG
CGCGCGCATTATGCCGCGGCATGCTGCGCACACACAGTACTATAGCATTAGTAAAAA
GGCCGCGTATATTTTACACGATAGTGCGGCGCGGCGCGTAGCTAGTGCTAGCTAGTC
TCCGGTTACACAGGTAGCTAGCTAGCTGCTAGCTAGCTGCTGCATGCATGCATTAGT
AGCTAGTGTAGCTAGCTAGCATGCTGCTAGCATGCAGCATGCATCGGGCGCGATGCT
GCTAGCGCTGCTAGCTAGCTAGCTAGCTAGGCGCTAATTATTTATTTTGGGGGGTTA
AAAAAAAAAATTTCGCTGCTTATACCCCCCCCCACATGATGATCGTTAGTAGCTACT
AGCTCTCATCGCGCGGGGGGATGCTTAGCGTGGTGTGTGTGTGTGGTGTGTGTGGTC
CTATAATTAGTGCATCGGCGCATCGATGGCTAGTCGATCGATCGATTTTATATATCT
AAAGACCCCATCTCTCTCTCTCTTTTCCCTTCTCTCGCTAGCGGGCGGTACGATTTACC
GGCCGCGTATATTTTACACGATAGTGCGGCGCGGCGCGTAGCTAGTGCTAGCTAGTC
AGCTCTCATCGCGCGGGGGGATGCTTAGCGTGGTGTGTGTGTGTGGTGTGTGTGGTC
TGCATCGATCGATGCATGCTAGCTAGCTAGCTAGCATGCTAGCTAGCTAGCTATTGG
CTATAATTAGTGCATCGGCGCATCGATGGCTAGTCGATCGATCGATTTTATATATCT
CGCTAGCTAGCATGCATGCATGCATCGATGCATCGATTATAAGCGCGATGACGTCAG
TCCGGTTACACAGGTAGCTAGCTAGCTGCTAGCTAGCTGCTGCATGCATGCATTAGT
```
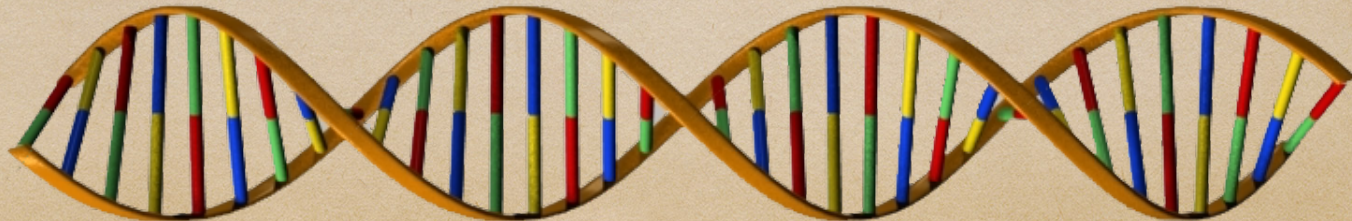
Infer this

"The double helix is indeed a remarkable molecule. Modern man is perhaps 50,000 years old, civilization has existed for scarcely 10,000 years and the United States for only just over 200 years; but DNA and RNA have been around for at least several billion years. All that time the double helix has been there, and active, and yet we are the first creatures on Earth to become aware of its existence."
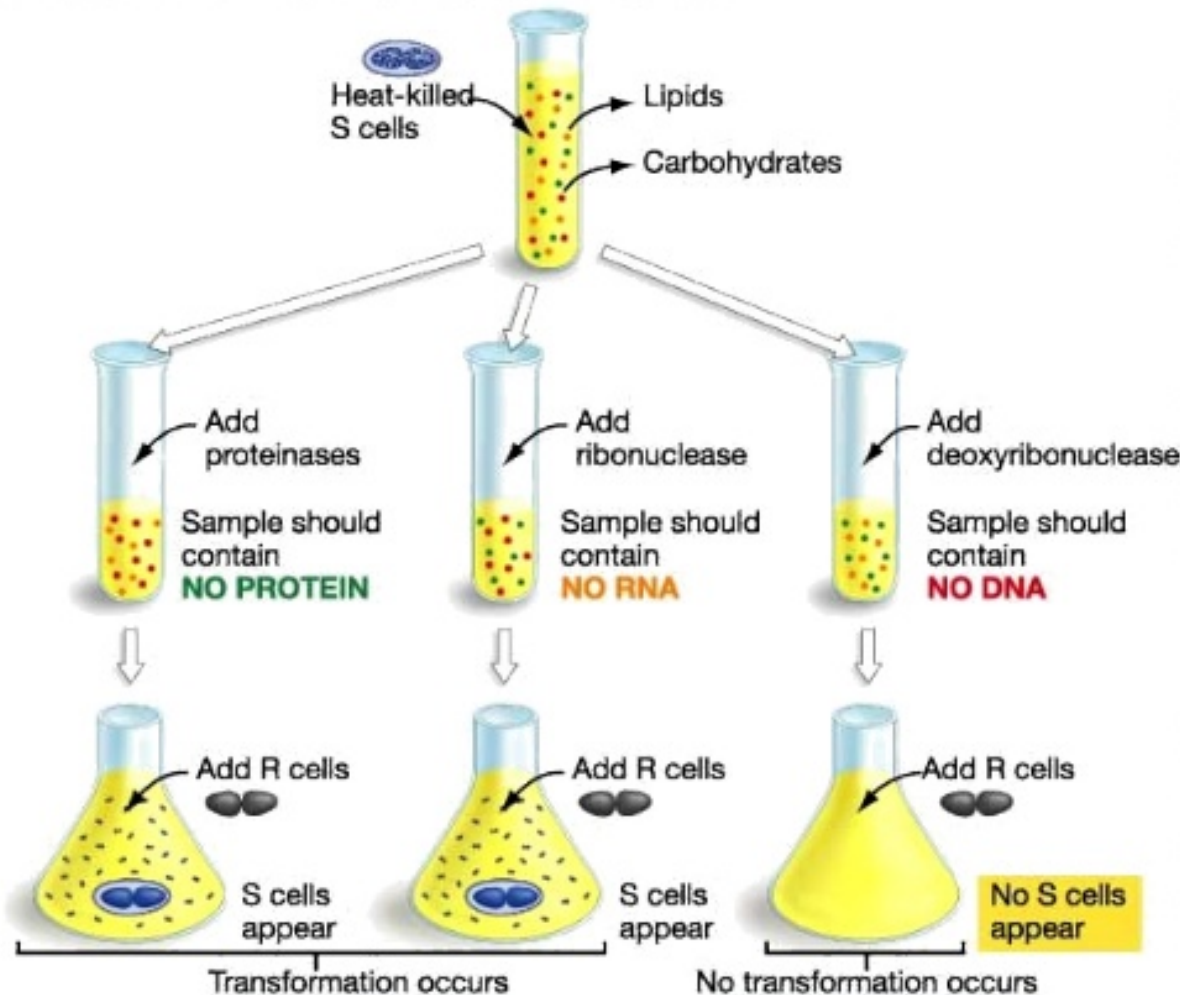
Francis Crick (1916–2004)

# DNA story

1870 Friedrich Miescher discovers DNA

1944 Oswald Avery proves that DNA is a genetic material

# DETERMINING THAT DNA IS THE HEREDITARY MATERIAL



Heat-killed S cells → Lipids
→ Carbohydrates

**1.** Remove the lipids and carbohydrates from a solution of heat-killed S cells. Proteins, RNA, and DNA remain.

Add proteinases — Sample should contain **NO PROTEIN**

Add ribonuclease — Sample should contain **NO RNA**

Add deoxyribonuclease — Sample should contain **NO DNA**

**2.** Subject the solution to treatments of enzymes to destroy either the proteins, RNA, or DNA.

Add R cells — S cells appear

Add R cells — S cells appear

Add R cells — No S cells appear

**3.** Add a small portion of each sample to a culture containing R cells. Observe whether transformation has occurred by testing for the presence of virulent S cells.

Transformation occurs

No transformation occurs

# DNA story

1953 James Watson and
Francis Crick discover
DNA structure

("Double Helix")

# Sequencing: beginnings

1964 Robert W. Holley determines nucleotide sequences (77 nt) of the yeast Alanine tRNA
J. Biol. Chem. 240: 2122-2128

1968 Ray Wu and A. Dale Kaiser sequenced 12 bases (!) of λ phage's 5' cohesive ends of its DNA, using radioactively labeled nucleotides and polyacrylamide gel electrophoresis
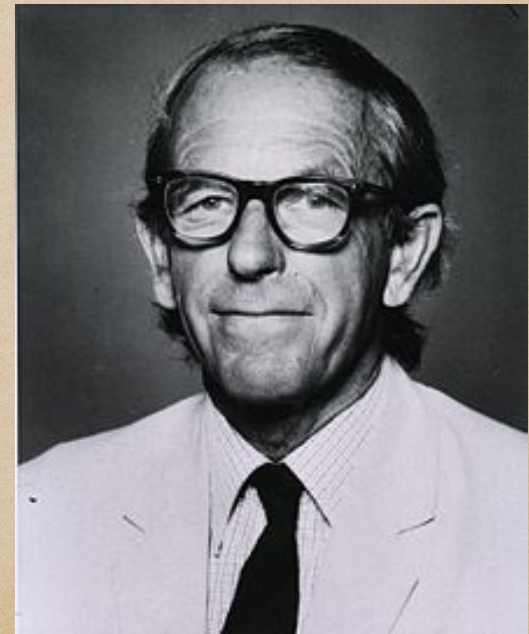J. Mol. Biol. 35: 523-537

# Sequencing:

## 1st generation sequencing

1977 – Allan Maxam and Walter Gilbert develop DNA sequencing method by chemical degradation

1977 Fred Sanger develops 2',3'-dideoxy chain termination method

# Chemical degradation sequencing

(Maxam & Gilbert)

DNA extraction

↓

DNA fragmentation

↓

Strand dissociation

↓

Radioactive labeling of 5' end

↓

Nucleotide-specific chemical reaction

↓

DNA cleavage at modified site



**The G reaction**

Molecule to be sequenced
(many copies)

Dimethyl sulfate

Piperidine

Figure 4.8  *Genomes 3* (© Garland Science 2007)

# Chemical degradation sequencing

(Maxam&Gilbert)

$5'$ $^{32}$PGCTACGTA $3'$

Cleavage at: A+G     G     C     C+T

**Four different reactions to detect four different nucleotides**

$^{32}$PGCT    $^{32}$PGCTAC   $^{32}$PG    $^{32}$PG
$^{32}$PGCTAC      $^{32}$PGCTA   $^{32}$PGC
$^{32}$PGCTACGT           $^{32}$PGCTA
$^{32}$PGCTACG

**Polyacrylamide gel electrophoresis can resolve single-stranded DNA molecules that differs in length by just one nucleotide and a sequence is read from an autoradiograph**

| | A+G | G | C | C+T | |
|---|---|---|---|---|---|
| 7 | ▬ | | | | A |
| 6 | | | | ▬ | T |
| 5 | ▬ | ▬ | | | G |
| 4 | | | ▬ | ▬ | C |
| 3 | ▬ | | | | A |
| 2 | | | | ▬ | T |
| 1 | | | ▬ | ▬ | C |

Sequencing Gel

# Chain termination DNA sequencing
## (Sanger)



**(A) Initiation of strand synthesis**

**(B) A dideoxynucleotide**

\* Position where the −OH of a dNTP is replaced by −H

**(C) Strand synthesis terminates when a ddNTP is added**

THE 'A' FAMILY

- use of DNA polymerase

- need for primers

- for each nucleotide a different analog

- similarly to M&G method separation of DNA fragments on polyacrylamide gel

- for each nucleotide a separate reaction

- sequence reading from an autoradiograph

Figure 4.2 *Genomes 3* (© Garland Science 2007)

# Sequencing: maturation

- 1983 – Marvin Caruthers developed a method to construct fragments of DNA of predetermined sequence from five to about 75 base pairs long. He and Leroy Hood invented instruments that could make such fragments automatically.

- 1983 – Kary Mullis invented the polymerase chain reaction (PCR) technique

- 1987 – ABI 370; first fully automated sequencing machine by Leroy Hood

- 1995 – Craig Venter uses whole-genome shotgun sequencing technique to determine complete genome of bacterium Haemophilus influenzae

- 2005 – introduction of GS-20 sequencing machine; first in the line of "Next Generation Sequencing", allowing hihg-throughput production

# Sequencing: maturation



**Chromatogram of a DNA sequence generated by ABI sequencing machine (https://www.dnalc.org/view/15912-Sequencing-DNA.html )**

# Sequencing: maturation



**SLAB GEL**

**CAPILLARY GEL**

Wells for samples

Back plate

Front plate

Capillary, 50–80 cm in length

Gel, <1 mm thick

Gel, 0.1 mm diameter

# Sequencing: maturation

- 1983 –  Marvin Caruthers developed a method to construct fragments of DNA of predetermined sequence from five to about 75 base pairs long. He and Leroy Hood invented instruments that could make such fragments automatically.

- 1983 – Kary Mullis invented the polymerase chain reaction (PCR) technique

- 1987 – ABI 370; first fully automated sequencing machine

- 1995 – Craig Venter uses whole-genome shotgun sequencing technique to determine complete genome of bacterium Haemophilus influenzae

- 2005 – introduction of GS20 sequencing machine (454 Lige Sciences); first in the line of "Next Generation Sequencing"

# Next Generation Sequencing

- Massive parallelization of the sequencing process

- Relatively short reads

- Different approaches from improving Sanger's technique to direct "observation" of DNA through a microscope

- Attempts to sequence single molecules without amplification step

# Next Generation Sequencing

- 1 – Pyrosequencing (Roche454)

- 2 – Ion torrent (Thermo Fisher)

- 3 – Illumina

# NGS – pyrosequencing

## library preparation

# NGS - pyrosequencing

## sample preparation

# NGS – pyrosequencing

- After the emulsion PCR has been performed, the oil is removed, and the beads are put into a "picotiter" plate. Each well is just big enough to hold a single bead.

- The pyrosequencing enzymes are attached to much smaller beads, which are then added to each well.

- The plate is then repeatedly washed with the each of the four dNTPs, plus other necessary reagents, in a repeating cycle.

- The plate is coupled to a fiber optic chip. A CCD camera records the light flashes from each well.

# NGS ~ pyrosequencing

Extension with individual dNTPs gives a readout. The readout is recorded by a detector that measures position of light flashes and intensity of light flashes.

