**Introduction to Bioinformatics**

# Next Gen Sequencing

Stevan Radanović

# Resources

- Illumina inc., https://www.youtube.com/channel/UCxWMU29FF4kIG8YmQf6Zv0g
- AWS https://aws.amazon.com/documentation/ec2/

# Refresher

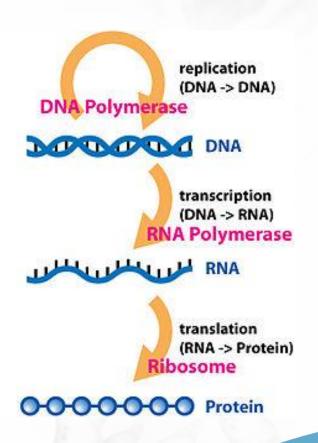- Central dogma: DNA makes RNA makes Proteins
- PCR



Polymerase chain reaction - PCR

original DNA to be replicated

DNA primer

nucleotide

5′ 3′
3′ 5′

1 Denaturation at 94-96℃
2 Annealing at ~68℃
3 Elongation at ca. 72 ℃



replication
(DNA -> DNA)

**DNA Polymerase**

**DNA**

transcription
(DNA -> RNA)

**RNA Polymerase**

**RNA**

translation
(RNA -> Protein)

**Ribosome**

**Protein**

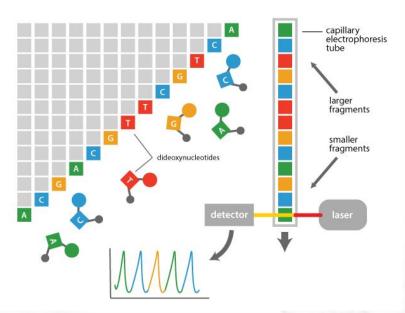# Refresher
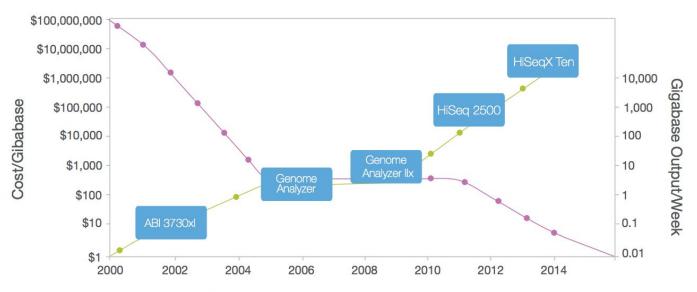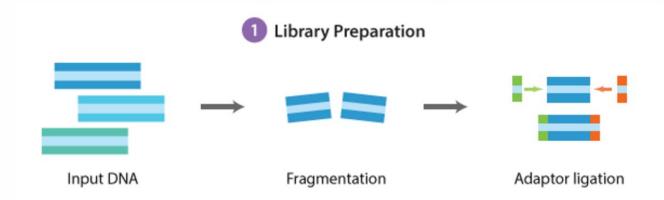
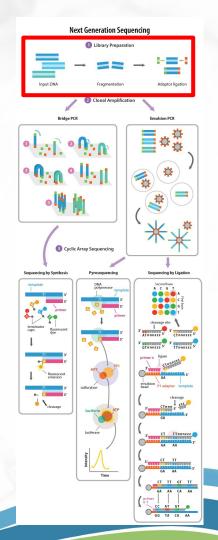- Sanger sequencing

