

# Introduction to Bioinformatics

Instructor: Sakhaa Alsaedi

TA: Ebtihal Hani

[Sakhaa.Alsaedi@kaust.edu.sa](mailto:Sakhaa.Alsaedi@kaust.edu.sa)

---

Day 1: Sequence Analysis  
9<sup>th</sup> January 2024

# Course Schedule and Content Overview



Day	Content Overview
Day 1	<ul style="list-style-type: none"><li>Introduction to Bioinformatics</li><li>Sequence analysis</li></ul>
Day 2	<ul style="list-style-type: none"><li>Genetic variant analysis</li></ul>
Day 3	<ul style="list-style-type: none"><li>Diagnosing human genetic disease</li></ul>
Day 4	<ul style="list-style-type: none"><li>Genome Assembly</li></ul>
Day 5	<ul style="list-style-type: none"><li>Transcriptomics + Exam</li></ul>

## Course Timing and Duration

- Training Sessions:** Total of 6 hours,
  - Morning (9:00 AM - 12:00 PM)
  - Afternoon (2:00 PM - 5:00 PM)
- Breaks:** 12:00 PM to 2:00 PM (Prayer and Lunch)
- Operational Days:** 5 consecutive days, (9th - 13th of Jan.)

**GitHub:** <https://github.com/Sakhaa-Alsaedi/Bioinformatics-/blob/main/README.md>

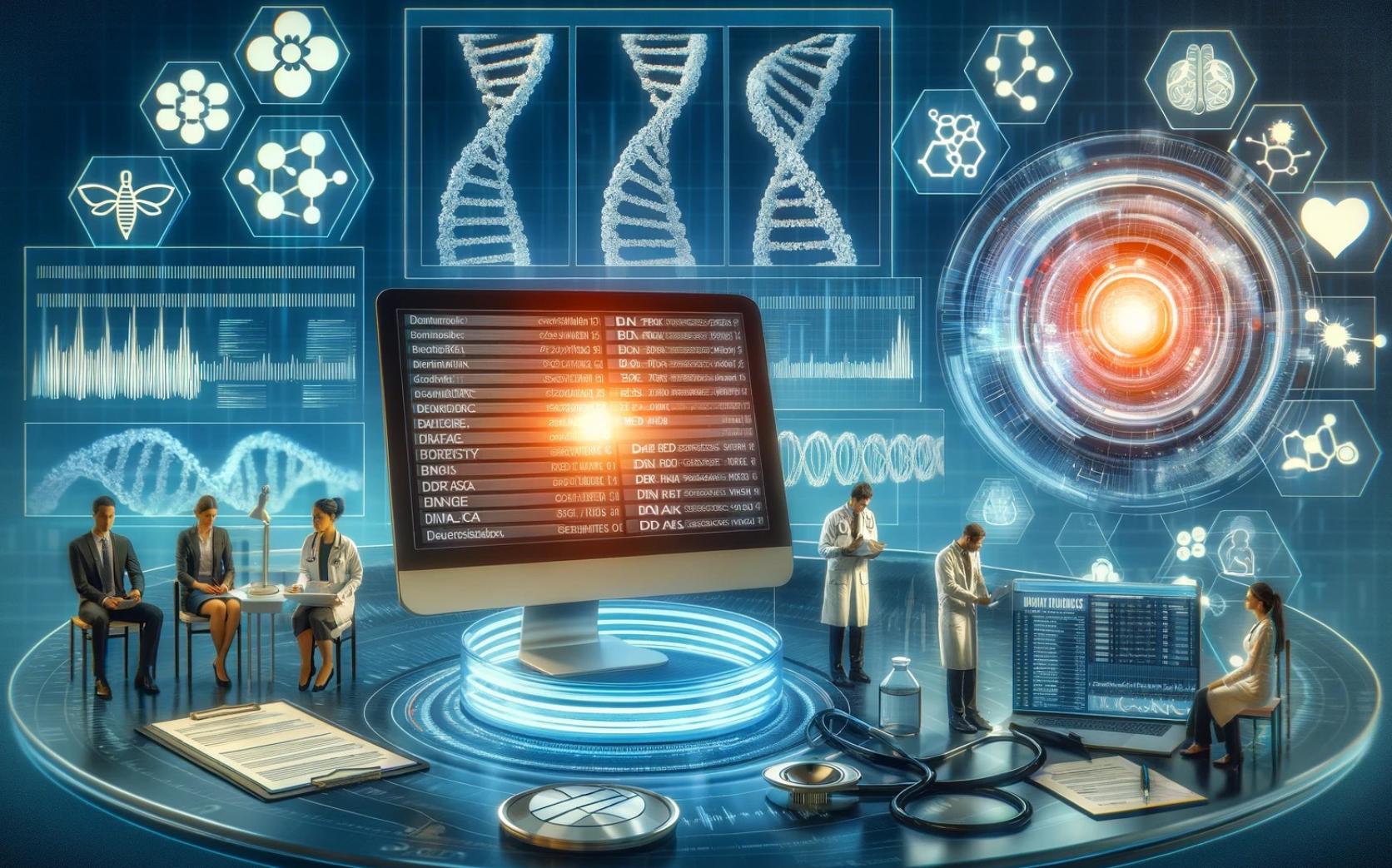


# Introduction to Bioinformatics

## Understanding the Digital Frontier in Biomedicine

---

### Part 1: Introduction

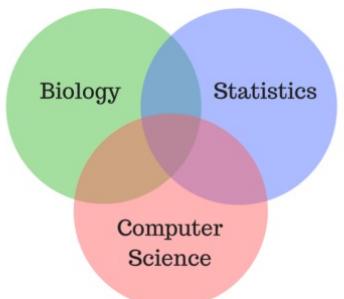


# Bioinformatics

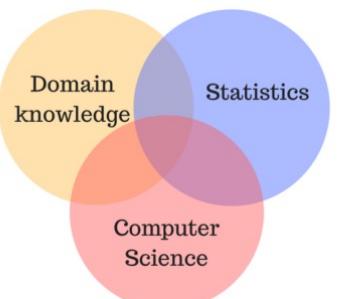
The science of collecting and analysing complex and digital **biological data**

→ Understanding the biological significance behind vast amounts of data in different applications

## Bioinformatics



## Data Science



# Biological Data and Complexity



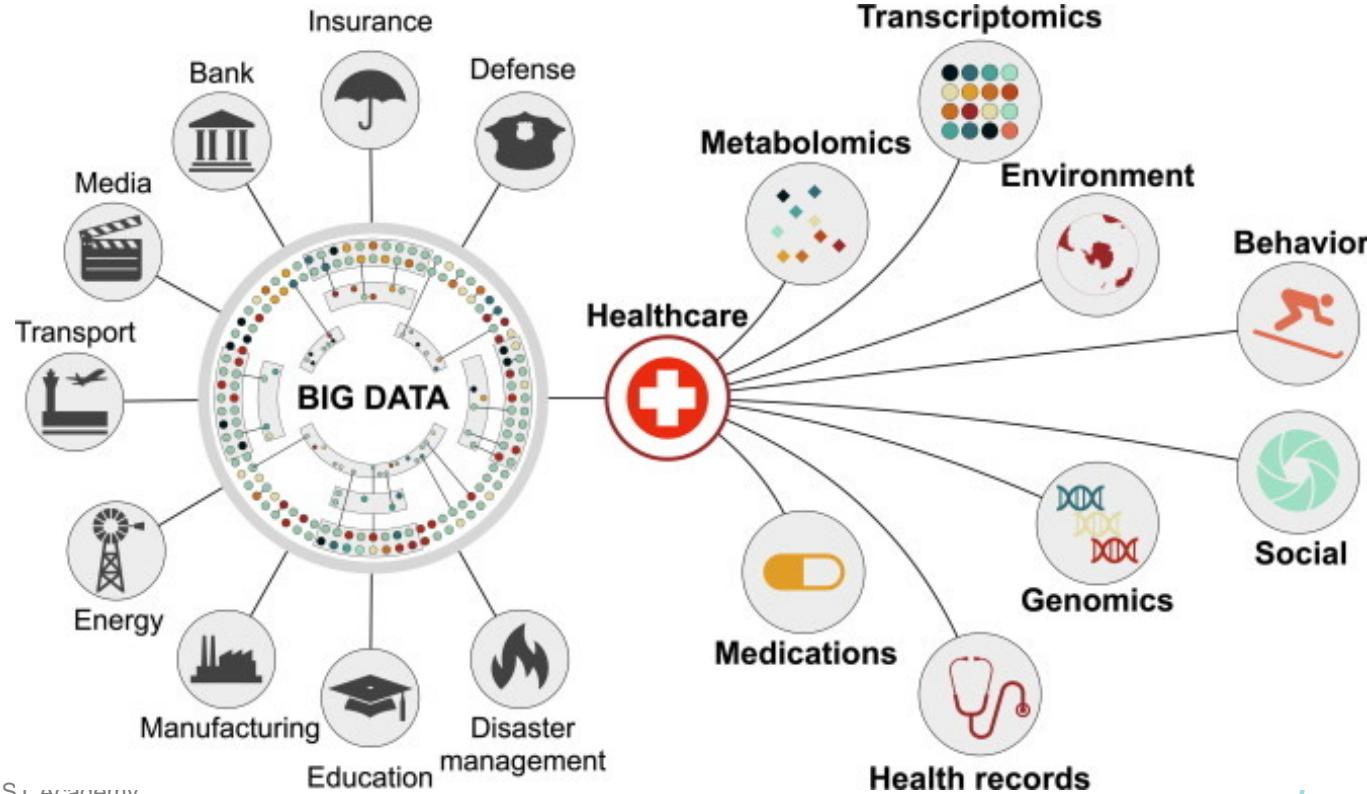
- Complex biological data refers to the intricately detailed, highly interconnected, and voluminous datasets that are characteristic of biological research and studies.

## Dimensions of Biological Data Complexity

- Diversity of data Types, high dimensionality, and noise and uncertainty.
- Heterogeneity, Computational Intensity, and Rapid Data Generation.



# Example: Biomedical Data → Big Data



# Example: Complex Molecular Networks



- Protein-Protein interactions
- Protein-DNA interactions
- Genetic interactions
- Metabolic reactions
- Co-expression interactions
- Text mining interactions
- Association Networks
- Etc.

# Digitizing Biomedical Data

- Converting traditional biological and medical data into digital formats for computational analysis.



## Why study computational algorithms?

- What is possible?
- What is practical?
- Where will the next contribution be?

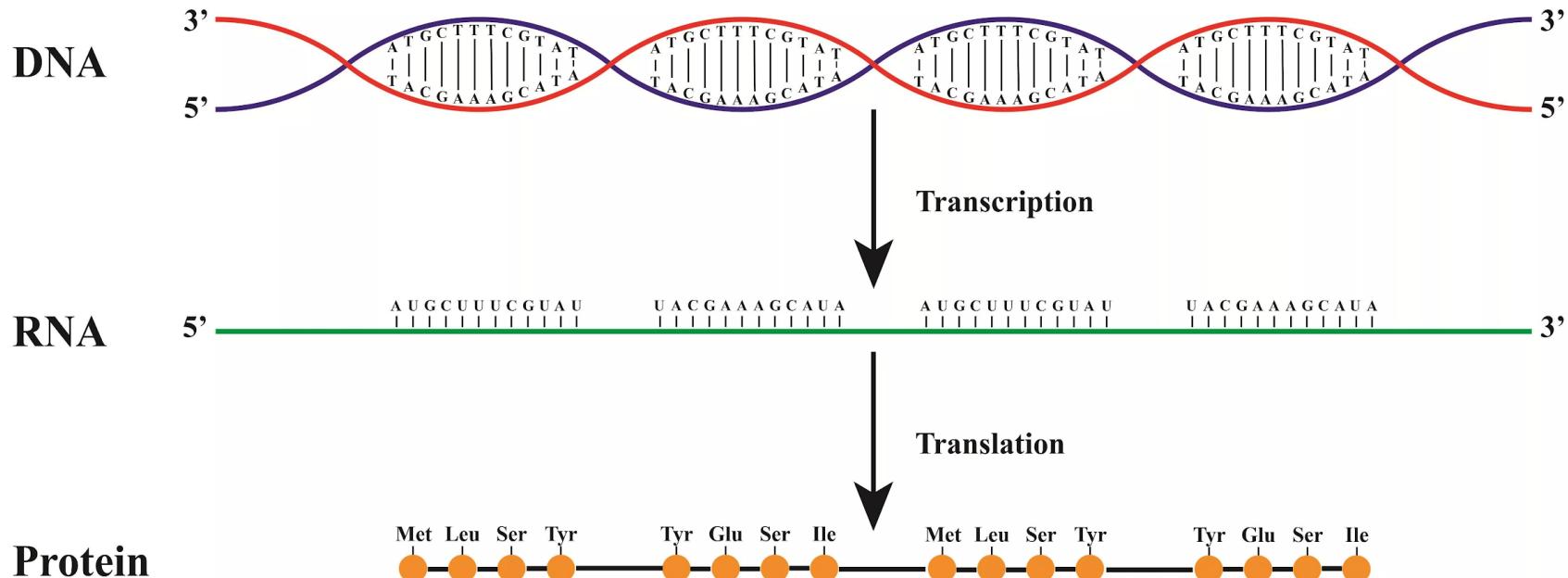
## The Need for Digitizing Medical Data

- Enhances accessibility, storage, and analysis capabilities.
- Improved diagnosis, treatment planning, and predictive medicine.

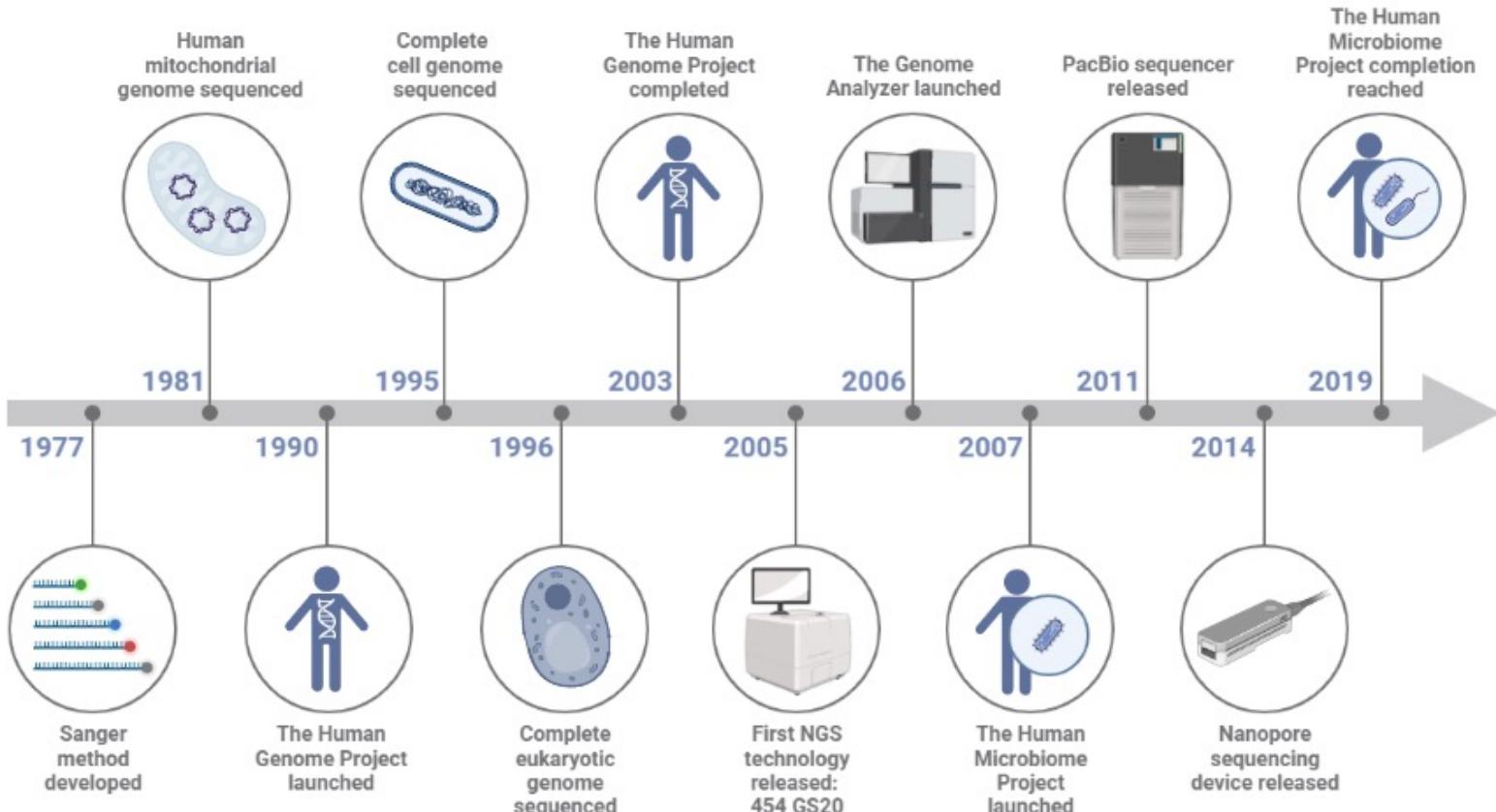
# DNA Sequencing



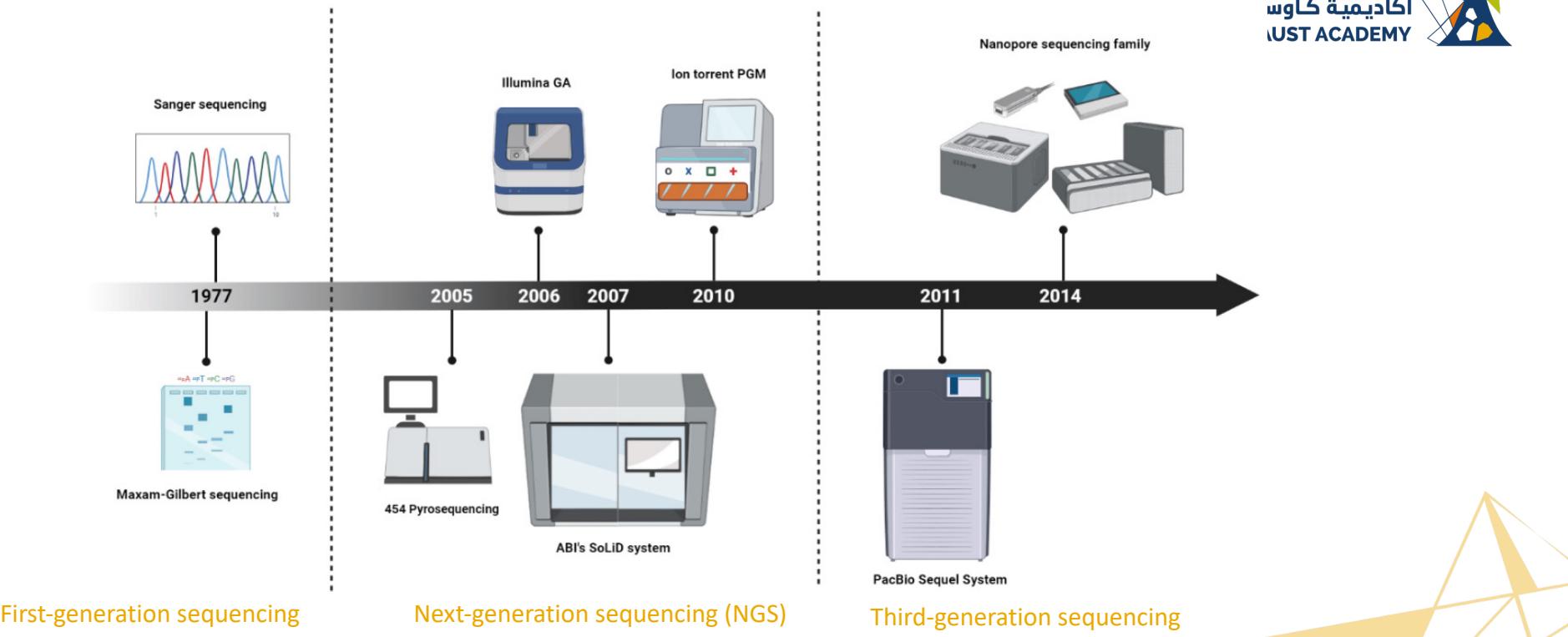
أكاديمية كاوهست  
KAUST ACADEMY



# History of Sequencing Technology



# History of Sequencing Technology



First-generation sequencing

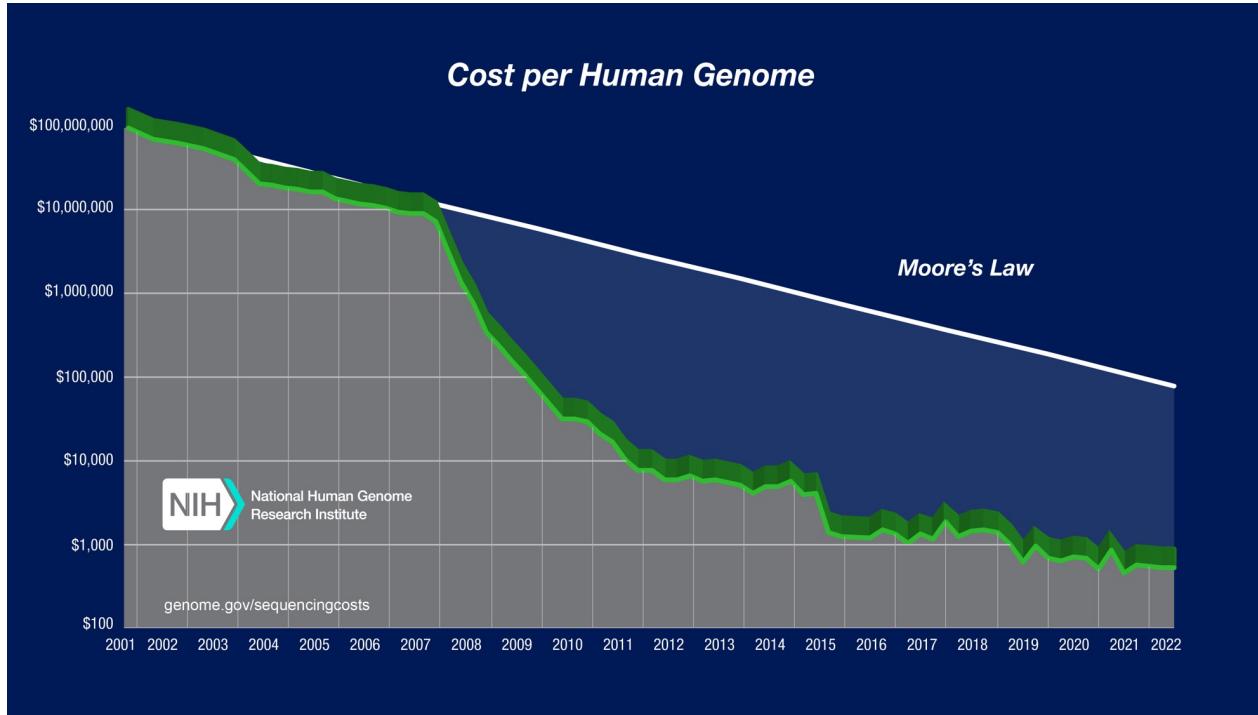
Next-generation sequencing (NGS)