# NGS - pyrosequencing



Example of a Flowgram

# NGS ~ion torrent

◆ Ten times faster workflow than other NGS systems

◆ ~2 hour sequencing runs (real-time detection of sequence extension)

◆ Batch sample preparation (six samples in six hours)

◆ Capable of six samples/day on two PGM Systems

https://www.youtube.com/watch?v=DyijNS0LWBY

# NGS –ion torrent
# Simple Natural Chemistry

## Sequencing by synthesis



Eliminate source of sequencing errors:
- Modified bases
- Fluorescent bases
- Laser detection
- Enzymatic amplification cascades

Eliminate source of read length limitations:
- Unnatural bases
- Faulty synthesis
- Slow cycle time

# NGS –ion torrent

# Fast Direct Detection



**DNA → Ions → Sequence**

- Nucleotides flow sequentially over Ion semiconductor chip
- One sensor per well per sequencing reaction
- Direct detection of natural DNA extension
- Millions of sequencing reactions per chip
- Fast cycle time, real time detection

# NGS -ION TORRENT

Four nucleotides flow sequentially

NGS –ION TORRENT

Base call

NGS –ION TORRENT

NGS –ION TORRENT

Base call

ATCGTGTTTTAGGGTCCCCGGGGGTTAAAA…
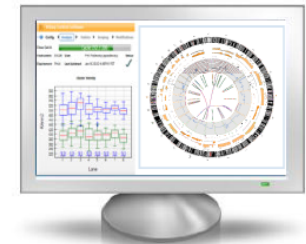
# NGS – Illumina

## Workflow



| SAMPLE PREP | cBot CLUSTER GENERATION | Genome Analyzer SEQUENCING | DATA PROCESSING & ANALYSIS |

# NGS ~ Illumina
## The flow cell ~ a core component

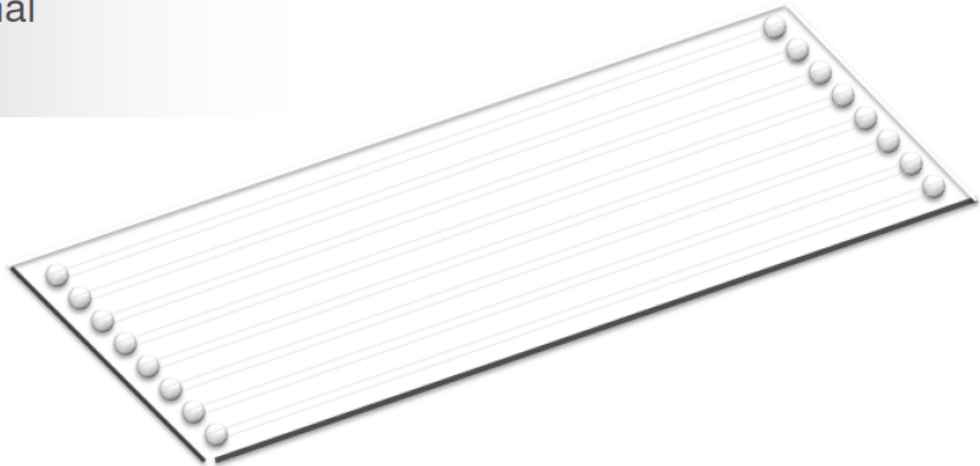**EVERYTHING EXCEPT SAMPLE PREPARATION IS COMPLETED ON THE FLOW CELL**

template annealing (1 - 96 samples)

template amplification

sequencing primer hybridization

Sequencing-by-synthesis reaction

generation of fluorescent signal

# NGS – Illumina
# Preparation of template

# NGS – Illumina
## Preparation of template



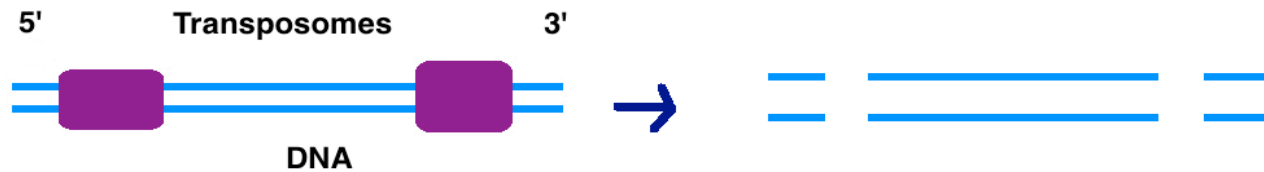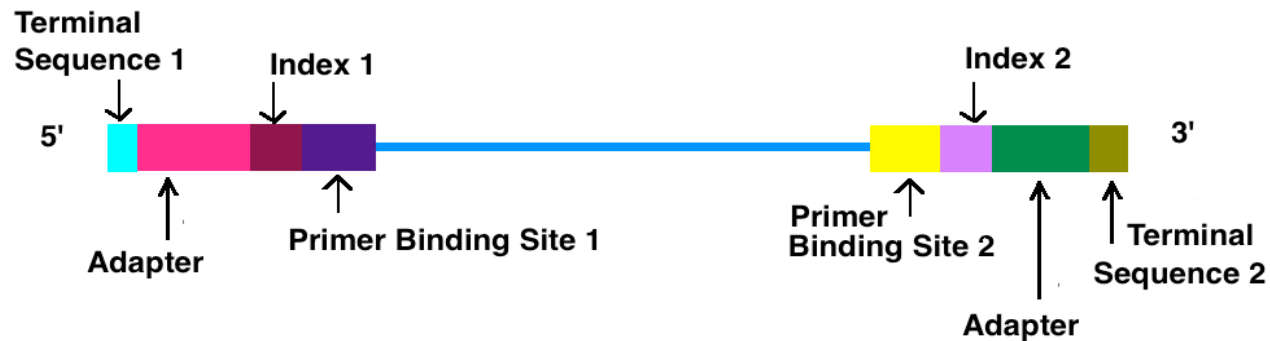Another sheering method: transposomes – enzymes for DNA cleavage

# NGS – Illumina
# The flow cell is mounted on the cBot

**AUTOMATICALLY**

loads library into the lanes of the flow cell

amplifies templates

anneals sequencing primer to templates

**FEATURES**
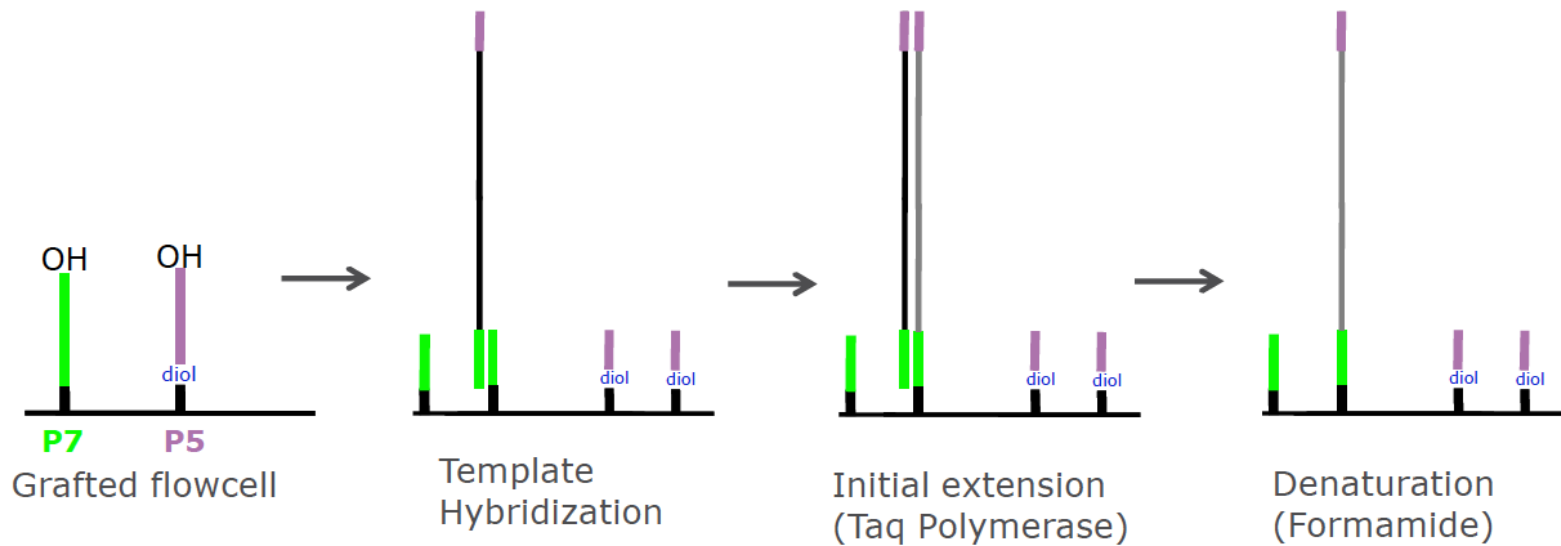
intervention-free clonal amplification in 4 hours

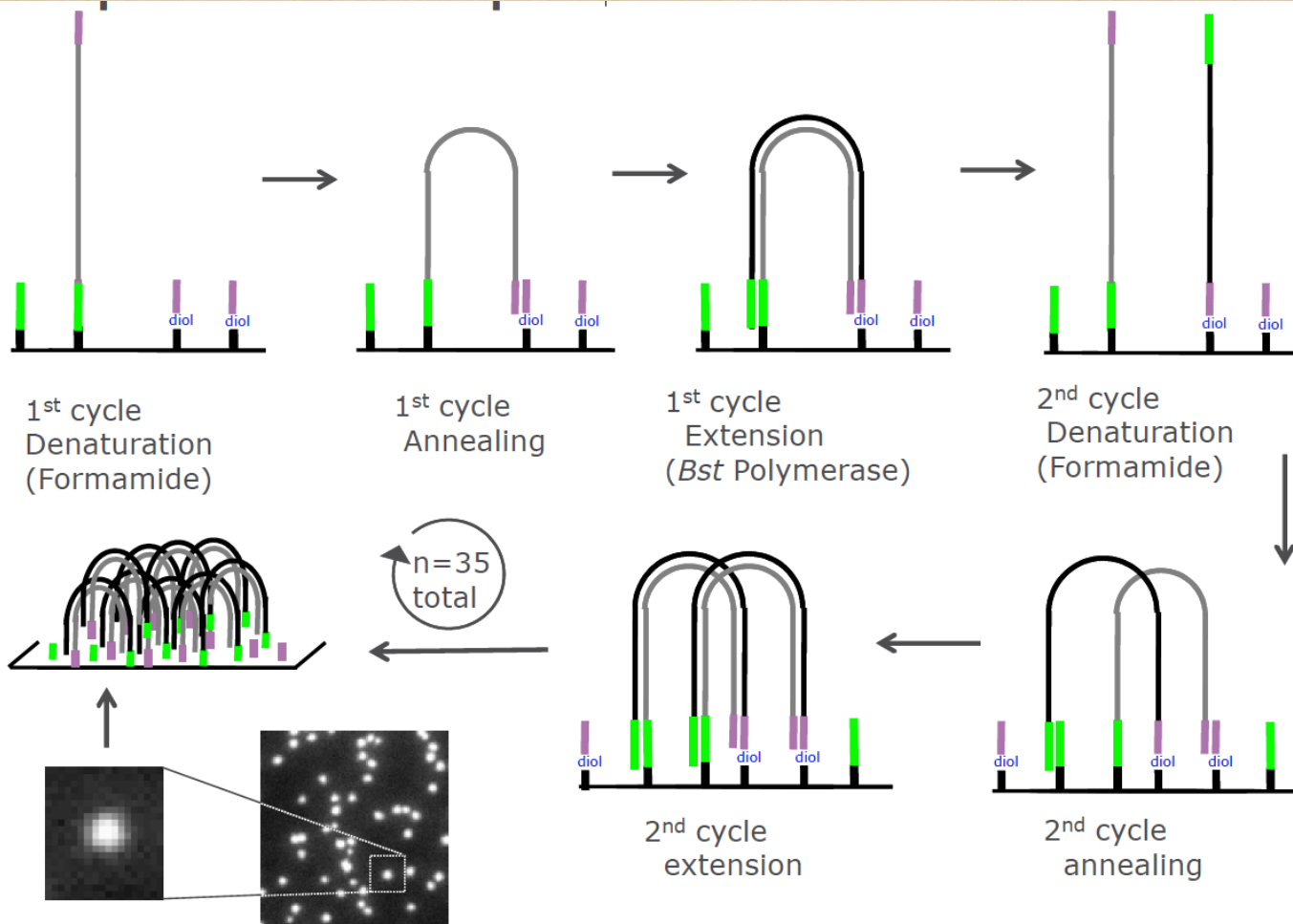simple touch screen operation

# NGS – Illumina
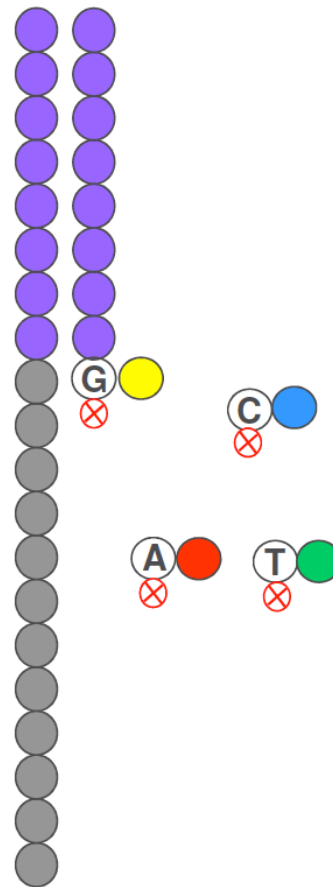## Hybridization of template
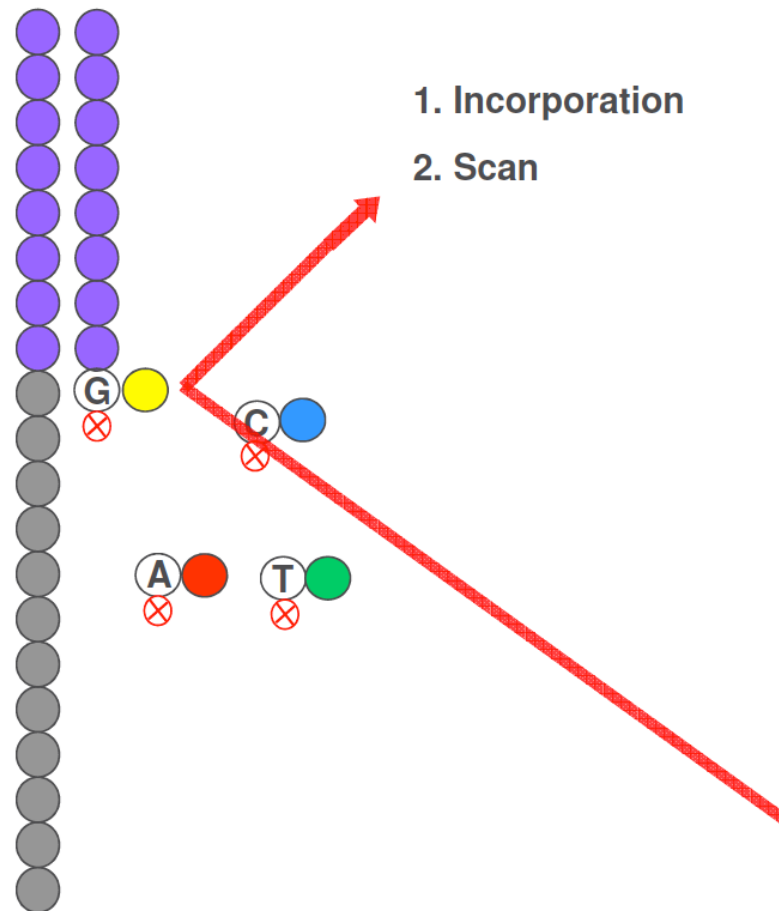
# NGS – Illumina
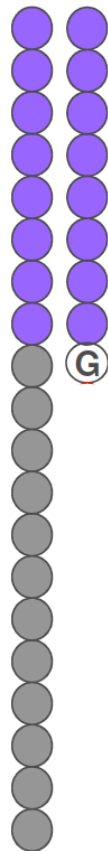
# Incorporation



1. Incorporation

# NGS – Illumina

## Scanning

# NGS ~ Illumina
# Cleavage



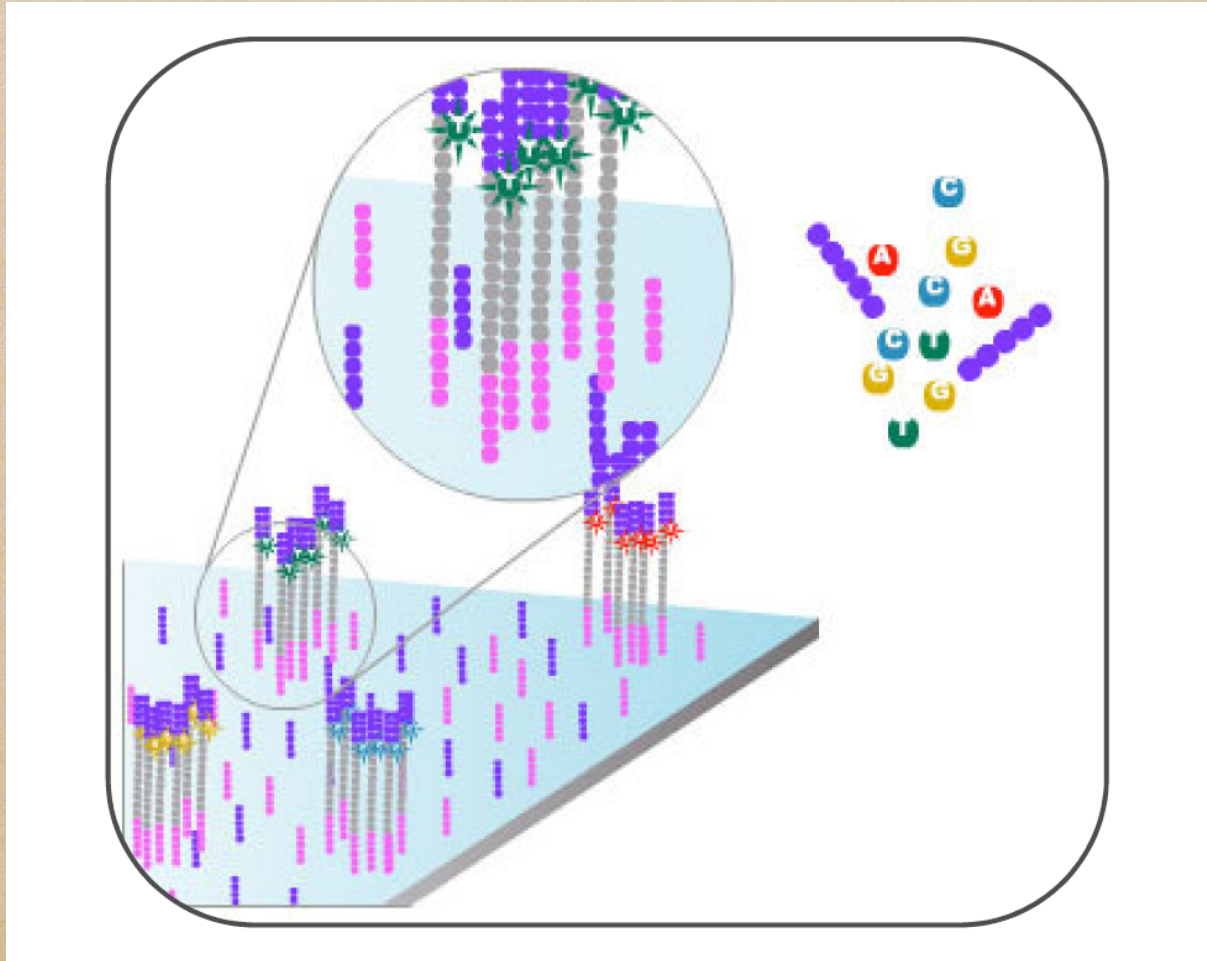1. Incorporation
2. Scan
3. Cleavage

# NGS – Illumina
## Millions of clusters are sequenced in parallel

# NGS – Illumina
## A picture is taken every time a new base is added

# NGS – Illumina
## The flow cell is mounted on the sequencer



CCD camera collects laser-excited fluorescence

sequencing reagents pass through the 8 lanes inside the flow cell

sequencing reaction is temperature controlled

# Third Generation Sequencing

1 – Pacific Bioscience (PacBio)

2 – MinIon (Oxford NanoTechnologies)

# PacBio



https://www.youtube.com/watch?v=_B_cUZ8hSYU

# MinIon

# MinIon: Sequencing using nanopores

- Nanopores as polymer sensors.

