

WHO Collaborating Centre
for Reference and
Research on Influenza
VIDRL



**World Health
Organization**
Philippines



KDCA
Korea Disease Control and
Prevention Agency

NGS DATA ANALYSIS

FRANCIS A. TABLIZO

Core Facility for Bioinformatics
Philippine Genome Center
University of the Philippines System

April 15, 2024



**CENTRE FOR
PATHOGEN
GENOMICS**



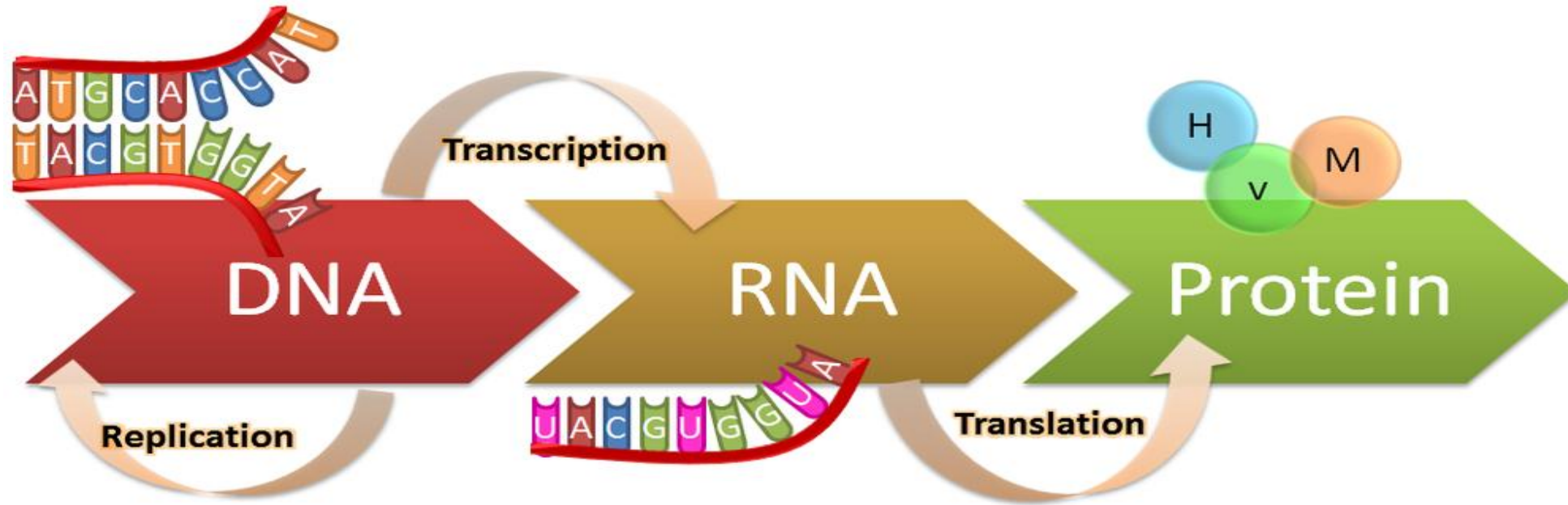
A joint venture between The University of Melbourne and The Royal Melbourne Hospital



Objectives

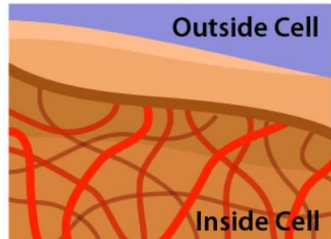
- Go through the currently available sequencing technologies
- Provide an overview on the handling and analysis of high-throughput sequencing data sets, particularly in the field of genomics

Central Dogma of Molecular Biology

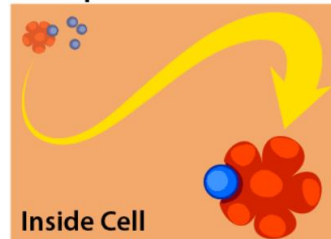


<https://genius.com/Biology-genius-the-central-dogma-annotated>

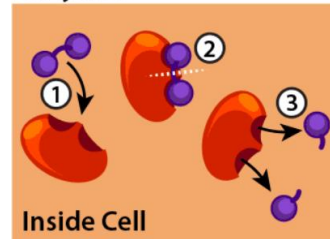
Structure Protein



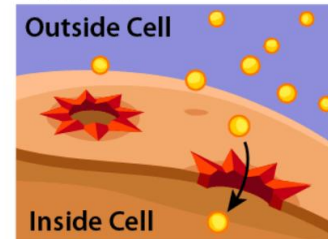
Transport Protein



Enzymes



Channels

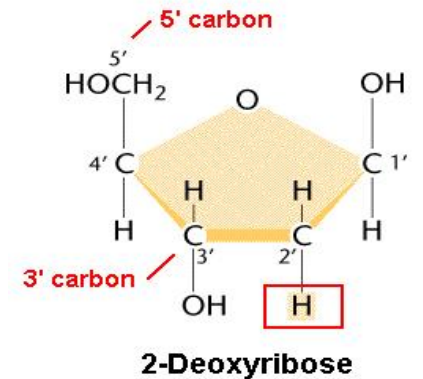
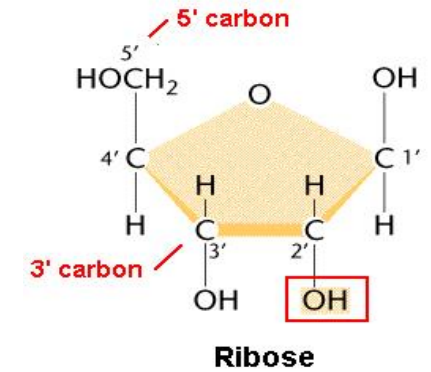
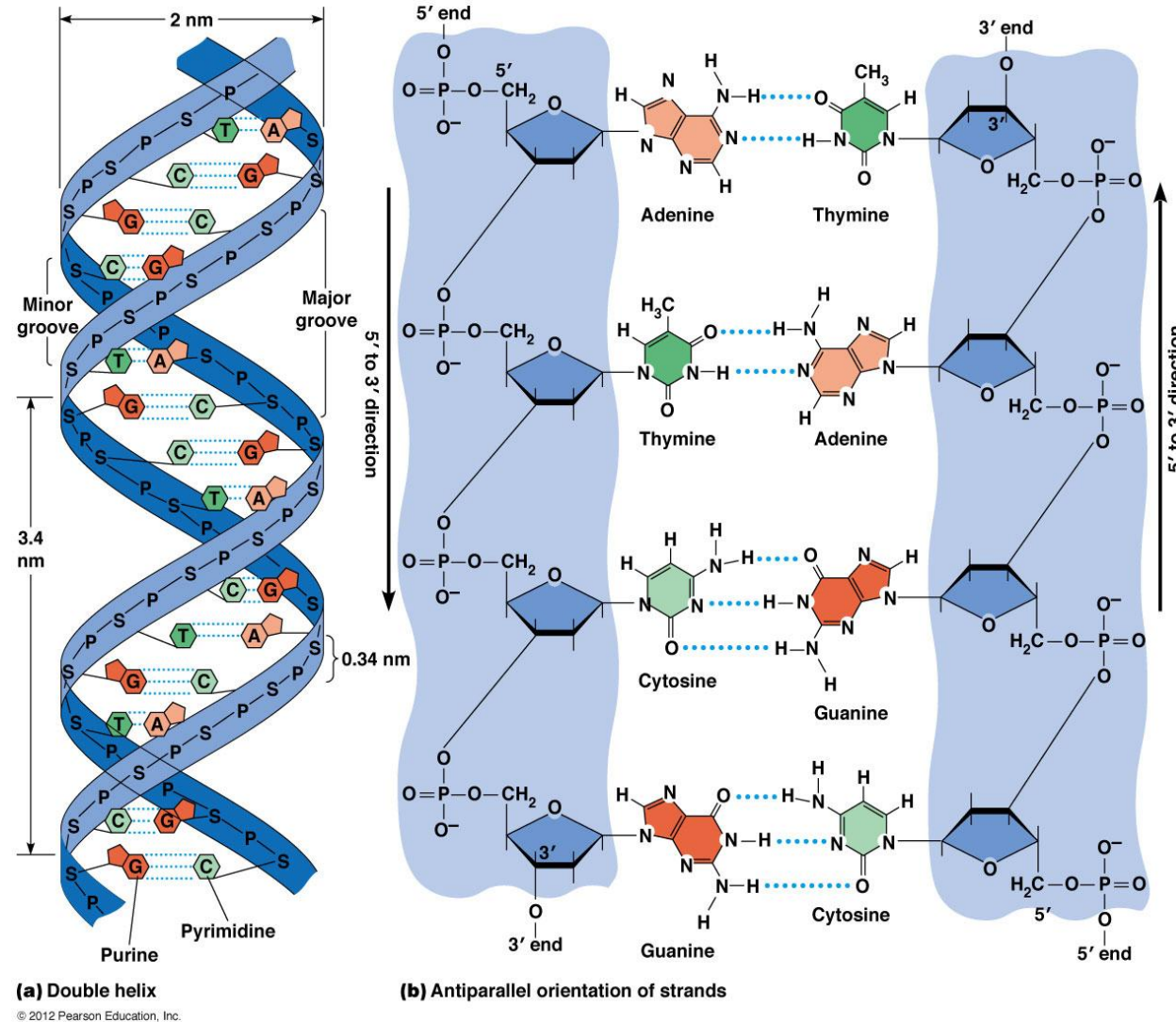


<https://askabiologist.asu.edu/venom/what-are-proteins>

<http://www.pixabay.com/>

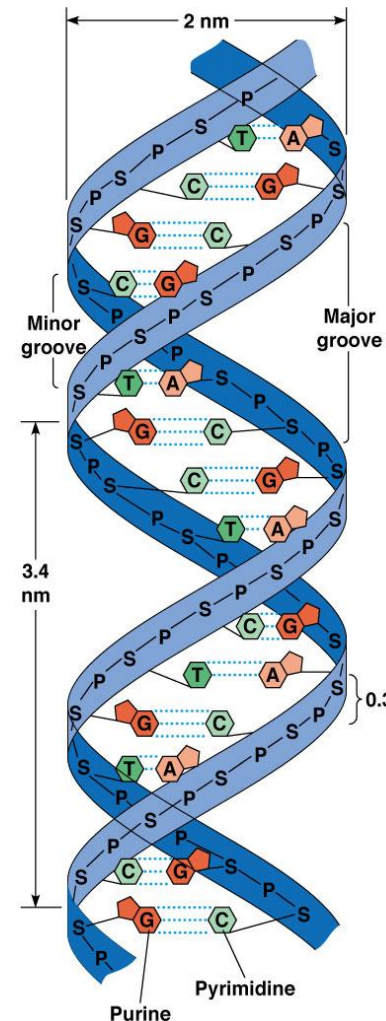
DNA Molecule

Deoxyribonucleic acid



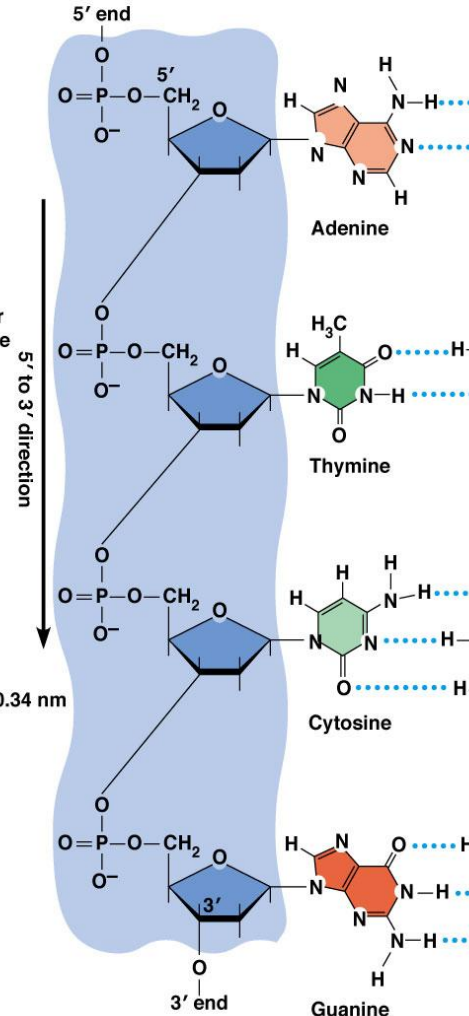
https://www.mun.ca/biology/scarr/Fg10_09b_revised.gif

DNA Sequencing



(a) Double helix

© 2012 Pearson Education, Inc.