Third-generation sequencing

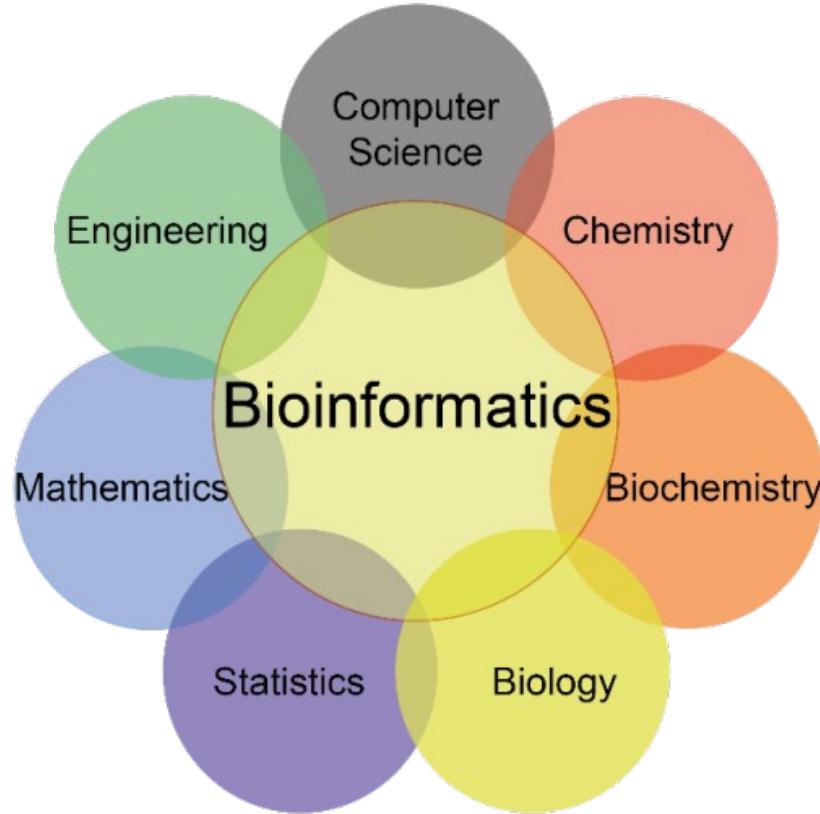
# Sequencing Technology Cost



أكاديمية كاوهست  
KAUST ACADEMY



# Interdisciplinary Foundations of Bioinformatics

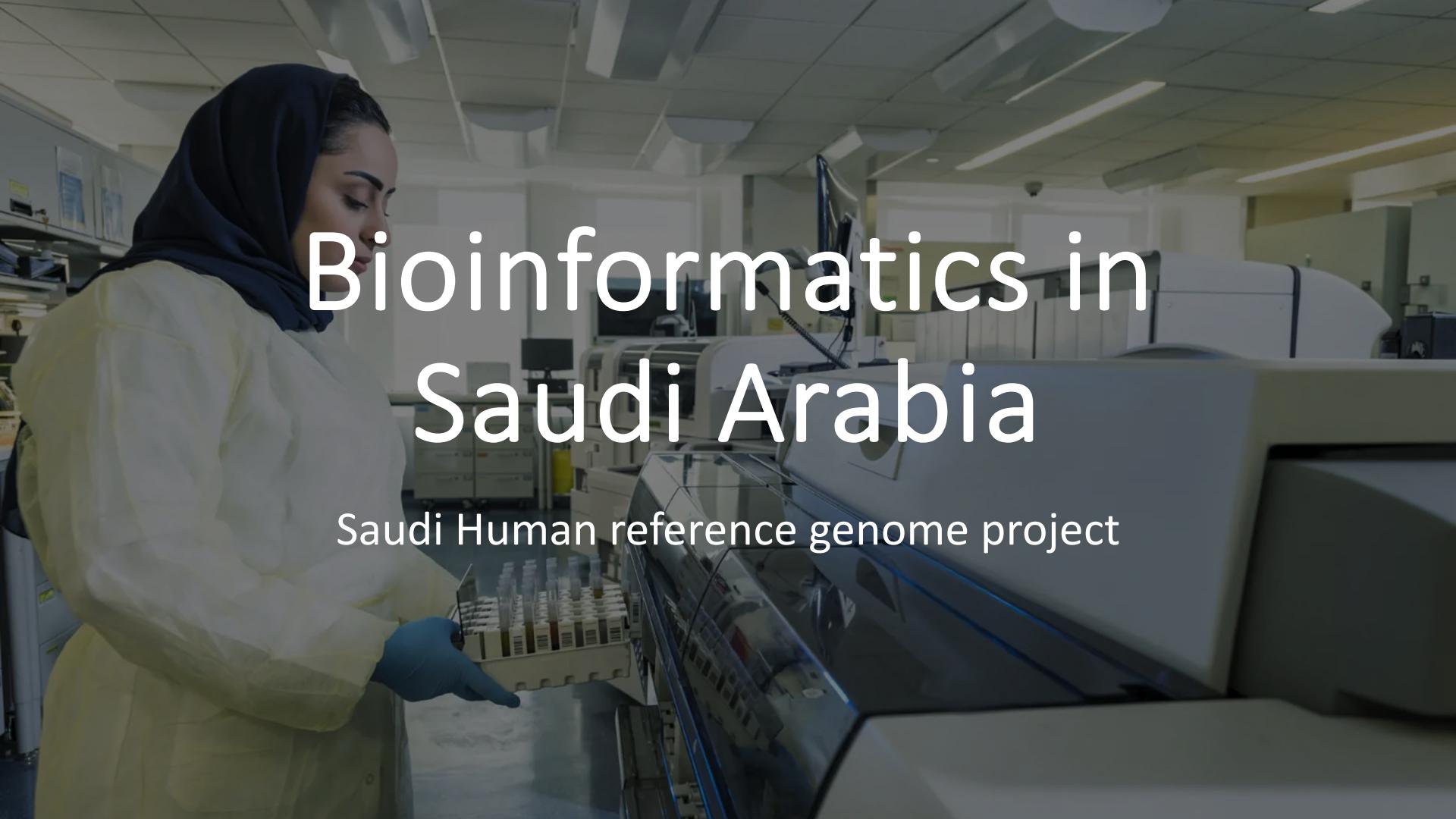


# Saudi gene hunters comb country's DNA to prevent rare diseases

Research could help prevent disorders that result from marriages between relatives

8 DEC 2016 • BY JOCELYN KAISER



A female scientist in a lab coat and blue gloves is shown in a laboratory setting, handling a tray of test tubes. The background shows various pieces of scientific equipment and shelving.

# Bioinformatics in Saudi Arabia

Saudi Human reference genome project



# Bioinformatics in Saudi Arabia

- SARS-COV2 comparative analysis

# NoorDX Startup

- Rapid COVID-19 PCR test
- Metagenomics analysis
- KAUST



# KMAP Platform



## KAUST Metagenomic Analysis Platform (KMAP), enabling access to massive analytics of re-annotated metagenomic data

Intikhab Alam <sup>1</sup>, Allan Anthony Kamau <sup>2</sup>, David Kamanda Ngugi <sup>3</sup>, Takashi Gojobori <sup>2</sup>,  
Carlos M Duarte <sup>2</sup> <sup>4</sup>, Vladimir B Bajic <sup>2</sup>

Affiliations + expand

PMID: 34075103 PMCID: PMC8169707 DOI: 10.1038/s41598-021-90799-y

Free PMC article



KAUST Academy



## Microbial Habitats

1

### Samples

- Metadata [Temp., Salinity]
- Shotgun metagenomic sequencing
- 40,000-60,000 free living or host associated shotgun metagenomic samples available at EBI.



### Microbes

Microbes are everywhere, an estimated number is ~1 trillion species.

2

### Assemble Metagenomes

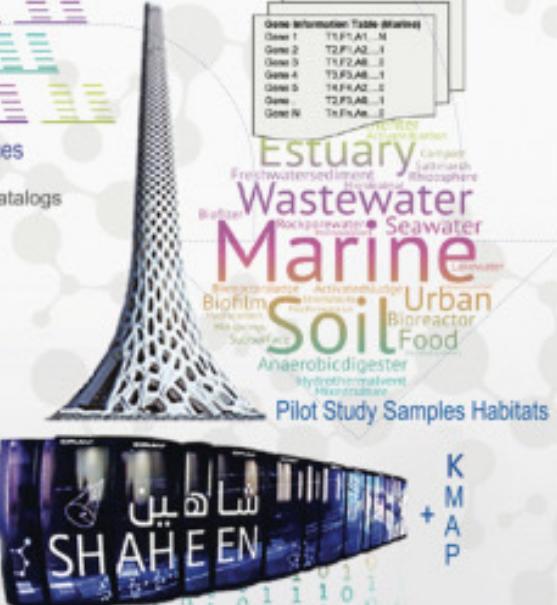
- Get full-length genes
- Make unique genes catalogs & map sample reads

Unique Genes Catalog (Metagenome)	
Gene 1	91.9, 50...54
Gene 2	0, 1, 0, 1
Gene 3	1, 2, 1, 5, 6
Gene 4	0, 1, 0, 5, 6
Gene 5	1, 3, 0, 1, 8
Gene 6	0, 1, 0, 5, 6
Gene N	1, 1, 9, 1, 8

4

### Data Sharing

Shaheen and KMAP were used to produce Pilot Study: 40 Gene Catalogs, 275 million genes.



### Research Groups

Gene information tables (AAMG TSV) available for research organizations with advanced computational resources and skills.



### Individuals

Graphic User Interface access to indexed gene information tables for browsing and comparisons. Available for individuals with less computational resources.





NEOM

ABOUT ▾

REGIONS ▾

OUR BUSINESS ▾

NEWS ▾

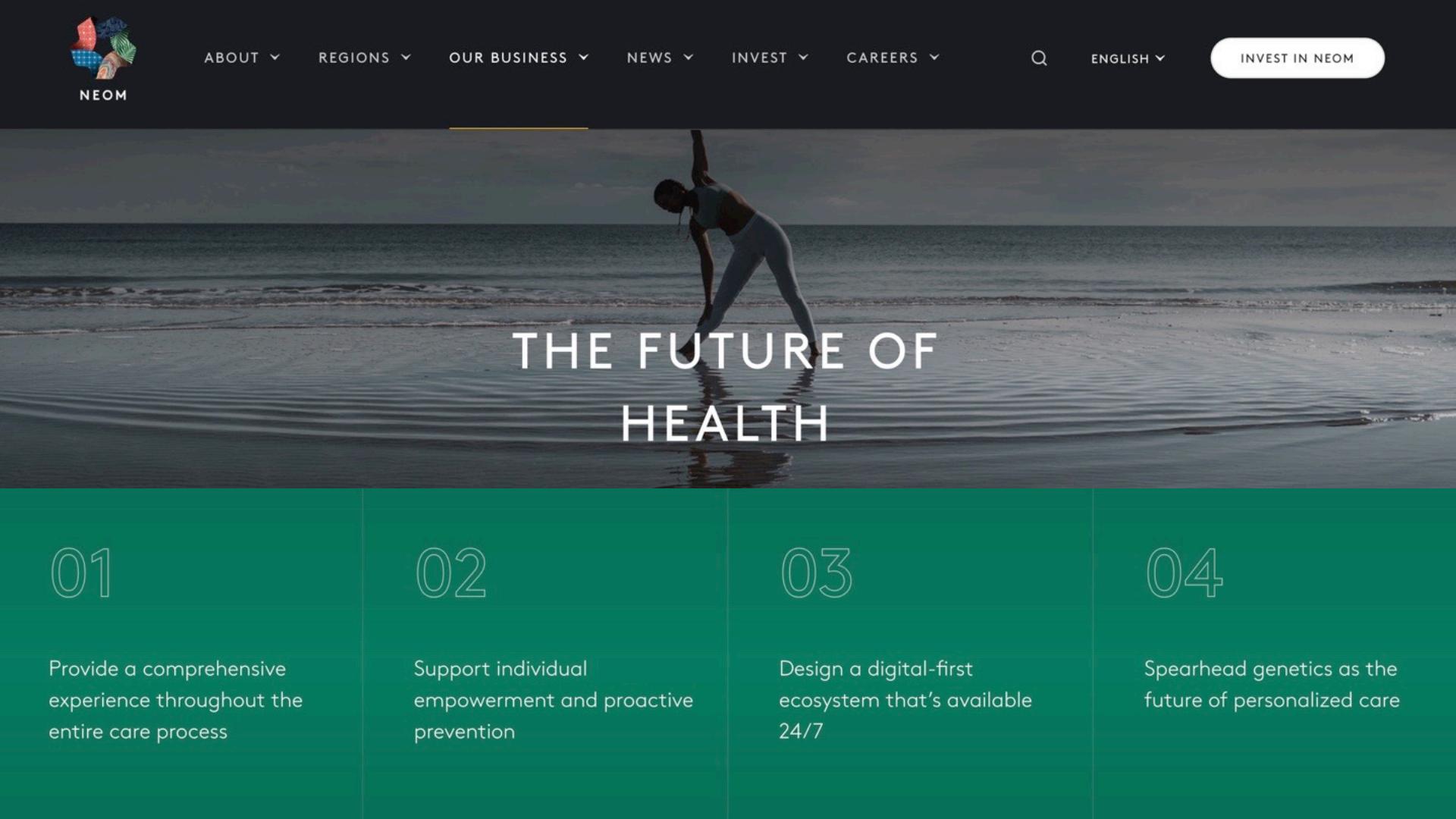
INVEST ▾

CAREERS ▾



ENGLISH ▾

INVEST IN NEOM



# THE FUTURE OF HEALTH

01

Provide a comprehensive experience throughout the entire care process

02

Support individual empowerment and proactive prevention

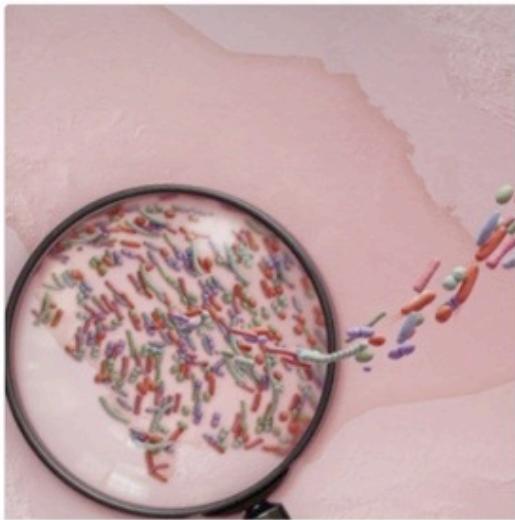
03

Design a digital-first ecosystem that's available 24/7

04

Spearhead genetics as the future of personalized care

## Health and Wellness: Accelerating Impact in KSA



04 December, 2023

### The weird and wonderful world of Saudi Arabia's microbiomes

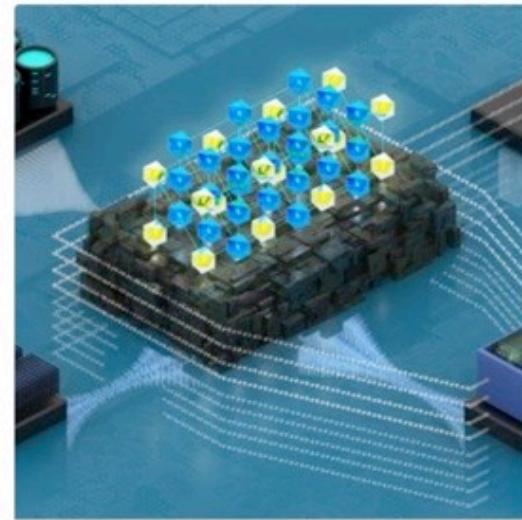
Saudi Arabia is home to microbial communities that can boost coral reef health, sequester carbon in the soil, and reduce desertification.



03 December, 2023

### The first ride across Saudi Arabia on a hand bike

To showcase the remarkable potential of people with physical disabilities, KAUST Professor Matteo Parsani will travel from the east to the west of Saudi Arabia by



22 November, 2023

### Safeguarding the right to be forgotten

An open-source software can help align artificial intelligence applications in healthcare with data privacy regulations.

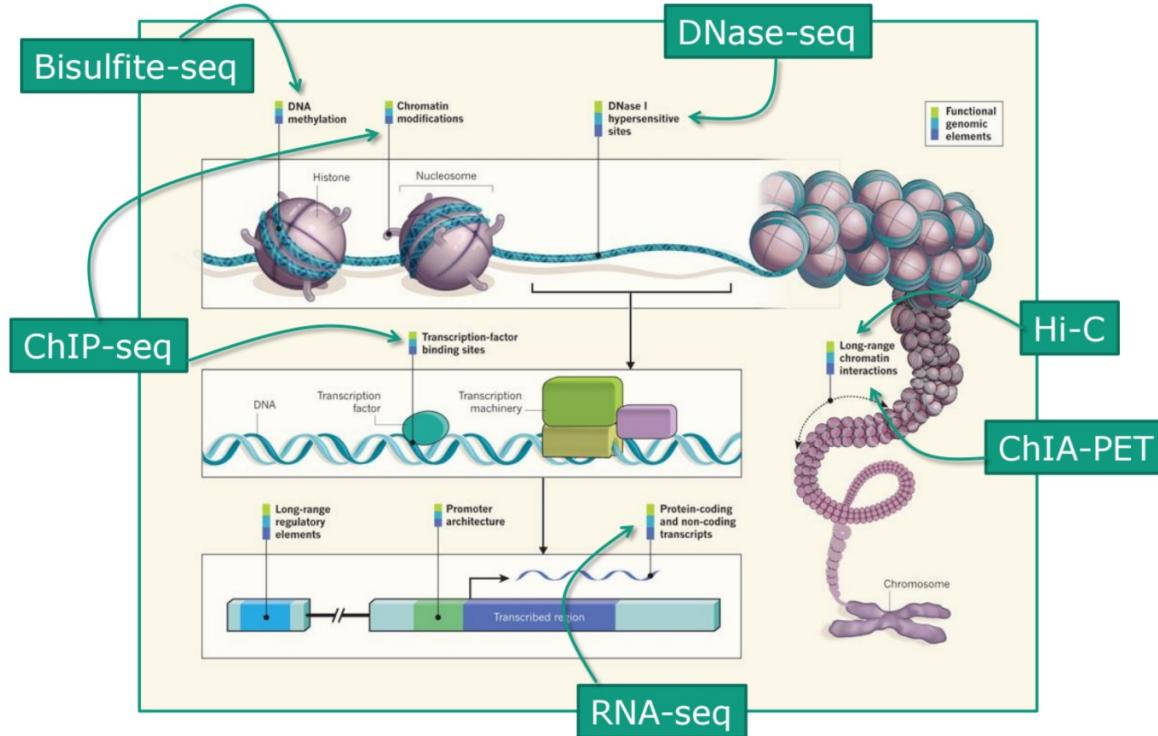
# Introduction to Bioinformatics

## How to Get Started in Bioinformatics

---

### Part 2: Sequence analysis

# Where is Data Coming From?



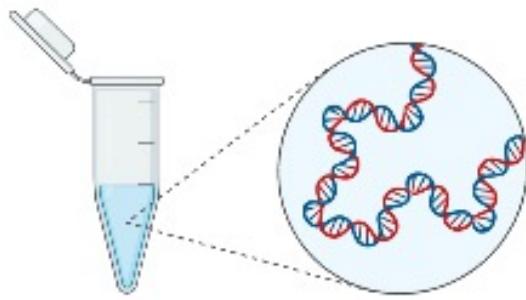
*Ecker et al, Nature, 2012*

# Sequencing Technologies

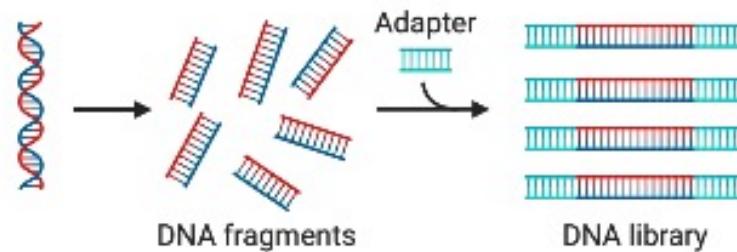
NGS Technology	Description
ChIP-seq	Analyzes protein-DNA interactions to identify binding sites of DNA-associated proteins
RNA-seq	Sequences RNA to examine transcriptome profiles, including gene expression and splicing patterns
Exome-seq	Focuses on sequencing exons, the coding regions of the genome, to find genetic variants
DNA-seq	Involves sequencing the whole DNA to identify genetic variations and mutations.
Metagenomic Sequencing	Studies genetic material directly from environmental samples.
Single-Cell Sequencing	Analyzes genomes or transcriptomes of individual cells for detailed cellular variation
Methyl-Seq	Studies DNA methylation, an important epigenetic modification
ChIP-exo	An enhanced version of ChIP-seq offering higher resolution.
ATAC-seq	Assesses chromatin accessibility, important for understanding transcriptional regulation.