- The idea emerged in early 1990s.

- Fundamental work done by David Deamer and Daniel Branton in collaboration with John Kasianowicz. (PNAS 1996 146:13770-13773)
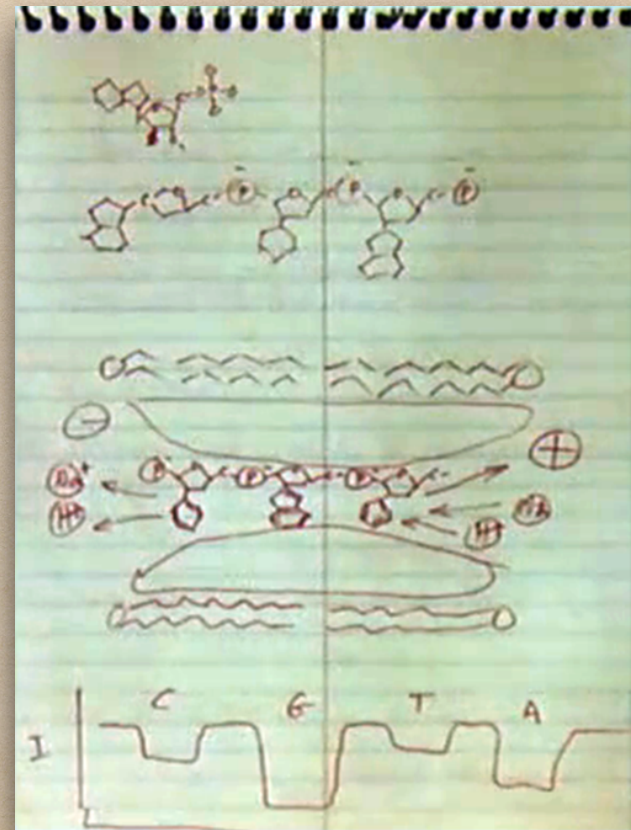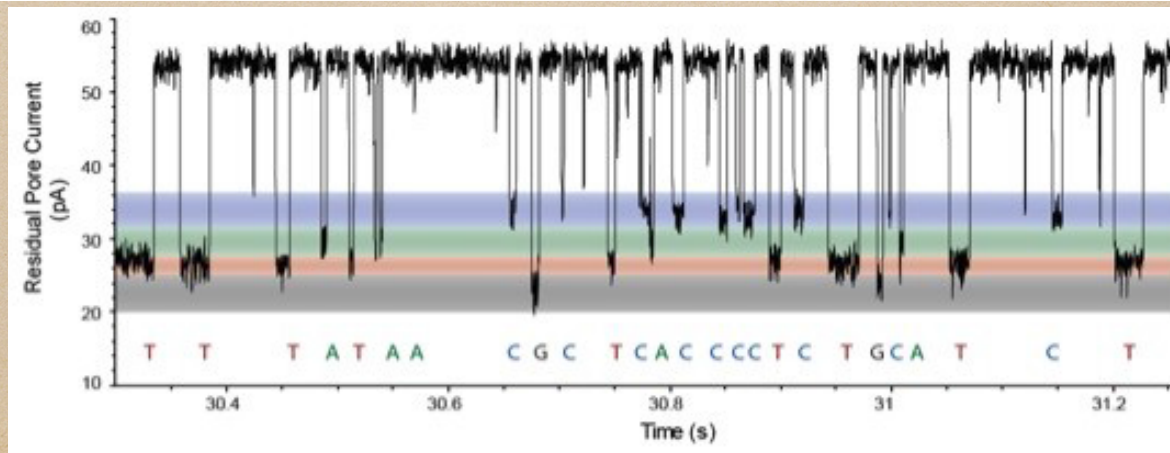
- Biologicaly relevant experiments – since 2010.

Current flow

current

Residual Pore Current (pA)

T  T     T  A T A A     C G C  T C A C  C C C T C  T  G C A  T       C     T

Time (s)

# MinION basics

- Synthetic membrane

- Nanopore (2) is created by modified protein pores: α-hemolysin, CsgG from E.coli

- Non-destructive motor protein (1) (actually serves as a break)

# MinION basics

https://nanoporetech.com/science-technology/introduction-to-nanopore-sensing/introduction-to-nanopore-sensing
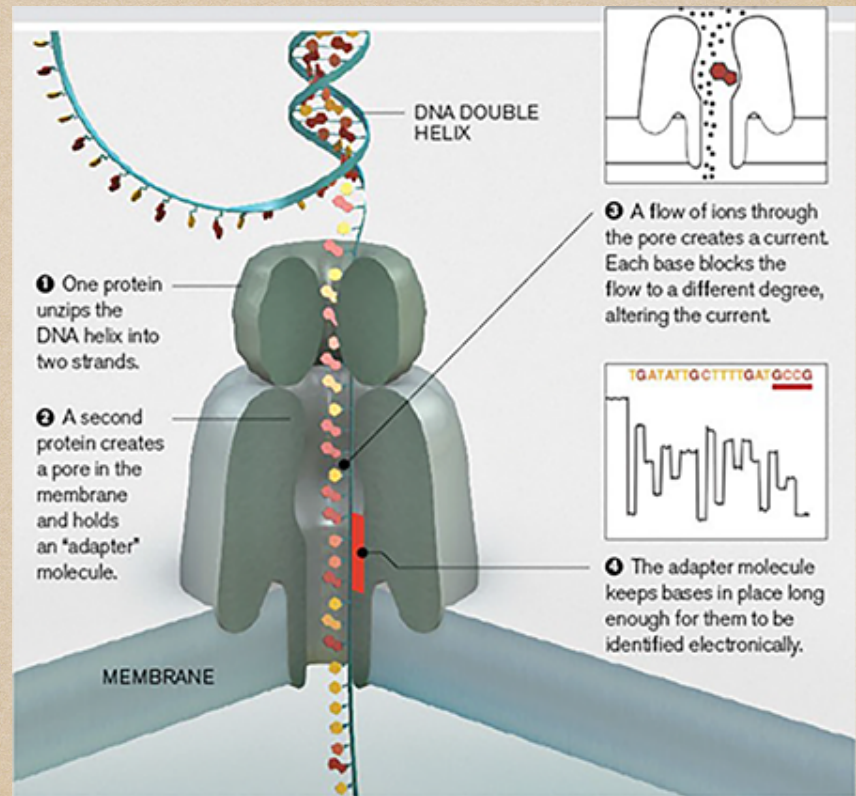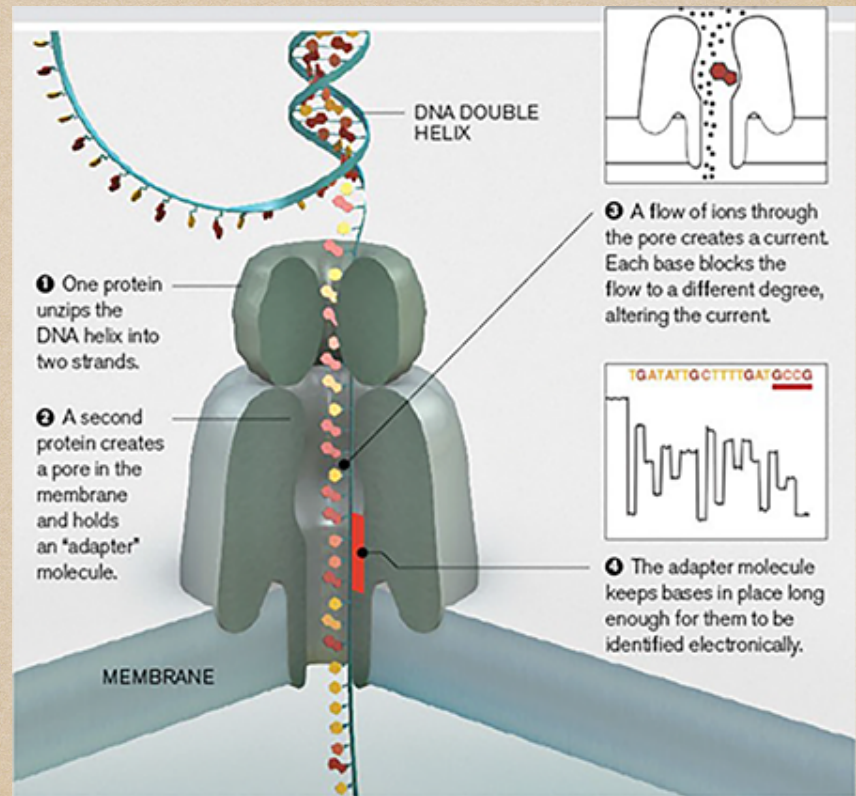
- 512 channels (pores) per flow cell. Usually about 90% are working.

- Read length: over a million of bp

- Read speed: 8 bases to 20 bases/sec

- Run time: max 48 hours

- Error rate = 5-10 %

- Sequence yield per flow cell: 15 Gb

- Complex algorithm for base calling using neural network approach



DNA DOUBLE HELIX

❶ One protein unzips the DNA helix into two strands.

❷ A second protein creates a pore in the membrane and holds an "adapter" molecule.

❸ A flow of ions through the pore creates a current. Each base blocks the flow to a different degree, altering the current.

TGATATTGCTTTTGATGCCG

❹ The adapter molecule keeps bases in place long enough for them to be identified electronically.

MEMBRANE

Easy, standard template preparation

Time of library preparation:
1D – about ten minutes
2D – up to two hours

Cost of a single run:
reagents $200
flow cell   $1000

# MinION dataflow

## MinION – the device

Nanopore sensing is carried out on the sensor chip, contained in the flow cell inside the MinION device. Data is processed by an Application-Specific Integrated Circuit (ASIC) also in the flow cell and processed in real time by the MinKNOW software

## MinKNOW – the software

MinKNOW is the software that controls the MinION. It carries out several core data tasks and can be used to change experimental workflows or parameters. MinKNOW runs on the user's computer.

## ALBACORE – base calling

Albacore is a command-line (some programming skills are required) base-calling software, developed for MinIon and accounts for specific sequencing errors

# Numerous applications explored by MinION Access Program (MAP)

- Genomic DNA sequencing

- Metagenomic analysis

- Medical diagnostics (in development)

- Species identification in the field

- Splice variants identification

- Virus detection in the field

- Sequencing in space, etc … ☺

# Comparison table

|  | 454 | Illumina | Ion Torrent | PacBio | MinIon |
|---|---|---|---|---|---|
| **Method** all sequence by synthesis | Pyrosequencing: pyrophosphates detection by chemoluminicent reaction (luciferase enzyme). Detector: CCD camera | Bridge amplification; detection of fluorescently labeled nucleotides. Detector: CCD camera | Ion semiconductor: label free detection of released protons. Detector: ion sensor | Single-molecule in real-time: detection of fluorescently labeled cleaved pyrophosphates. Detector: ZMW camera (sensitive!) | Nanopores: modified pore proteins detect current change when different nucleotides pass the pore. Detector: ASIC -measures ionic current flow |