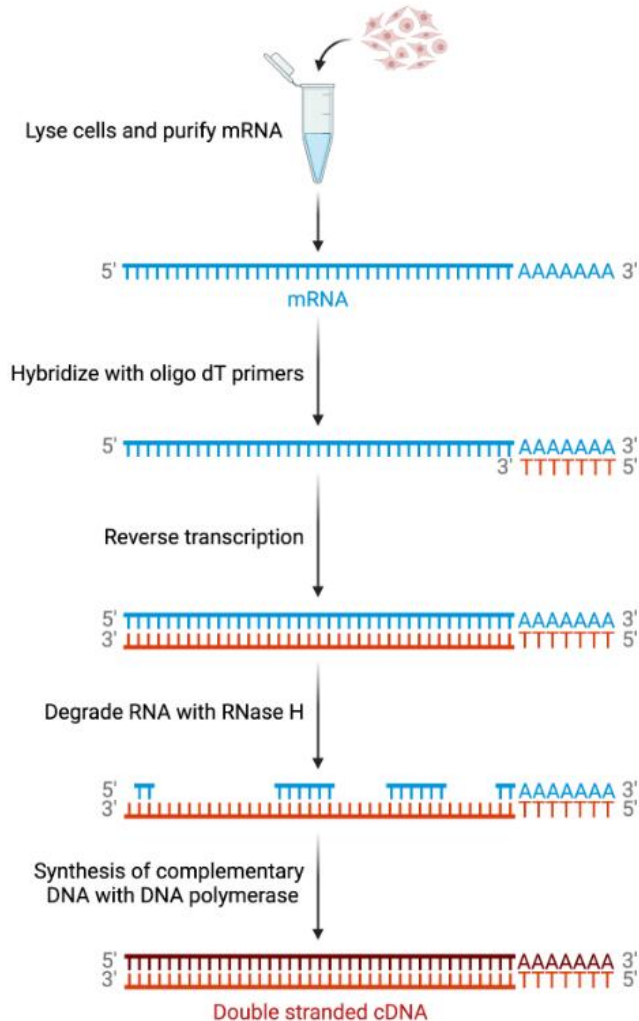


(b) Antiparallel orientation of strands

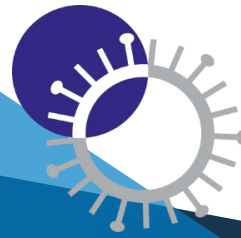
Refers to the process of determining the order or sequence of nucleotides along a DNA strand

DNA Sequencing

cDNA Synthesis



- Most current sequencing technologies only work with **double-stranded DNA** molecule as starting material.
- Usually, if the starting material is RNA, there is a need to **convert the RNA to its complementary DNA (cDNA)** form prior to sequencing.
- Viruses have **other forms of genetic material** (e.g., ssDNA, dsRNA, etc.). Make sure to take this in consideration before sequencing.



WHO Collaborating Centre
for Reference and
Research on Influenza
VIDRL



**World Health
Organization**
Philippines



KDCA
Korea Disease Control and
Prevention Agency

CURRENT SEQUENCING TECHNOLOGIES



**CENTRE FOR
PATHOGEN
GENOMICS**

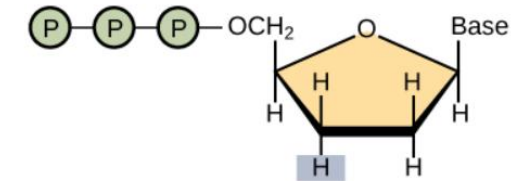
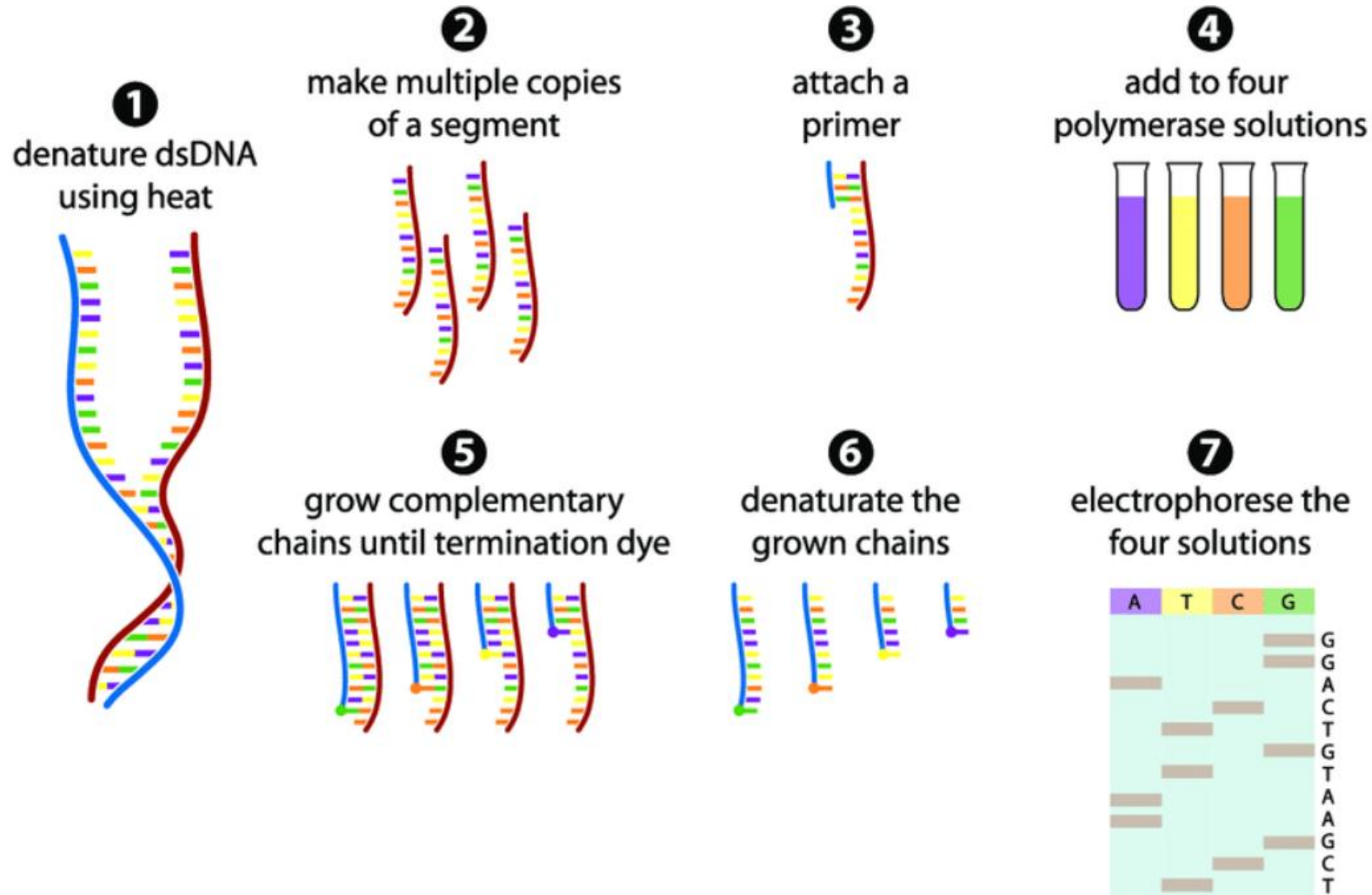


A joint venture between The University of Melbourne and The Royal Melbourne Hospital

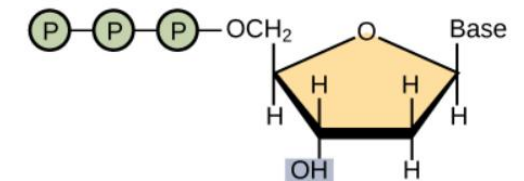


Sanger Sequencing

Dideoxy chain termination method



Dideoxynucleotide (ddNTP)

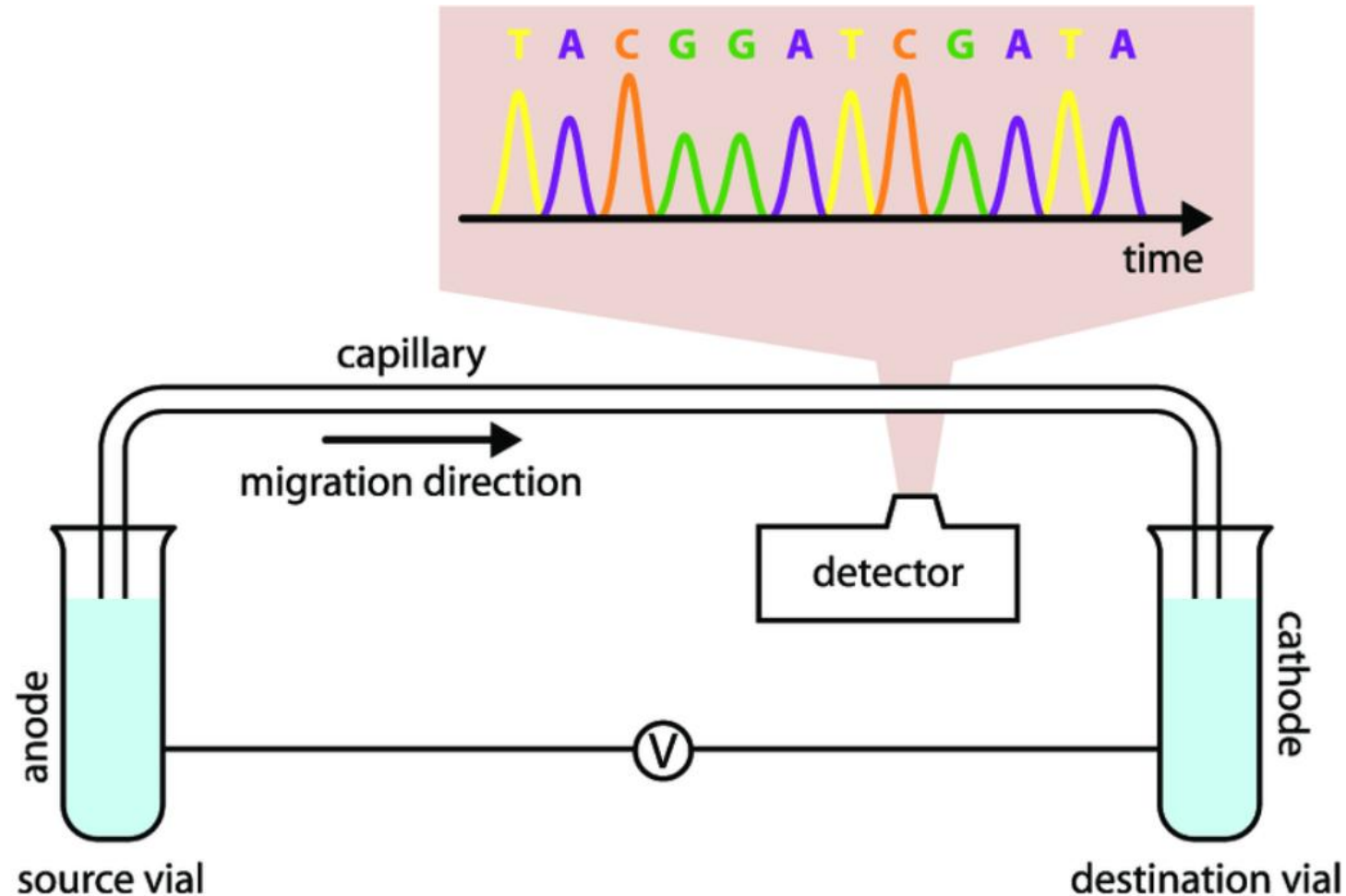


Deoxynucleotide (dNTP)

<https://openstax.org/books/biology/pages/17-3-whole-genome-sequencing>

Capillary Sequencing

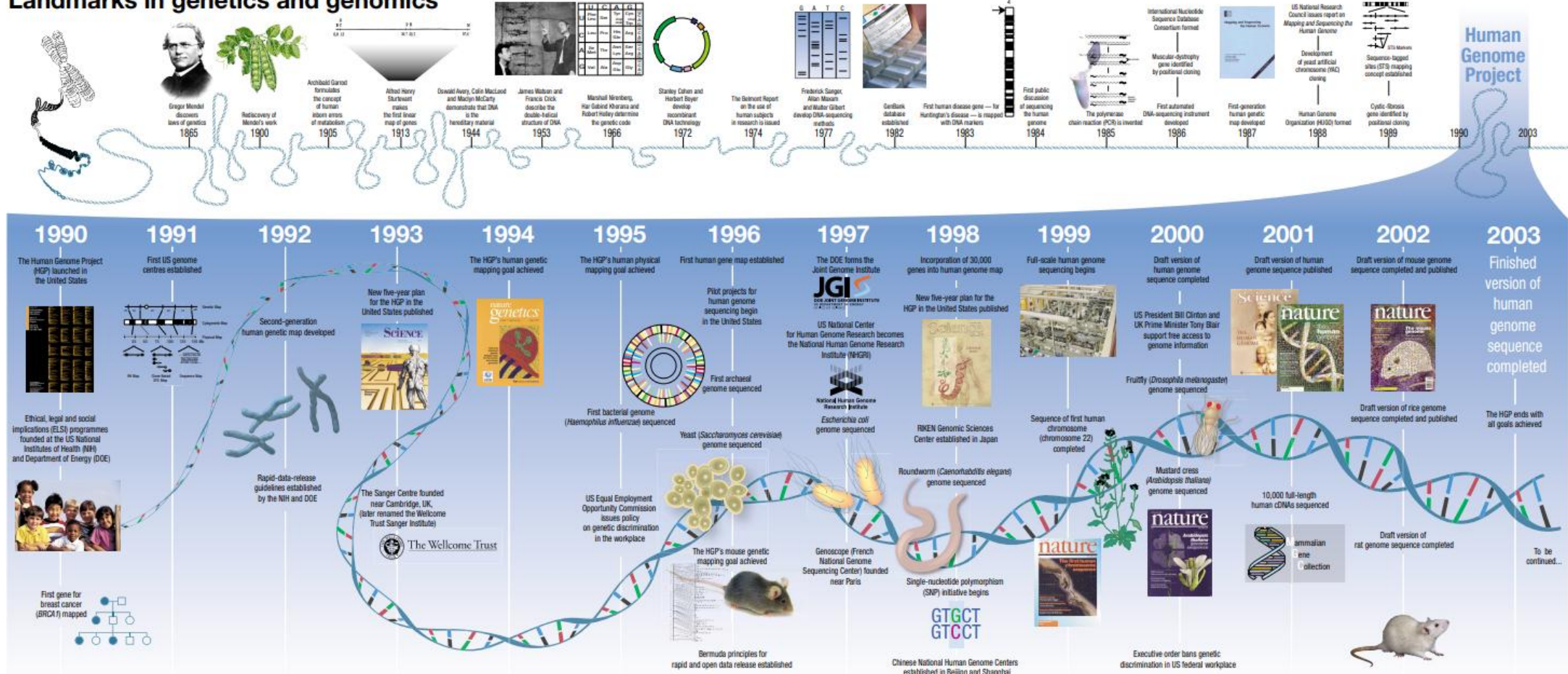
Dideoxy chain termination method



Capillary Sequencing

Main technology available during the Human Genome Project

Landmarks in genetics and genomics



PEAS COURTESY J. BLAMBLE, CITY UNIV, NEW YORK; WATSON & CRICK COURTESY A. BARRINGTON-BROWN/NIH; SCIENCE CLIPARTS COURTESY AAAA