# Wet Lab

## Step 1: DNA extraction



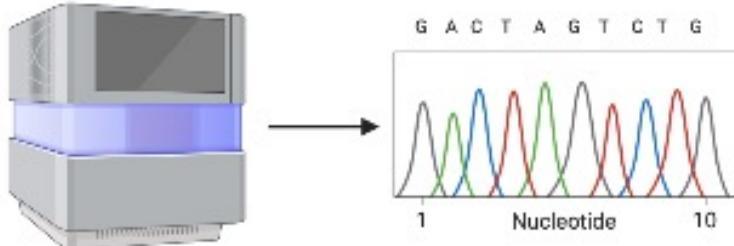
## Step 2: Library preparation



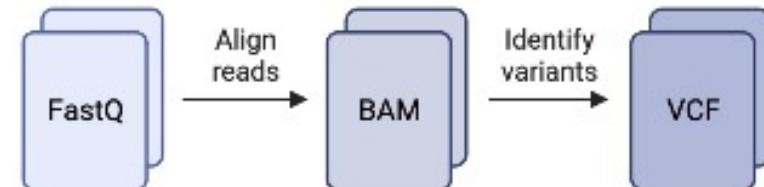
## Next Generation Sequencing Workflow

# Dry Lab

## Step 3: Sequencing



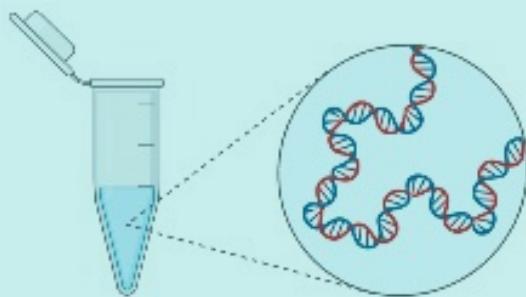
## Step 4: Analysis



# Wet Lab

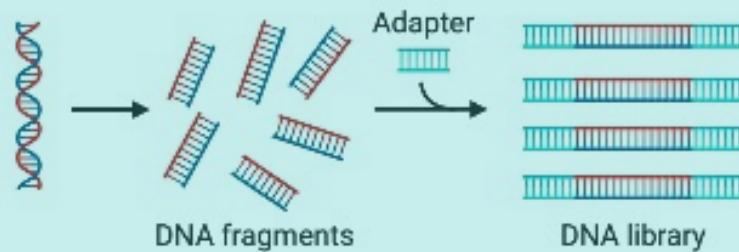
## Step 1:

### DNA extraction



## Step 2:

### Library preparation

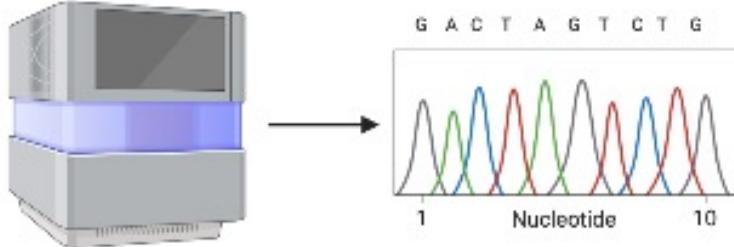


## Next Generation Sequencing Workflow

# Dry Lab

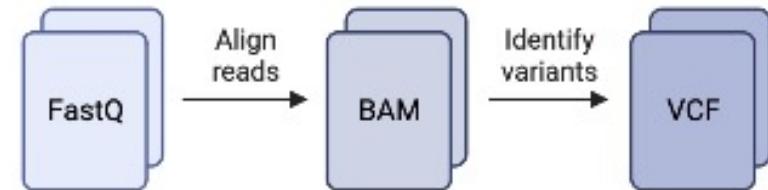
## Step 3:

### Sequencing

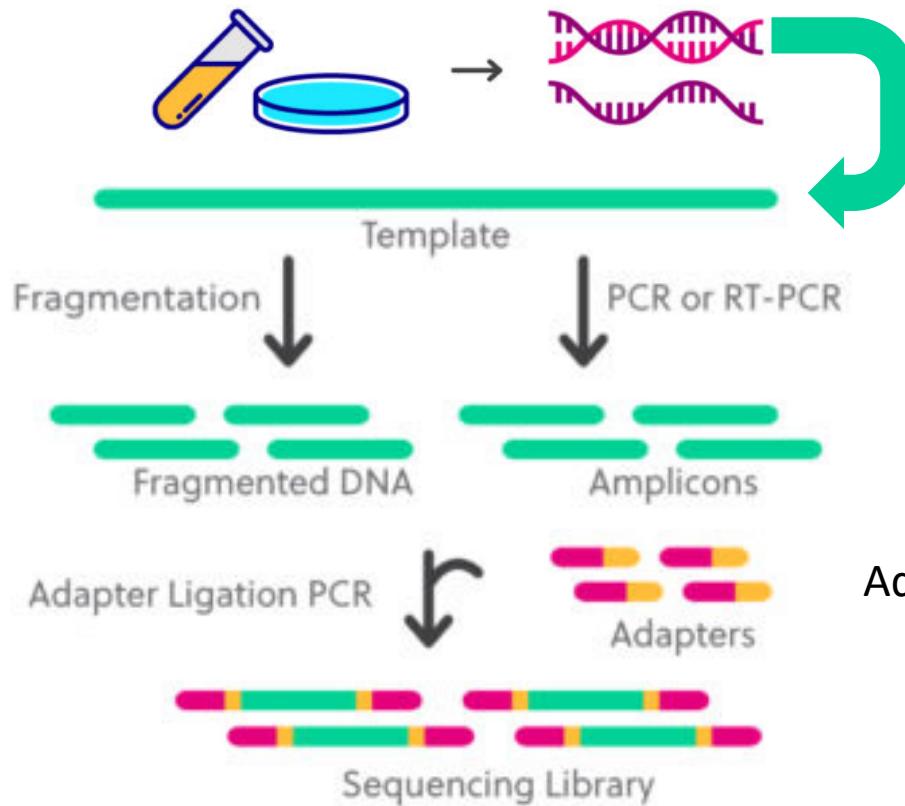


## Step 4:

### Analysis

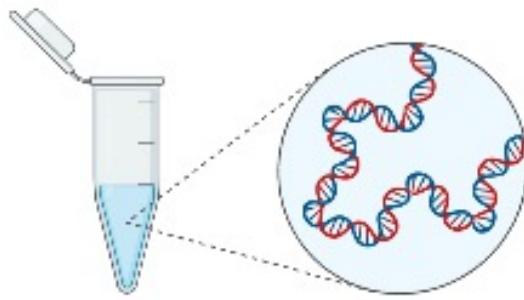


# Wet-Lab: Library Preparation

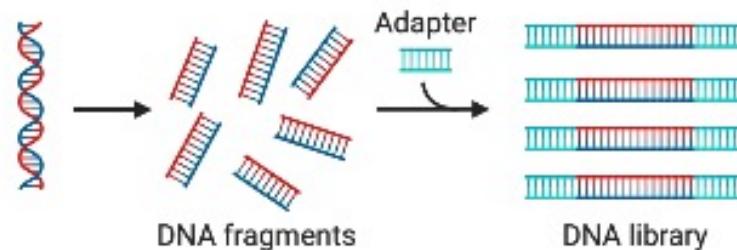


# Wet Lab

## Step 1: DNA extraction



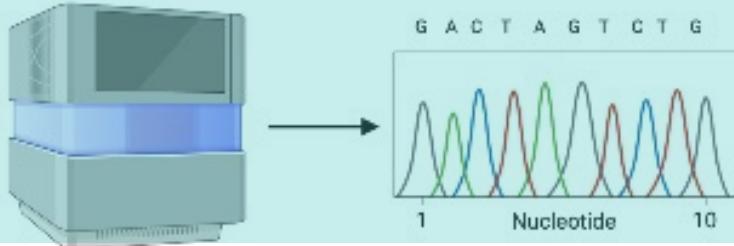
## Step 2: Library preparation



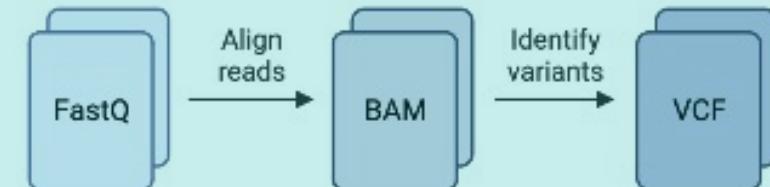
## Next Generation Sequencing Workflow

# Dry Lab

## Step 3: Sequencing



## Step 4: Analysis



# DNA Sequencer



Input DNA

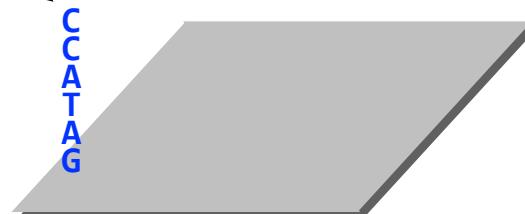
CCATAGTATATCTGGCTCTAGGCCCTCATTTTT  
CCATAGTATATCTGGCTCTAGGCCCTCATTTTT  
CCATAGTATATCTGGCTCTAGGCCCTCATTTTT  
CCATAGTATATCTGGCTCTAGGCCCTCATTTTT

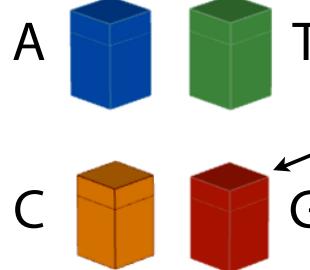


Cut into snippets

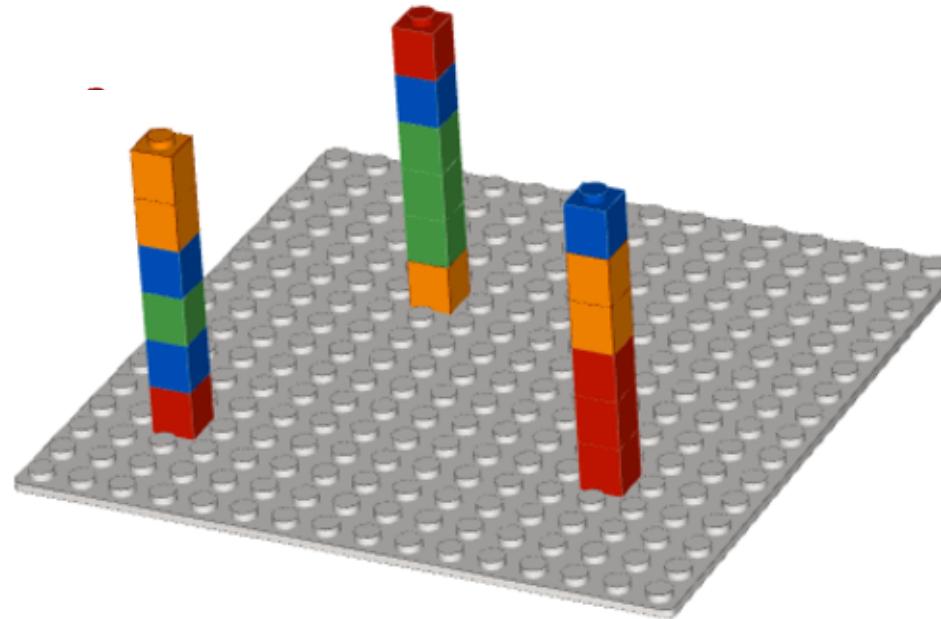
CCATAGTA TATCTGG CTCTAGGCCCTC ATTTTTT  
CCA TAGTATAT CTCGGCTCTAGGCCCTCA TTTTTT  
CCATAGTAT ATCTCGGCTCTAG GCCCTCA TTTTTT  
CCATAG TATATCT CGGCTCTAGGCCCT CATTTTT

Deposit on slide



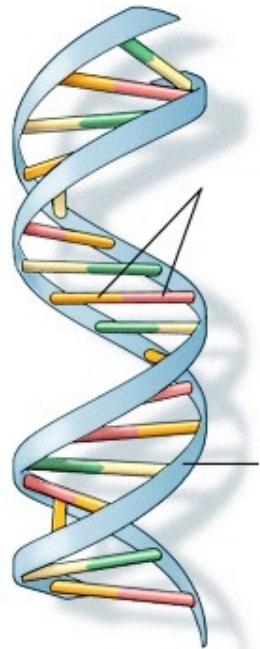


Template  
(billions of them!)



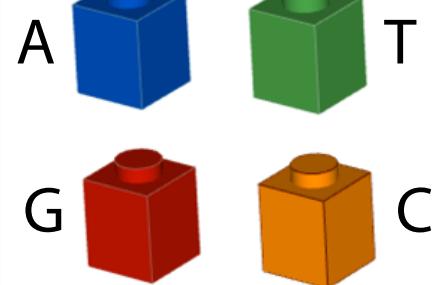
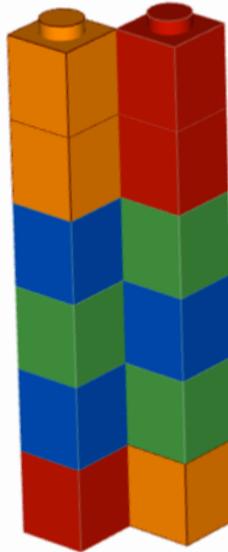
Slide

# How DNA is copied



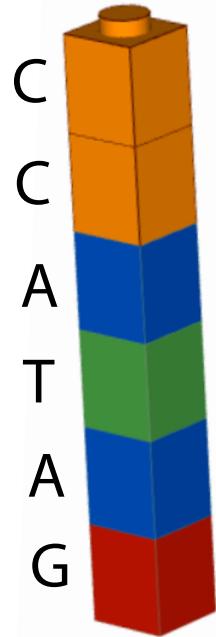
U.S. National Library of Medicine

Double stranded  
DNA (double helix)

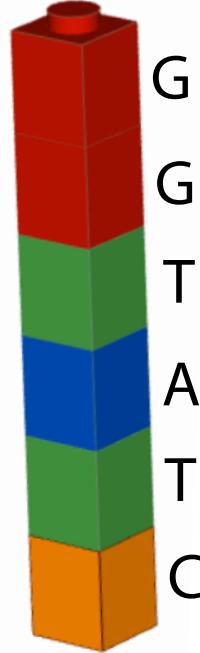


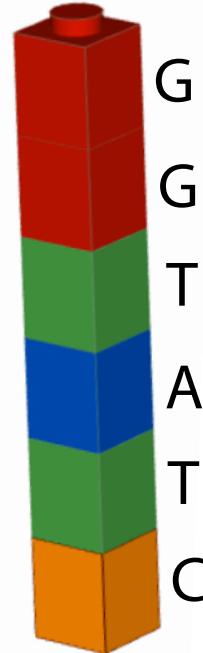
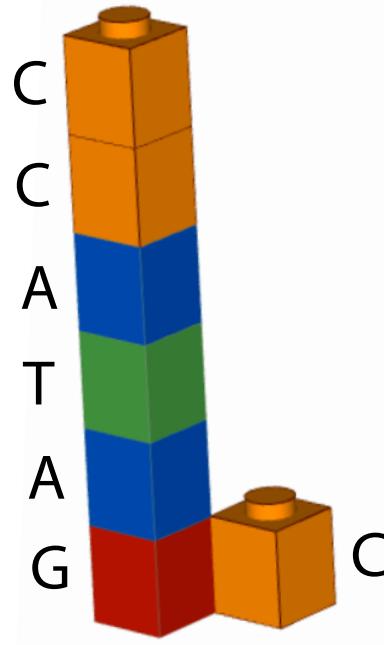
Double stranded  
DNA (lego version)

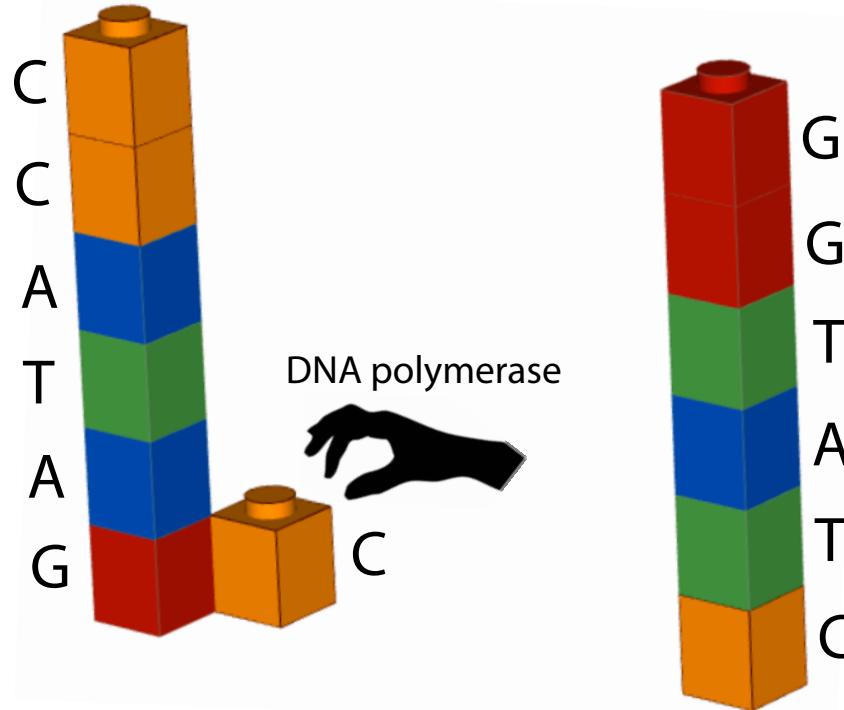


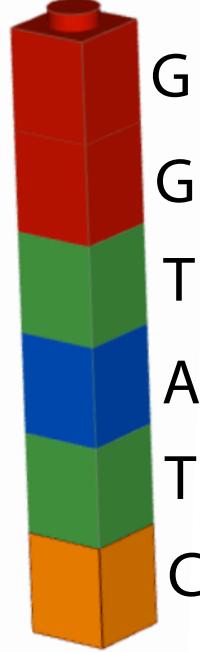
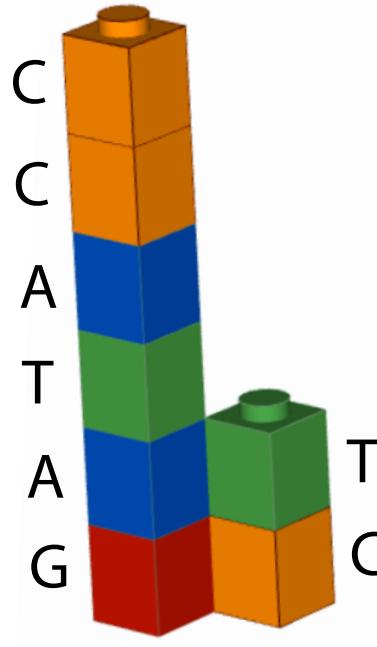


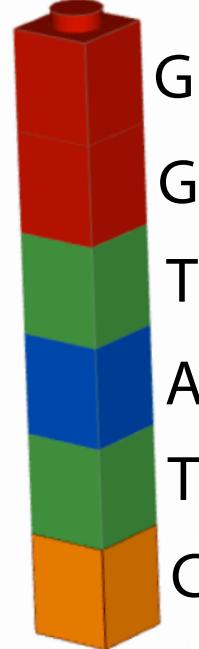
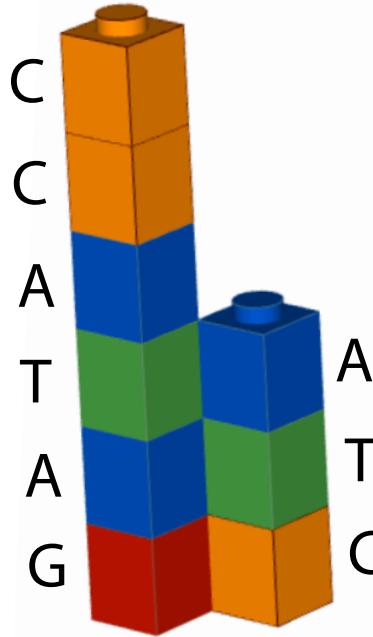
Single stranded  
templates

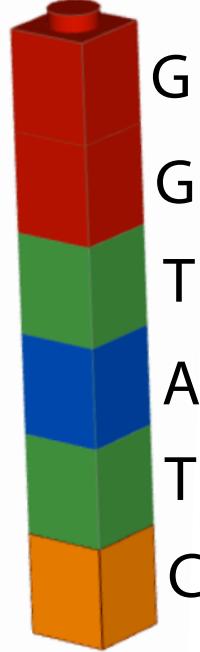
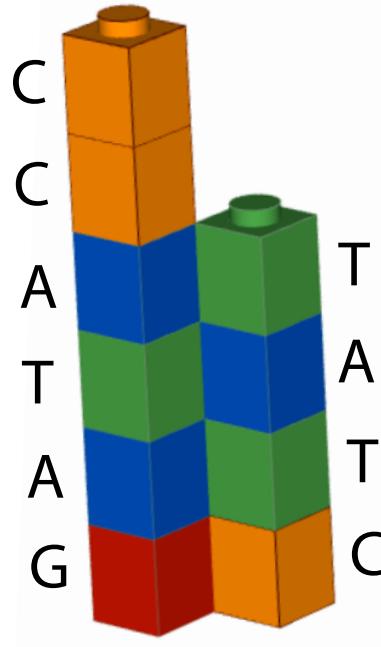


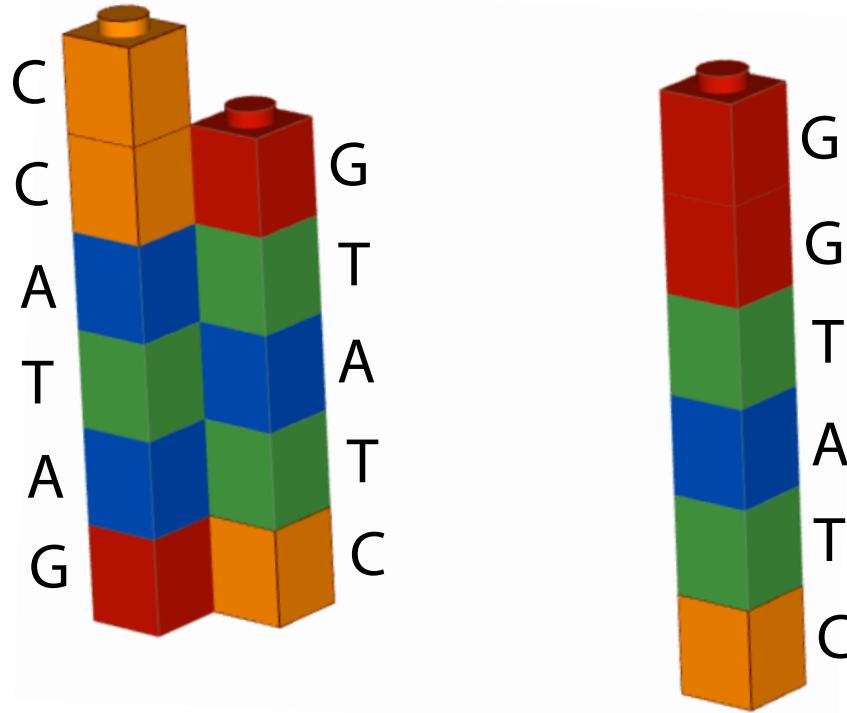


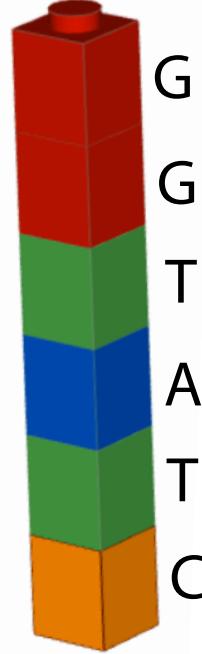
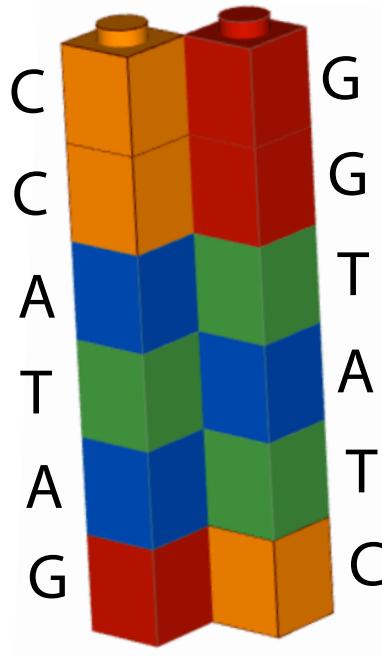