454: https://www.youtube.com/watch?v=nFfgWGFe0aA

Illumina: https://www.youtube.com/watch?v=fCd6B5HRaZ8

Ion Torrent: https://www.youtube.com/watch?v=WYBzbxIfuKs

PacBio: https://www.youtube.com/watch?v=_B_cUZ8hSYU

MioIon: https://nanoporetech.com/how-it-works

# Comparison table

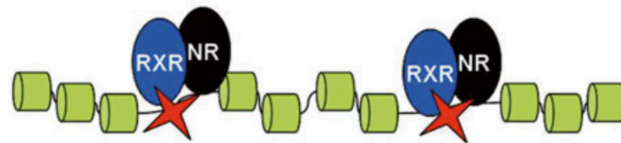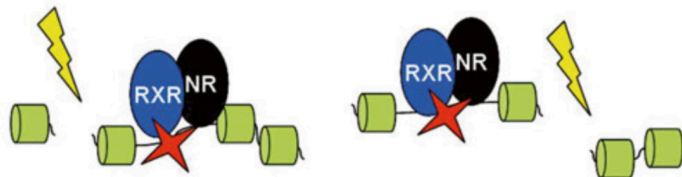|  | 454 | Illumina | Ion Torrent | PacBio | MinIon |
|---|---|---|---|---|---|
| **Read length** | 700 bp | 50-250 bp | 200 bp | 3000-15000 bp | 500-100000 |
| **Reads per run** | 1 million | up to 3 billion | up to 5 million | 35000-75000 | 30-400 million |
| **Time per run** | 24 hours | 1-10 days | 2 hours | 30 min – 2 hours | 6-48 hours |
| **Cost per million bases** | 10$ | 0.05-0.15$ | 1$ | 2$ | 2$ |
| **Machine cost** |  | 120.000-650.000$ | 80.000$ | 695.000$ | 1500$ |
| **Error rate** | 0.1-1% | 0.5-1% | 1-2% | 12% | 5-10% |

1.USA : 1,000,000 Genome @ Veterans Project & All of Us Reserach Program | 2.United Kingdom: 100,000 Genomes project | 3.China: 100,000 genomes project | 4.Saudi Arabia: 100,000 Genome Project (Saudi Genome) | 5.United Arab Emirates: 3,000,000 Genome Project | 6.Estonia: 100,000 Genome Project (Personalized Medicine Program) 7.France: 100,000 Genome Project (French Plan for Genomic Medicine 2025) | 8.Australia: 100,000 Genome Project (The Australian Genomics Health Futures Mission) 9.Japan: 2,000 Genome Project (Initiatives on Rare and Undiagnosed Diseases)
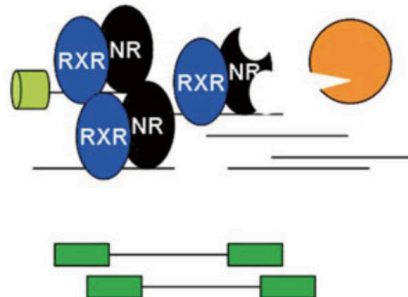
# chip-seq experiments



1. Chromatin crosslinking

2. Chromatin shearing

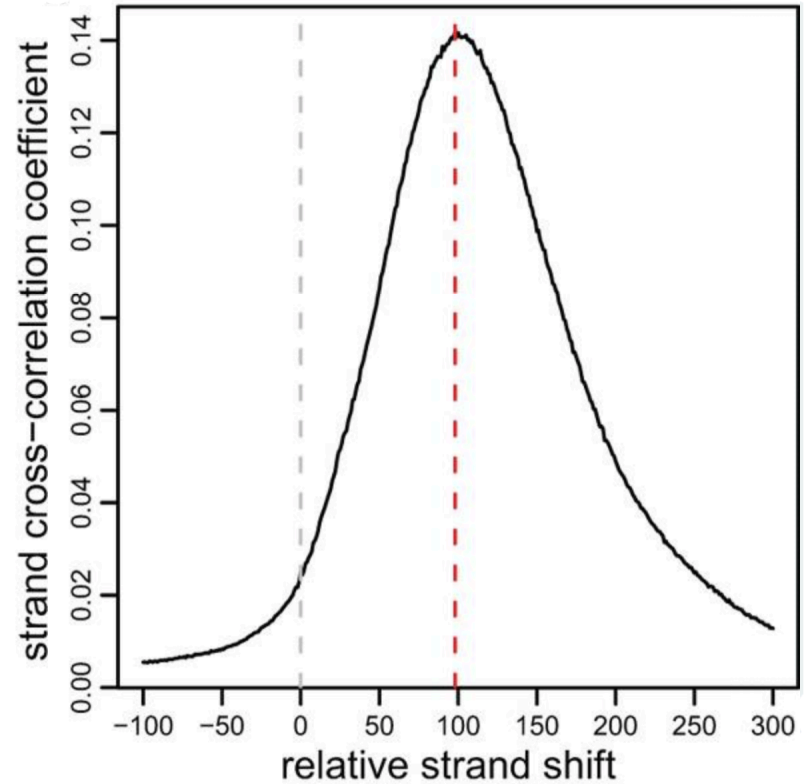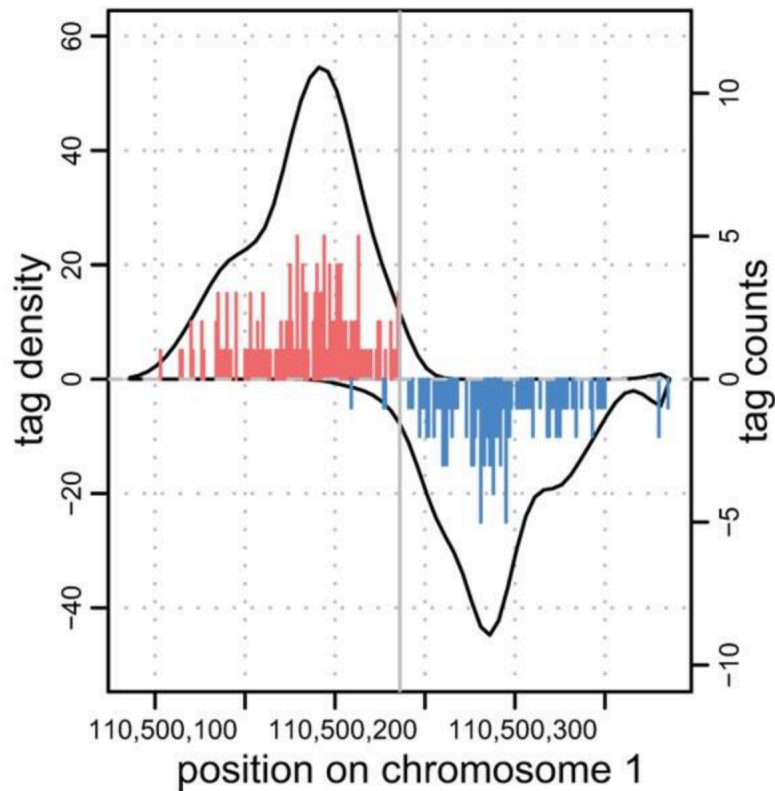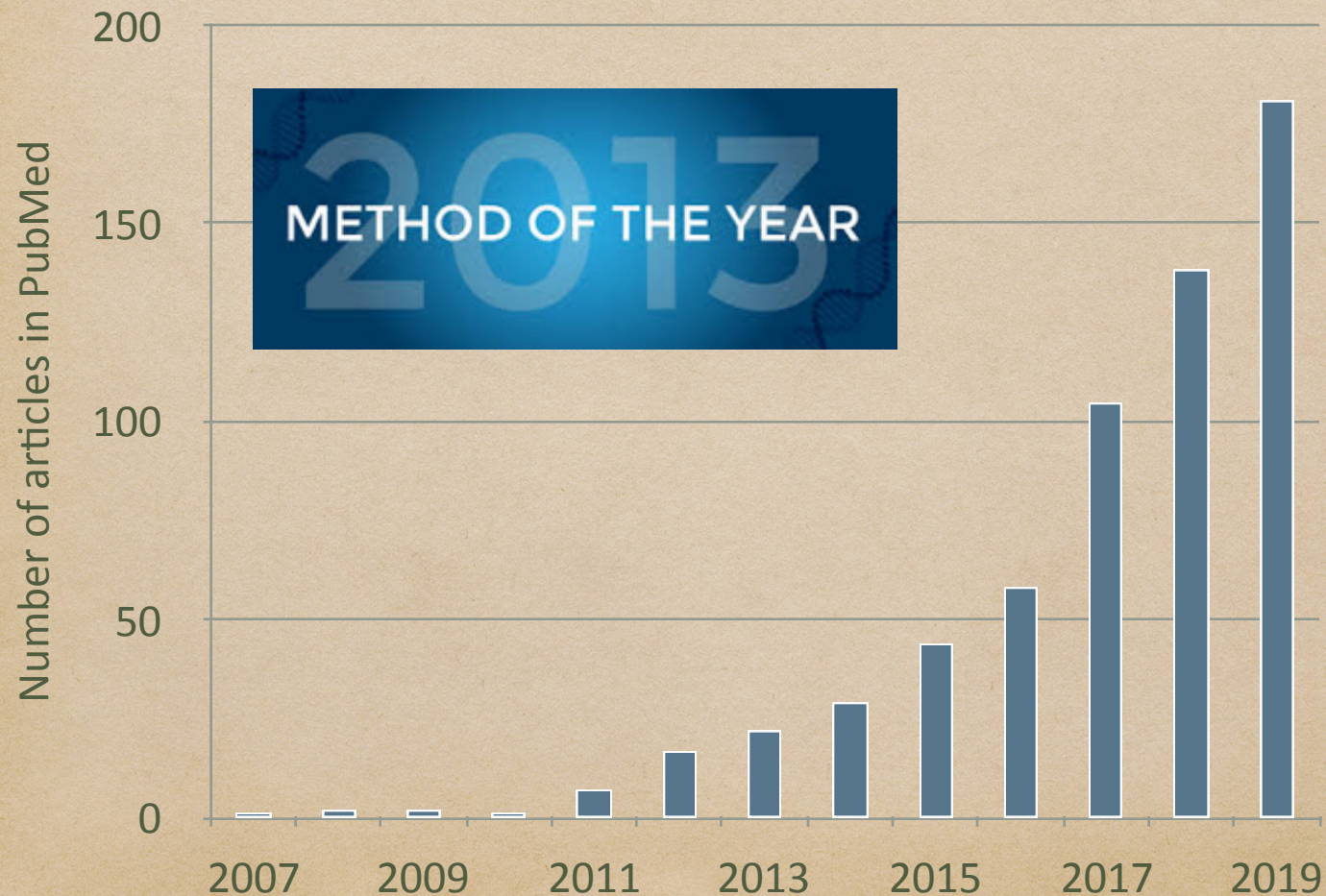3. Immunoprecipitation

4. De-crosslinking

5. Library preparation

# chip-seq experiments

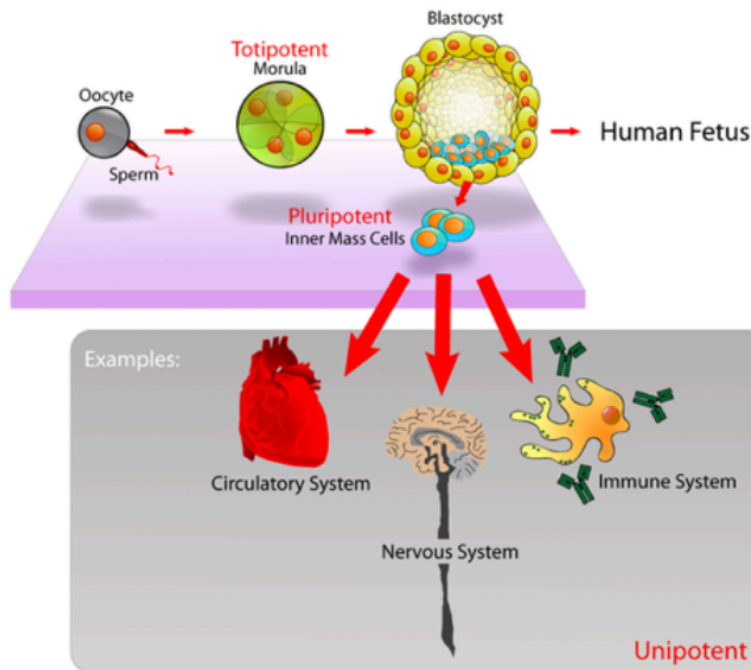# Single-cell sequencing

# Single-cell sequencing
## applications

- Developmental Biology
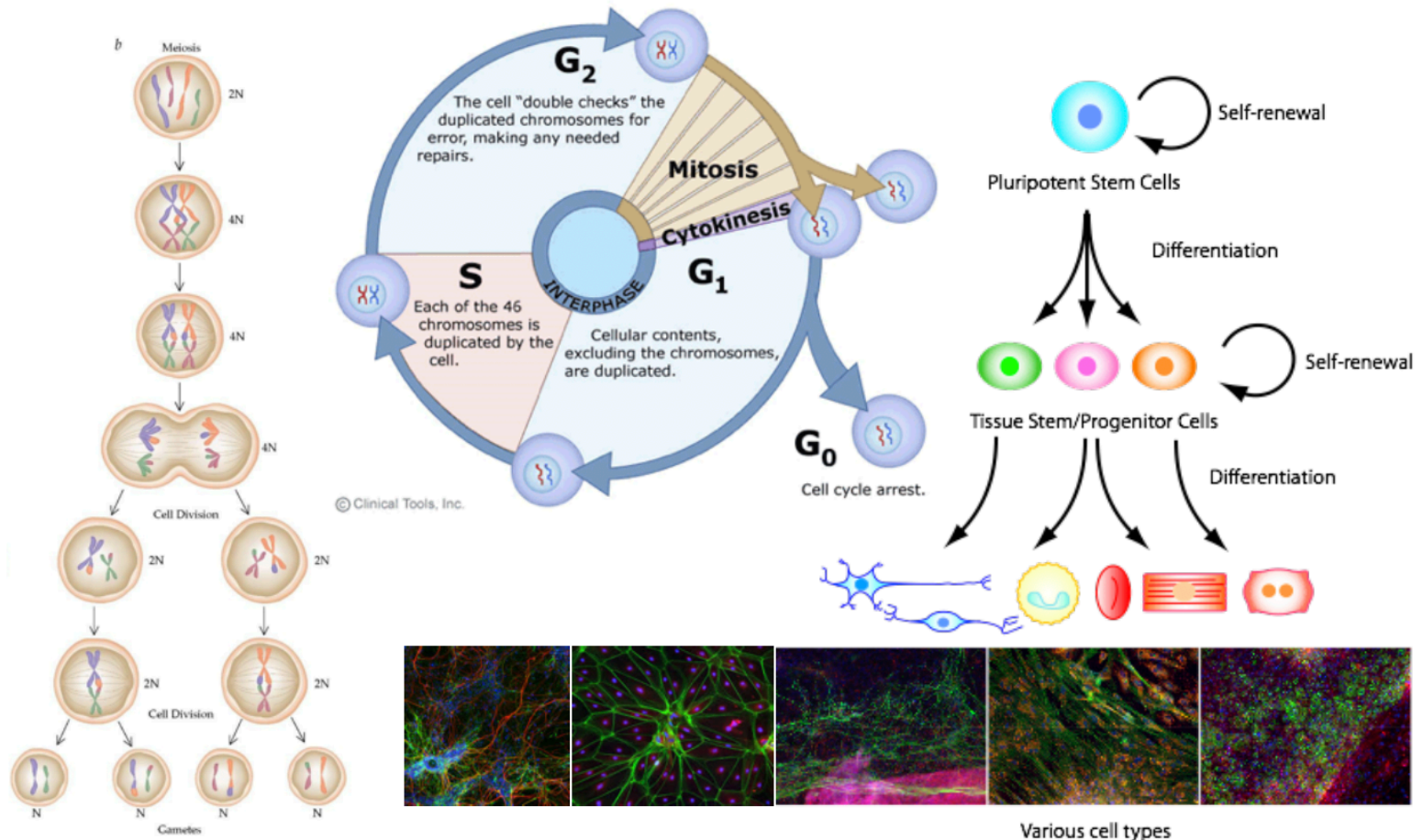
- Cancer Biology

- Microbiology

- Neurology

# Developmental Biology

How do animals grow and develop from a single cell?

# Developmental Biology

# Developmental Biology

- We need single-cell resolution to:

  - Discover more complicated mechanisms in cellular development

  - Confirm the distinct gene expression signatures across different cell types

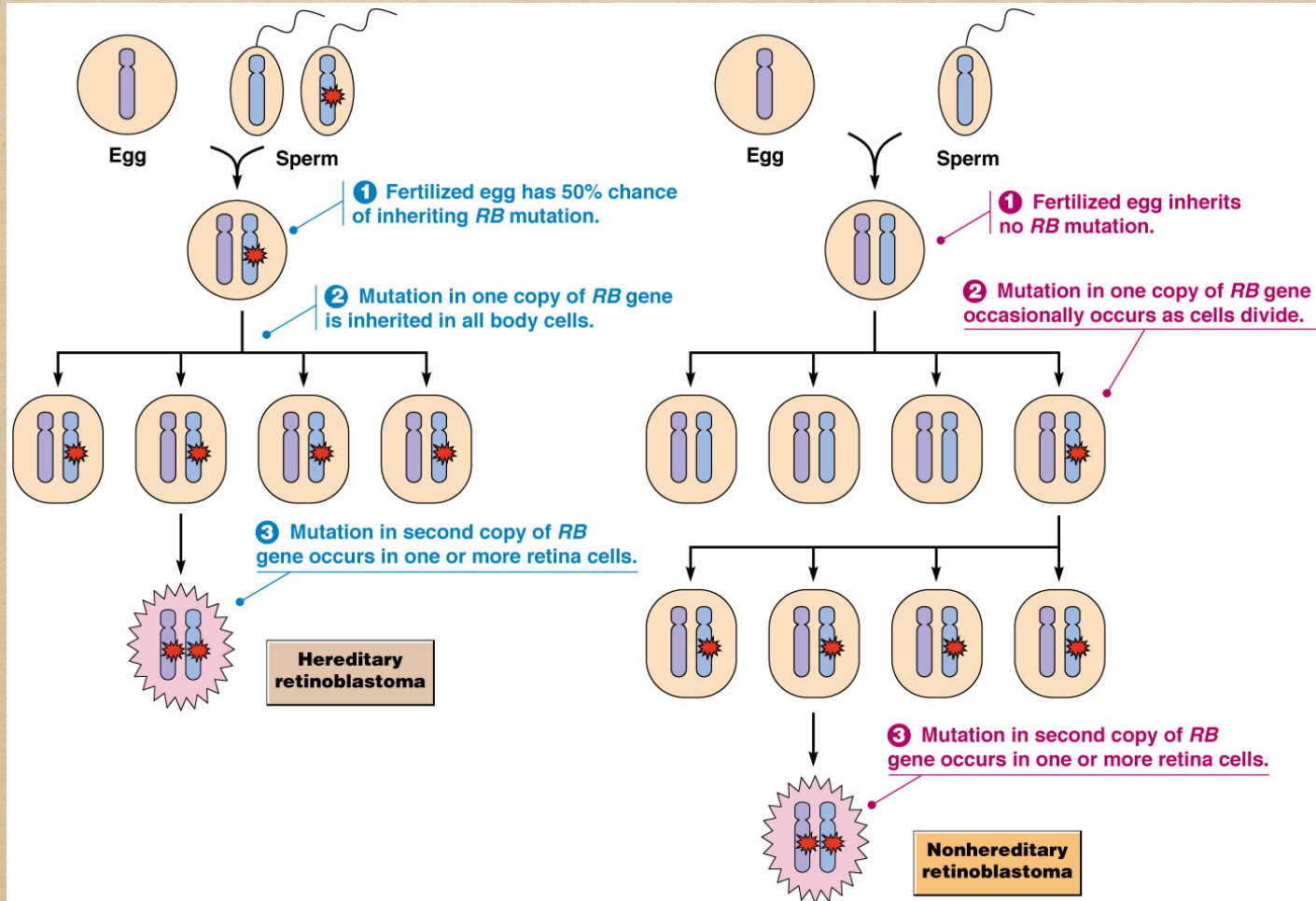  - Identify functional differences among the same cell cell type

# Single-cell sequencing
## applications

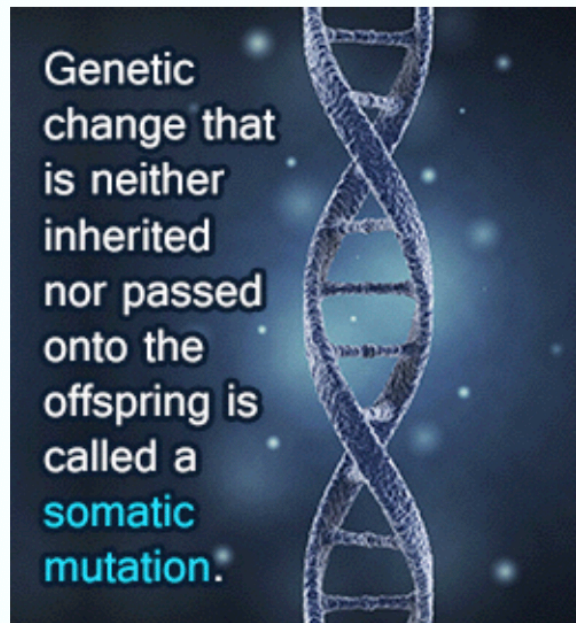- Developmental Biology

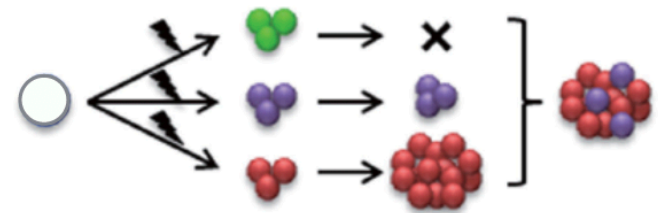- Cancer Biology
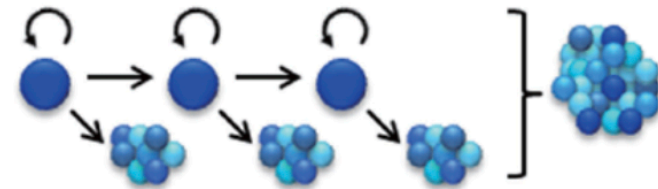
- Microbiology

- Neurology

# Cancer Biology



① Fertilized egg has 50% chance of inheriting *RB* mutation.

② Mutation in one copy of *RB* gene is inherited in all body cells.

③ Mutation in second copy of *RB* gene occurs in one or more retina cells.

**Hereditary retinoblastoma**

① Fertilized egg inherits no *RB* mutation.

② Mutation in one copy of *RB* gene occasionally occurs as cells divide.

③ Mutation in second copy of *RB* gene occurs in one or more retina cells.

**Nonhereditary retinoblastoma**

# Cancer Biology

Tumors are composed of genetically and phenotypically **heterogeneous** clones

Genetic change that is neither inherited nor passed onto the offspring is called a somatic mutation.