High-Throughput Sequencing

First NGS Platform on the Market: 454 Life Sciences, Pyrosequencing, 2004

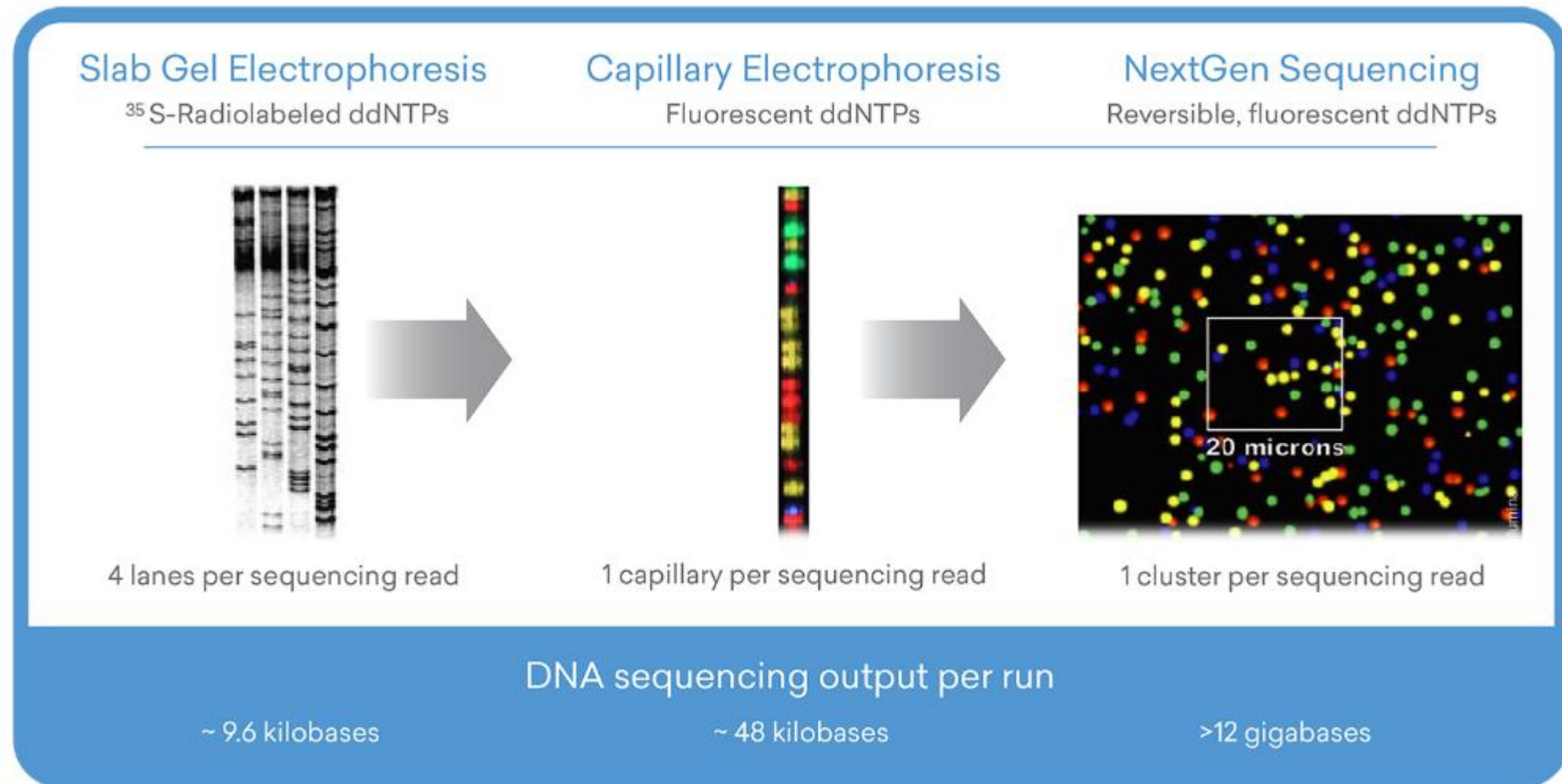
Proof of Concept: “Project Jim”, Published in 2008

Jim Watson	Human Genome Project
454 Life Sciences, Pyrosequencing Technology, analyzed by BCM-HGSC	Sanger Capillary Sequencing
2 months, 3 instruments	10-13 Years
\$1-2 million \$ 250,000 with GS Titanium FLX	\$ 100 million - \$ 2.7 billion
8x coverage	7.5x coverage
250 bp read length 400 bp with GS Titanium FLX	500-800 bp

Source: History of DNA Sequencing & Current Applications. Roche Applied Sciences.

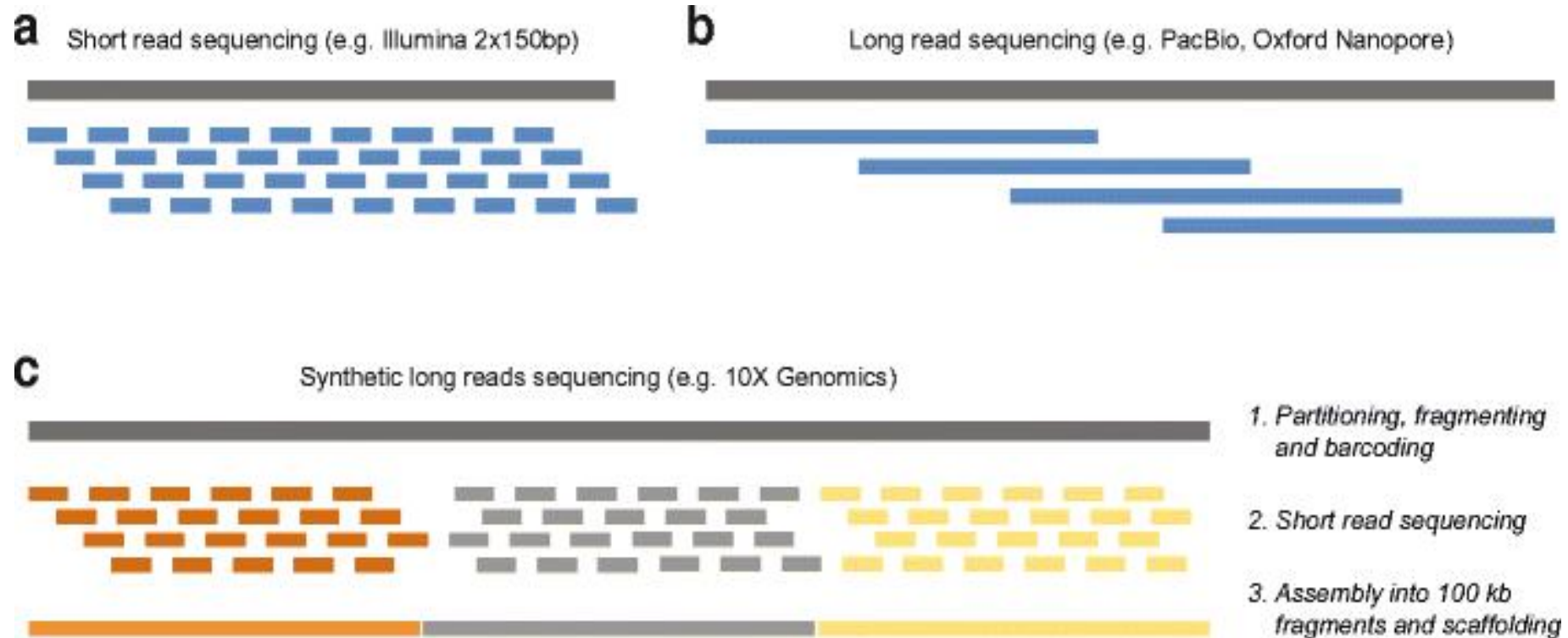
High-Throughput Sequencing

Simultaneous sequencing of thousands to even billions of DNA fragments in a single run



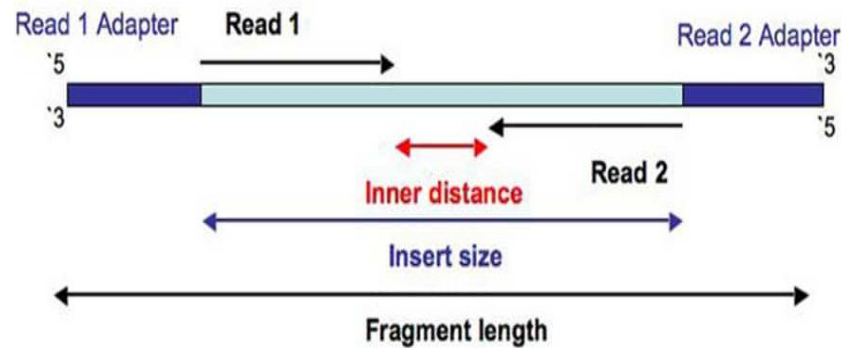
Important Concepts

Short vs. Long Read Sequencing



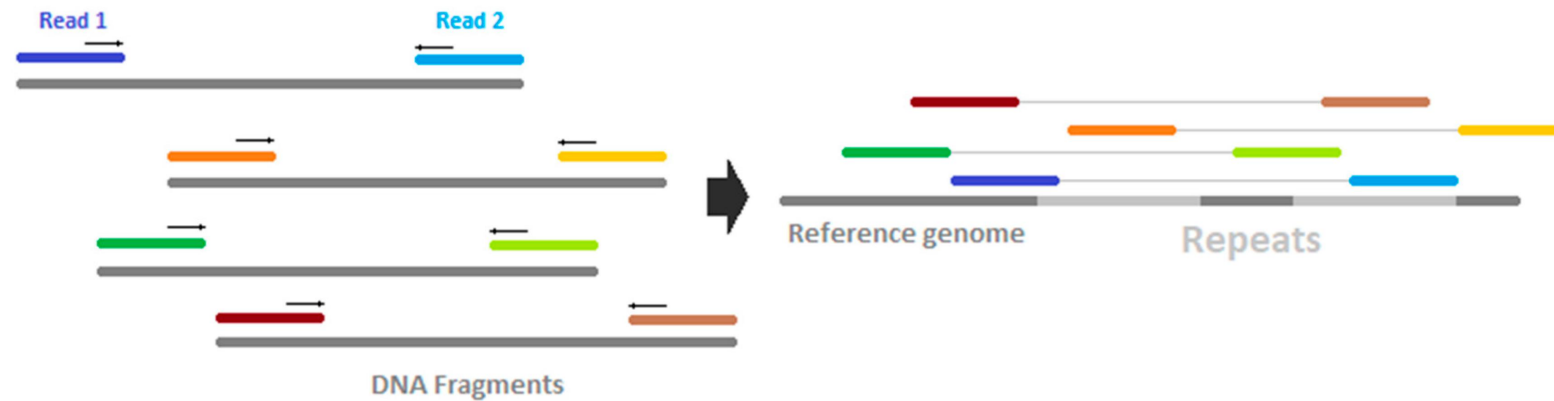
Important Concepts

Single-end vs. Paired-end Data (Short Reads)