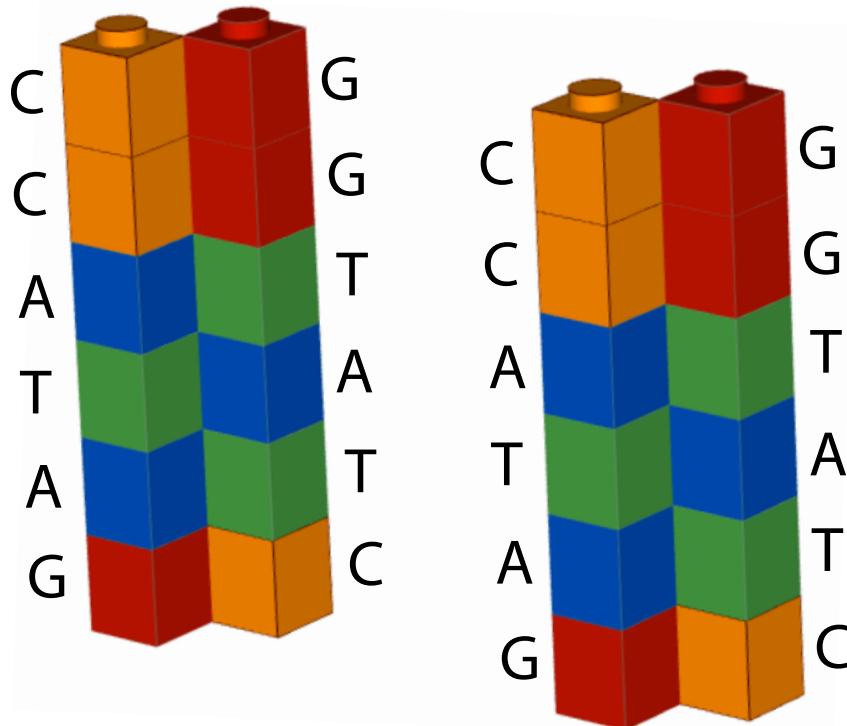


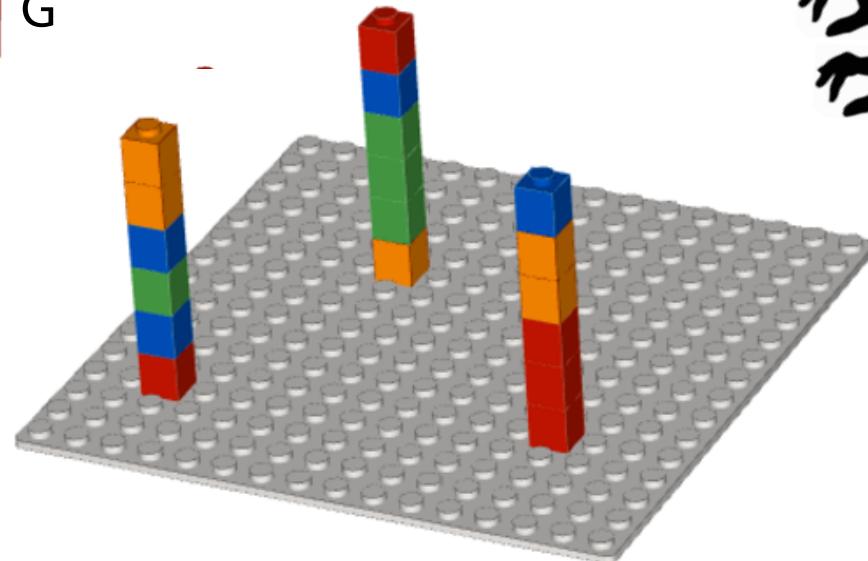
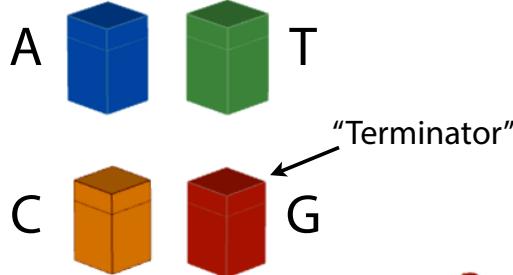