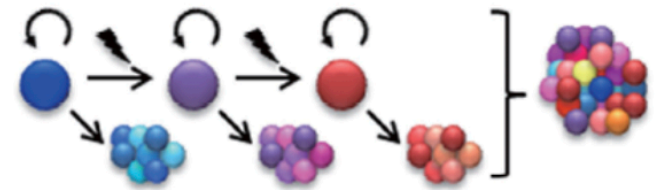
A   Stochastic model

B   Cancer stem cell model

C   Combination model

Major genetic/epigenetic events

# Cancer Biology



Cancer stem cell specific therapy → Tumor regression

Conventional cancer therapy → Tumor relapse

A Stochastic model

B Cancer stem cell model

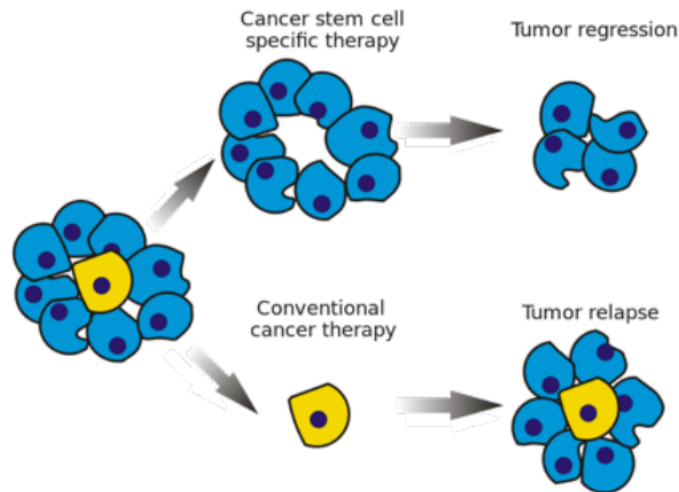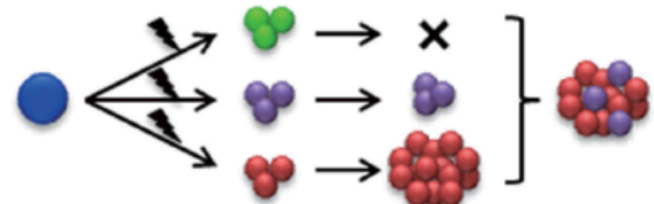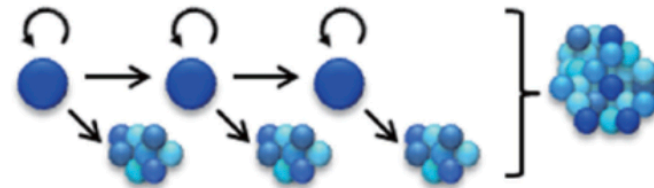C Combination model

Major genetic/epigenetic events

Deep (bulk) sequencing can only capture 1% of the cell population (excluding some types such as circulating tumor cells).

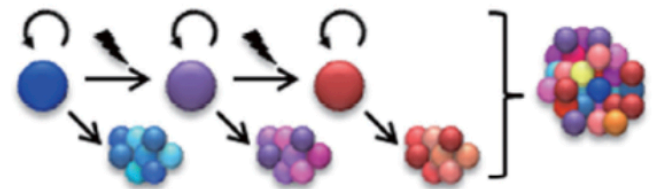http://www.thetcr.org/article/viewFile/1415/html/10439

# Caner Biology

- We need single-cell resolution to:

  - Find evidence for models of cancer

  - Infer timing of mutations and the drivers
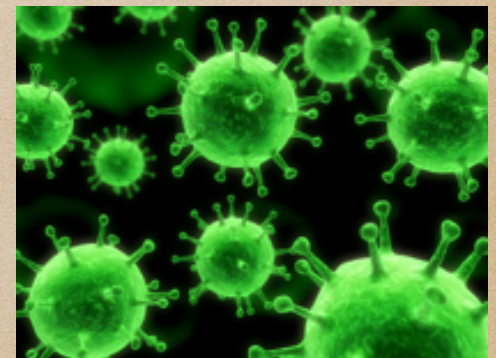
  - Evaluate effectiveness of targeted therapy

# Single-cell sequencing
## applications

- Developmental Biology

- Cancer Biology

- Microbiology

- Neurology

# Microbiology

# Microbiology

- We need single-cell resolution to:

  - Discover low-abundance species that are difficult to culture in vitro

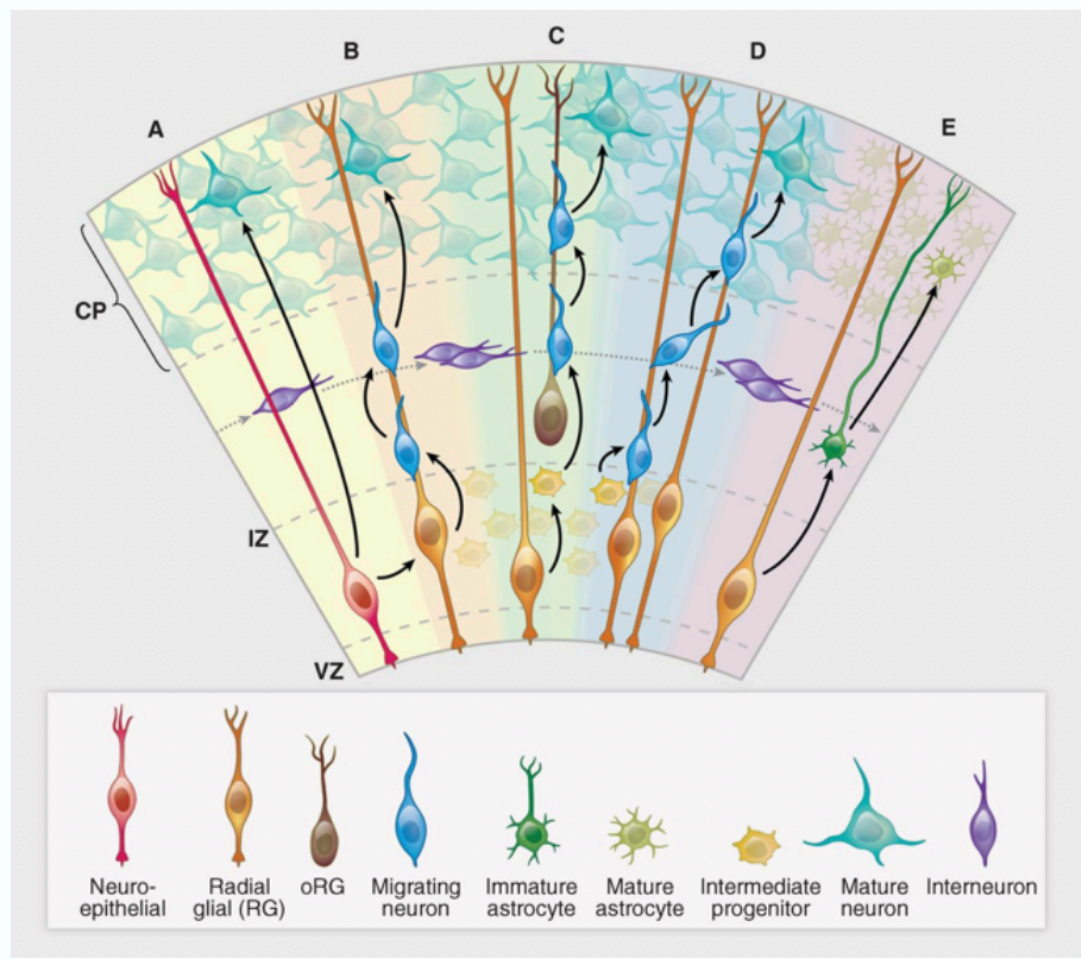  - Monitor transcriptional gene activation mechanisms for functional annotation

# Single-cell sequencing
## applications

- Developmental Biology

- Cancer Biology

- Microbiology

- Neurology
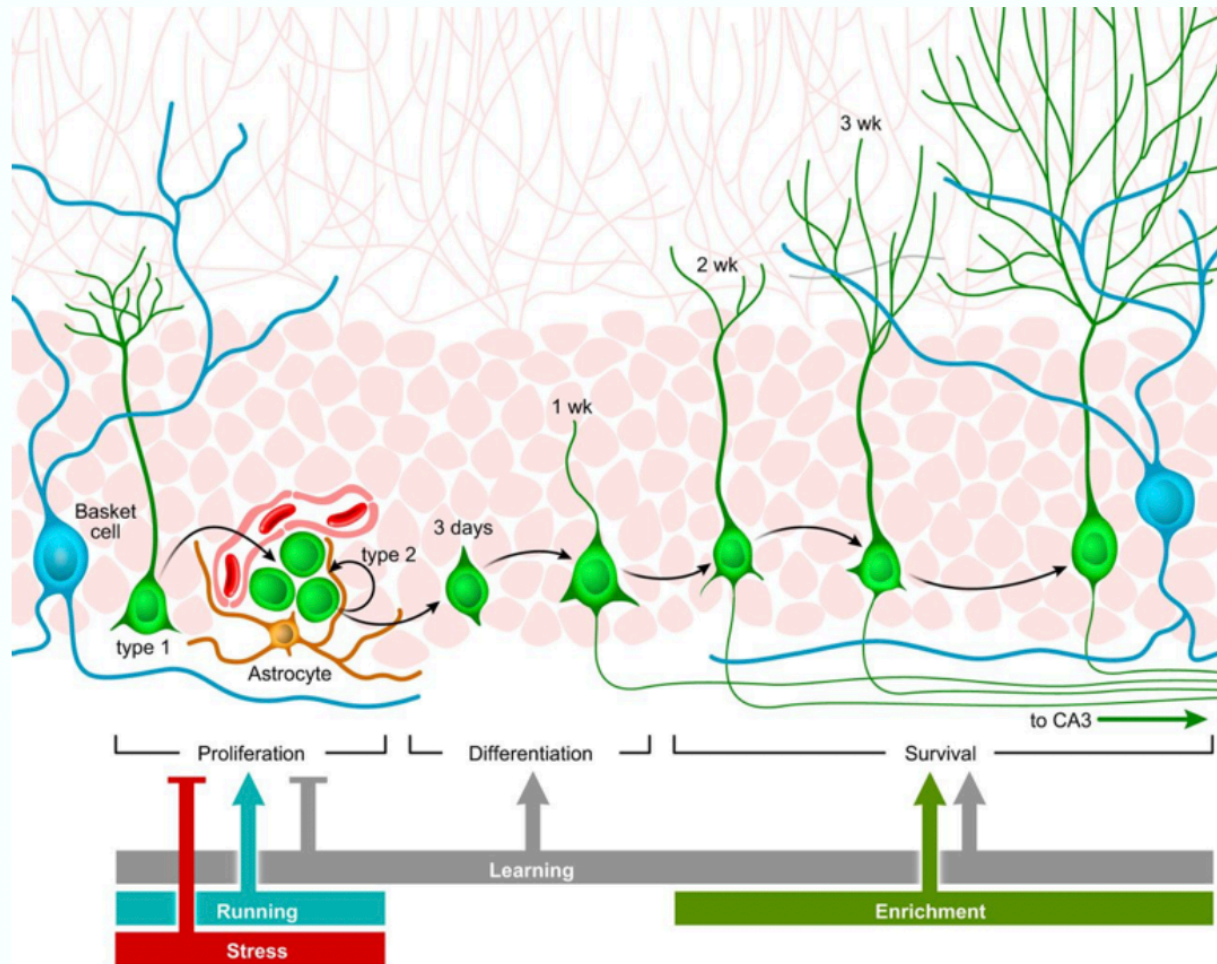
# Microbiology

# Microbiology

# Neurology

- We need single-cell resolution to:

  - Study the mosaic genomes of individual neurons and compositions in the brain

  - Follow genetic variations during fetal development

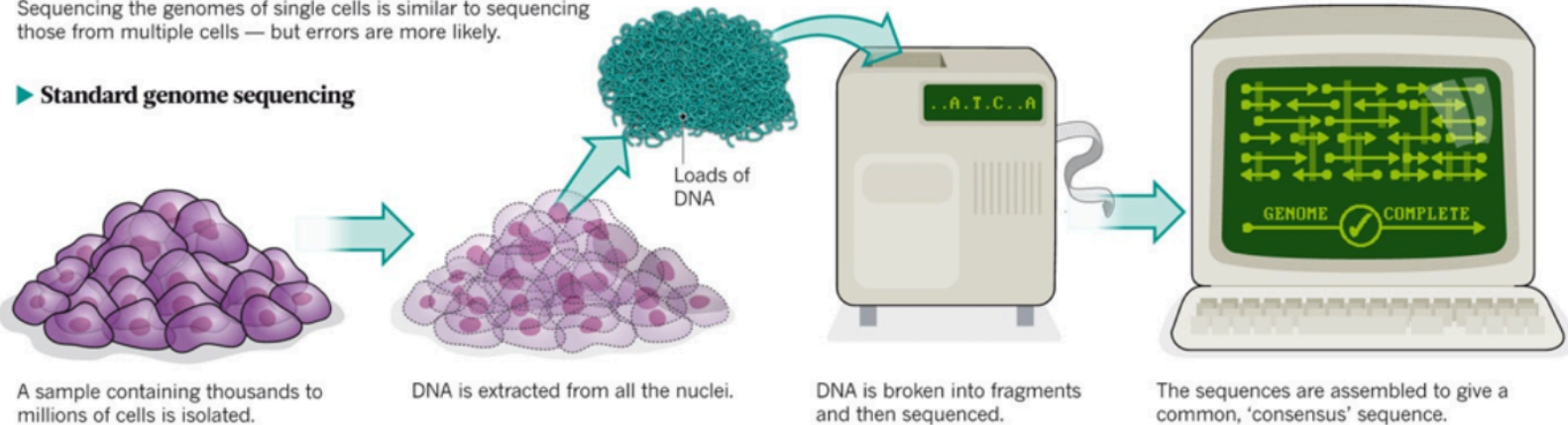  - Develop targeted therapy for neurological diseases for specific cell types
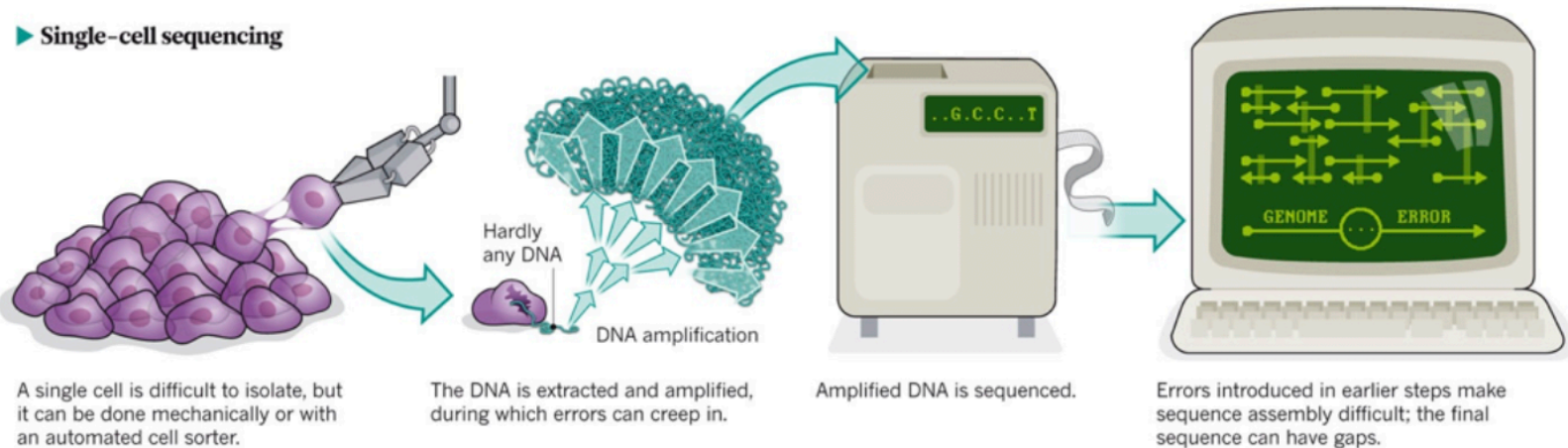
# Traditional vs. Single-cell sequencing



**ONE GENOME FROM MANY**

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.