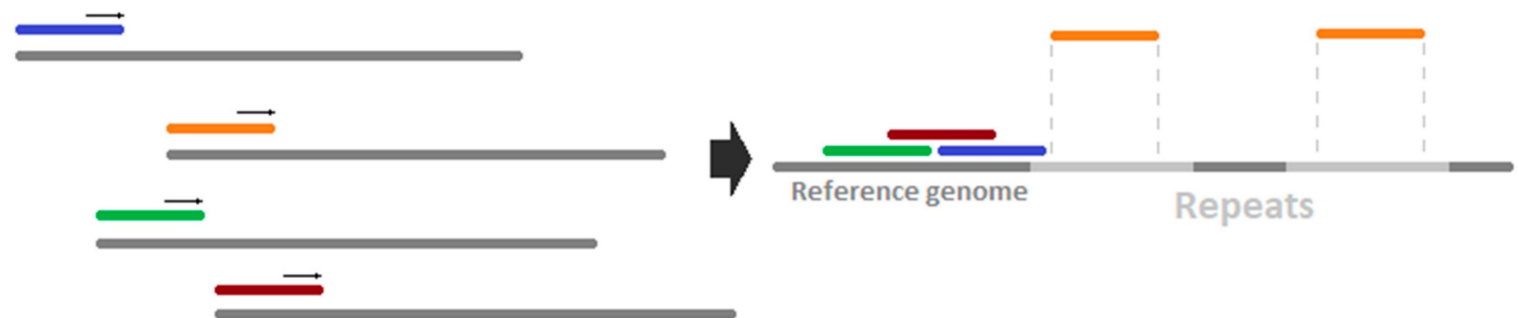


<https://thesequencingcenter.com/knowledge-base/what-are-paired-end-reads/>

A) Paired-end Sequencing

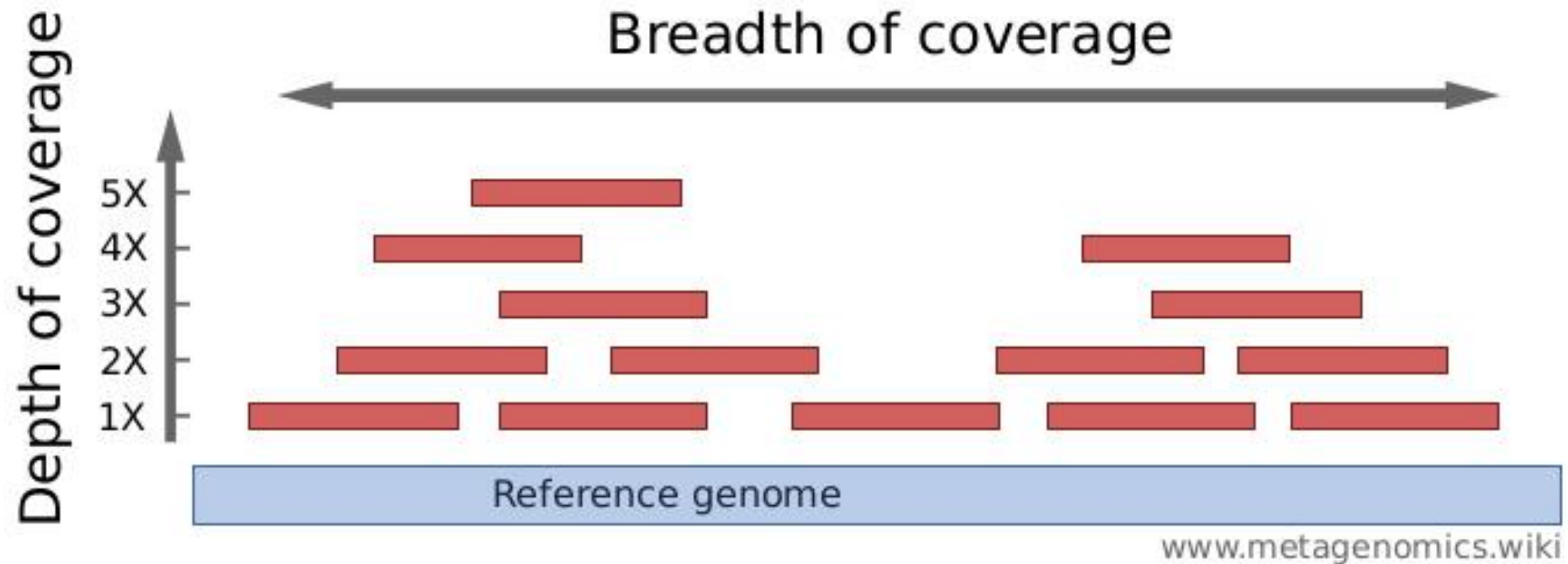


B) Single-end Sequencing



Important Concepts

Breadth vs. Depth of Coverage

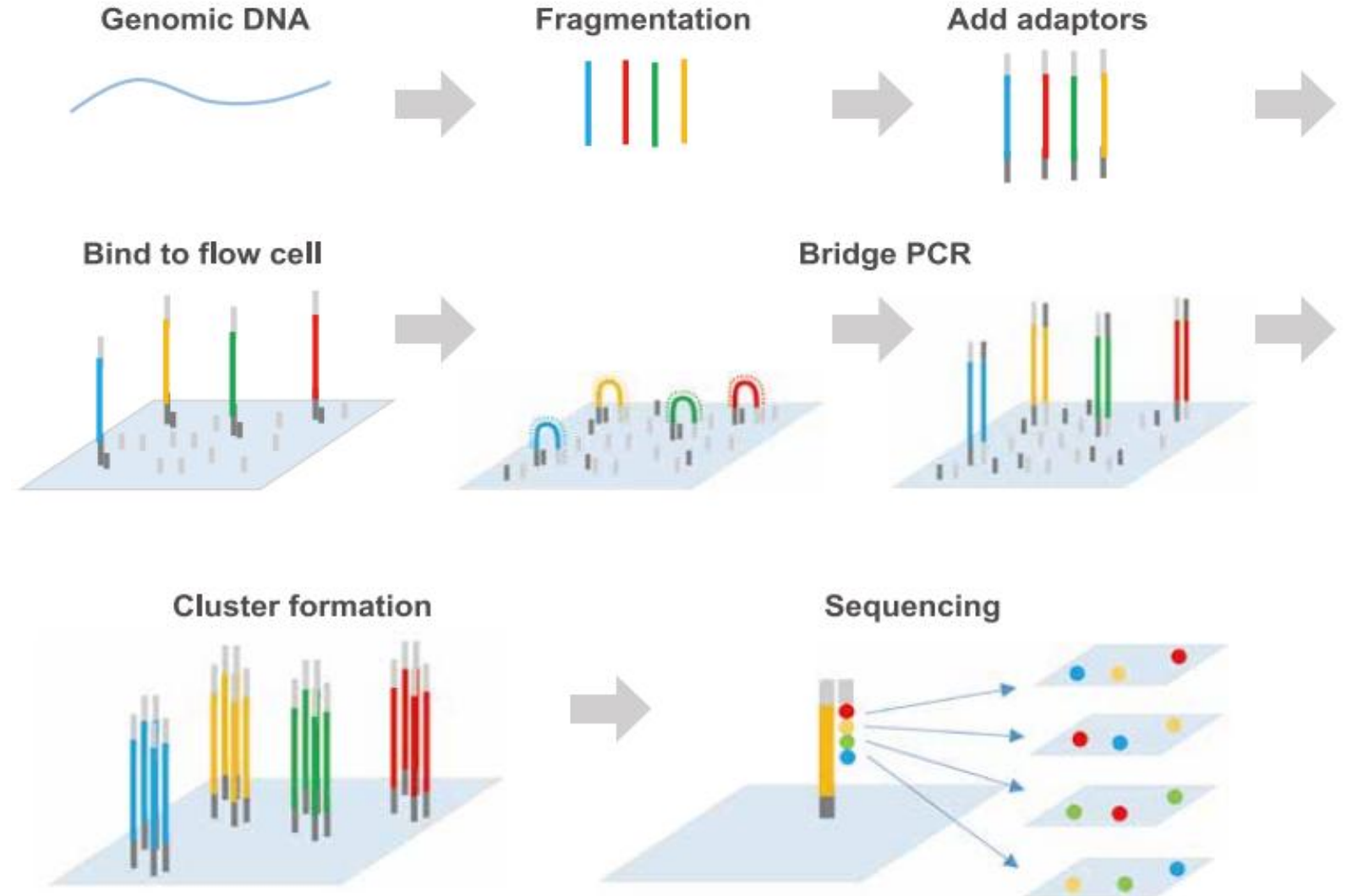


Next (2nd) Generation Sequencing

Platform: Illumina

Sequencing by synthesis

- Read Length: 50-300 bp
- Read Fragments: Paired-end, Single-end
- Throughput: 1.2GB to 8TB

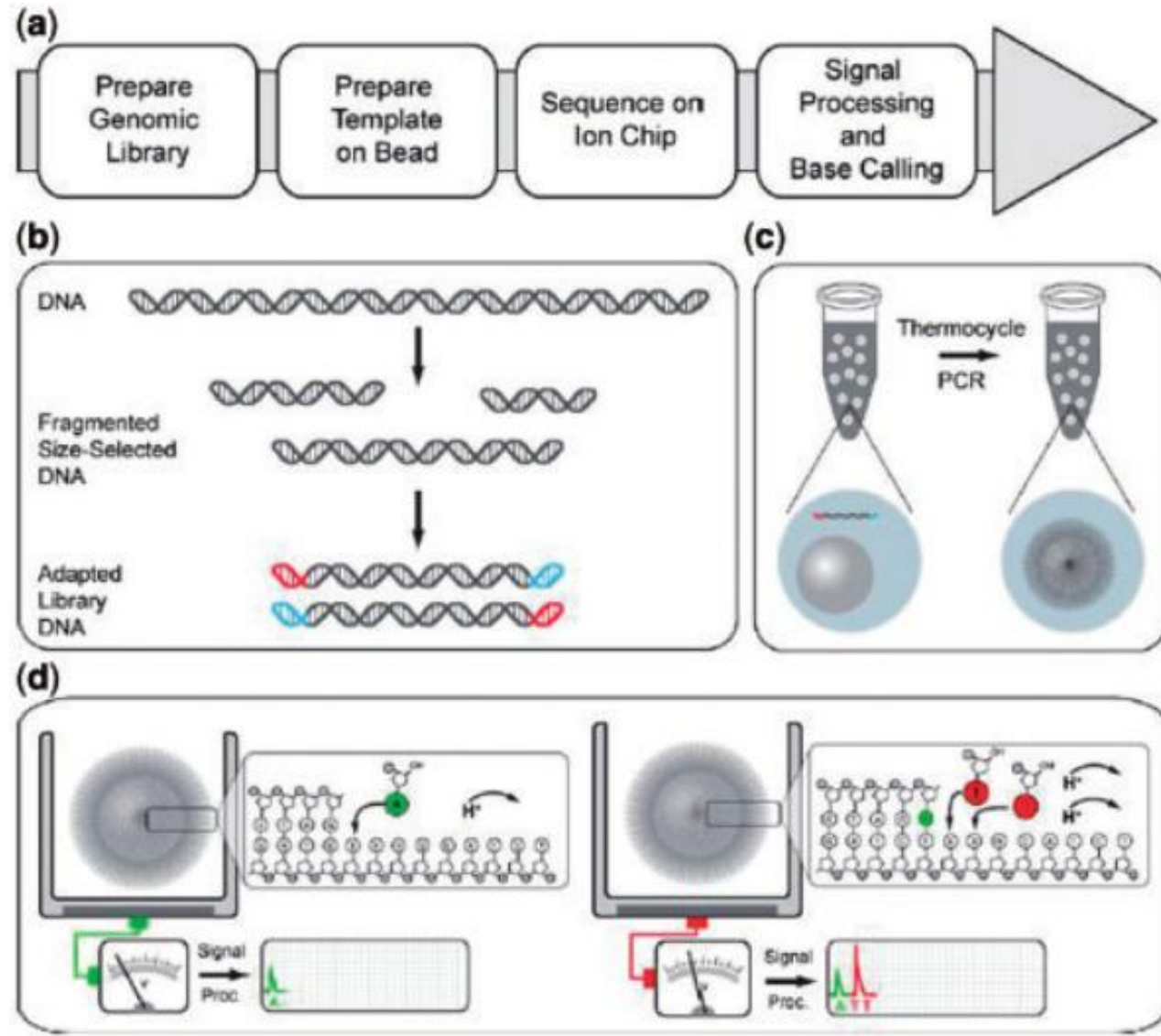


Next (2nd) Generation Sequencing

Platform: Ion

Semiconductor sequencing

- Read Length: 200-400 bp
- Read Fragments: Single-end, Paired-end (requires additional preparation steps)
- Throughput: 400MB to 26GB

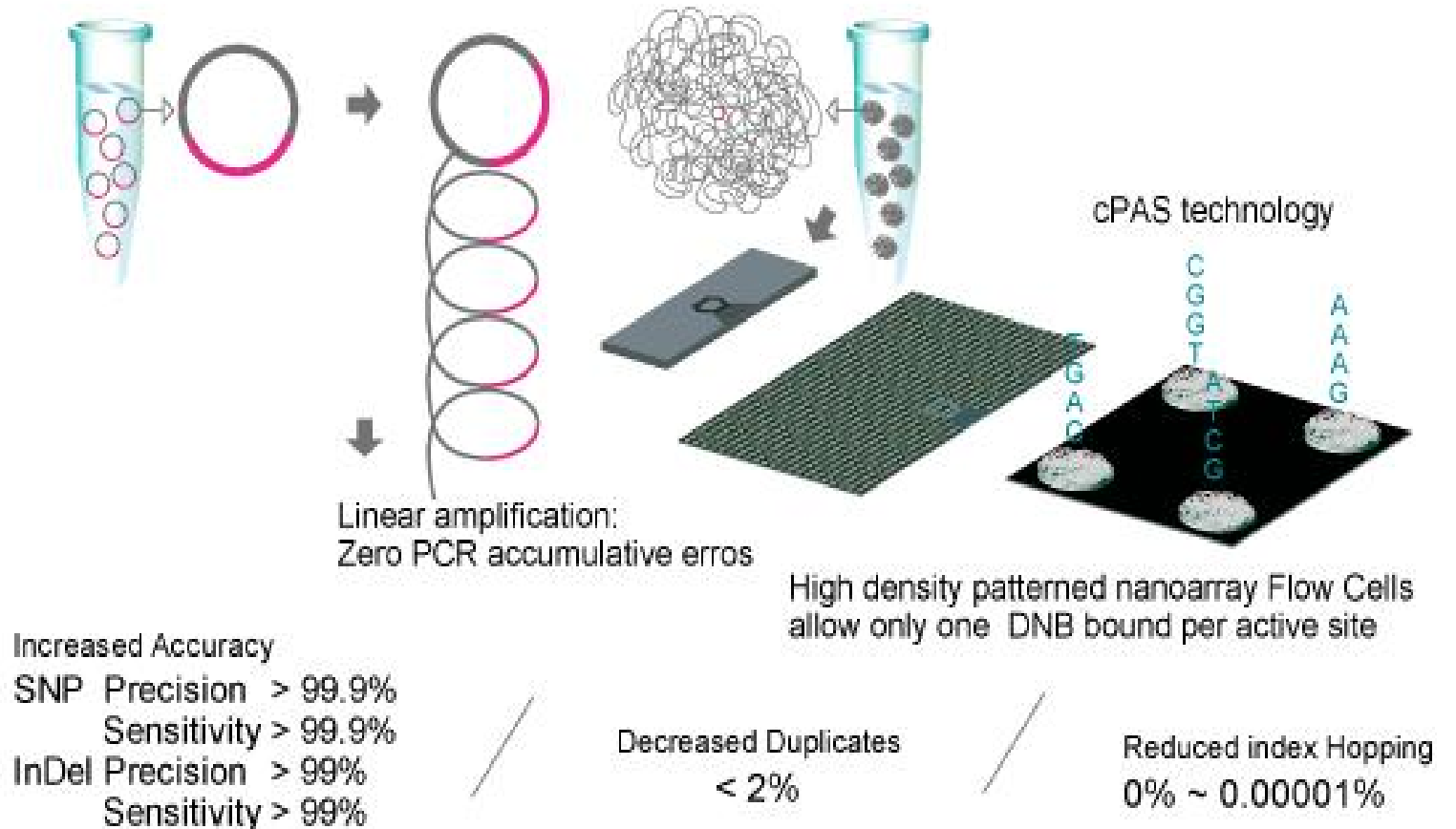


Next (2nd) Generation Sequencing

Platform: MGI / BGI

DNA nanoball sequencing (DNBSeq)