# Sequencing



**Figure 1: Sequencing Cost and Data Output Since 2000**—The dramatic rise of data output and concurrent falling cost of sequencing since 2000. The Y-axes on both sides of the graph are logarithmic.
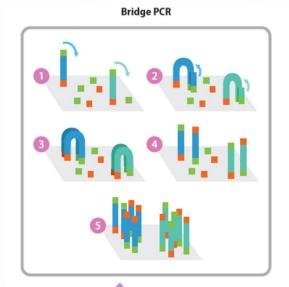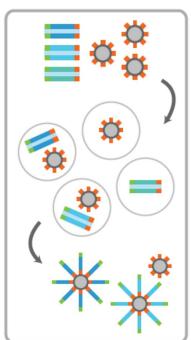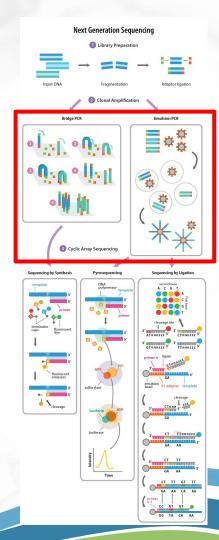
# Next Gen Sequencing

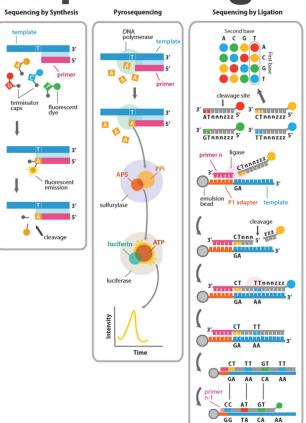# Next Gen Sequencing

# Next Gen Sequencing

# Illumina Sequencing



A. Library Preparation

Genomic DNA

Fragmentation

Adapters

Ligation

Sequencing Library

NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

# Illumina Sequencing



Transposome
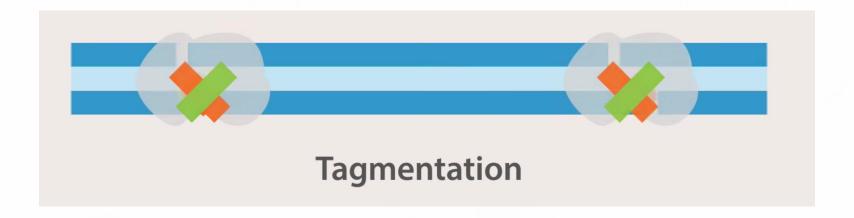
Makes cuts in DNA
_____
Inserts a portion of itself in DNA

# Illumina Sequencing



Tagmentation

# Illumina Sequencing



Tagmentation

# Illumina Sequencing



PCR Amplification

# Illumina Sequencing

# Illumina Sequencing



## A. Library Preparation

Genomic DNA

↓ Fragmentation

Adapters

↓ Ligation

Sequencing Library

NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

## A. Cluster Amplification

Flow Cell

Bridge Amplification Cycles

1 2 3 4

Clusters

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

# Illumina Sequencing

# Data Analysis

# Illumina Sequencing



**Figure 4: Paired-End Sequencing and Alignment**—Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.
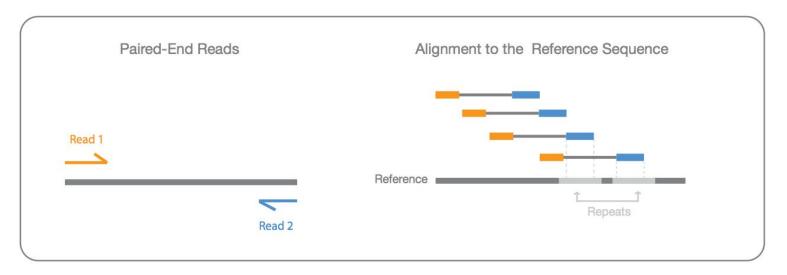
# Illumina Sequencing



**Figure 7: *De Novo* Assembly with Mate Pairs**—Using a combination of short and long insert sizes with paired-end sequencing results in maximal coverage of the genome for *de novo* assembly. Because larger inserts can pair reads across greater distances, they provide a better ability to read through highly repetitive sequences and regions where large structural rearrangements have occurred. Shorter inserts sequenced at higher depths can fill in gaps missed by larger inserts sequenced at lower depths. Thus a diverse library of short and long inserts results in better *de novo* assembly, leading to fewer gaps, larger contigs, and greater accuracy of the final consensus sequence.
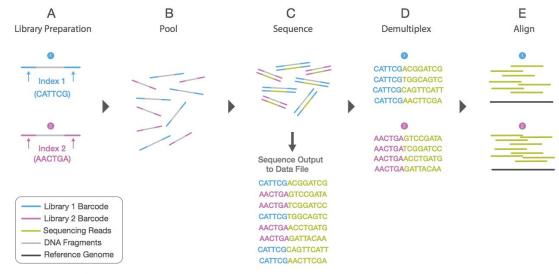
# Illumina Sequencing