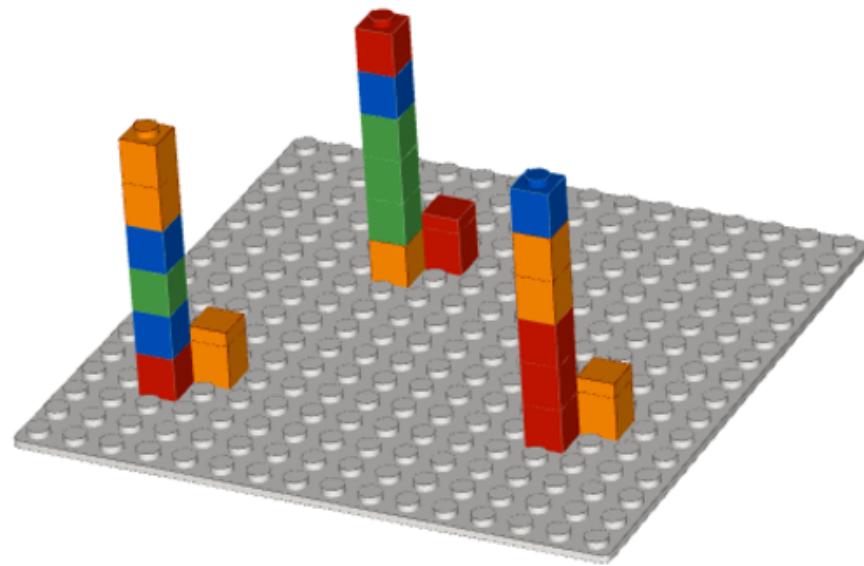


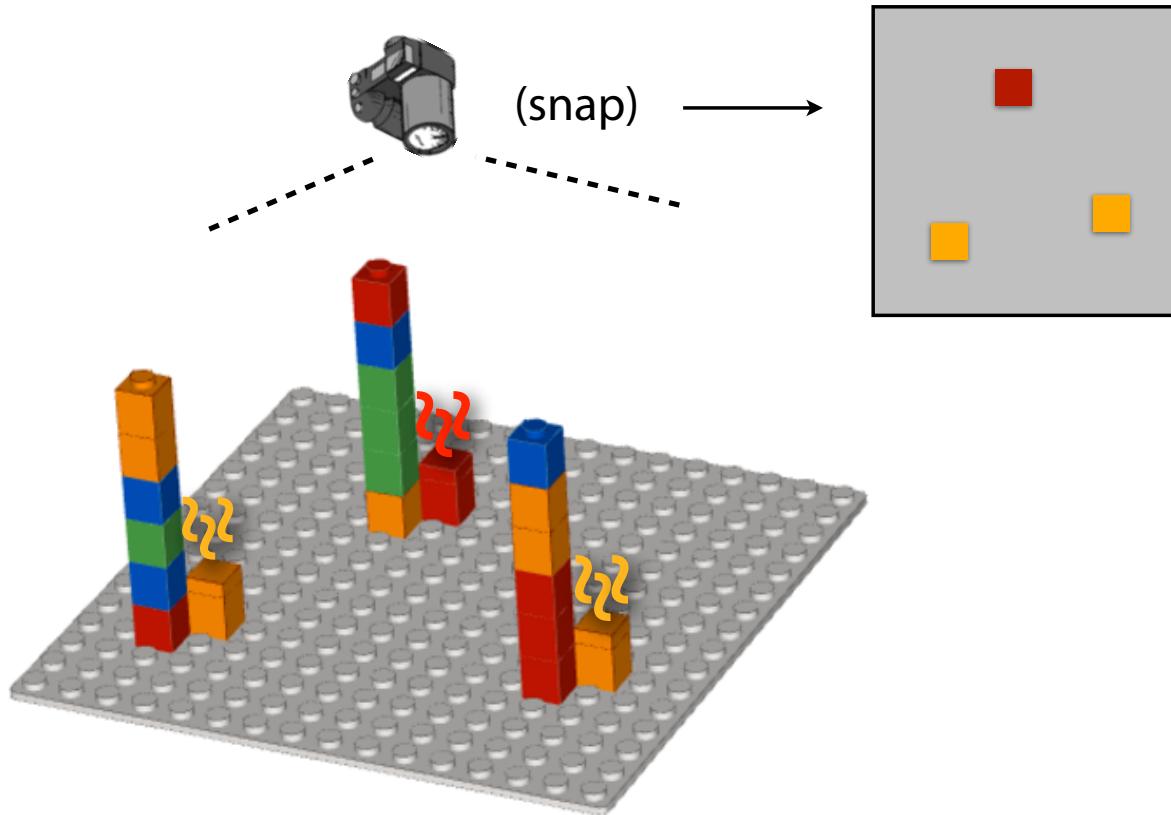




DNA polymerase

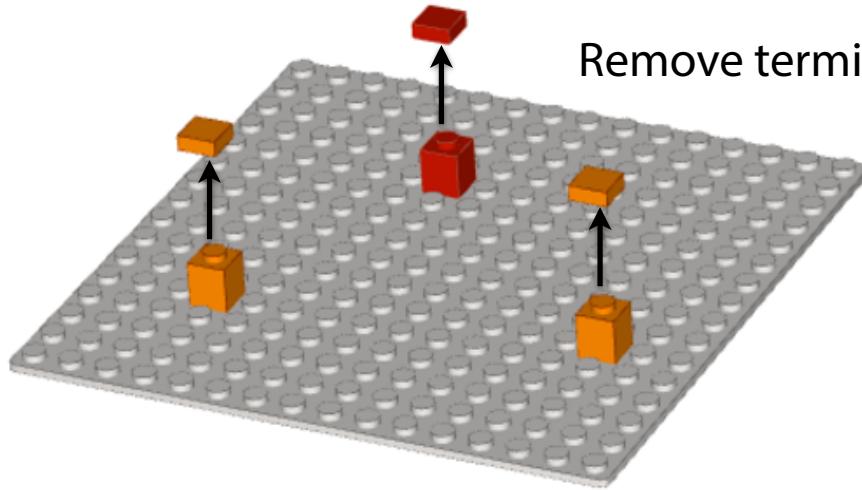






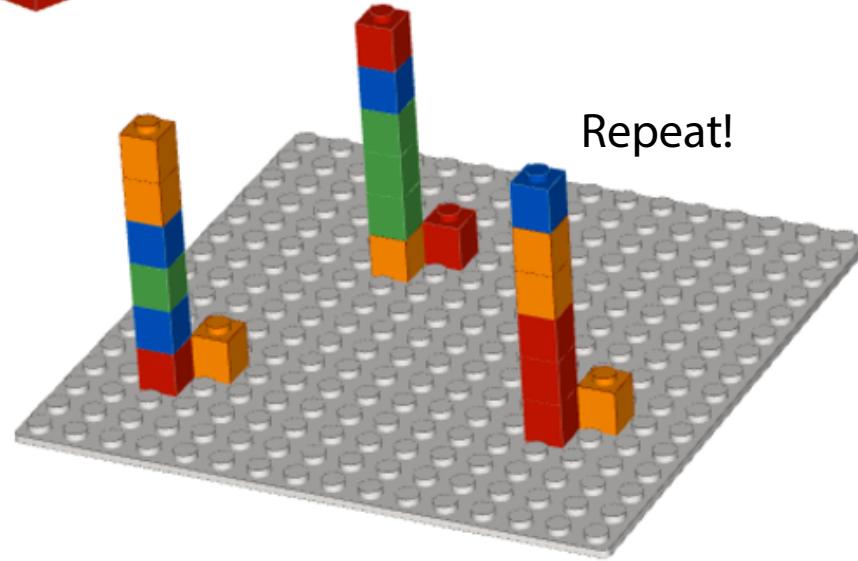


Remove terminators



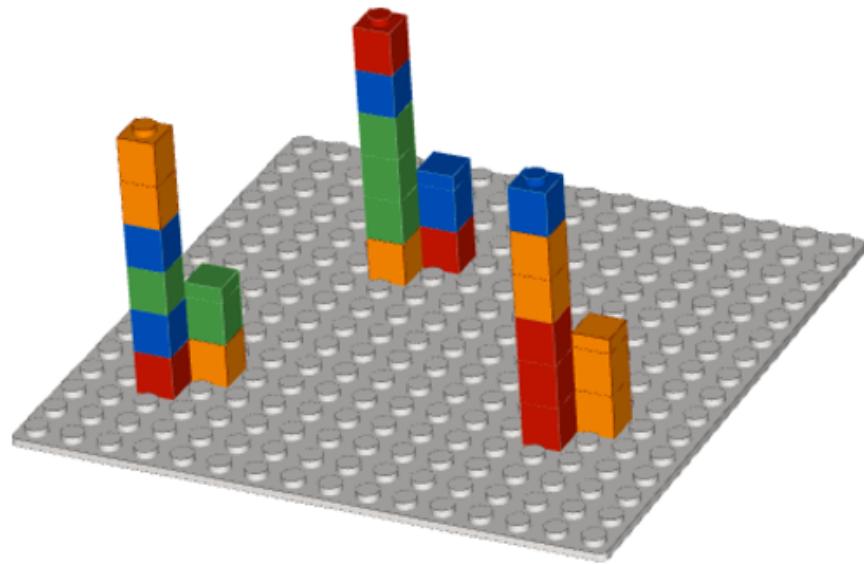
A T

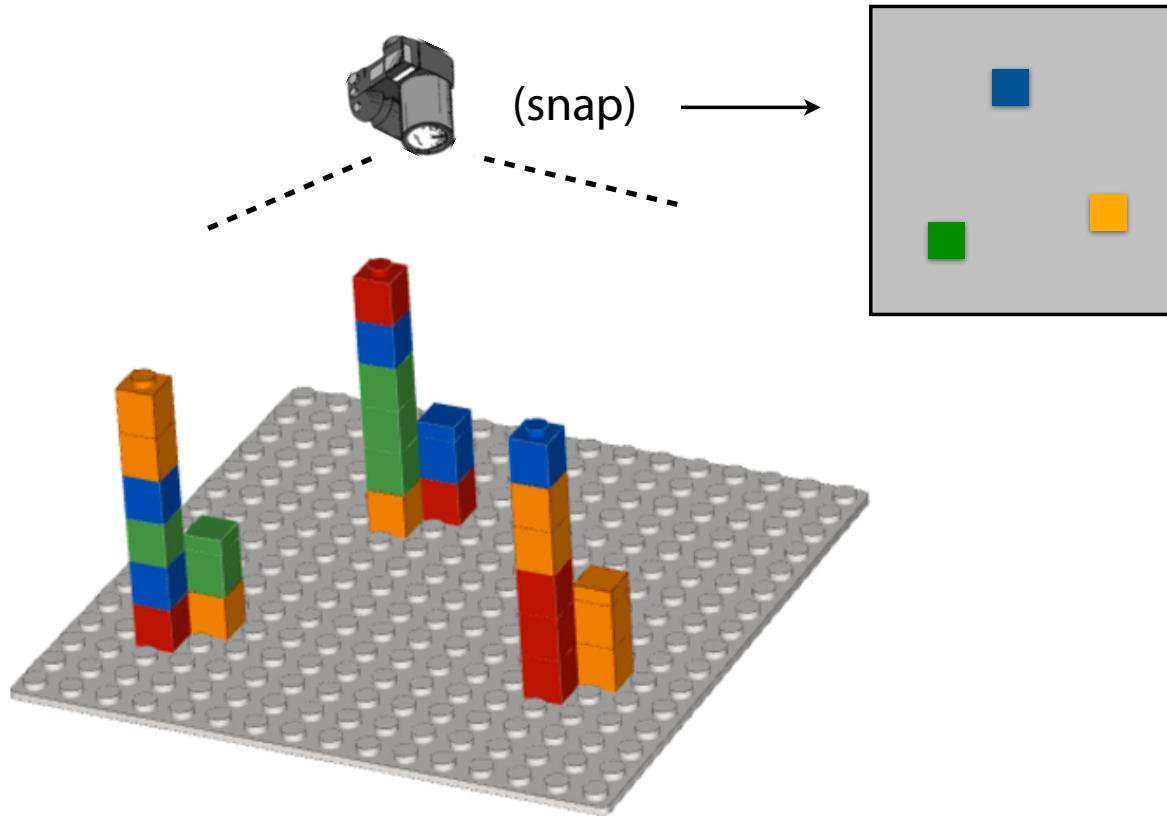
C G

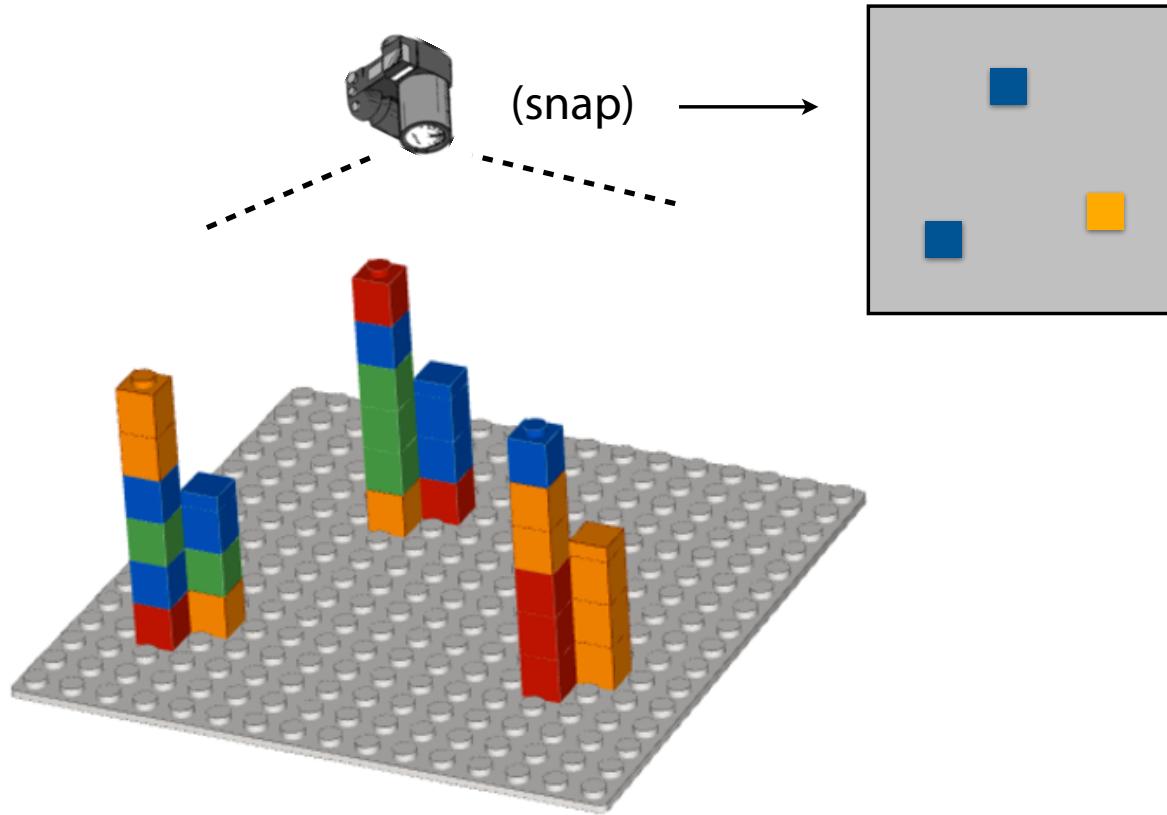


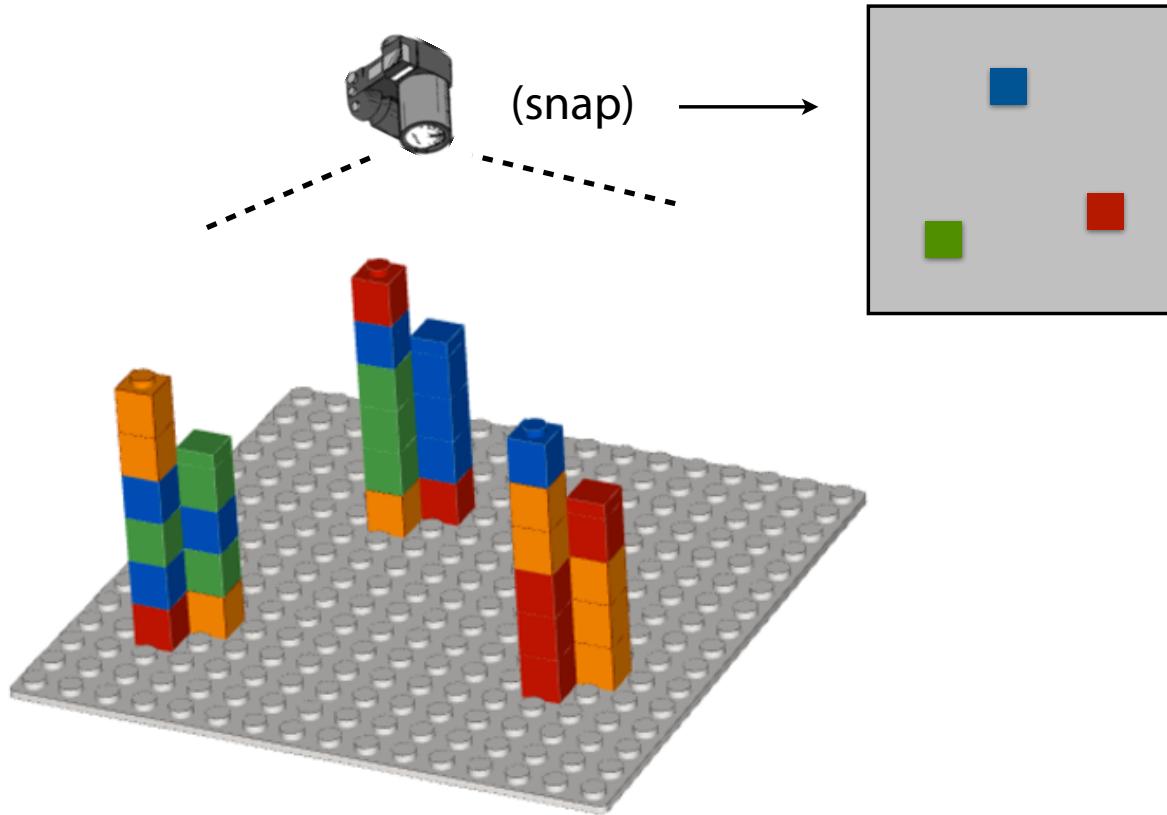
DNA polymerase



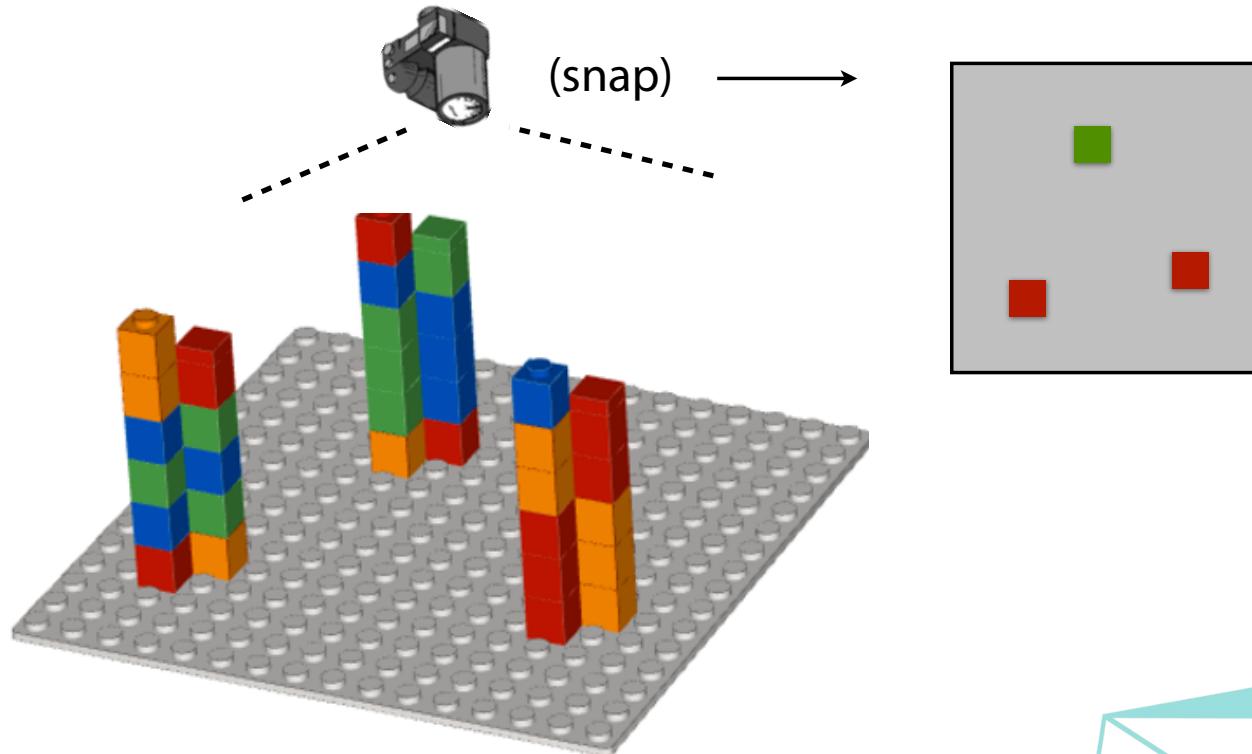






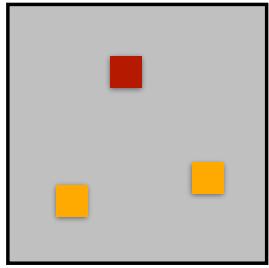


# Example

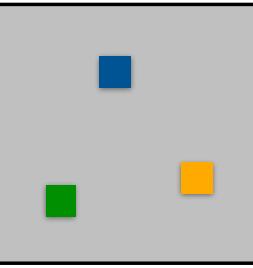


# Cycles

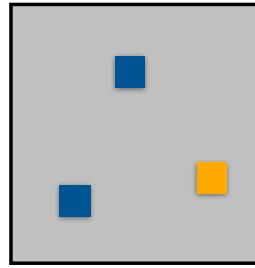
Cycle 1



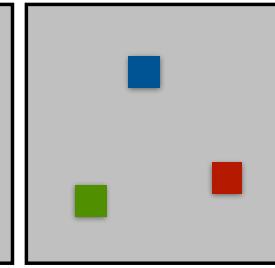
Cycle 2



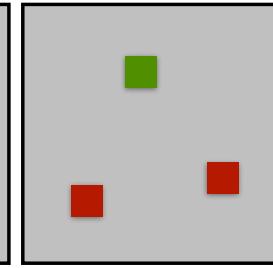
Cycle 3



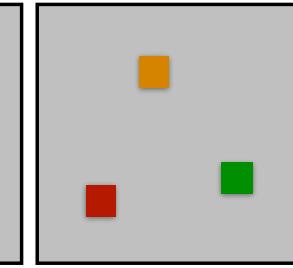
Cycle 4



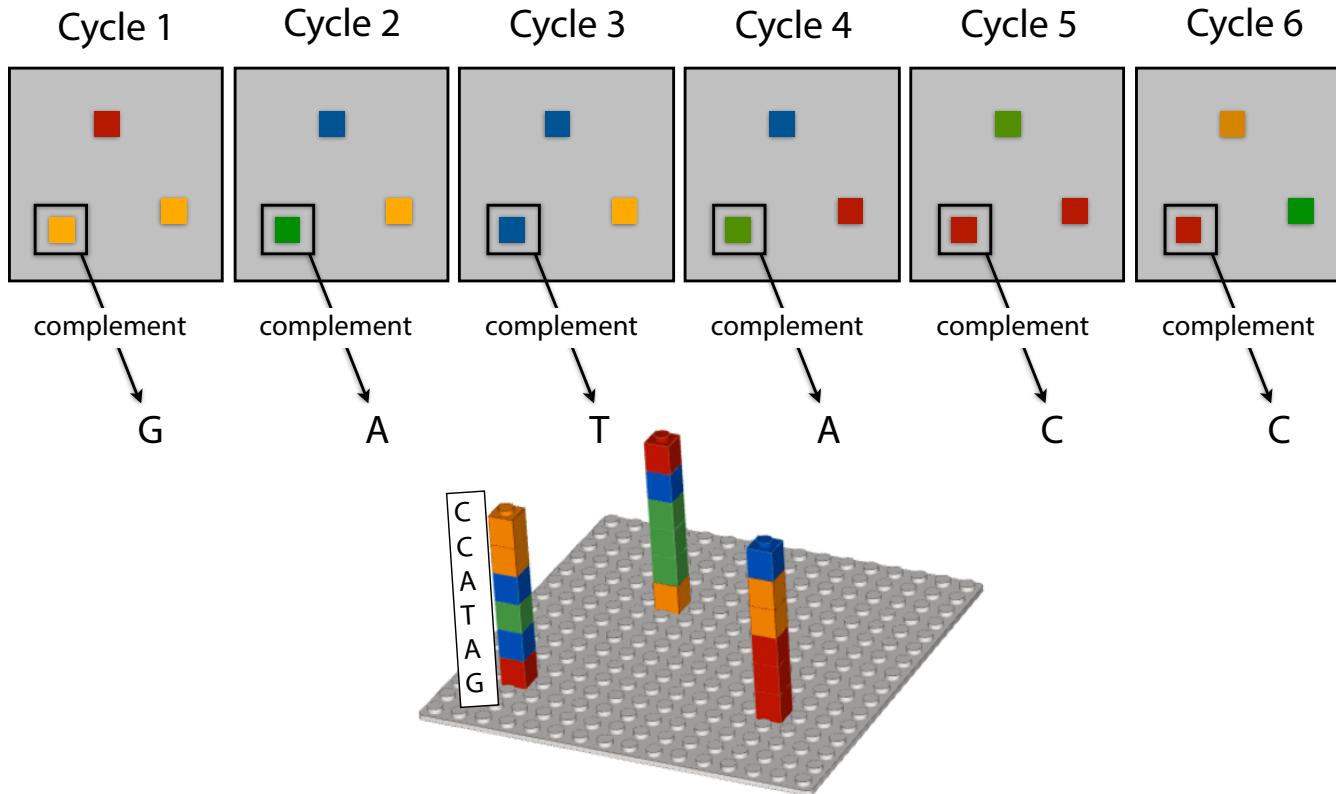
Cycle 5



Cycle 6



# Sequencing by Synthesis



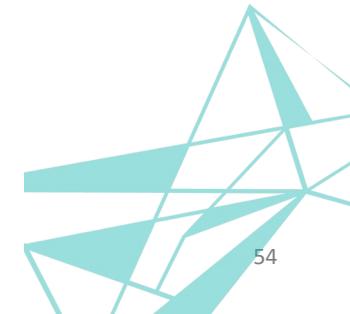
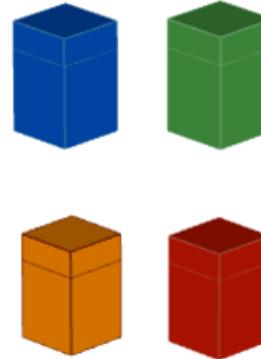
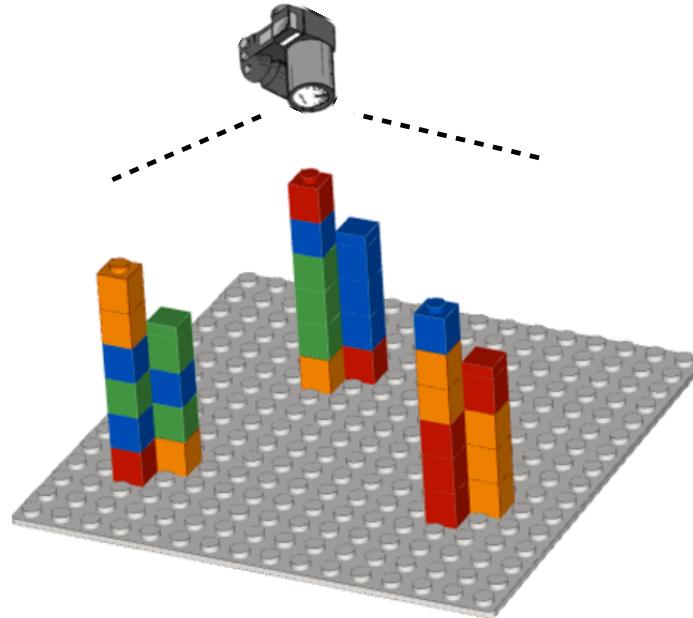
# Exan Sequencing by synthesis



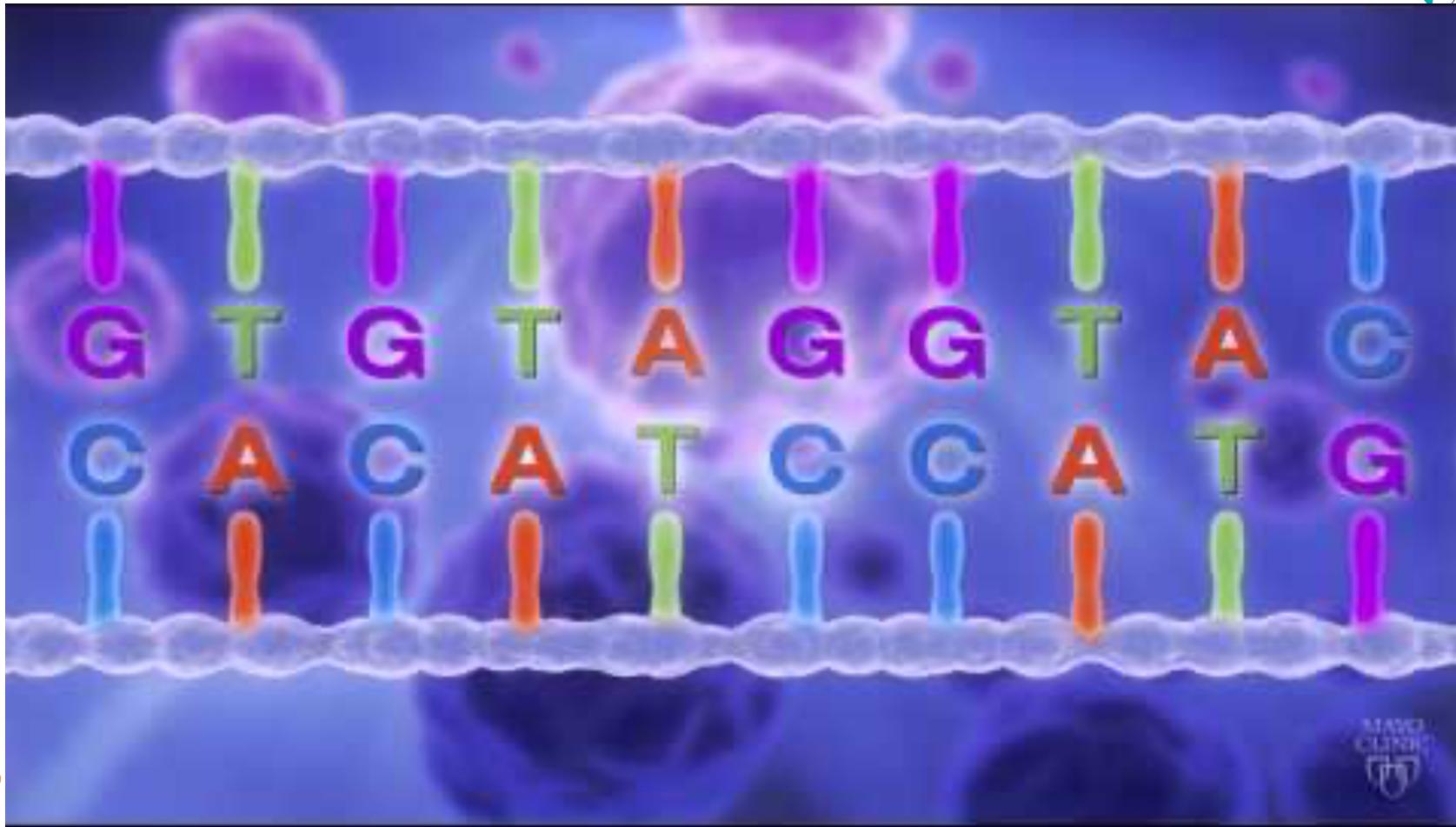
Billions of templates on a slide

Massively parallel: photograph captures all templates simultaneously

Terminators are “speed bumps,” keeping reactions in sync



# From experiments to data



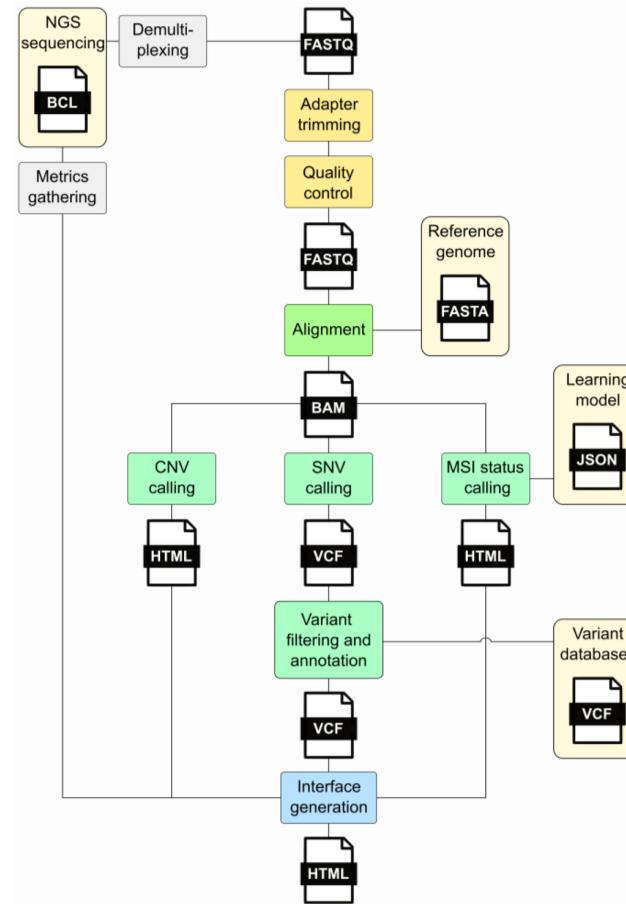
# Sequencing Data File Formats

File Format	Description
FASTA	Represents nucleotide or peptide sequences without quality scores
FASTQ	Stores raw sequence data and quality scores for each nucleotide
BCL	Binary Base Call format; raw base call data from Illumina sequencing platforms
SAM	Verbose text-based format for storing sequence alignment data
BAM	Binary version of SAM, compressed and indexed for efficient data processing
CRAM	Compressed file format for storing sequencing reads, optimized for space
VCF	Describes gene sequence variations like SNPs and indels
BED	Specifies genomic regions, such as those identified in sequencing analyses
GFF/ GTF	Detailed annotation of genes and other genomic features
JSON	JavaScript Object Notation; a lightweight data-interchange format, often used for APIs and config files.
HTML	HyperText Markup Language; used for creating web pages and web applications

# Sequencing Data Format Workflow



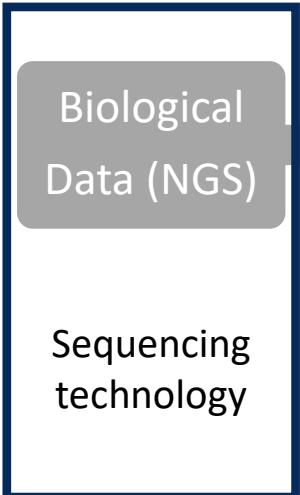
أكاديمية كاوهست  
KAUST ACADEMY



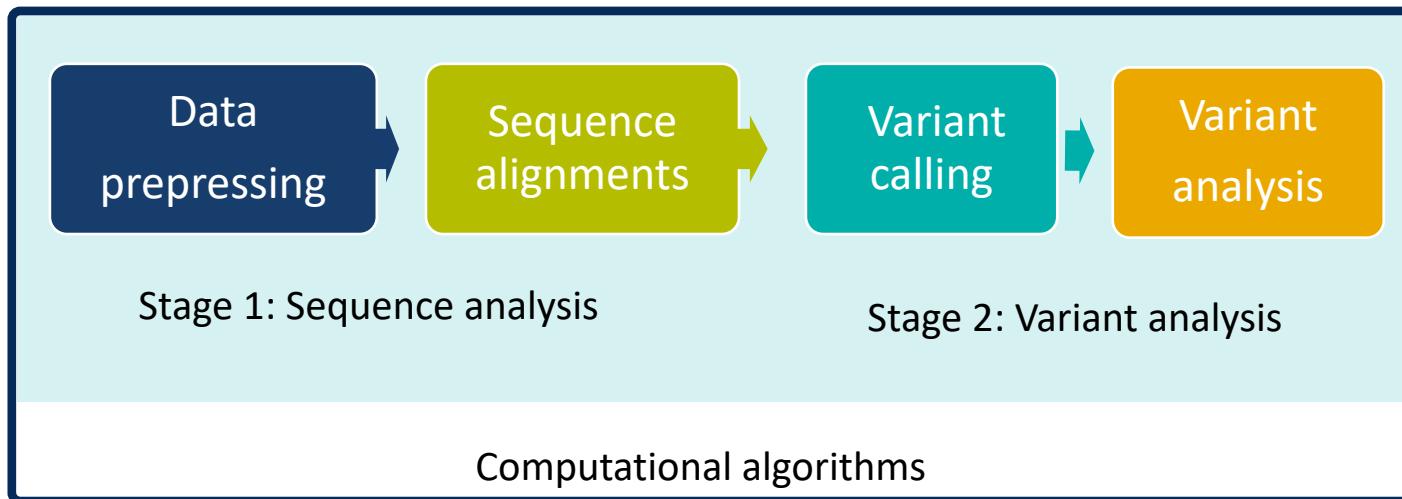
# How to Get Started in Bioinformatics



Wet-Lab



Dry-Lab



NGS Bioinformatics workflow

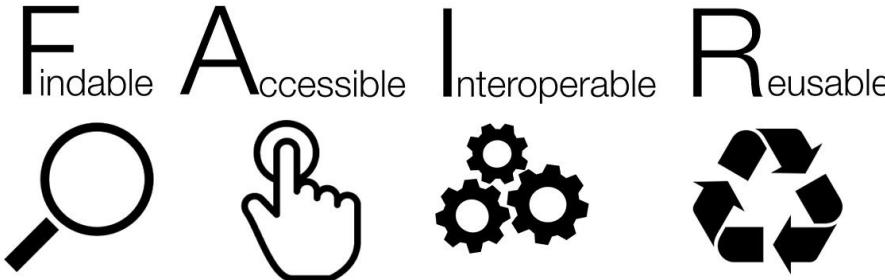


KAUST Academy

# Introduction of Galaxy Platform

Galaxy is an **open-source** platform for **FAIR** data analysis that enables users to:

- Use **tools** from various domains through its graphical web interface.
- Run code in **interactive environments** along with other tools or workflows.
- **Manage data** by sharing and publishing results, workflows, and visualizations.
- **Ensure reproducibility** by capturing the necessary information to repeat and understand data analyses.



# Discover galaxy platform!

The screenshot shows the Galaxy web interface with three main panels:

- Tools Panel (Left, Blue Border):** A sidebar menu listing various tools categorized under "GENERAL TEXT TOOLS", "GENOMIC FILE MANIPULATION", "COMMON GENOMICS TOOLS", and "GENOMICS ANALYSIS".
- Home Panel (Middle, Red Border):** The main content area featuring a banner about the James P. Taylor Foundation for Open Science, a "Learn More" button, and a callout for SARS-CoV-2 analysis.
- History Panel (Right, Green Border):** A list of datasets in the "Galaxy 101 History" session, including "2: SNPs" and "1: Exons".

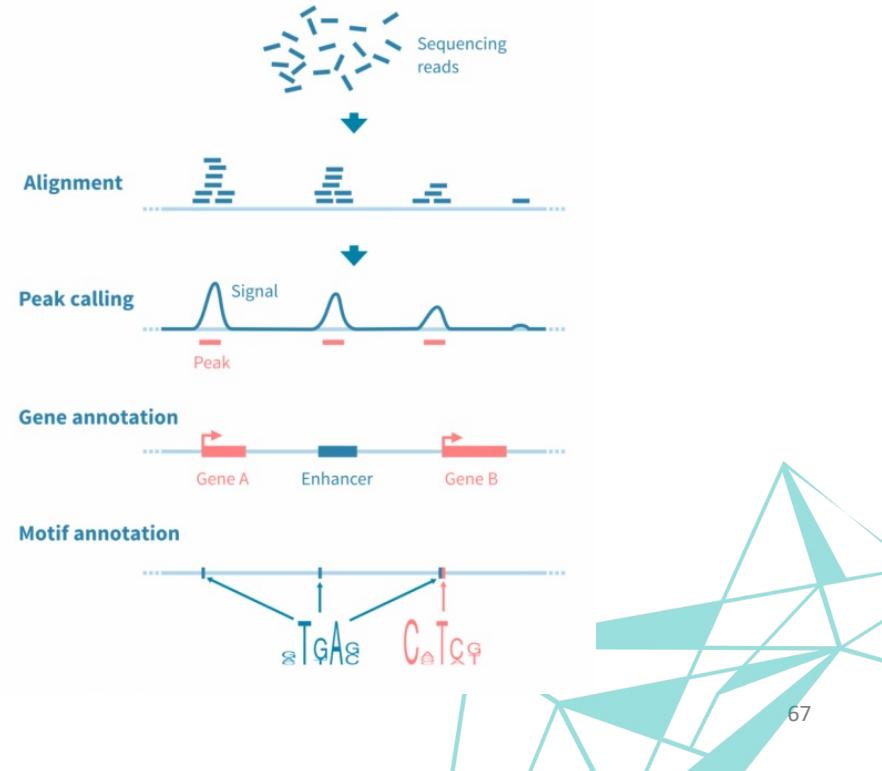
At the bottom, logos for Penn State, Johns Hopkins University, and Oregon Health & Science University are displayed, along with a note about the Galaxy Team being part of the Center for Comparative Genomics and Bioinformatics.



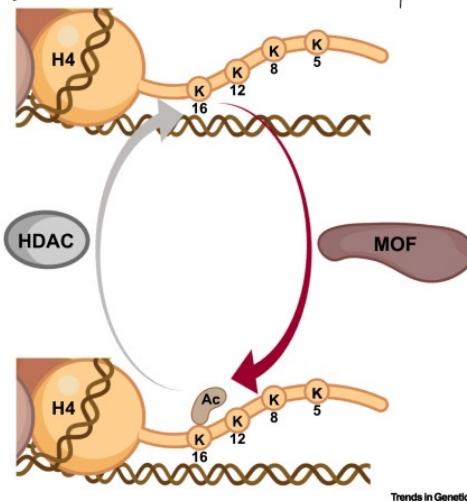
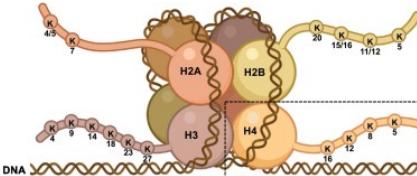
# From peaks to genes

## How to get from peak regions to a list of gene names?

- What is a peak?
- Peak regions in DNA, identified via sequencing techniques like ChIP-seq, are areas with high read enrichment, suggesting where proteins bind to the genome. These peaks are mapped to nearby genes to infer which genes may be regulated by these proteins.