- Read Length: 50-150 bp
- Read Fragments: Paired-end, Single-end
- Throughput: 7.5GB to 76.8TB

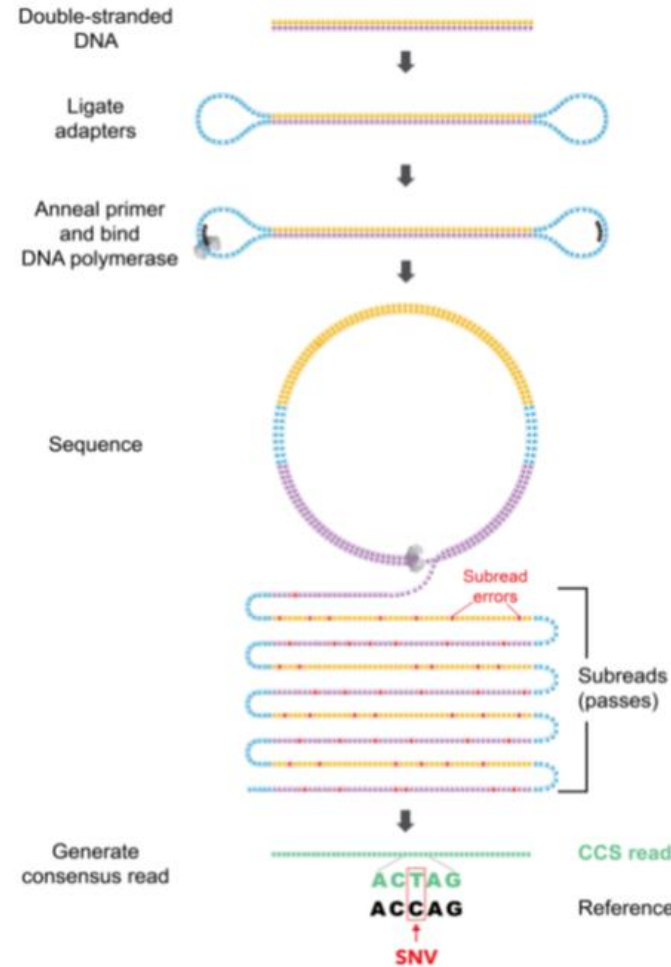


Third Generation Sequencing

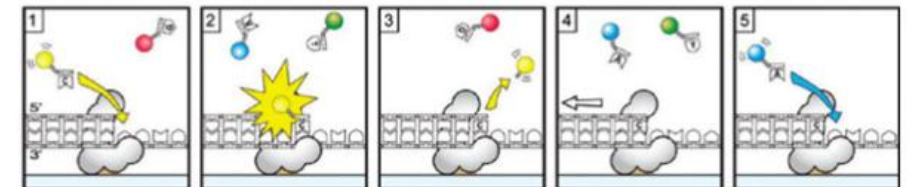
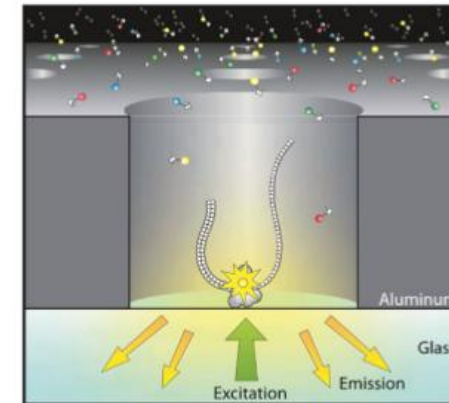
Platform: Pacific Biosciences (Pacbio)

Single Molecule, Real Time (SMRT) Sequencing

- Read Length: 15-25 kbp
- Throughput: 24Gb to 360Gb



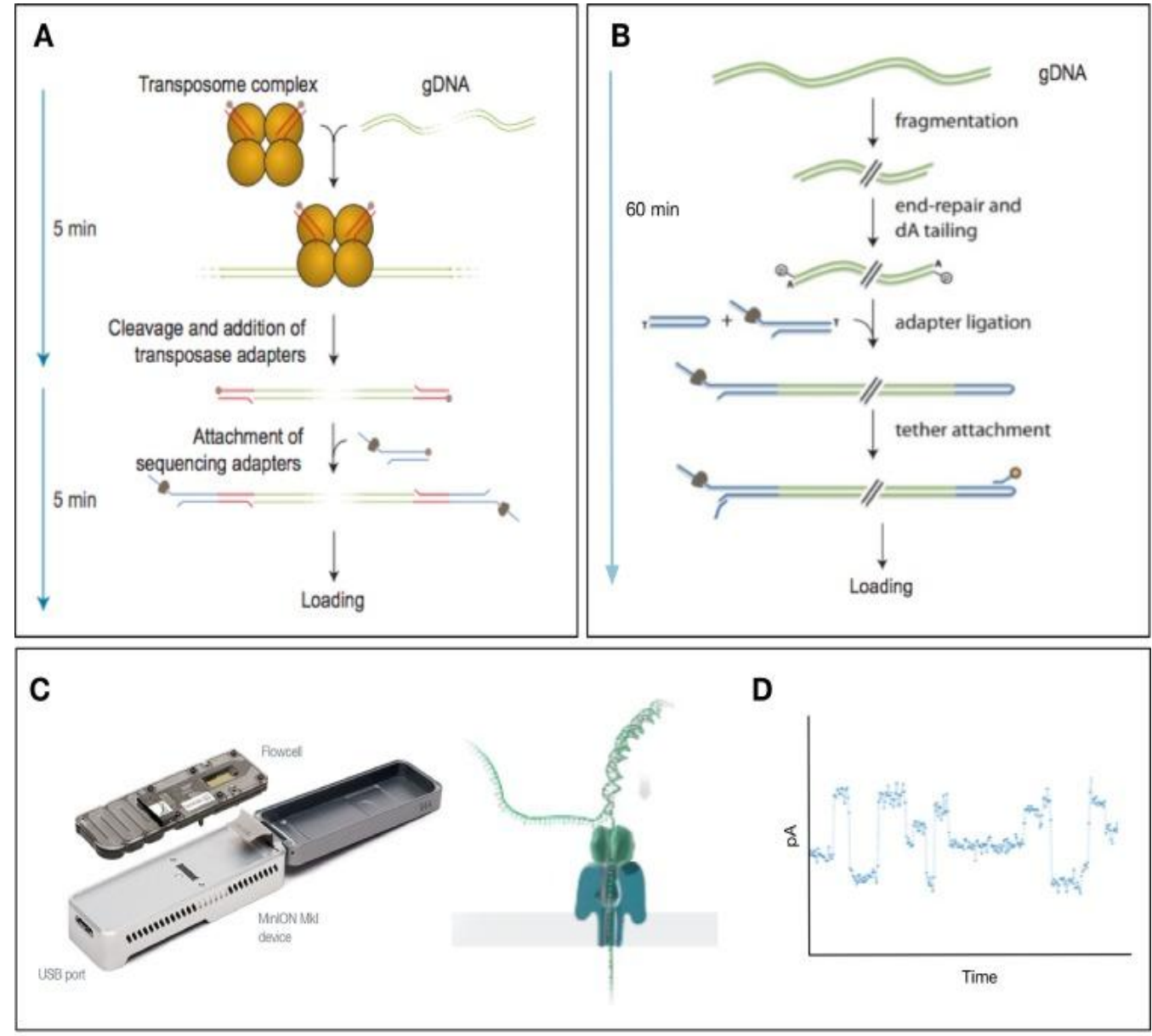
Single Molecule, Real-Time (SMRT) Sequencing technology delivers the highest consensus accuracy with unprecedented read lengths



Third Generation Sequencing

Platform: Oxford Nanopore Nanopore Sequencing

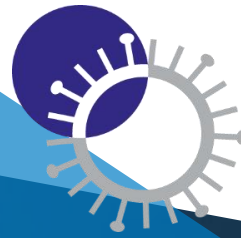
- Read Length:
Short Fragment: >20 bp
Standard: 5-30kbp
Ultra-long: >50kbp
- Throughput: 35Gb to 240Gb



Comparison of Technologies

	First generation	Second generation ^a	Third generation ^a
Fundamental technology	Size-separation of specifically end-labeled DNA fragments, produced by SBS or degradation	Wash-and-scan SBS	SBS, by degradation, or direct physical inspection of the DNA molecule
Resolution	Averaged across many copies of the DNA molecule being sequenced	Averaged across many copies of the DNA molecule being sequenced	Single-molecule resolution
Current raw read accuracy	High	High	Moderate
Current read length	Moderate (800–1000 bp)	Short, generally much shorter than Sanger sequencing	Long, 1000 bp and longer in commercial systems
Current throughput	Low	High	Moderate
Current cost	High cost per base Low cost per run	Low cost per base High cost per run	Low-to-moderate cost per base Low cost per run
RNA-sequencing method	cDNA sequencing	cDNA sequencing	Direct RNA sequencing and cDNA sequencing
Time from start of sequencing reaction to result	Hours	Days	Hours
Sample preparation	Moderately complex, PCR amplification not required	Complex, PCR amplification required	Ranges from complex to very simple depending on technology
Data analysis	Routine	Complex because of large data volumes and because short reads complicate assembly and alignment algorithms	Complex because of large data volumes and because technologies yield new types of information and new signal processing challenges
Primary results	Base calls with quality values	Base calls with quality values	Base calls with quality values, potentially other base information such as kinetics

^aThere are many TGS technologies in development but few have been reduced to practice. While there is significant potential of TGS to radically improve current throughput and read-length characteristics (among others), the ultimate practical limits of these technologies remain to be explored. Furthermore, there is active development of SGS technologies that will also improve read-length and throughput characteristics.



WHO Collaborating Centre
for Reference and
Research on Influenza
VIDRL



**World Health
Organization**
Philippines



KDCA
Korea Disease Control and
Prevention Agency

NGS DATA ANALYSIS AND CONSIDERATIONS



**CENTRE FOR
PATHOGEN
GENOMICS**



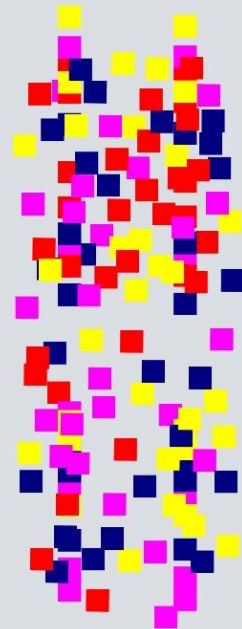
A joint venture between The University of Melbourne and The Royal Melbourne Hospital



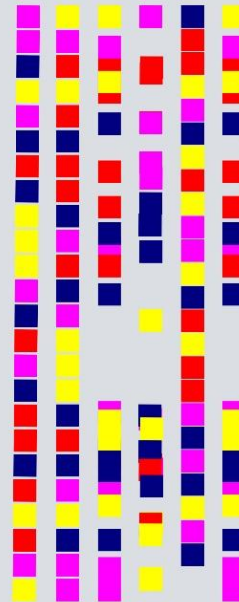
Big Data

Big data is high-volume, high-velocity and/or high-variety information assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and process automation.