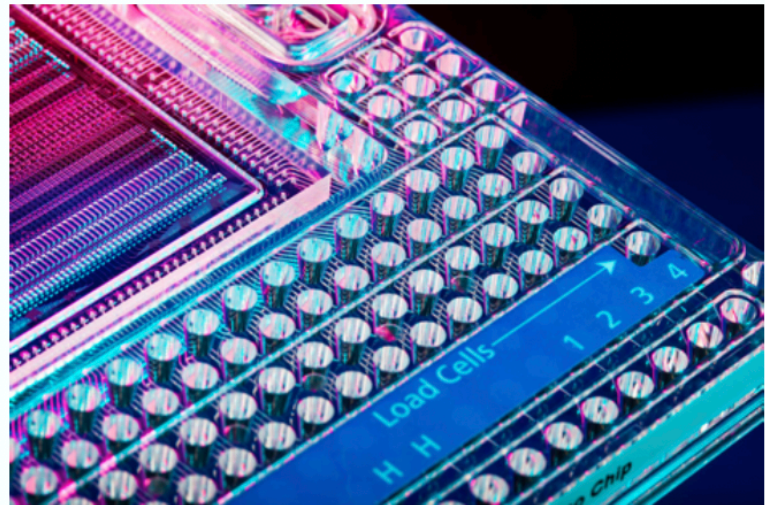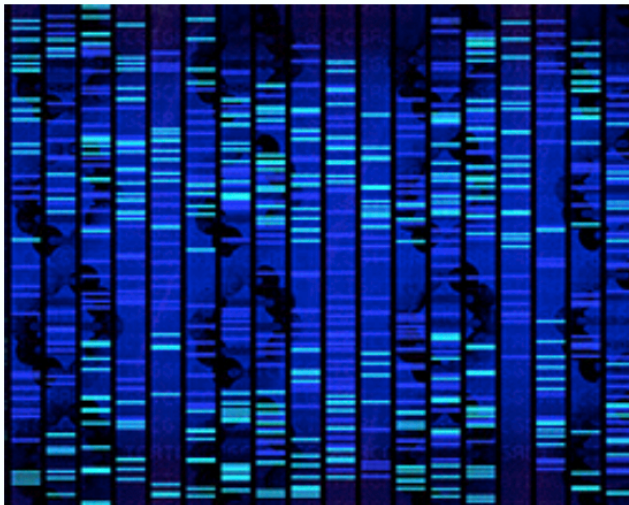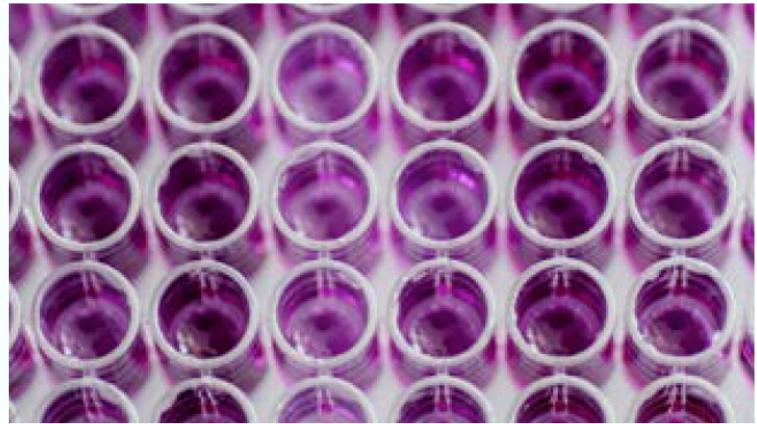
▶ **Standard genome sequencing**

Loads of DNA

..A.T.C..A

GENOME ✓ COMPLETE

A sample containing thousands to millions of cells is isolated.

DNA is extracted from all the nuclei.

DNA is broken into fragments and then sequenced.

The sequences are assembled to give a common, 'consensus' sequence.

▶ **Single-cell sequencing**

Hardly any DNA

DNA amplification

..G.C.C..T

GENOME ☺ ERROR

A single cell is difficult to isolate, but it can be done mechanically or with an automated cell sorter.

The DNA is extracted and amplified, during which errors can creep in.

Amplified DNA is sequenced.

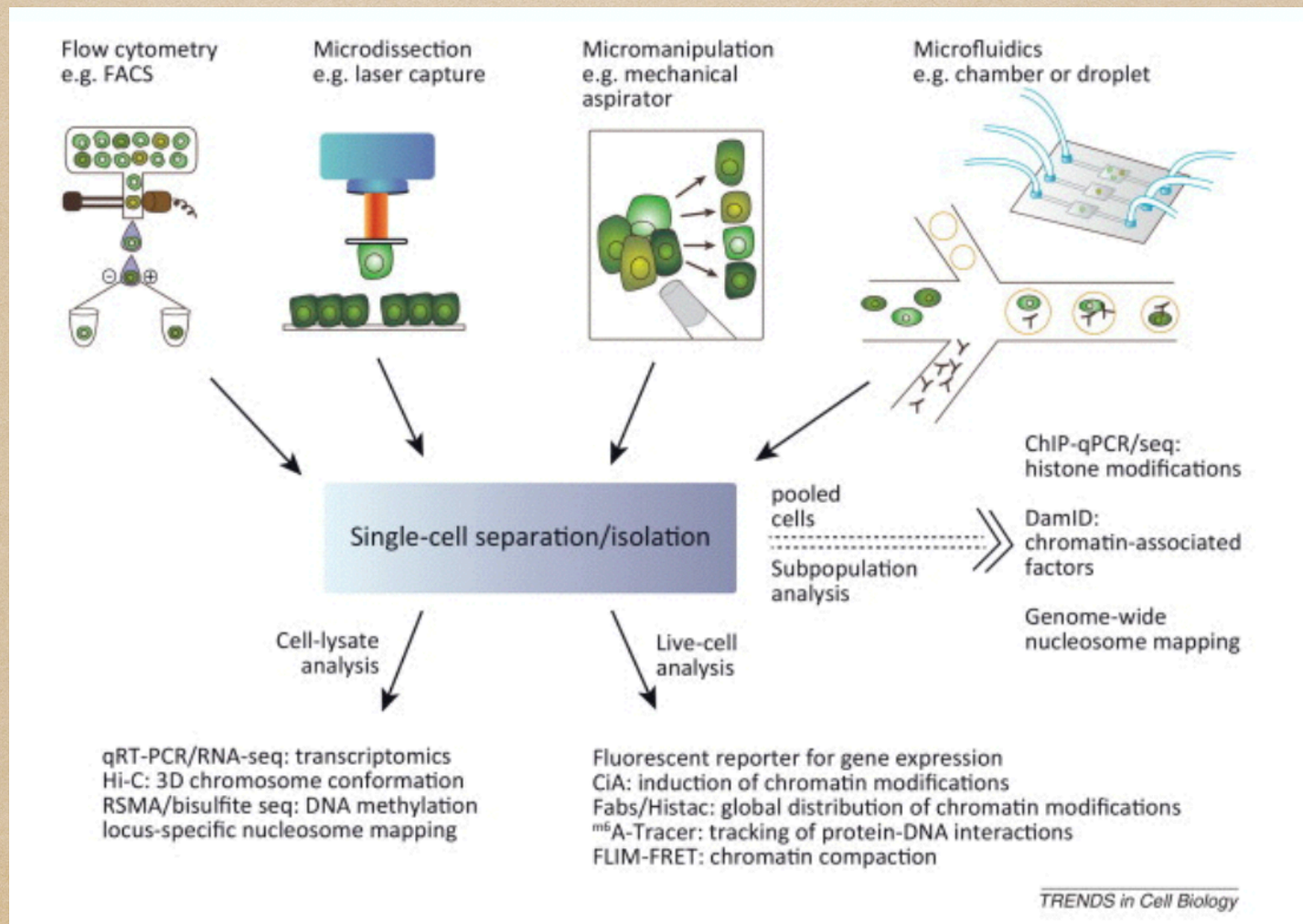Errors introduced in earlier steps make sequence assembly difficult; the final sequence can have gaps.

# Single-Cell Technologies

(i) isolate single cells

(ii) amplify genome efficiently

(iii) sequence DNA

# Single-Cell Technologies



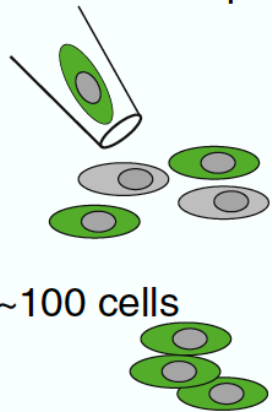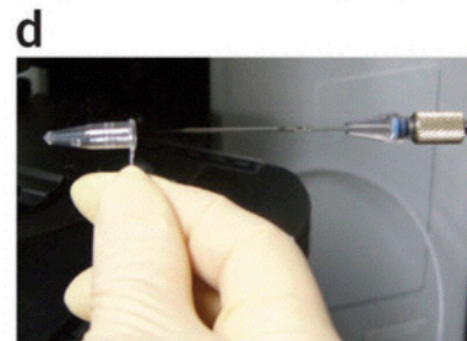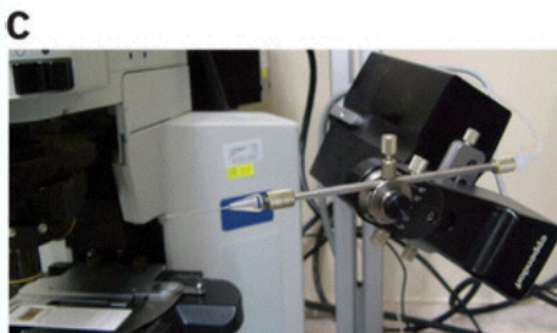Flow cytometry e.g. FACS

Microdissection e.g. laser capture

Micromanipulation e.g. mechanical aspirator

Microfluidics e.g. chamber or droplet

Single-cell separation/isolation

pooled cells

Subpopulation analysis

ChIP-qPCR/seq: histone modifications

DamID: chromatin-associated factors

Genome-wide nucleosome mapping

Cell-lysate analysis

Live-cell analysis

qRT-PCR/RNA-seq: transcriptomics
Hi-C: 3D chromosome conformation
RSMA/bisulfite seq: DNA methylation
locus-specific nucleosome mapping

Fluorescent reporter for gene expression
CiA: induction of chromatin modifications
Fabs/Histac: global distribution of chromatin modifications
$^{m6}$A-Tracer: tracking of protein-DNA interactions
FLIM-FRET: chromatin compaction

*TRENDS in Cell Biology*

# Cell Sorting

# Cell Sorting



FACS: fluorescence activated cell sorting

**FACS**

dissociated cells

Laser

FACS machine

>10 000 cells

nozzle

laser beam

flow cyto-meter

sort electro-nics

1. cell suspension

2. staining of membrane, cytoplasm, and nucleus

3. cell sorting

4. cell culture, analysis, etc.

Ellwart JW, 2000jan02

# Cell Sorting



LCM

cryosectioned tissue

UV laser

IR laser

>1000 cells

LCM: laser capture microdissection

# Cell Sorting



https://media.nature.com/full/nature-assets/nprot/journal/v8/n5/images_article/nprot.2013.046-F4.jpg
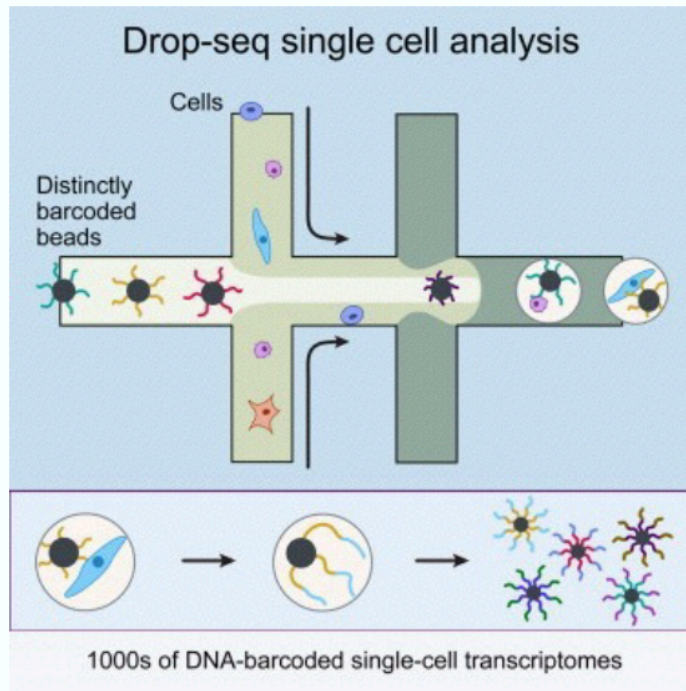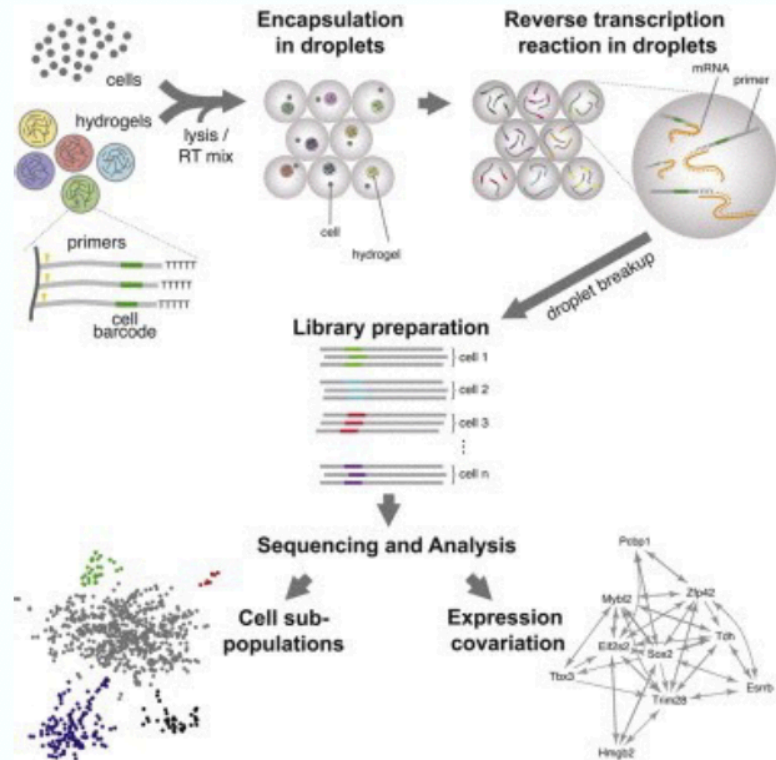