# GSE37268 Dataset for MOF Gene



- MOF(MOF histone acetyltransferase) is a protein that plays a crucial role in regulating the core transcriptional network of embryonic stem cells.
- Li et al. (2012) conducted a study to identify the target genes of MOF using ChIP-seq in mice.
- The raw data from the study is available through GEO, but the list of MOF target genes is not included in the paper's supplement or the GEO submission.

# GSE37268 Dataset for MOF Gene



أكاديمية كاوهست  
KAUST ACADEMY

## ChIP-seq analysis result

Chromosome Number	Start	End	Peak Length	Signal Strength	Score	Statistical Data
1	3660676	3661050	375	210	62.08762504	2.003293867

- **Chromosome Number:** The first column (e.g., '1') indicates the chromosome on which the region is located.
- **Start:** The second column (e.g., '3660676') shows the starting position of the peak region on the chromosome.
- **End:** The third column (e.g., '3661050') indicates the ending position of the peak region.
- **Peak Length:** The fourth column (e.g., '375') might represent the length of the peak region.
- **Signal Strength/Score:** The next columns, like '210' and '62.0876250438913', could represent the signal location or a score that quantifies the enrichment.
- **Statistical Data:** The last column (e.g., '-2.003') is statistical data related to the peak, such as a p-value or fold-change.



# Hands-on and Practical Part



## Part 1: Introduction of galaxy platform

- Collecting data
- Gene Identification

# Introduction to Bioinformatics

## Understanding the Digital Frontier in Biomedicine

---

### Part 3: Quality Control

# Expected Errors



أكاديمية كاوهست  
KAUST ACADEMY

## Nucleic acid extraction



## Library preparation



## Sequencing and analysis



- Samples collection
- Nucleic acid extraction
- Quality Control
- Adaptor ligation/barcoding
- Size selection
- Amplification/purification
- Quality control

- Sequencing
- Data analysis
  - Base calling
  - Read alignment
  - Variant calling
  - Variant annotation



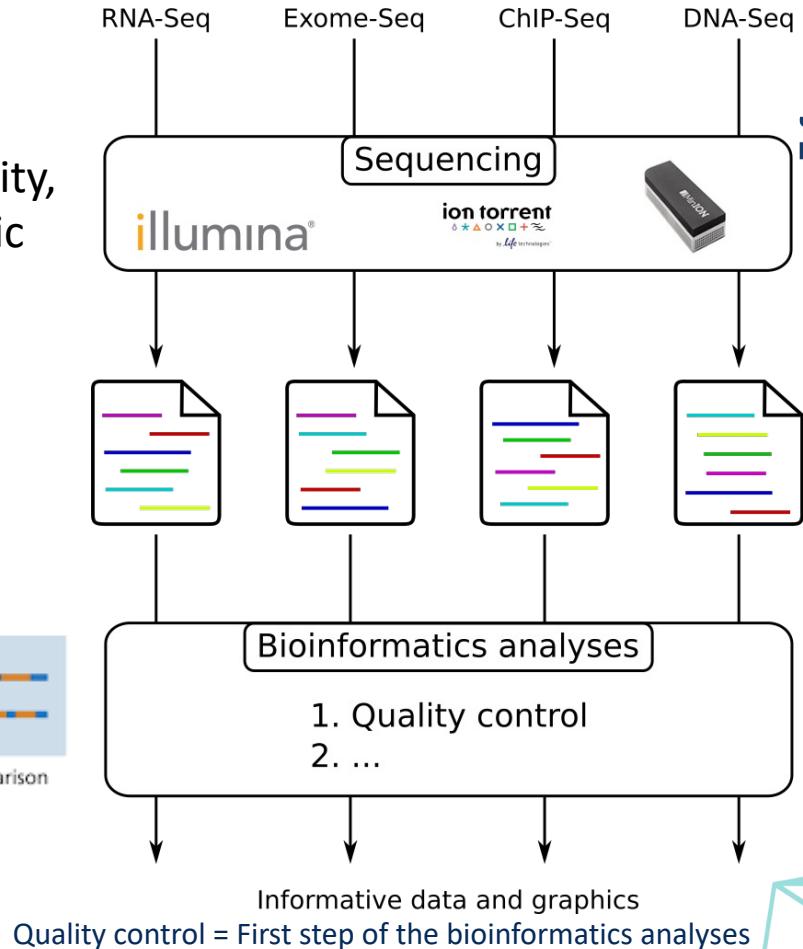
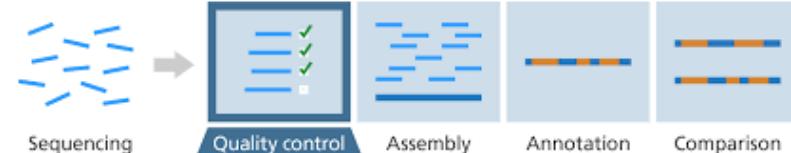
# Expected Errors

- **Sample Preparation:** Contamination with other DNA, RNA degradation, or insufficient quantity of starting material can lead to errors.
- **Library Construction:** Inaccurate size selection, adapter ligation problems, or PCR amplification errors can introduce biases.
- **Sequencing Process:** Issues like sequencing instrument calibration, cluster generation failures on flow cells, or phasing/pre-phasing problems can cause errors.
- **Base Calling:** Incorrect interpretation of fluorescent signals can lead to incorrect nucleotide identification.
- **Read Quality:** Low-quality reads can result from sequencing reaction issues or signal decay over time, especially at the end of reads.
- **Alignment:** Incorrectly mapped reads due to repetitive regions or sequence similarity can misrepresent the genomic location.
- **Variant Calling:** Misidentification of true genetic variants versus sequencing errors.

# Quality Control

The process of evaluating the quantity, purity, and intactness of the genomic DNA

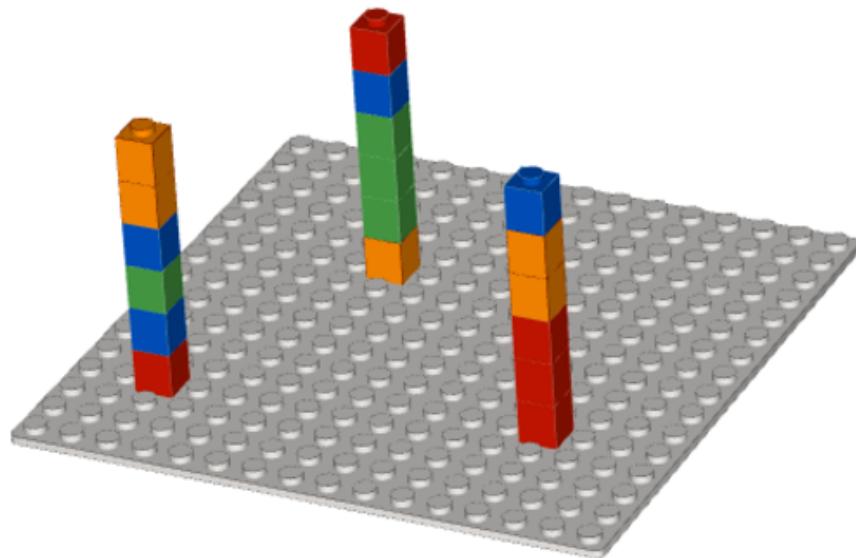
- Pre-sequencing Errors
- Sequencing Errors
- Post-sequencing Errors



# Example: Sequencing Errors



أكاديمية كاوهست  
KAUST ACADEMY

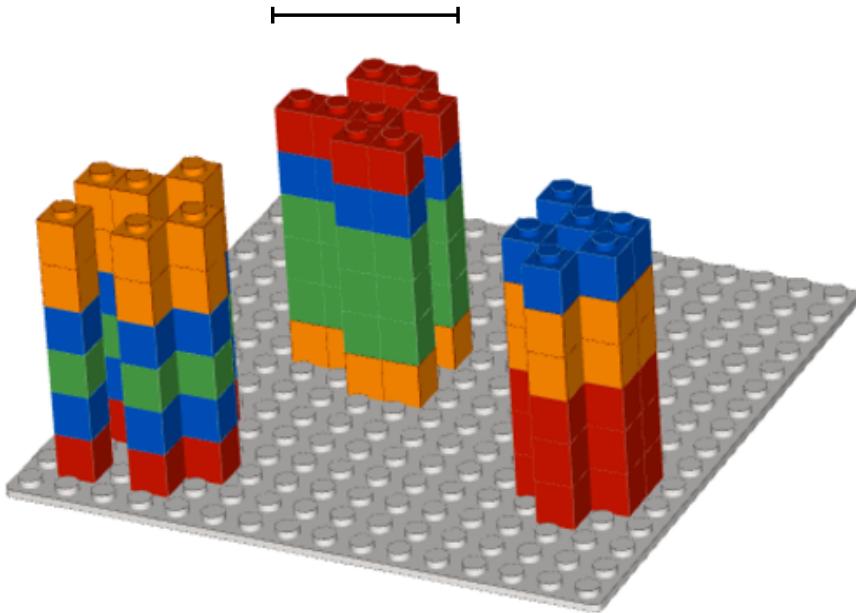


KAUST Academy

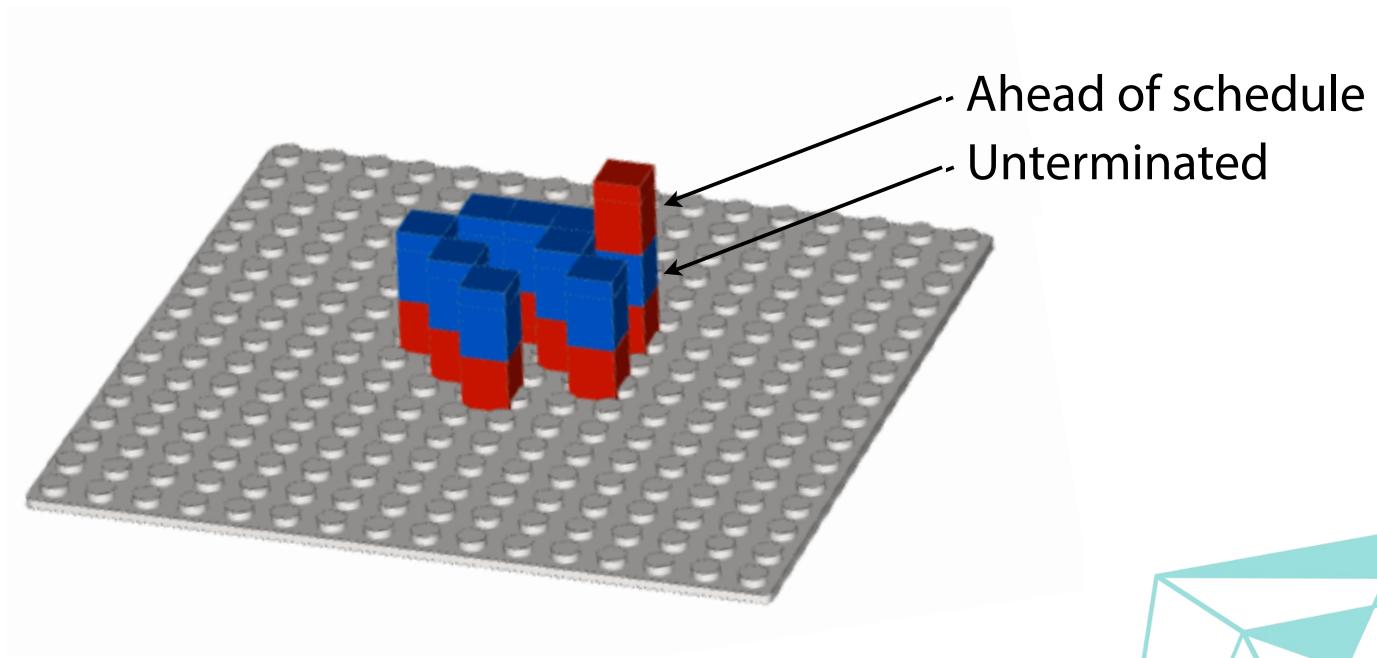


# Example: Sequencing Errors

Cluster of clones



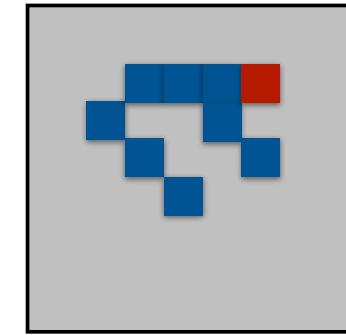
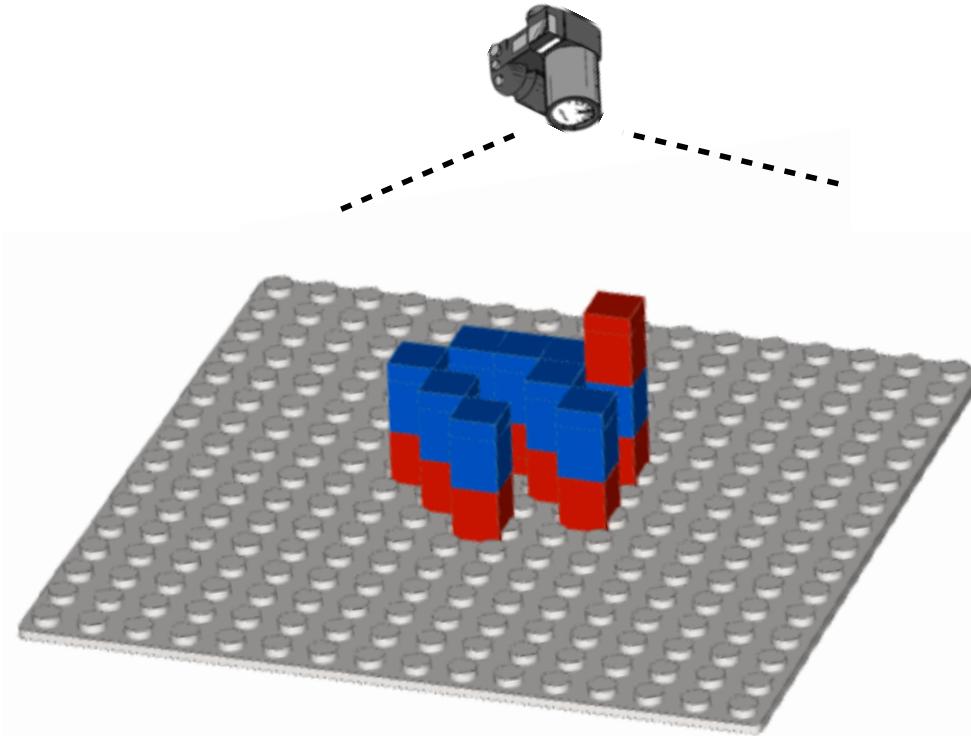
# Example: Sequencing Errors



# Example: Sequencing Errors



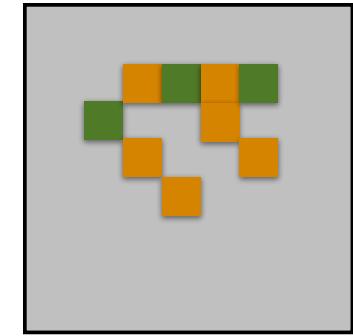
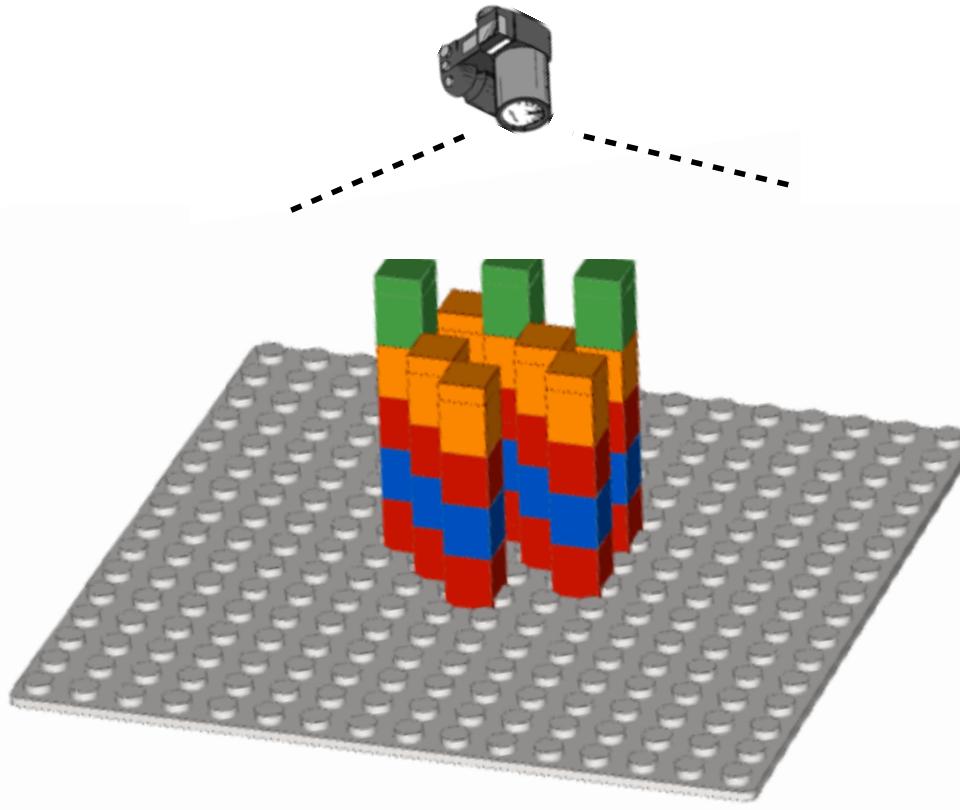
أكاديمية كاوهست  
KAUST ACADEMY



# Example: Sequencing Errors

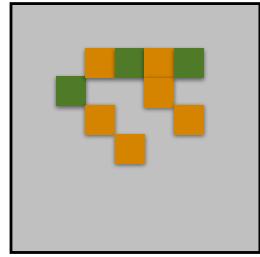


أكاديمية كاوهست  
KAUST ACADEMY



KAUST Academy

# Example: Sequencing Errors



Call: orange (C)

Estimate  $p$ , probability incorrect:  
non-orange light / total light

$$p = 3 \text{ green} / 9 \text{ total} = 1/3$$

$$Q = -10 \log_{10} 1/3 = 4.77$$



# Quality Score Calculation

$$Q = -10 \cdot \log_{10} p$$

Base quality                              Probability that base call is incorrect

$Q = 10 \rightarrow 1$  in 10 chance call is incorrect

$Q = 20 \rightarrow 1$  in 100

$Q = 30 \rightarrow 1$  in 1,000

# Explor FASTA File



أكاديمية كاوهست  
KAUST ACADEMY

- Sequences: FASTA

>Copy>Identifier1 (comment)

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

>Identifier2 (comment)

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
XX



# Explor FASTQ File



- Sequences: FASTAQ

>Copy>Identifier1 (comment)

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

+

QQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQ

QQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQ



# A read in FASTQ format

Name	@ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
Sequence	ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCTTAAAT
(ignore)	+
Base qualities	?@@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G



## FASTQ



# Base qualities

Bases and qualities line up:

AGCTCTGGTGACCCATGGGCAGCTGCTAGGGA  
||||| | | | | | | | | | | | | | | | | | | |  
HHHHHHHHHHHHHHHHHGCGC5FEFFFHBBBBB

Base quality is ASCII-encoded version of  $Q = -10 \log_{10} p$

# ASCII Code

0	<NUL>	32	<SPC>	64	@	96	'	128	Ä	160	†	192	¿	224	‡
1	<SOH>	33	!	65	A	97	a	129	Å	161	°	193	i	225	.
2	<STX>	34	"	66	B	98	b	130	Ç	162	¢	194	¬	226	,
3	<ETX>	35	#	67	C	99	c	131	É	163	£	195	✓	227	"
4	<EOT>	36	\$	68	D	100	d	132	Ñ	164	§	196	f	228	%
5	<ENQ>	37	%	69	E	101	e	133	Ö	165	•	197	≈	229	Â
6	<ACK>	38	&	70	F	102	f	134	Ü	166	¶	198	Δ	230	Ê
7	<BEL>	39	'	71	G	103	g	135	á	167	ß	199	«	231	Á
8	<BS>	40	(	72	H	104	h	136	à	168	®	200	»	232	É
9	<TAB>	41	)	73	I	105	i	137	â	169	©	201	...	233	È
10	<LF>	42	*	74	J	106	j	138	ã	170	™	202		234	Í
11	<VT>	43	+	75	K	107	k	139	ã	171	'	203	À	235	Î
12	<FF>	44	,	76	L	108	l	140	â	172	"	204	Ã	236	Ï
13	<CR>	45	-	77	M	109	m	141	ç	173	#	205	Õ	237	Ì
14	<SO>	46	.	78	N	110	n	142	é	174	Æ	206	Œ	238	Ó
15	<SI>	47	/	79	O	111	o	143	è	175	Ø	207	œ	239	Ô
16	<DLE>	48	0	80	P	112	p	144	ê	176	∞	208	-	240	apple
17	<DC1>	49	1	81	Q	113	q	145	ë	177	±	209	-	241	ò
18	<DC2>	50	2	82	R	114	r	146	í	178	≤	210	"	242	ú
19	<DC3>	51	3	83	S	115	s	147	ì	179	≥	211	"	243	û
20	<DC4>	52	4	84	T	116	t	148	î	180	¥	212	,	244	û
21	<NAK>	53	5	85	U	117	u	149	ï	181	µ	213	,	245	í
22	<SYN>	54	6	86	V	118	v	150	ñ	182	ð	214	÷	246	^
23	<ETB>	55	7	87	W	119	w	151	ó	183	Σ	215	◊	247	~
24	<CAN>	56	8	88	X	120	x	152	ò	184	Π	216	ÿ	248	-
25	<EM>	57	9	89	Y	121	y	153	ô	185	∏	217	ÿ	249	·
26	<SUB>	58	:	90	Z	122	z	154	ö	186	ƒ	218	/	250	.
27	<ESC>	59	;	91	[	123	{	155	õ	187	ä	219	€	251	°
28	<FS>	60	<	92	\	124		156	ú	188	ø	220	<	252	„
29	<GS>	61	=	93	]	125	}	157	ù	189	Ω	221	>	253	„
30	<RS>	62	>	94	^	126	~	158	û	190	æ	222	fi	254	„
31	<US>	63	?	95	_	127	<DEL>	159	ü	191	ø	223	fl	255	„