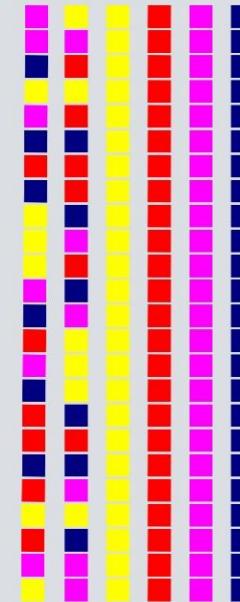
BIG DATA



ANALYTICS



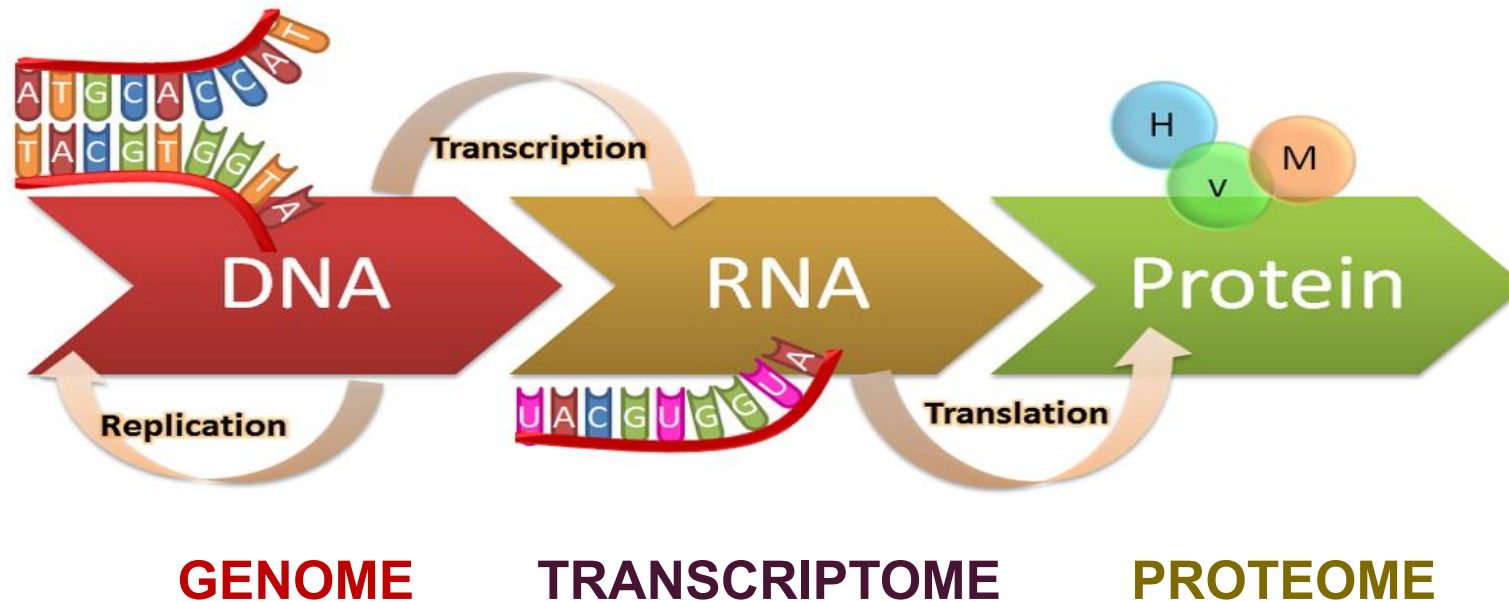
DECISIONS



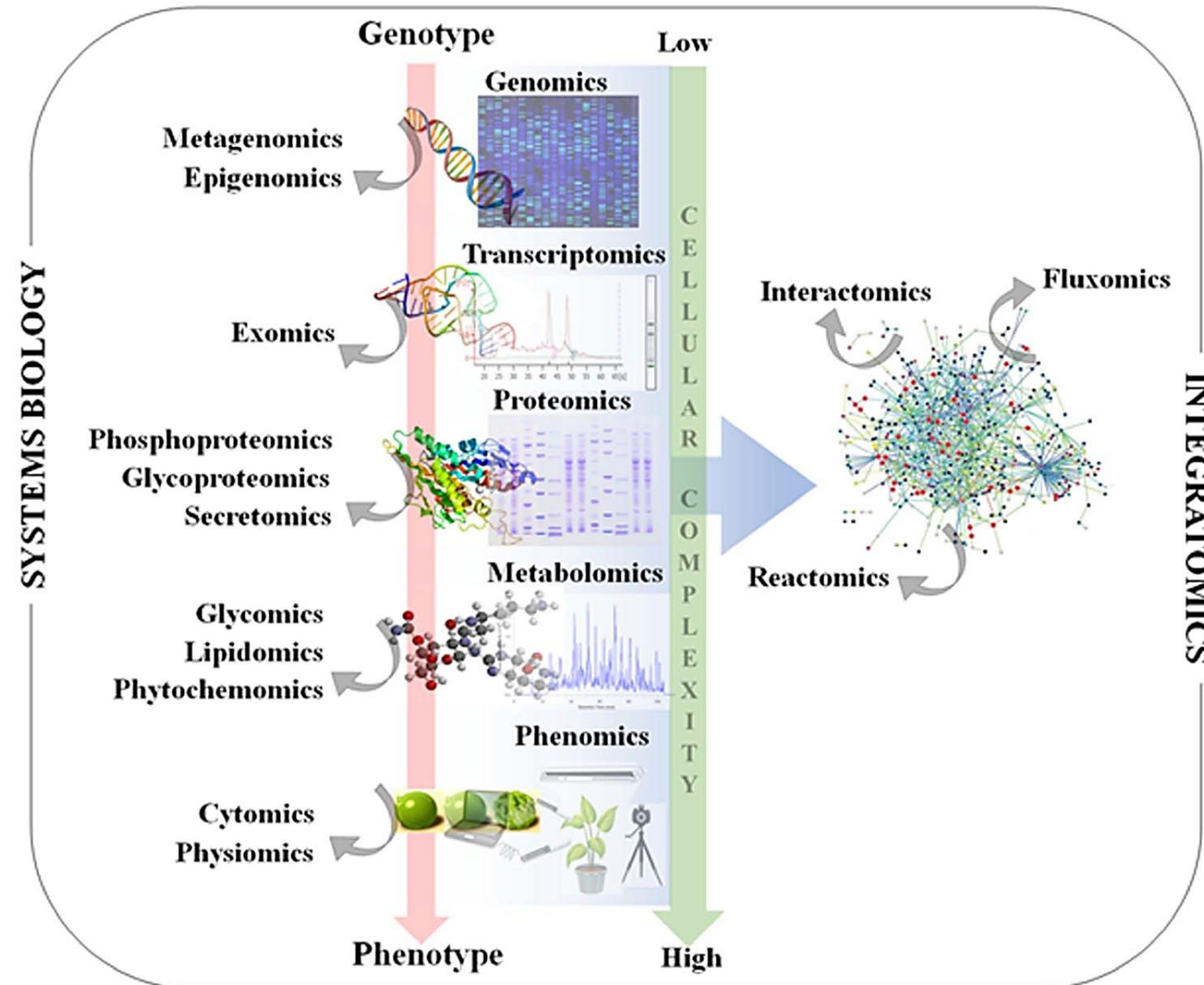
Omics - The Big Data Biology

-ome, -omics

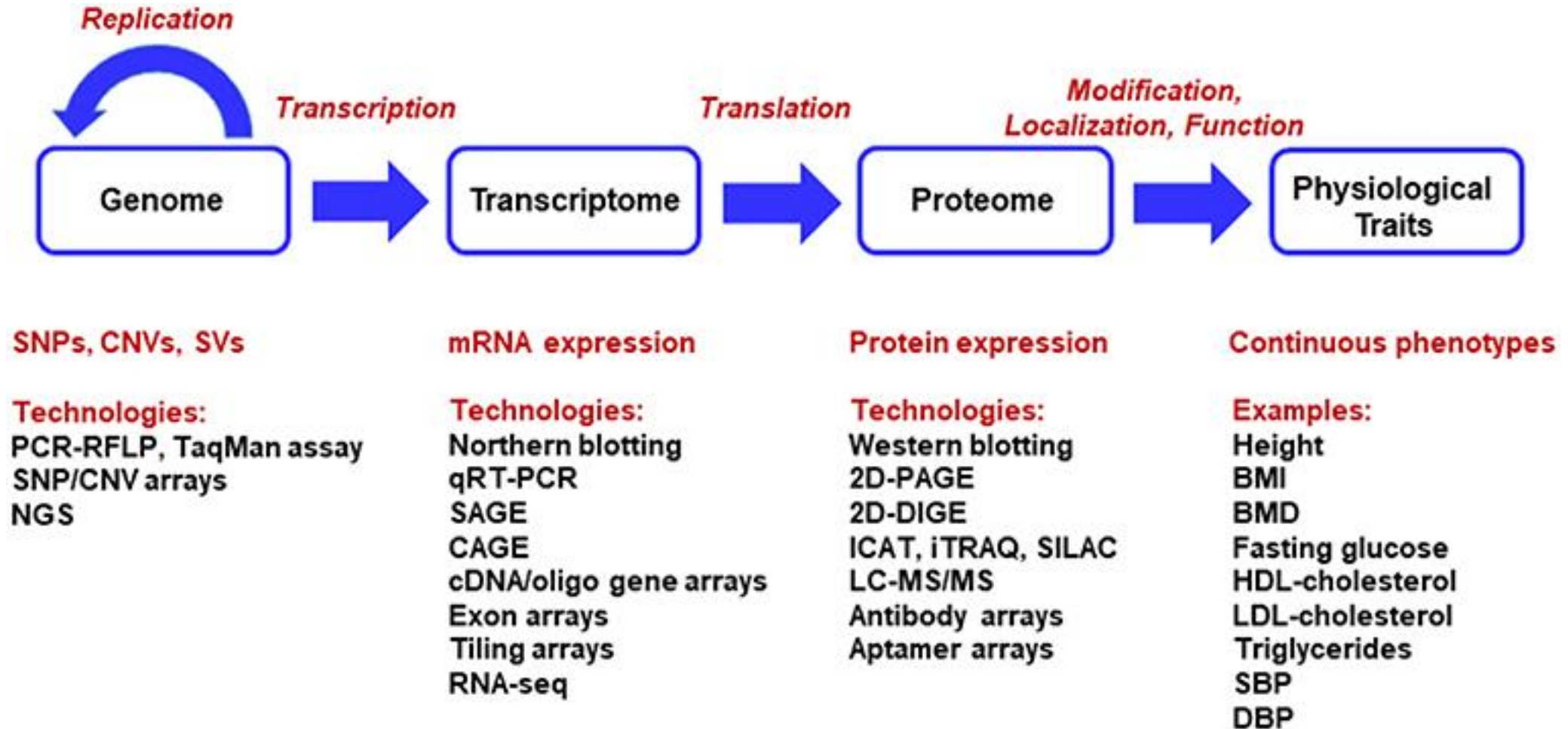
In molecular biology, suffixes used to refer to a ***totality*** of some sort



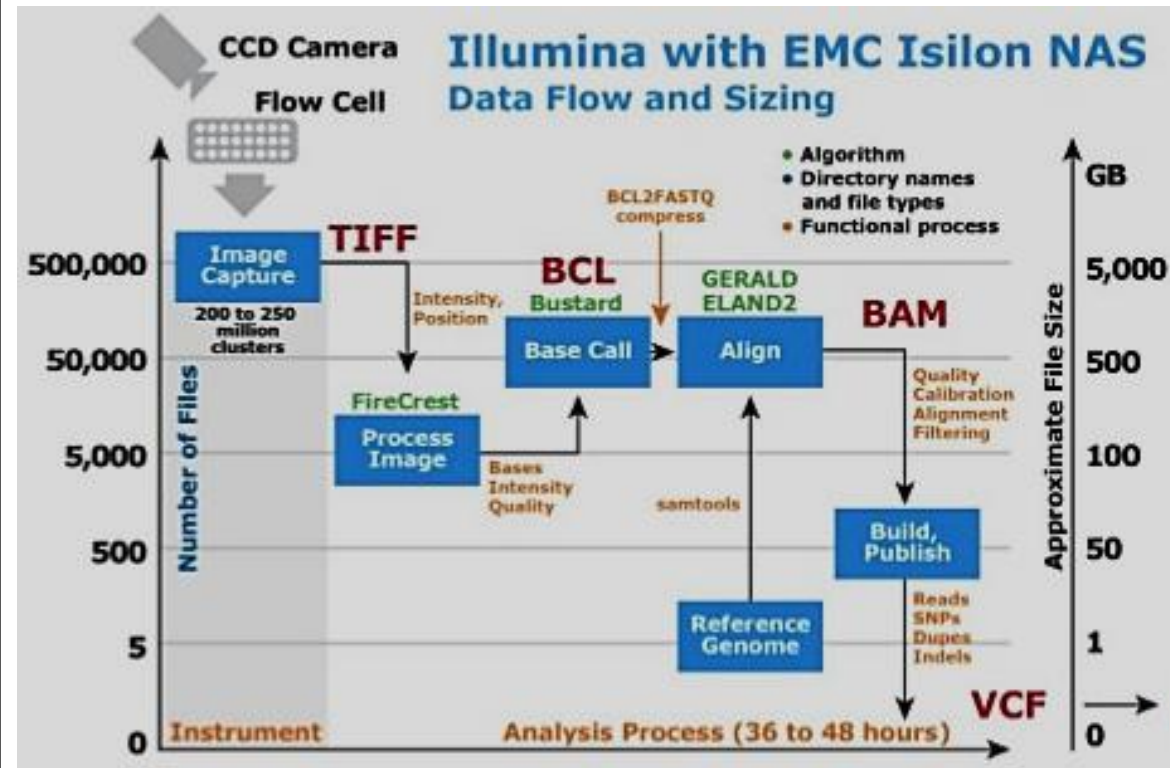
Omics - The Big Data Biology



Omics - The Big Data Biology



How huge are NGS data sets?



“... the process generates over **500,000 [files]** having aggregate size of greater than **5 Tb** over the course of the **48-hour run**.”

Figure 4. Data flow using Illumina NGS process

NGS production processes generate potentially millions of files with terabytes of aggregate storage impacting the capacity and manageability limits of existing file server structures.