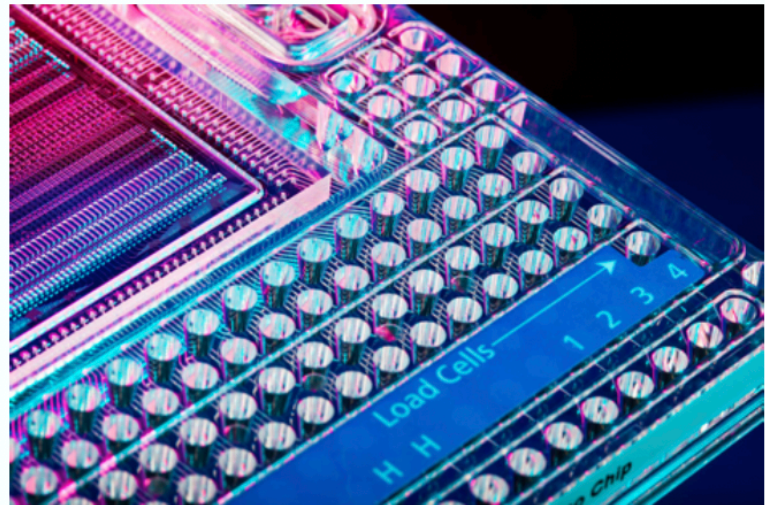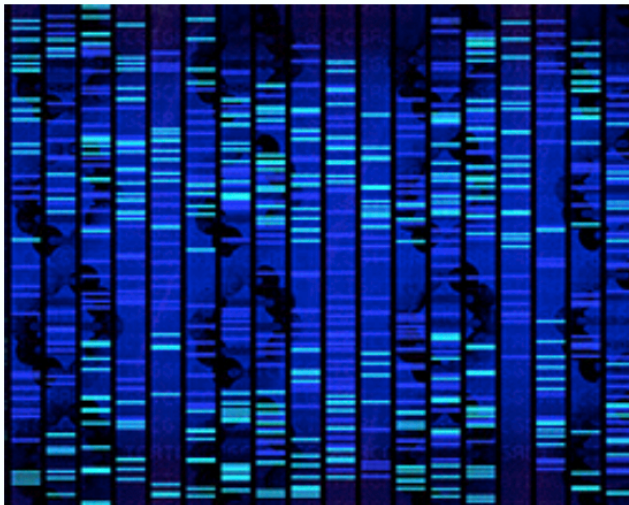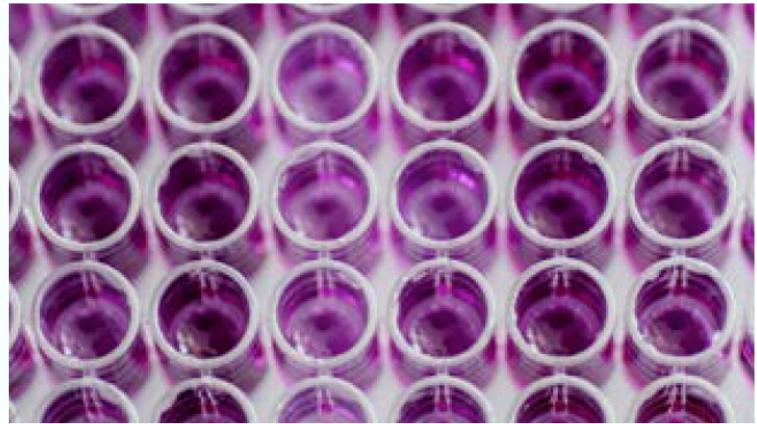# Cell Sorting



High-throughput (~100,000 cells)

Drop-seq

inDrop

# Single-Cell Technologies

(i) isolate single cells

(ii) amplify genome efficiently

(iii) sequence DNA

# SEQUENCING INFORMATICS ASSEMBLY AS A GIANT PUZZLE

# Sequencing informatics



**Clone contig approach**

**Whole-genome shotgun approach**

500 kb

Markers

A B    C    DE    F   G  H    Genome map

Mapped segment of DNA →

Shotgun sequencing of entire genome

A          B

Shotgun sequencing of mapped segment

**Assembled sequence**

A          B

D          E

H          **Assembled sequences**

**Position of sequence is already known**

**Markers used to anchor assembled sequences on to the map**
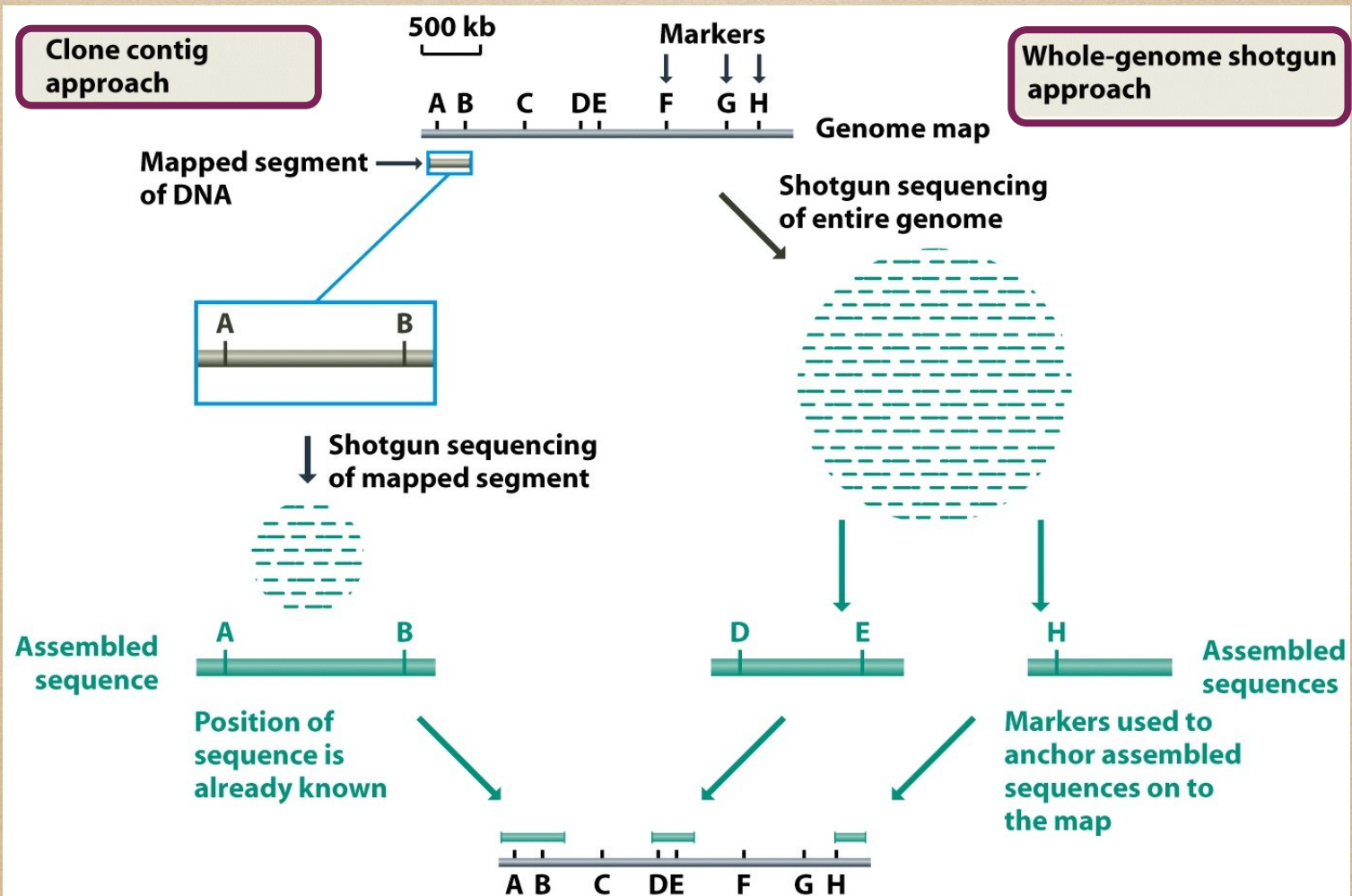
A B    C    DE    F    G  H

Figure 3-3  Genomes 3 (© Garland Science 2007)

# Sequencing informatics



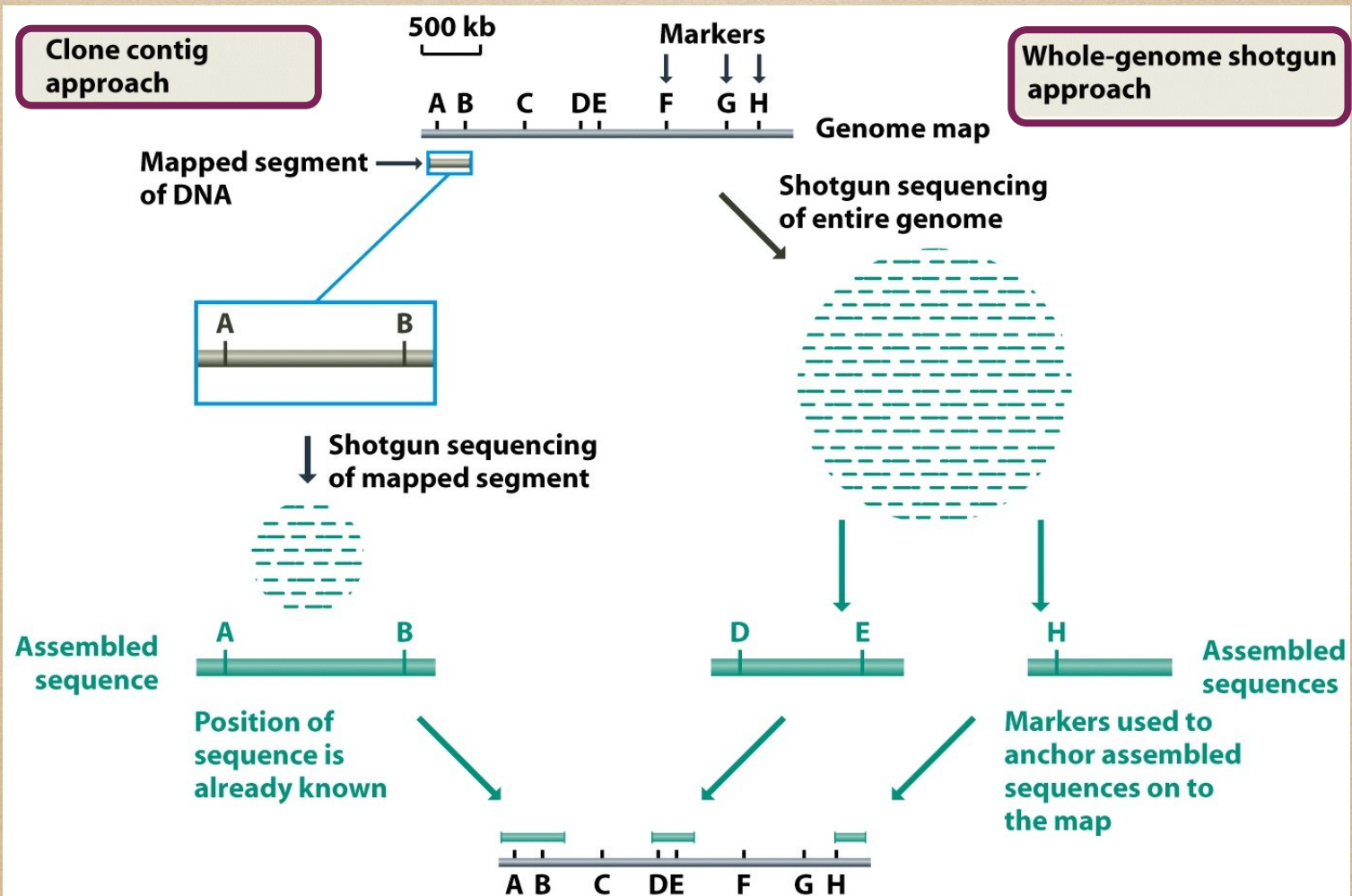Figure 3-3 Genomes 3 (© Garland Science 2007)

# Sequence assembly

- A fundamental goal of DNA sequencing has been to generate large, continuous regions of DNA sequence – CONTIGS

- In principle, assembling a sequence is just a matter of finding overlaps and combining them.

- In practice:

    - most genomes contain multiple copies of many sequences,
    - there are random mutations (either naturally occurring cell-to-cell variation or generated by PCR or cloning),
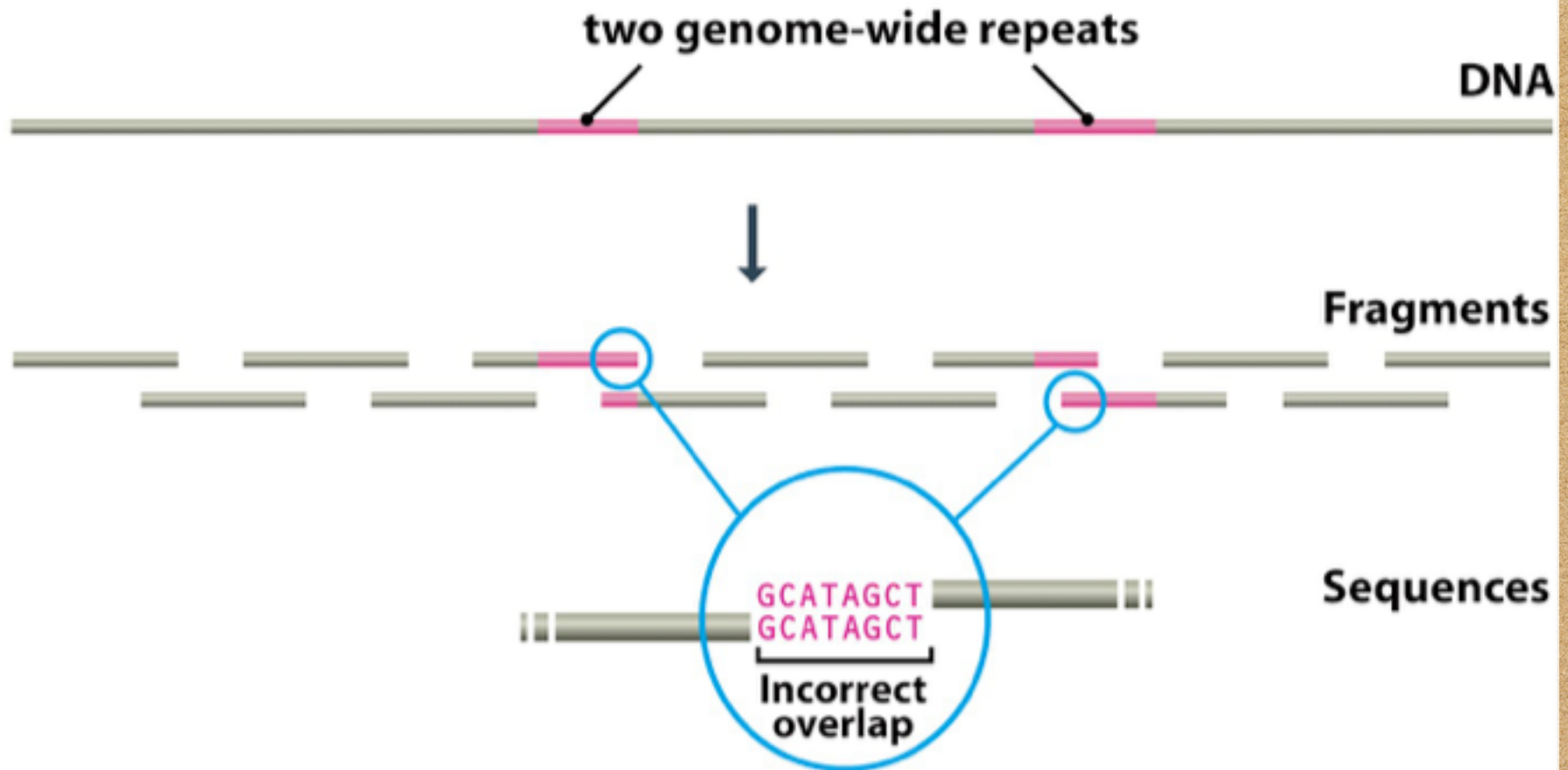    - there are sequencing errors

DNA

500 bp

Fragments

Sequences

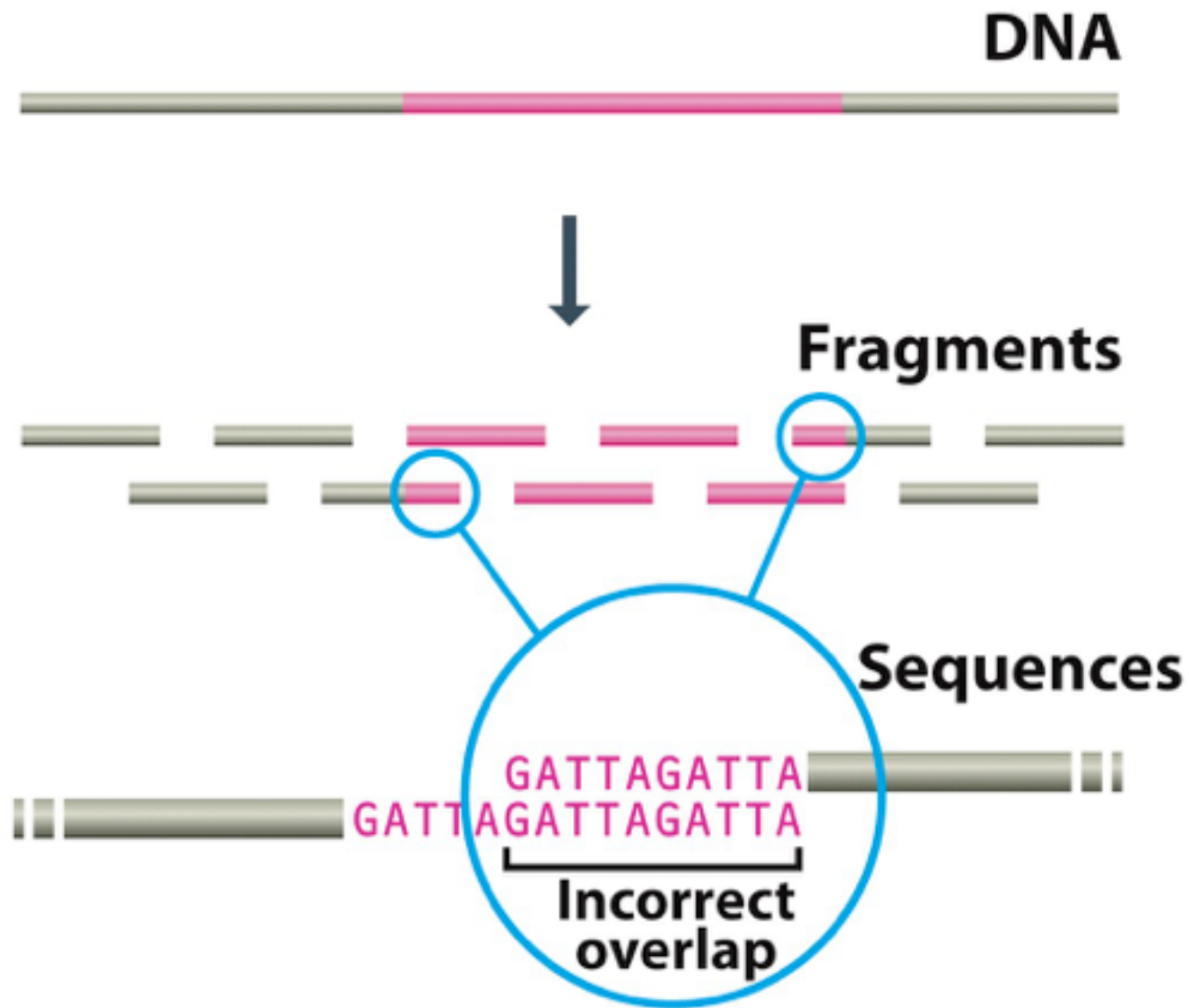CAATGCATTA

GCAGCCAATGC

Overlap

# Assembly problems



**Problems with genome-wide repeats**

# Problems with tandemly repeated DNA

**DNA**

**Fragments**

**Sequences**

GATTAGATTA
GATTAGATTAGATTA

**Incorrect overlap**

# Assembly problems: sequencing gaps



Genome sequence

Mini-sequences
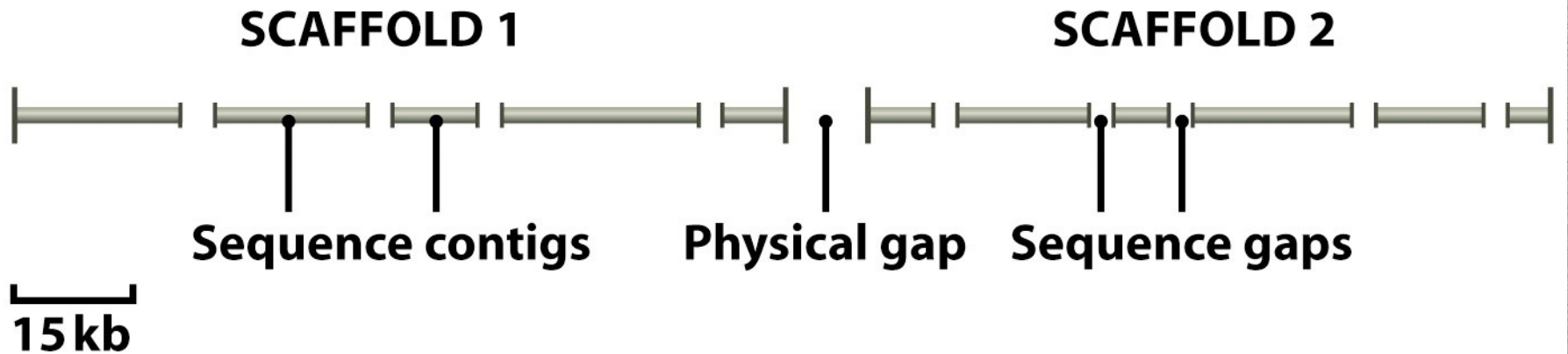
1 kb

**Scaffolds**

SCAFFOLD 1

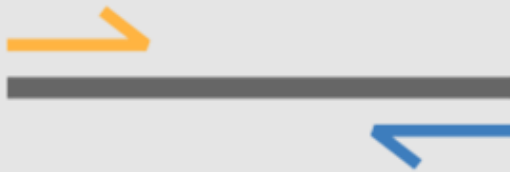SCAFFOLD 2

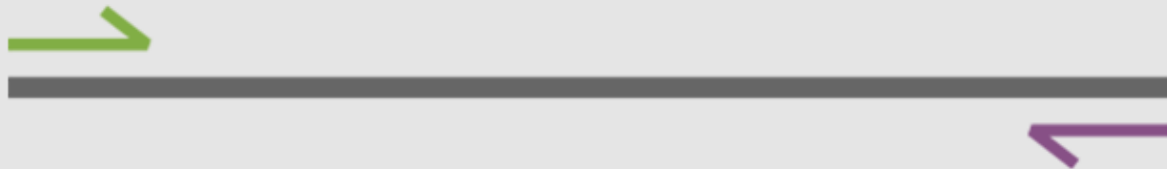Sequence contigs      Physical gap    Sequence gaps

15 kb

# Sequencing gaps - pair end reads to the rescue



Short-Insert Paired End Reads
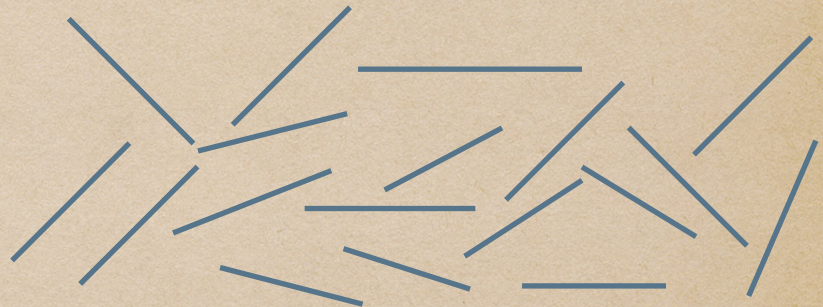
Long-Insert Paired End Reads (Mate Pair)

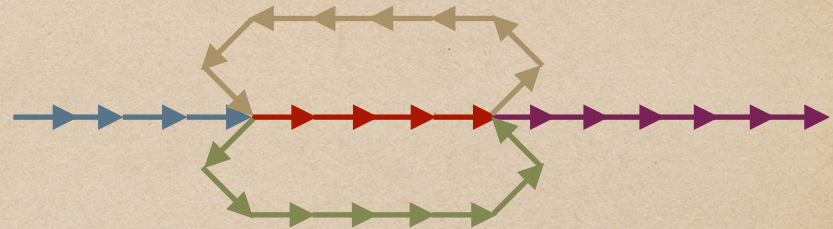# Overview of genome assembly (1)

Sample collection



DNA sequencing



Pairwise read overlaps

...AGCTTTAGGCTAGCAATGC
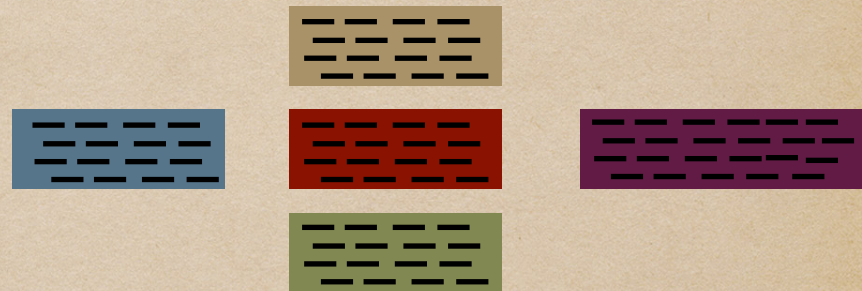GCAATGCTATAGGCCT...