# Quality Scores



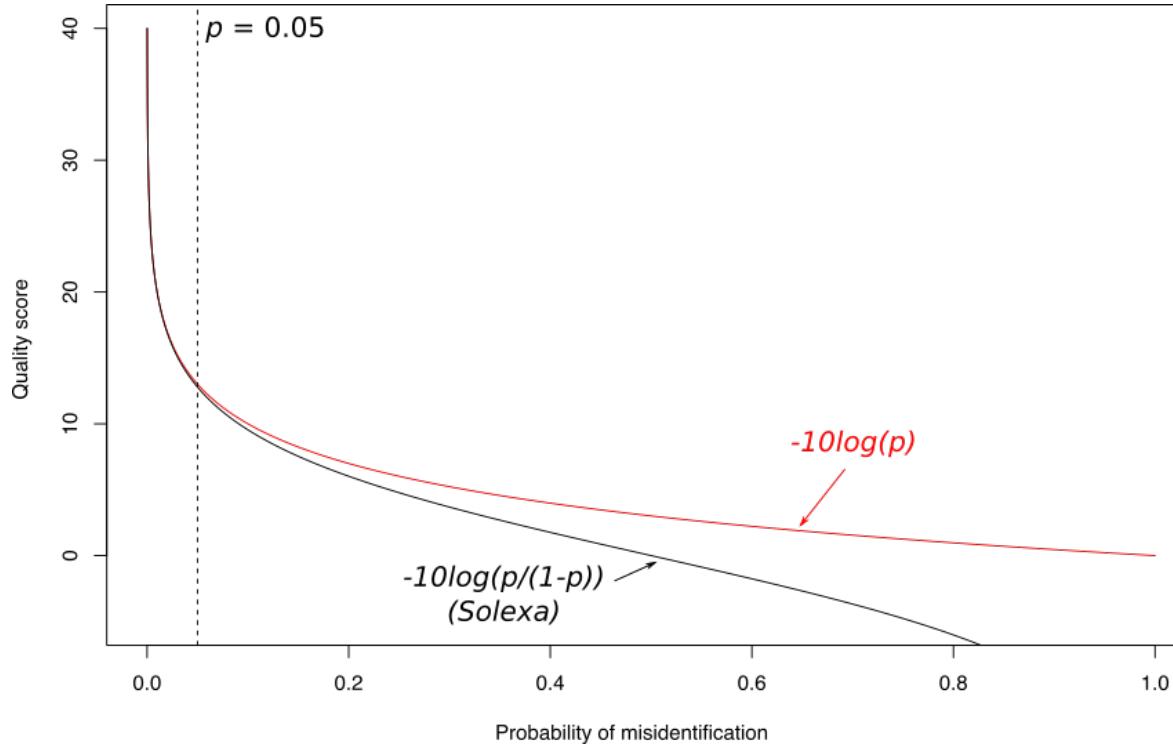
أكاديمية كاوهست  
KAUST ACADEMY

Measure of the quality of the identification of the nucleobases generated by automated DNA sequencing

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%



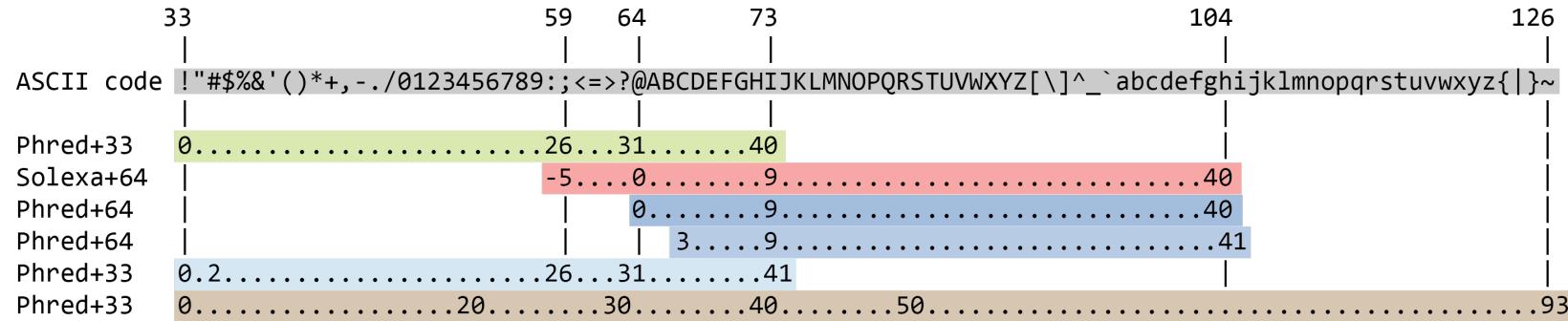
# Quality Score



# Quality score encoding



أكاديمية كاوهست  
KAUST ACADEMY

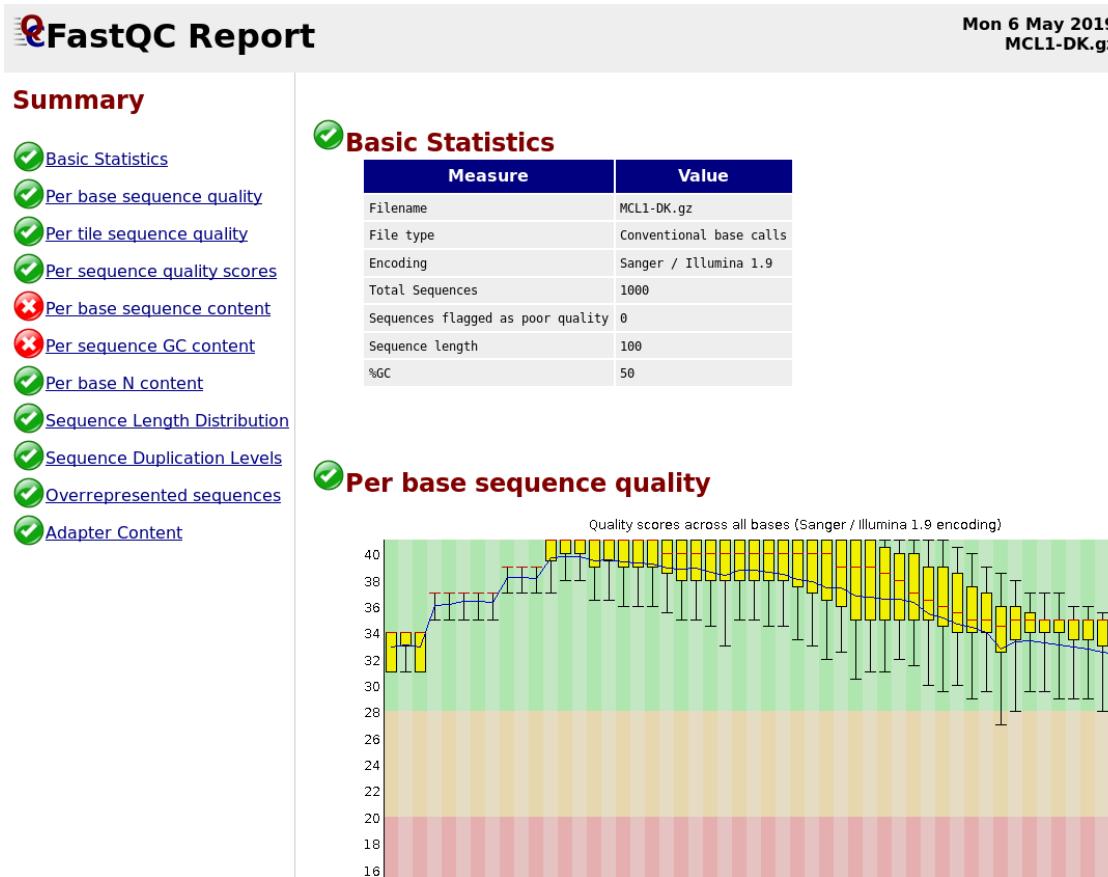


# Identifying Potential Quality Issues

- Per-base quality
- Pre-sequence quality
- Per-tile quality
- Per-base sequence content
- Per-sequence GC content

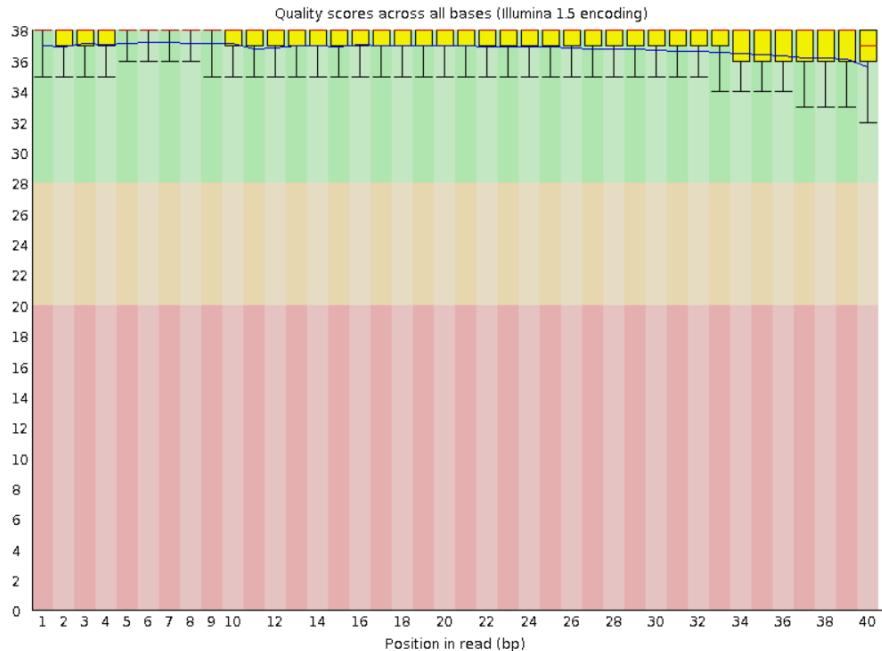


- FastQC tool : a versatile tool for short and long reads quality control

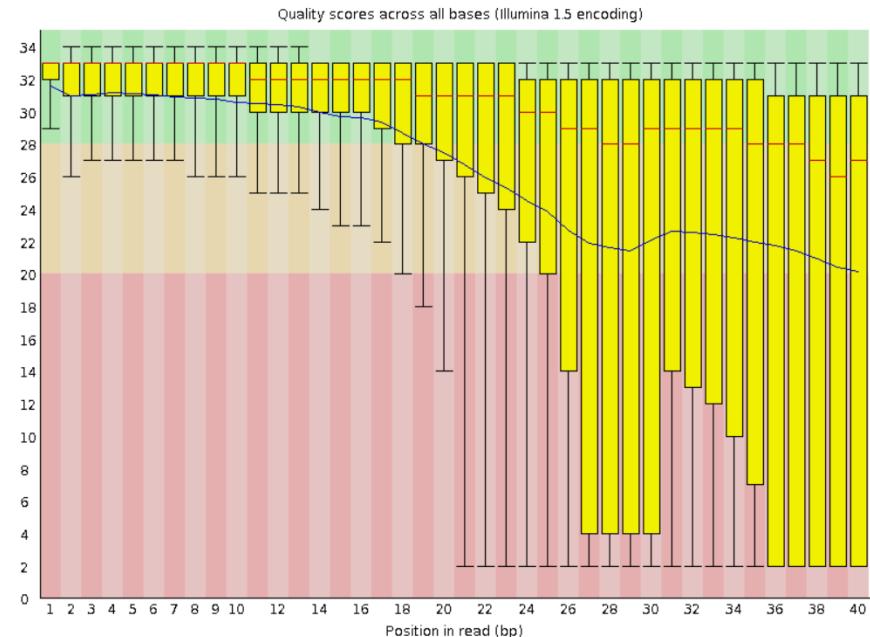


# Per-base Quality

## ✓ Per base sequence quality

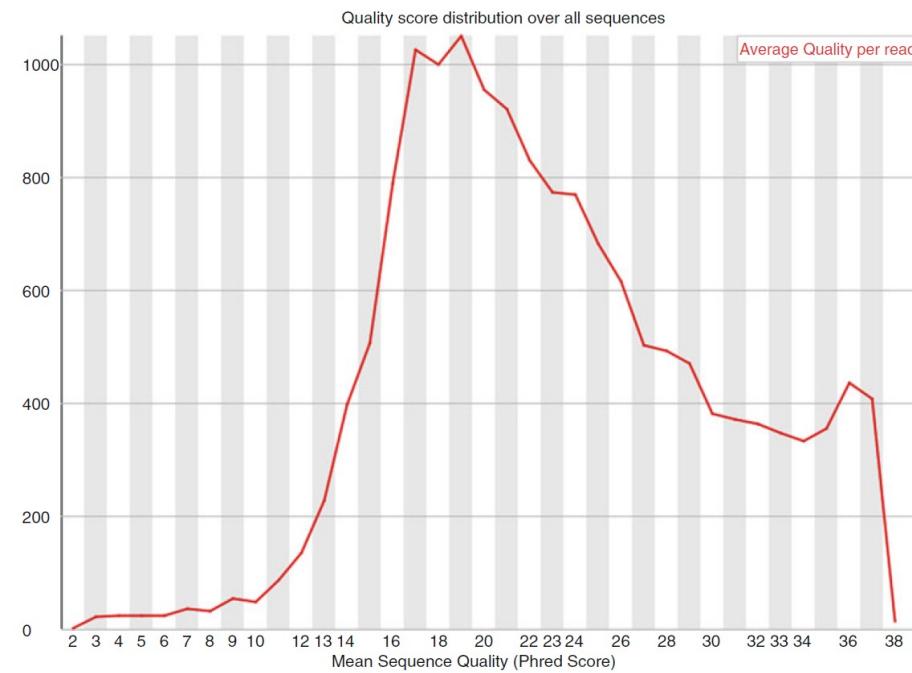
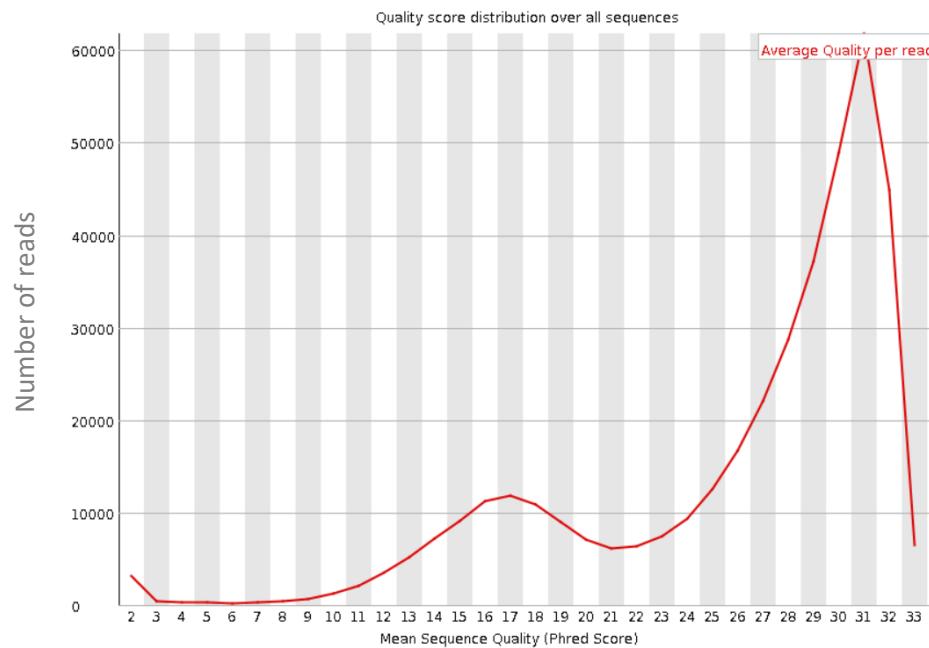


## ✗ Per base sequence quality

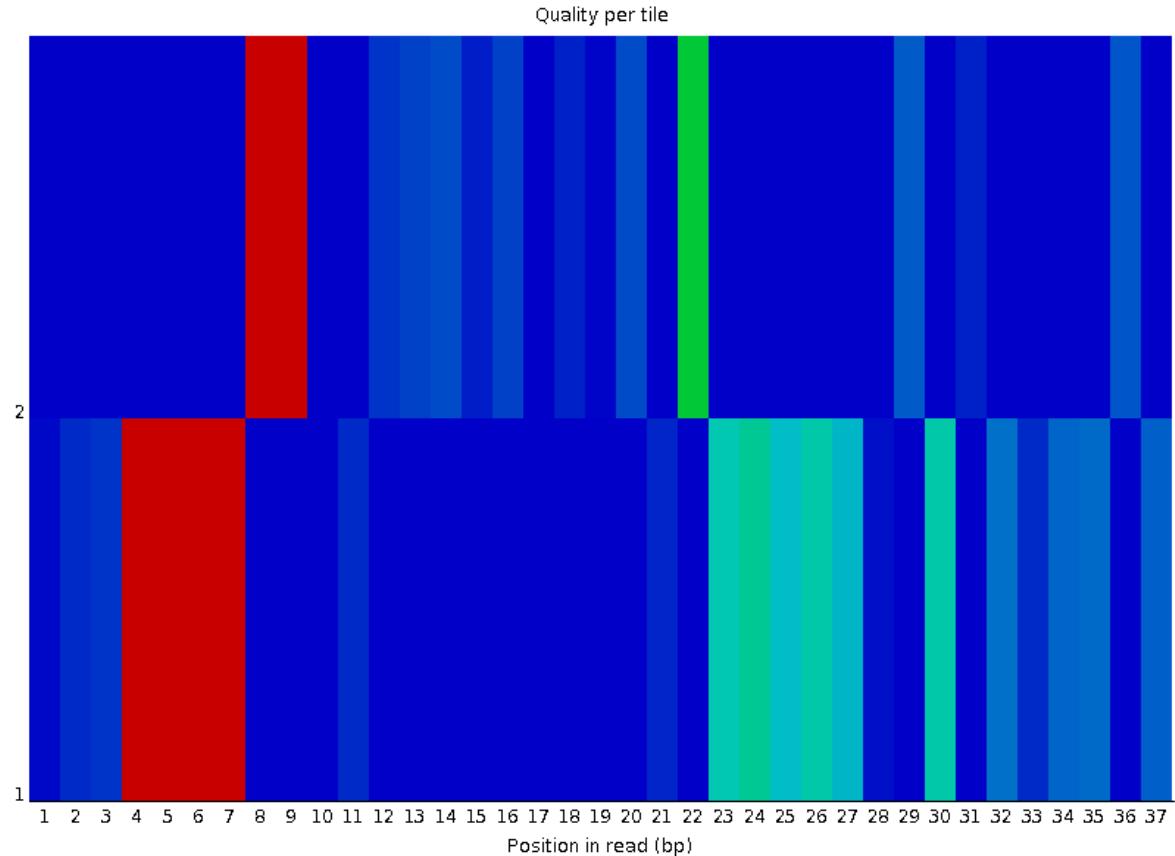
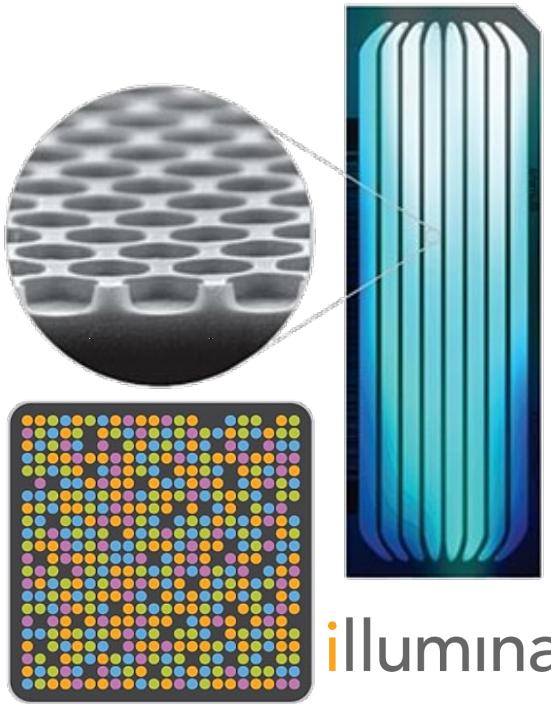