Figure 4 shows the data flow including a file number and capacity summary of an actual NGS process using Illumina sequencer and Isilon scale-out NAS storage. As can be seen, the process generates over 500,000 having aggregate size of greater than 5 TB over the course of the 48-hour run.

How huge are NGS data sets?

[IEEE Spectr.](#) Author manuscript; available in PMC 2014 Jun 9.

Published in final edited form as:

IEEE Spectr. 2013 Jul; 50(7): 26–33.

doi: [10.1109/MSPEC.2013.6545119](https://doi.org/10.1109/MSPEC.2013.6545119)

PMCID: PMC4048922

NIHMSID: NIHMS563699

PMID: [24920863](https://pubmed.ncbi.nlm.nih.gov/24920863/)

The DNA Data Deluge

Fast, efficient genome sequencing machines are spewing out more data than geneticists can analyze

[Michael C. Schatz](#) and [Ben Langmead](#)

The roughly 2000 sequencing instruments in labs and hospitals around the world can collectively sequence 15 quadrillion nucleotides per year, which equals about 15 petabytes of compressed genetic data. A petabyte is 2^{50} bytes, or in round numbers, 1000 terabytes. To put this into perspective, if you were to write this data onto standard DVDs, the resulting stack would be more than 2 miles tall. And with sequencing capacity increasing at a rate of around three- to fivefold per year, next year the stack would be around 6 to 10 miles tall. At this rate, within the next five years the stack of DVDs could reach higher than the orbit of the International Space Station.

Handling NGS data sets

Bioinformatics (Applied)

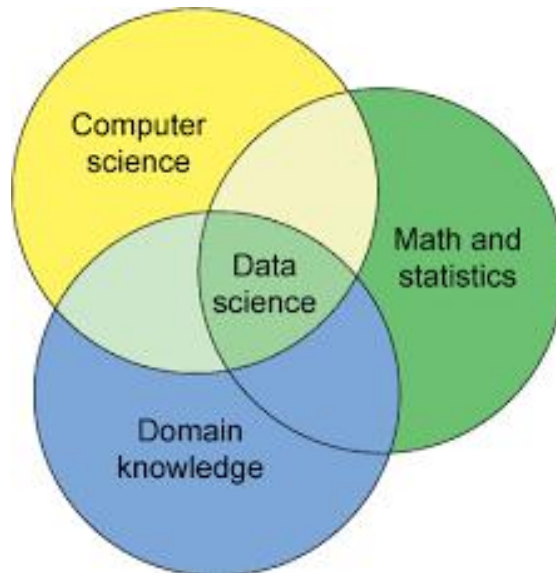
“Research, development, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze, or visualize such data.”

Computational Biology (Theoretical)

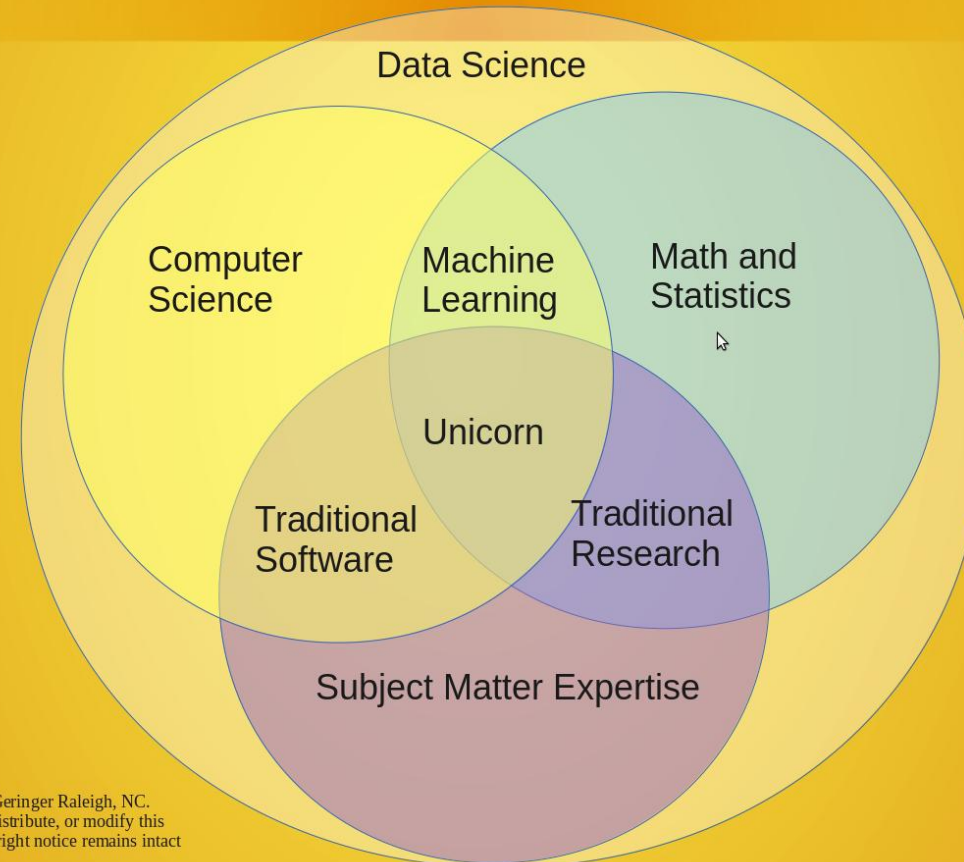
“The development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems.”

Handling NGS data sets

Data Science



Data Science Venn Diagram v2.0



Copyright © 2014 by Steven Geringer Raleigh, NC.
Permission is granted to use, distribute, or modify this
image, provided that this copyright notice remains intact

NGS Analogy

Multiple copies of
template genome



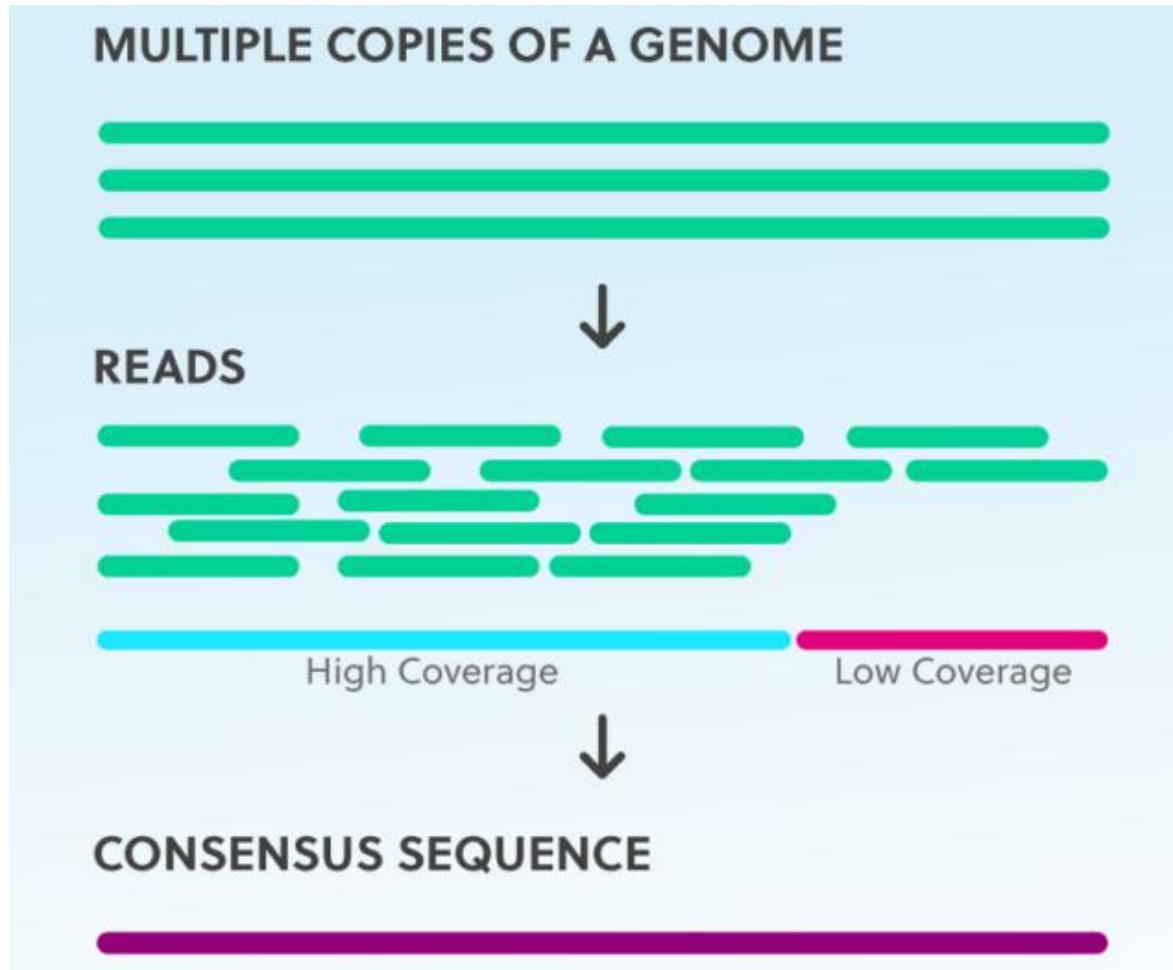
Sequence reads /
fragments



Reconstructed sequence with
readable features



NGS Analysis



Sequence Reads

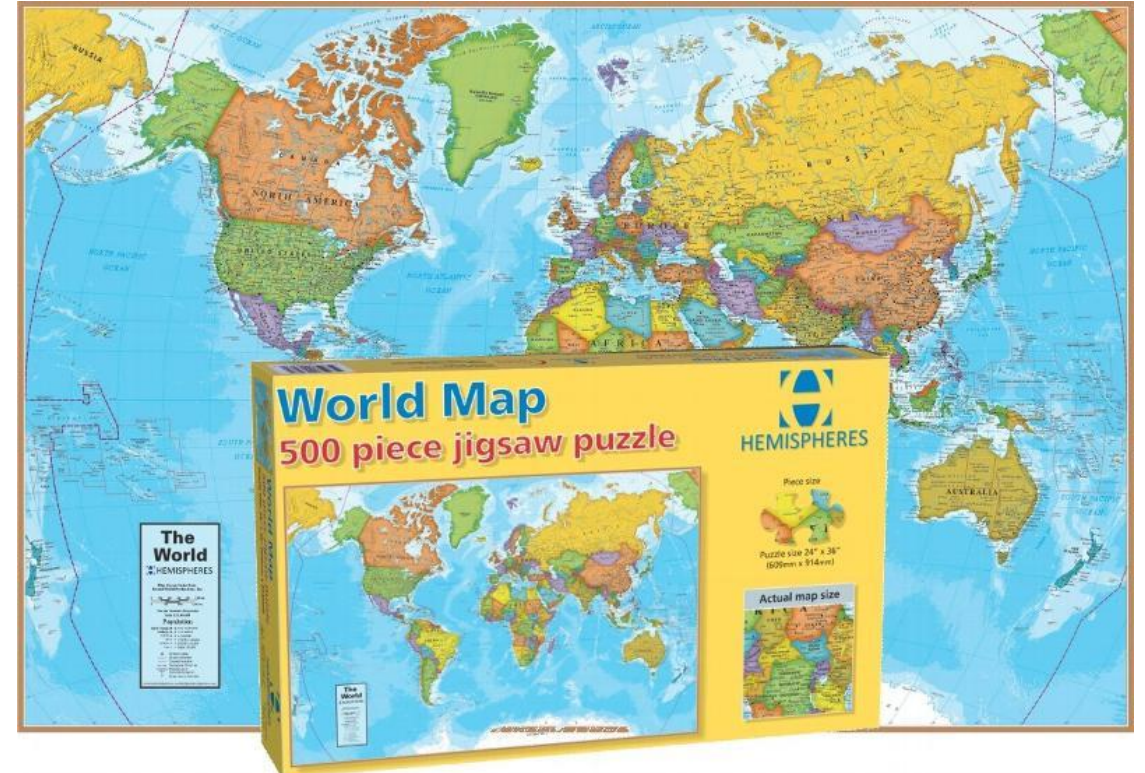
- Fragments of the original genome or any template sequence (FASTQ format)
- One of the initial goals of bioinformatics is to reconstruct the template sequence based on the read fragment information (sequence assembly)

NGS Analysis



https://www.123rf.com/photo_104112068_stock-illustration-missing-jigsaw-puzzle-pieces-in-unfinished-work-concept-white-pattern-texture-background-3d-illustration.html