Figure 5: Library Multiplexing Overview.

a. Two distinct libraries are attached to unique index sequences. Index sequences are attached during library preparation.

b. Libraries are pooled together and loaded into the same flow cell lane.

c. Libraries are sequenced together during a single instrument run. All sequences are exported to a single output file.

d. A demultiplexing algorithm sorts the reads into different files according to their indexes.

e. Each set of reads is aligned to the appropriate reference sequence.

# Illumina Sequencing



**MiSeq Series**

Small genome, amplicon and targeted gene panel sequencing.

**NextSeq Series**

Everyday genome, exome transcriptome sequencing, and more.
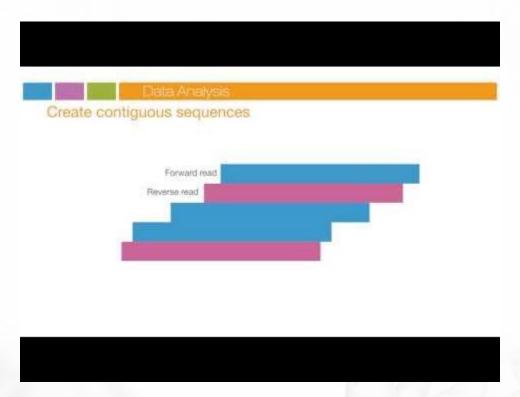
**HiSeq Series**

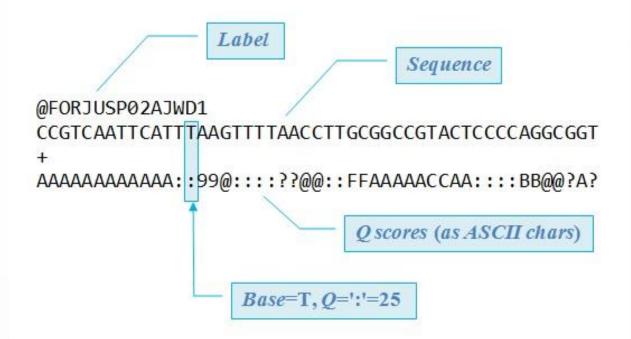Production-scale genome, exome, transcriptome sequencing and more.

**HiSeq X Series**

Population- and production-scale human whole-genome sequencing.
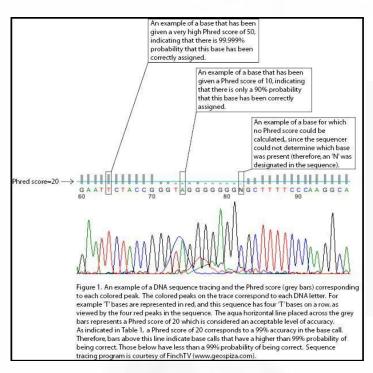
SevenBridges
genomics

# Illumina Sequencing

# FASTQ File Format

# FASTQ File Format



Figure 1. An example of a DNA sequence tracing and the Phred score (grey bars) corresponding to each colored peak. The colored peaks on the trace correspond to each DNA letter. For example 'T' bases are represented in red, and this sequence has four 'T' bases on a row, as viewed by the four red peaks in the sequence. The aqua horizontal line placed across the grey bars represents a Phred score of 20 which is considered an acceptable level of accuracy. As indicated in Table 1, a Phred score of 20 corresponds to a 99% accuracy in the base call. Therefore, bars above this line indicate base calls that have a higher than 99% probability of being correct. Those below have less than a 99% probability of being correct. Sequence tracing program is courtesy of FinchTV (www.geospiza.com).