# Overview of genome assembly (2)

String graph construction
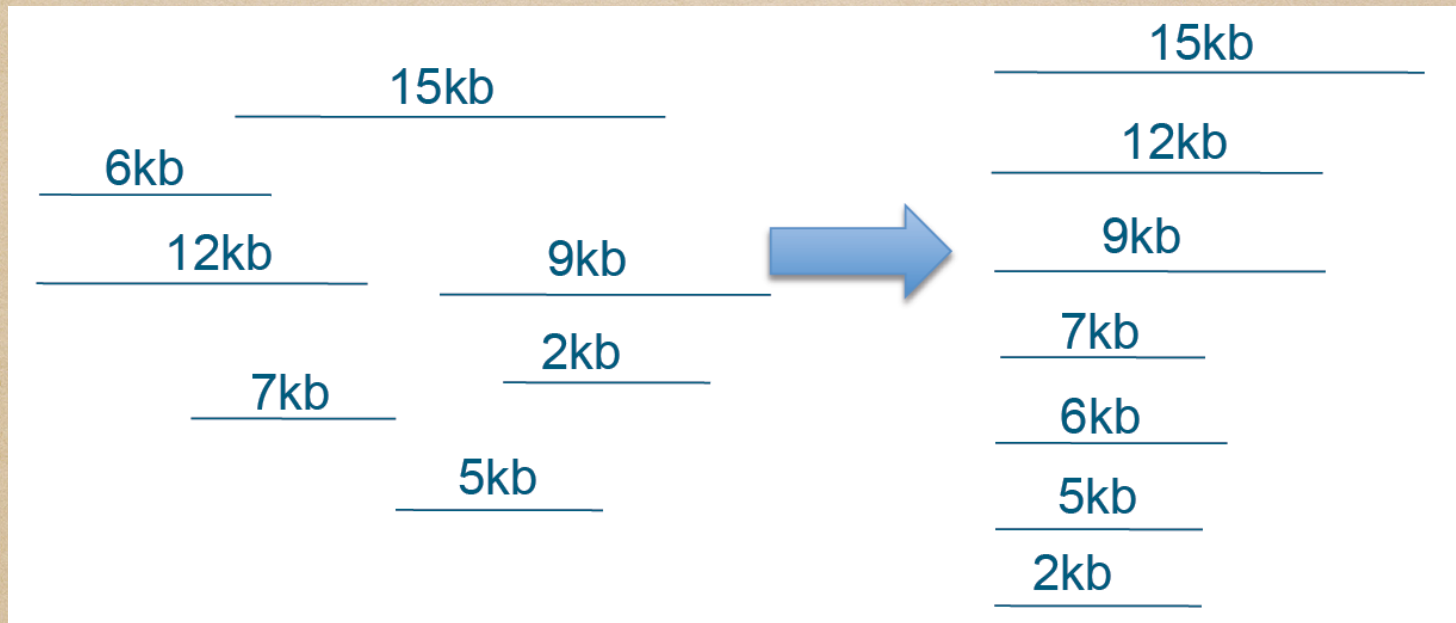
Contig construction

Scaffold construction

# Assembly evaluation ~ N50



If one orders the set of contigs produced by the assembler by size, then N50 is the size of the contig such that 50% of the total bases are in contigs of equal or greater size.

15+12+9+7+6+5+2 = 56.
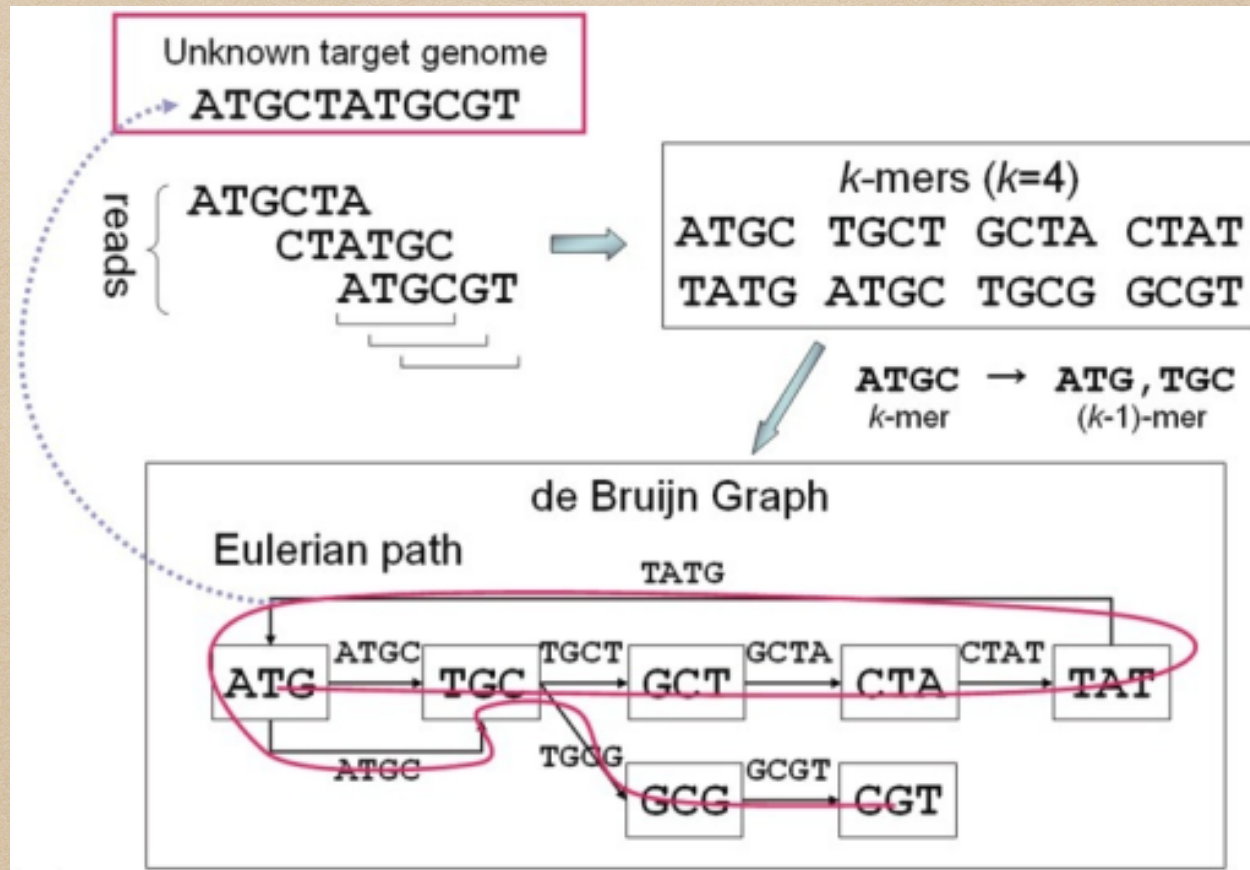
56/2 = 28   ->  N50 is 9kb (15+12 = 27 is less than 50%)

# Sequence assembly
## NGS case

- Volume and read length of data from next-gen sequencing machines meant that the read-centric overlap approaches were not feasible

- Already in 1980's Pevzner et al. introduced an alternative assembly framework based on de Bruijn graph

- Based on a idea of a graph with fixed-length subsequences (k-mers)

- Key is that not storing read sequences – just k-mer abundance information in a graph structure
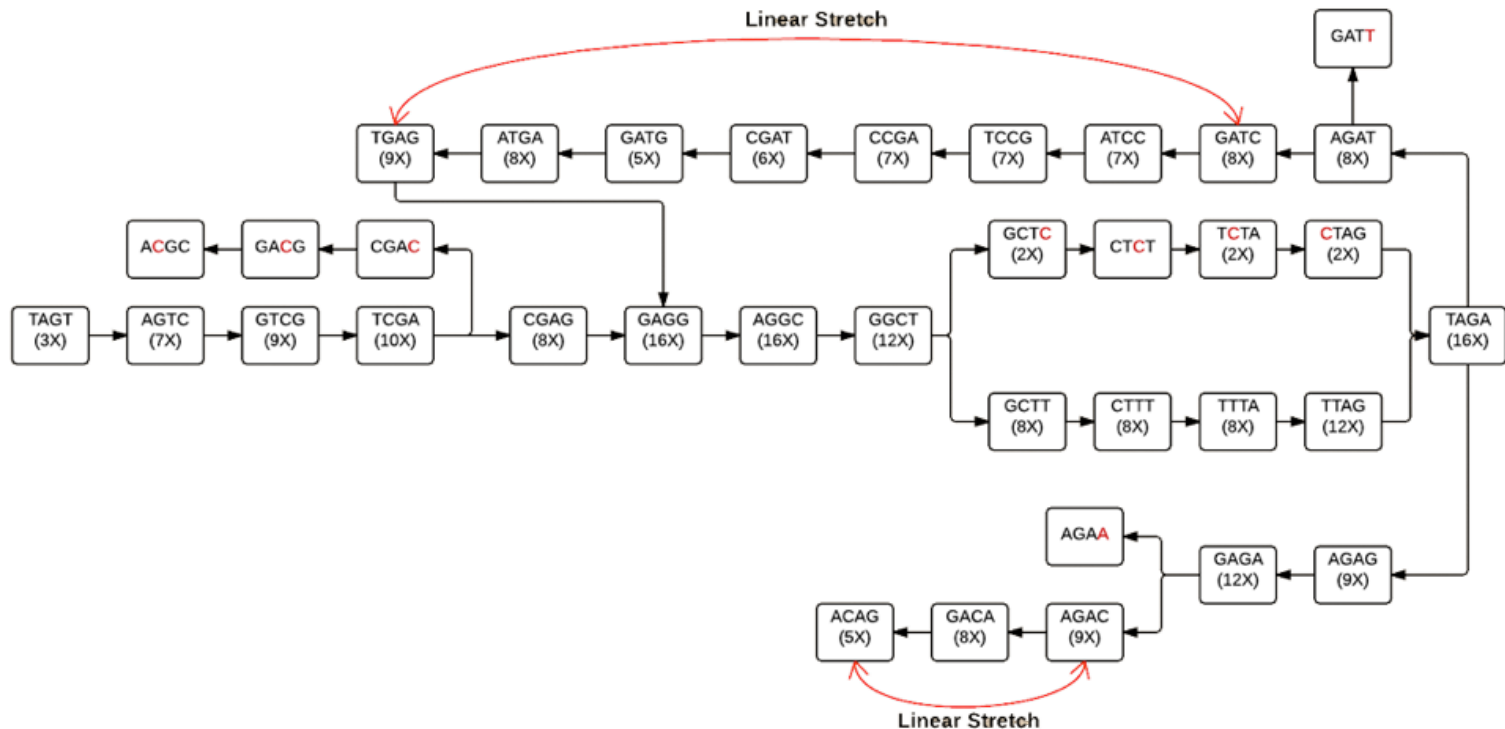
# De bruijn graph construction



- continuous linear stretches within the graph
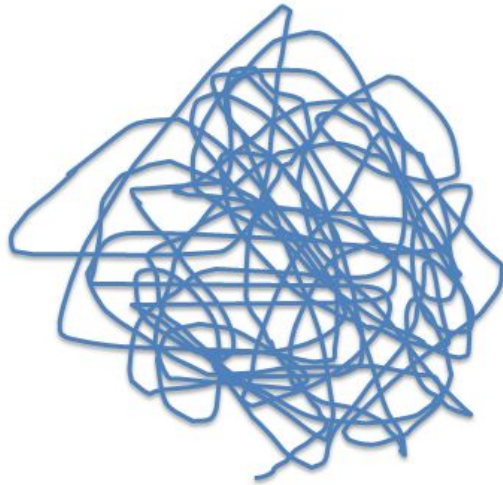- assembler keeps information about reads coverage for each k-mer/node.

Graph is simplified to combine nodes that are associated with the continuous linear stretches into single, larger nodes of various k-mer sizes.
Error correction removes the tips and bubbles that result from sequencing errors.
Sequencing errors are low frequency tips in the graph.

# Sequence assembly: genome or transcriptome



**Genome Assembly**
Single Massive Graph

Entire chromosomes represented.

**Trinity Transcriptome Assembly**
Many Thousands of Small Graphs

Ideally, one graph per expressed gene.

# Next-gen assemblers

- First de Bruijn based assembler was Newbler developed by 454 Life Sciences
  - Adapted to handle main source of error in 454 data – indels in homopolymer tracts
- Many de Bruijn assemblers subsequently developed
  - SHARCGS, VCAKE, VELVET, EULER-SR, EDENA, ABySS and ALLPATHS, SOAP
  - Most can use pair-mate information
- Slightly different approach to transcriptome assembly:
  - It has to allow many discontinuous graphs representing single transcript, including paralogs and alternatively spliced ones.
  - SOAP-Trans, Trinity

# BIOINFORMATICS CREED

- Remember about biology

- Do not trust the data

- Use comparative approach

- Use statistics

- Know the limits

- Remember about biology!!!