De novo



<https://www.puzzlewarehouse.com/World-Map-hmp01.html>

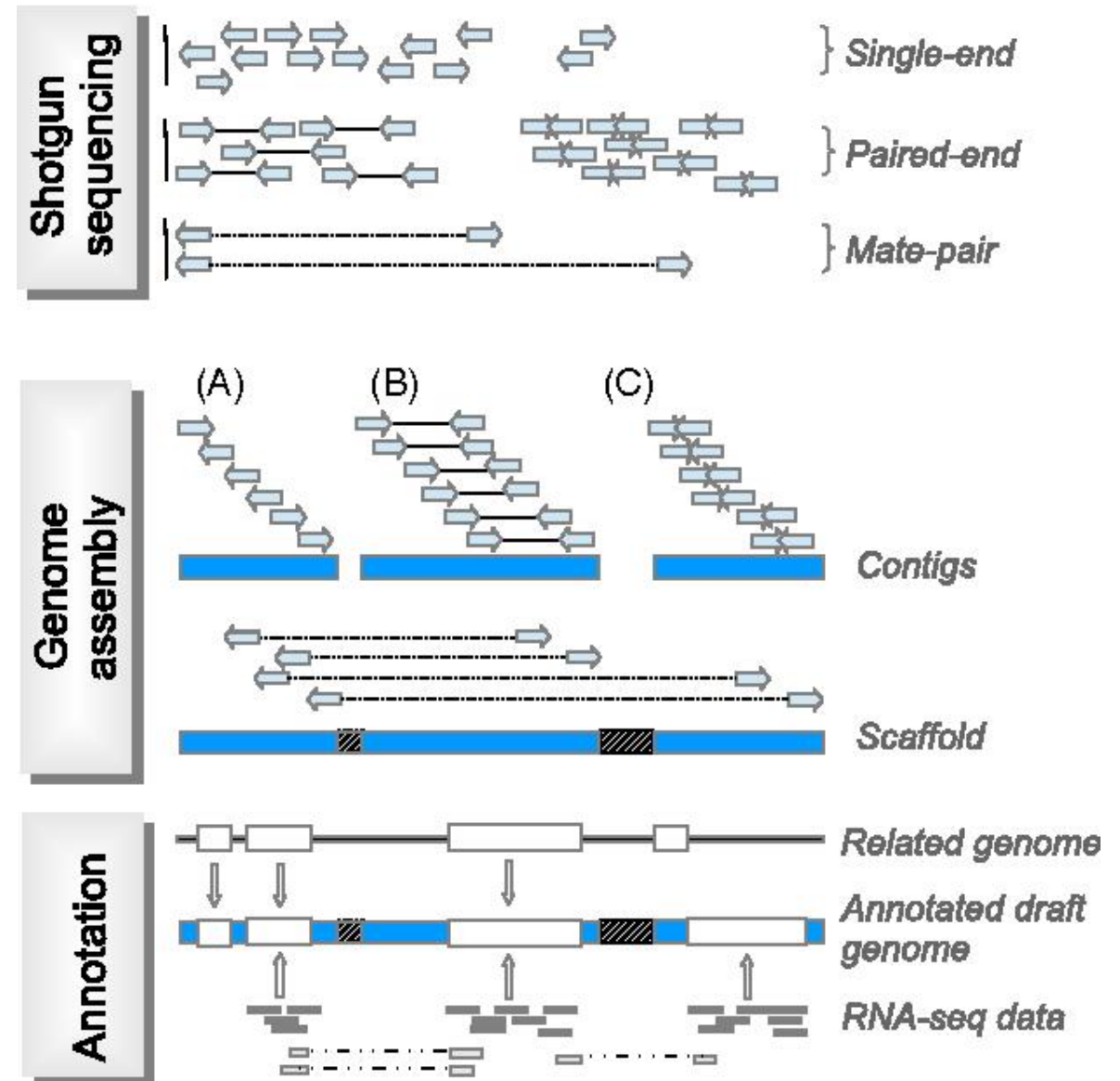
Reference-guided

NGS Analysis

De novo Workflow

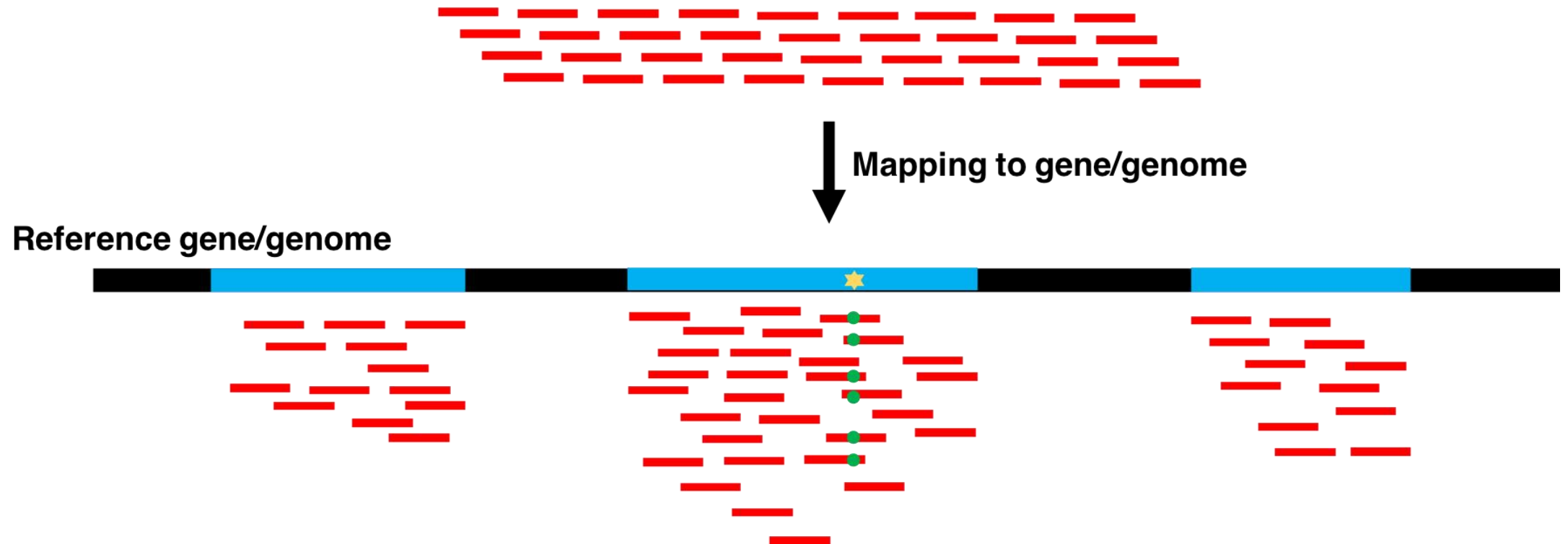
Output: FASTA file

Output: GFF, GenBank, EMBL



NGS Analysis

Reference-Guided Workflow



Output: SAM, BAM

NGS Analysis

Note on file format standardization

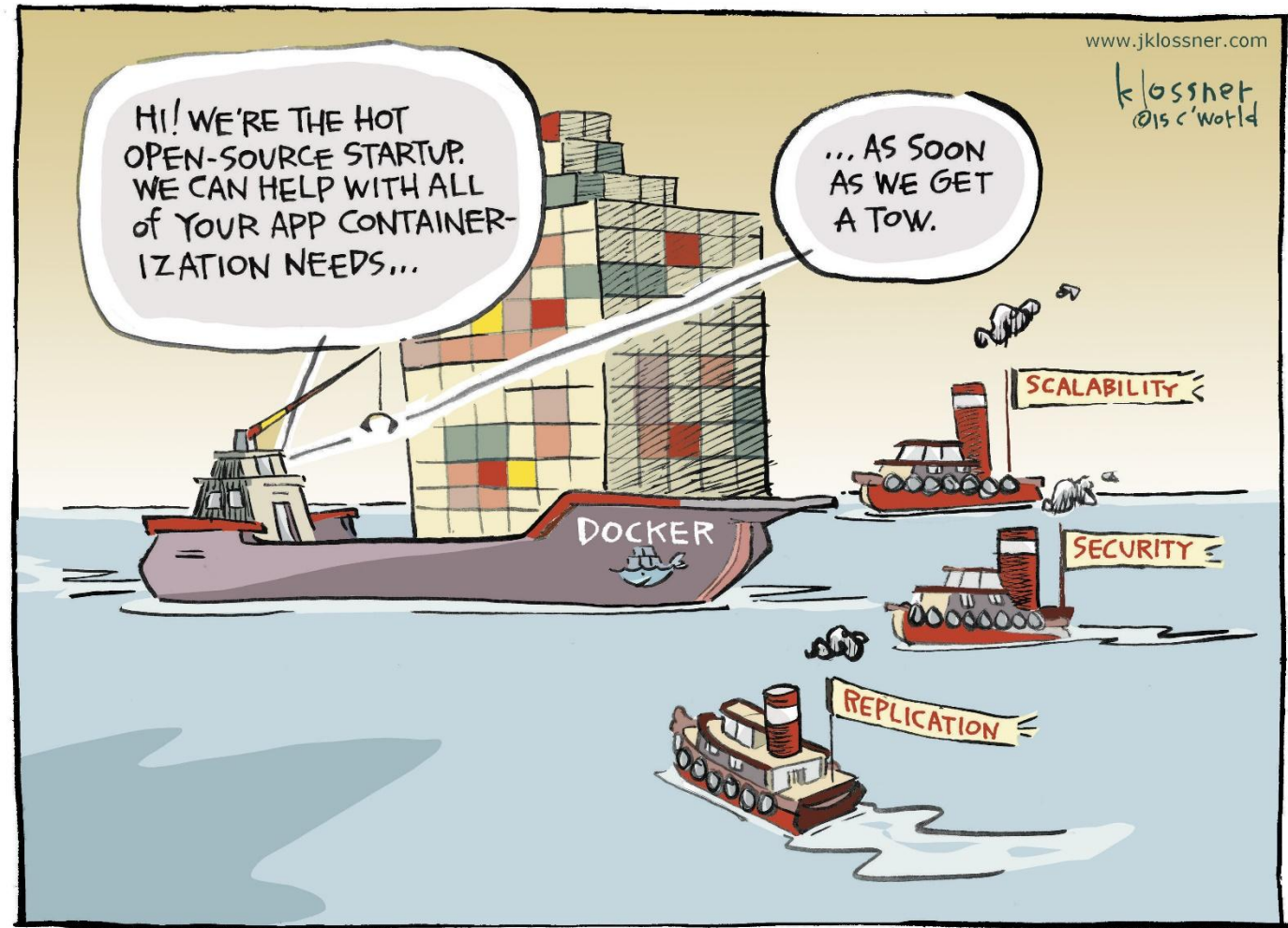
- In bioinformatics, many of the processes involve input and output files with standard formats
- File format standardization enables process automation

Process	Input	Output
Quality Control	FASTQ	FASTQ
Sequence Assembly	FASTQ	FASTA
Feature Annotation	FASTA	GFF, GenBank, EMBL
Read Mapping	FASTQ	SAM / BAM
Variant Calling	SAM / BAM	VCF
Alignment	Multi-FASTA	Multi-FASTA / Phylip
Phylogenetics	Multi-FASTA / Phylip	Newick / Nexus

Considerations in NGS Data Analysis

Infrastructure

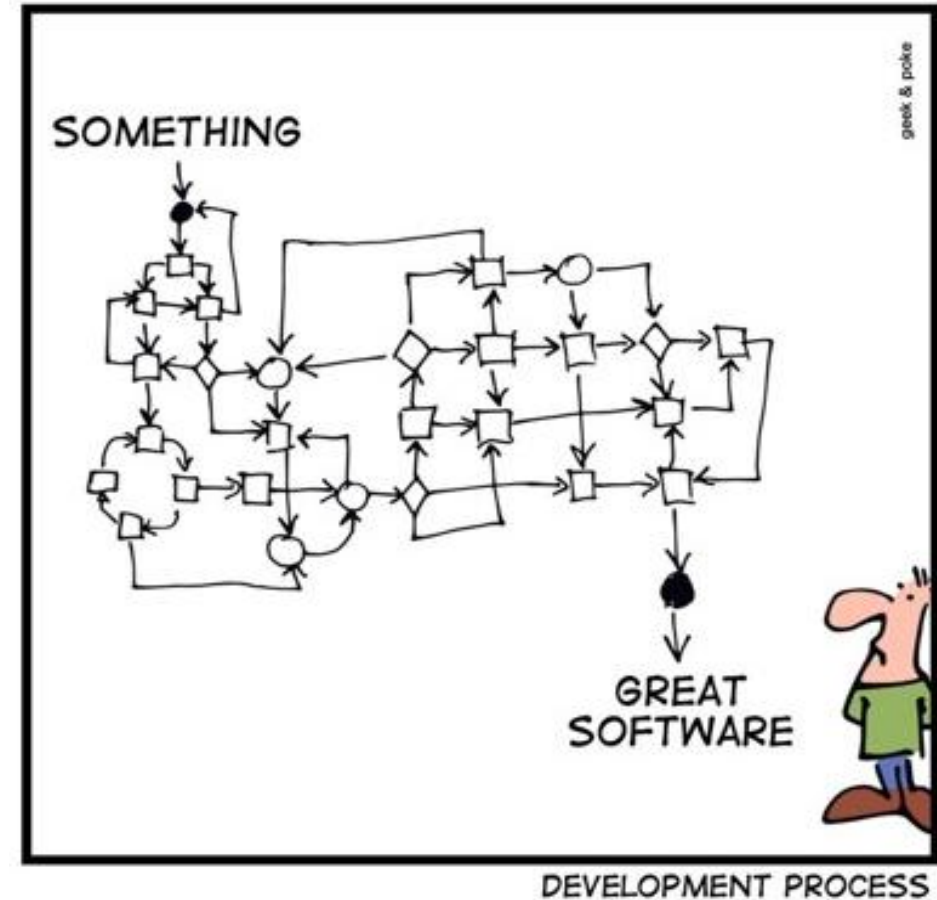
- Computing power
- Storage
- Network
- Security



Considerations in NGS Data Analysis

Algorithms / Implementation

- Statistical / Mathematical models (accuracy)
- Databases and Data Management
- Speed
- Usage of resources



Considerations in NGS Data Analysis

Skillset / Learning Curve / Manpower

- Biologists trying to understand computational concepts
- Computational scientists trying to understand biology

Biologists be like ...



Computational scientists be like ...



QUESTIONS?

fatablizo@up.edu.ph

bioinformatics@pgc.up.edu.ph



KDCA

Korea Disease Control and
Prevention Agency



A joint venture between The University of Melbourne and The Royal Melbourne Hospital