# Per-sequence Quality

## ✓ Per sequence quality scores

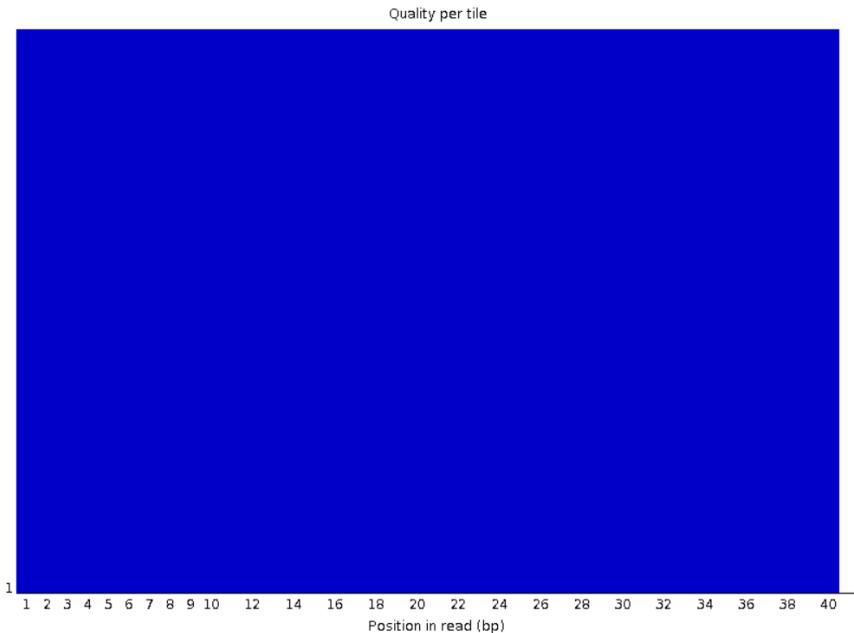


# Per-tile Quality

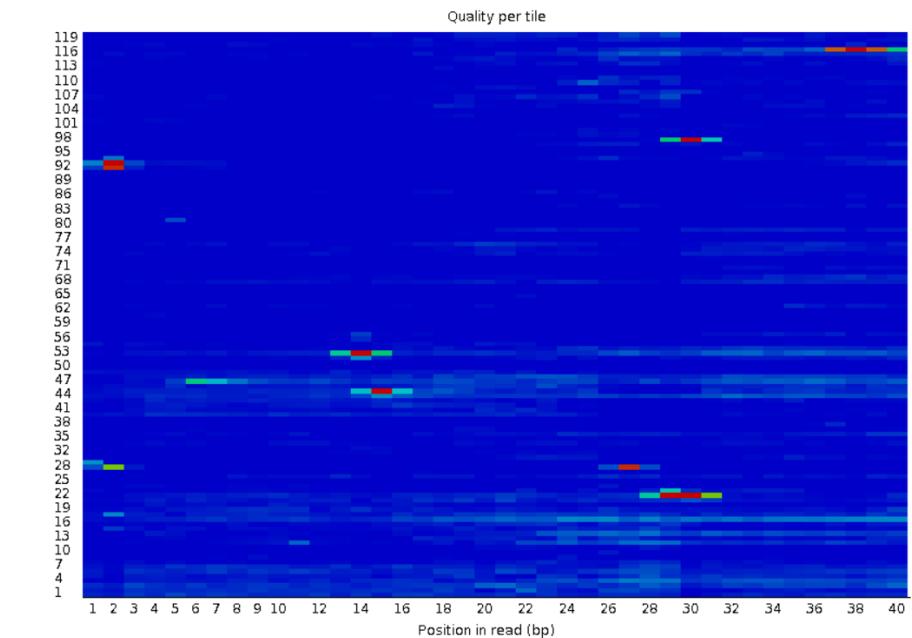


# Per-tile Quality

✓ Per tile sequence quality

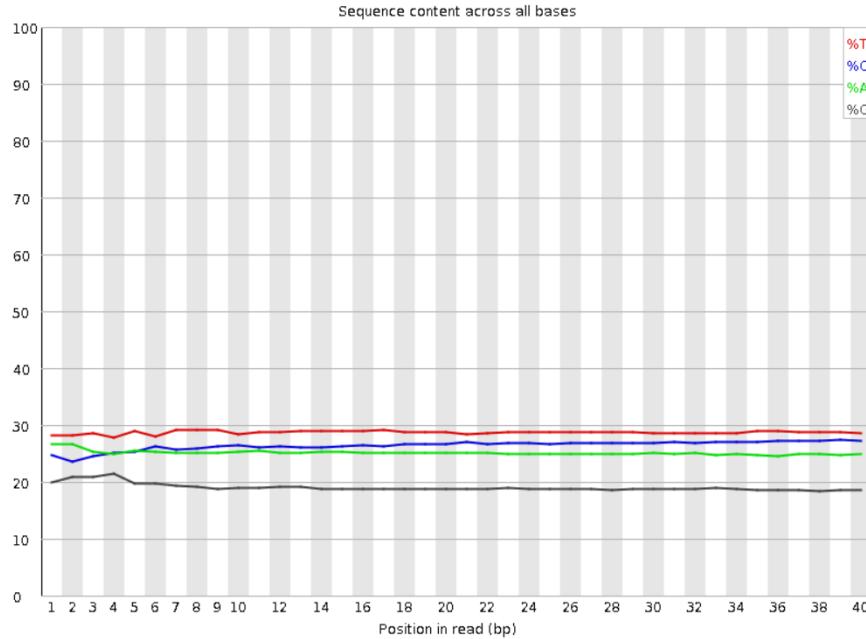


✗ Per tile sequence quality

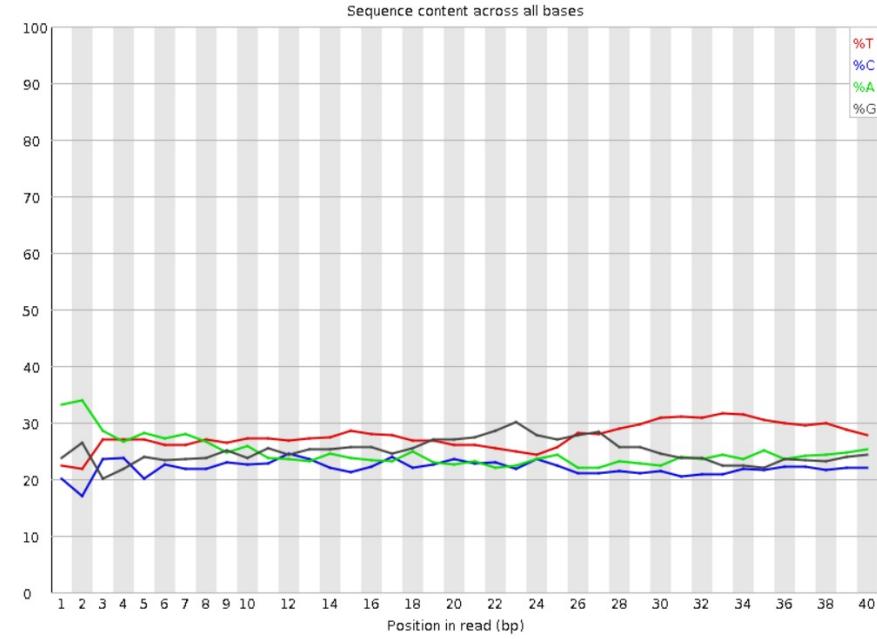


# Per-base Sequence Content

## ✓ Per base sequence content

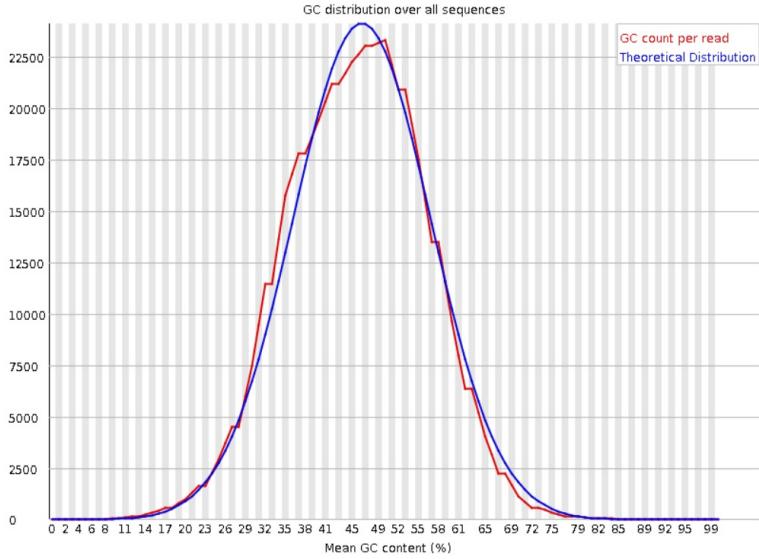


## ⚠ Per base sequence content

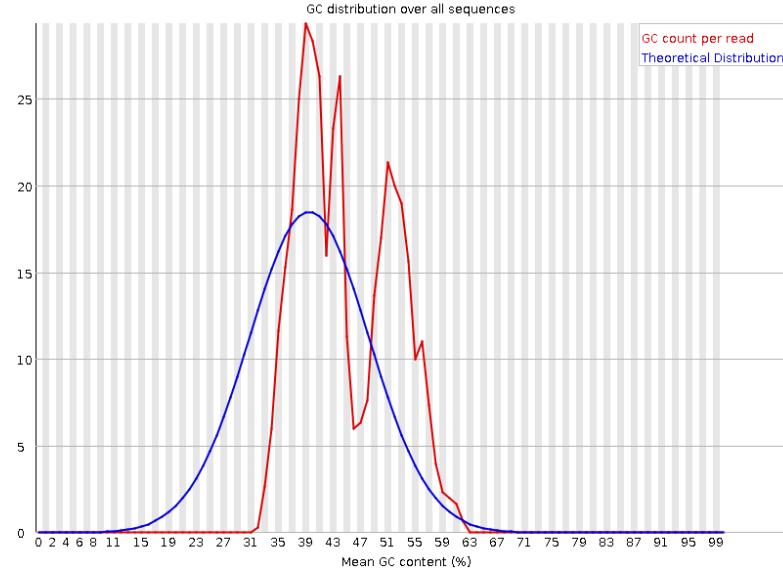


# Per-sequence GC Content

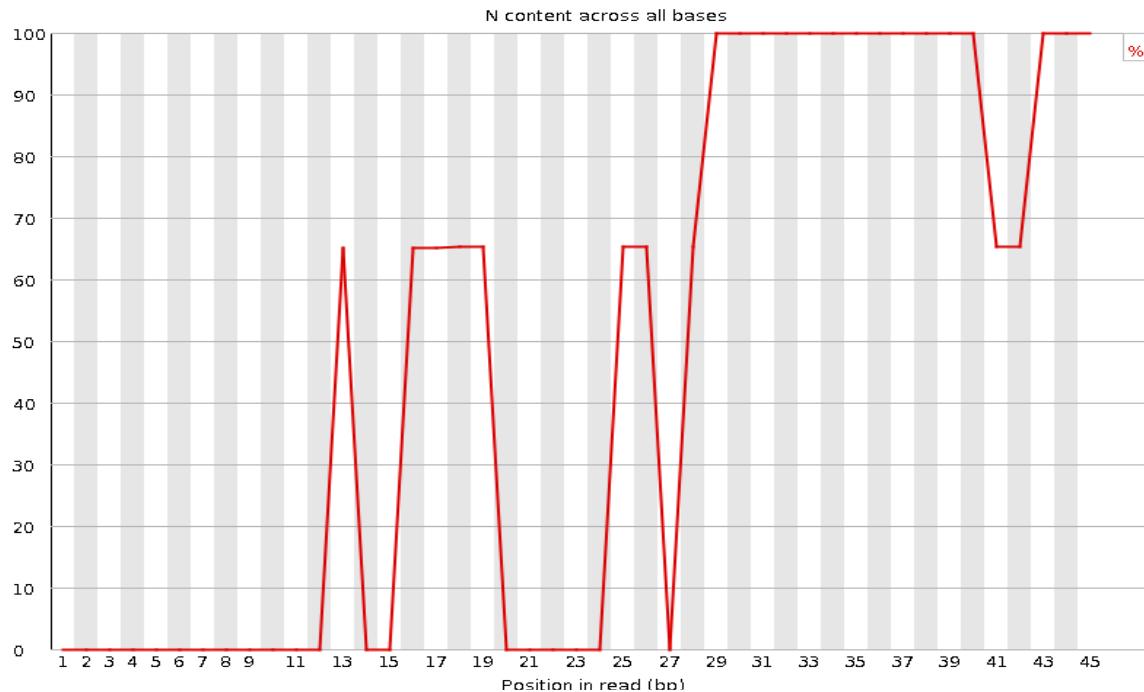
✓ Per sequence GC content



⚠ Per sequence GC content



# Per-base N Content

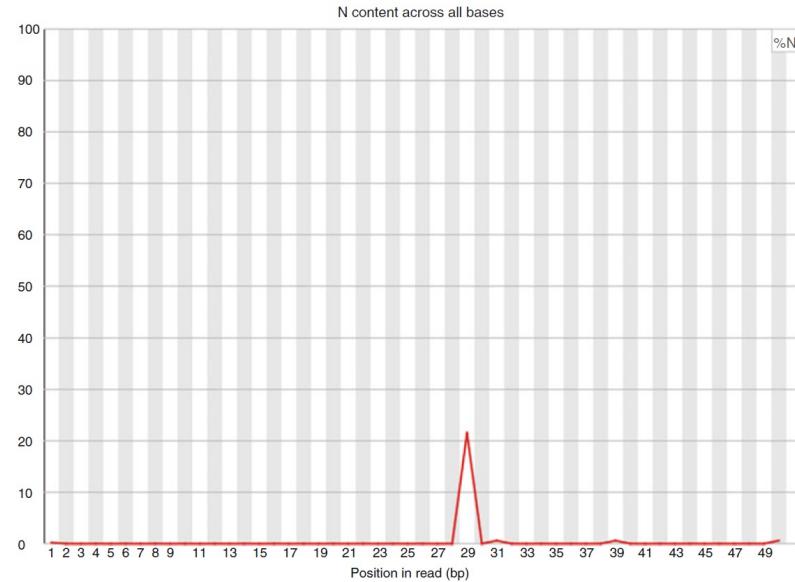
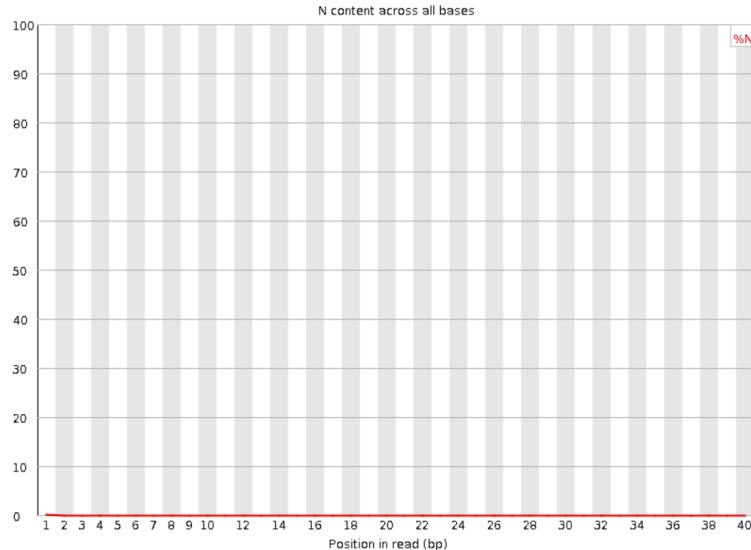


Name @ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1

Sequence ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAT

# Per-base N content

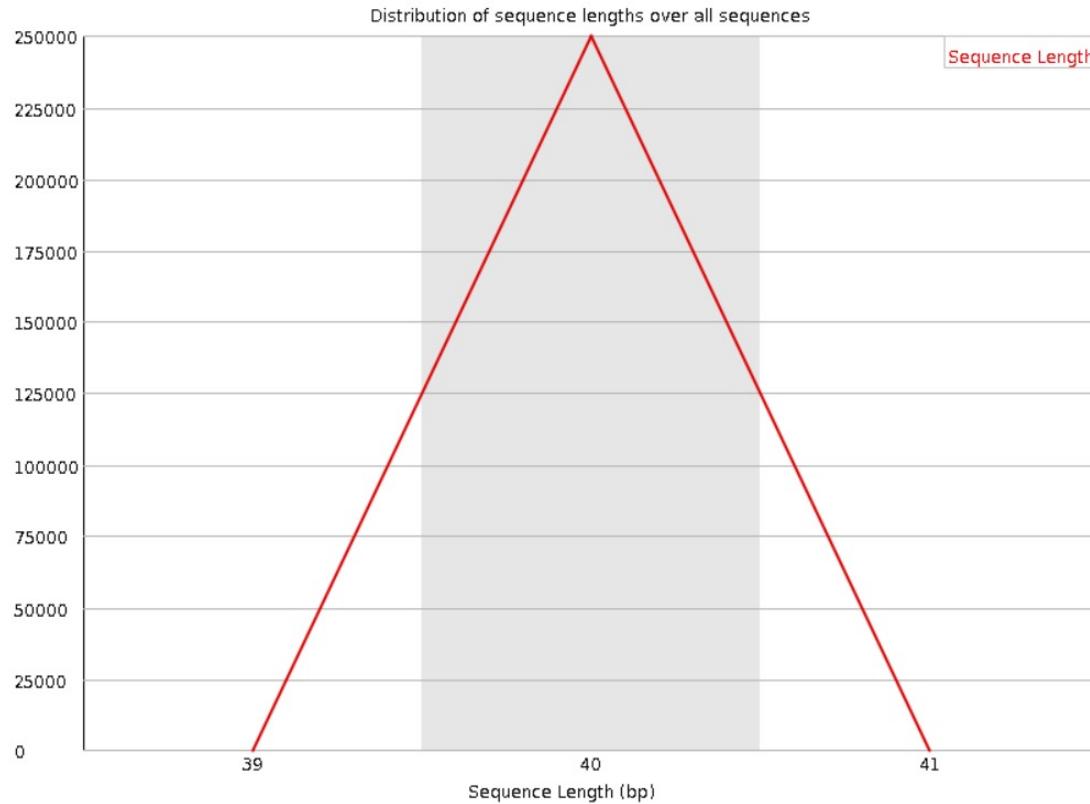
## ✓ Per base N content



# Sequence Length Distribution

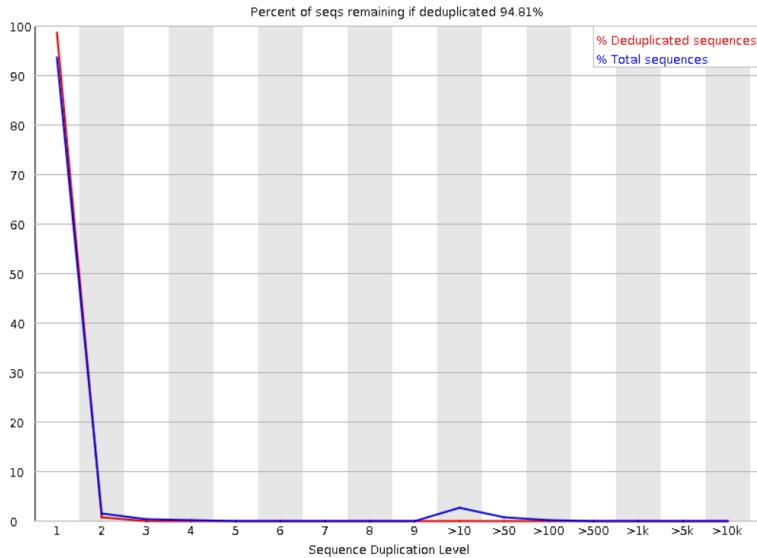


## Sequence Length Distribution

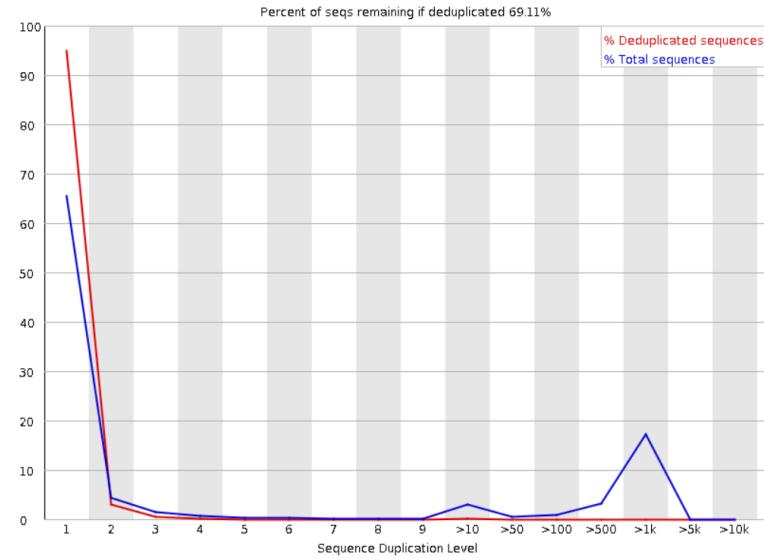


# Per-sequence Duplication Quality

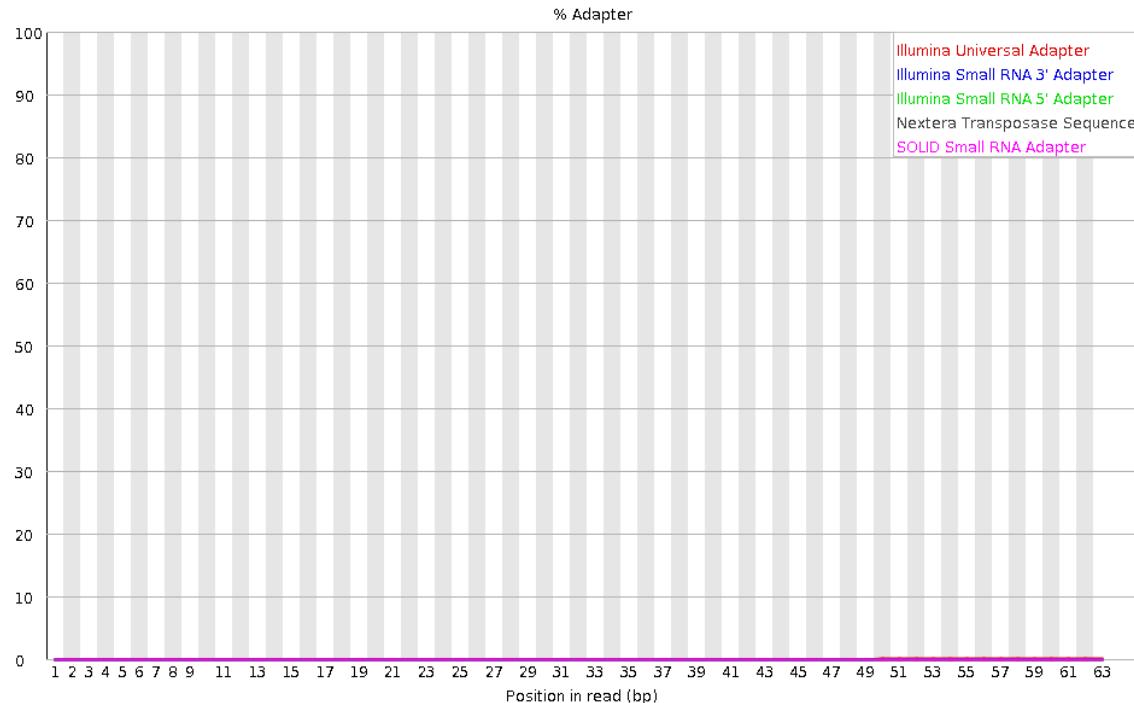
## Sequence Duplication Levels



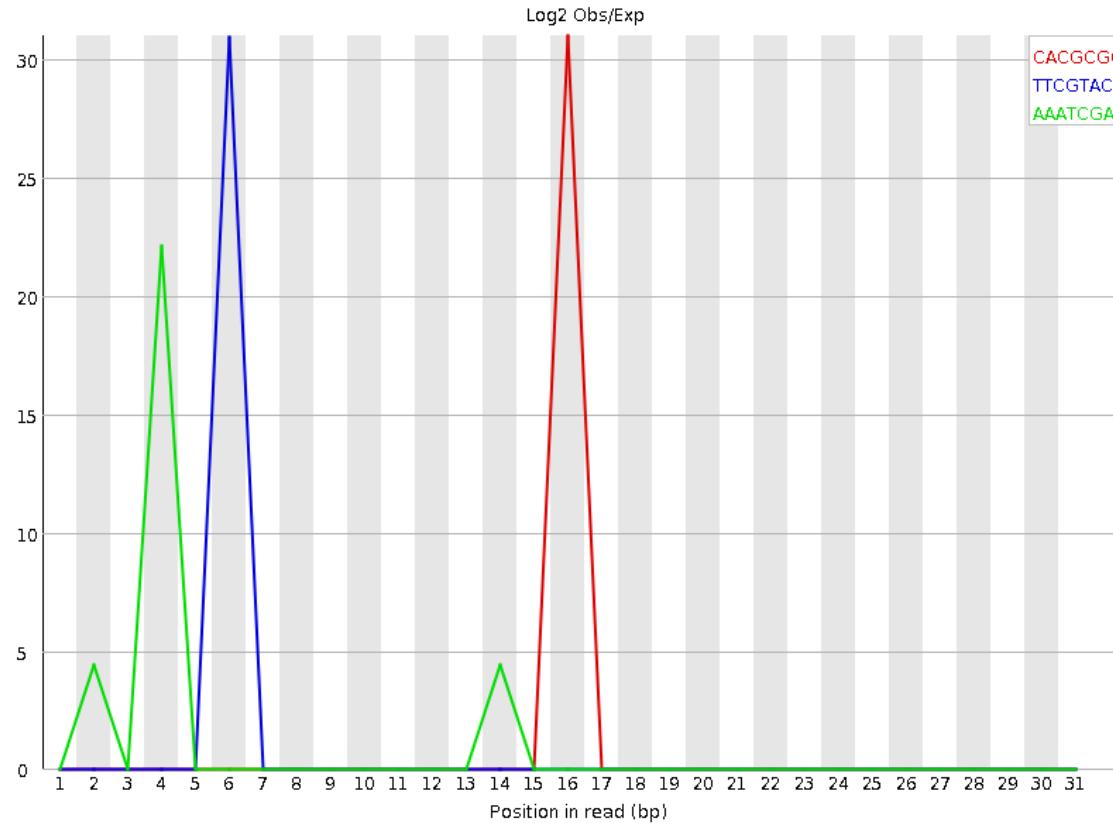
## Sequence Duplication Levels



# Tag sequences: Adapter Contamination



# Tag Sequences: K-mer Content



# Key Points



- Run quality control on every sequencing dataset before any other analyses
- Choose QC parameters carefully
- Re-run FastQC to check the impact of the quality control



# Hands-on and Practical Part



## Part 2: QC

- [Inspect a raw sequence file](#)
- [Assess quality with FASTQE](#)
- [Assess quality with FastQC](#)

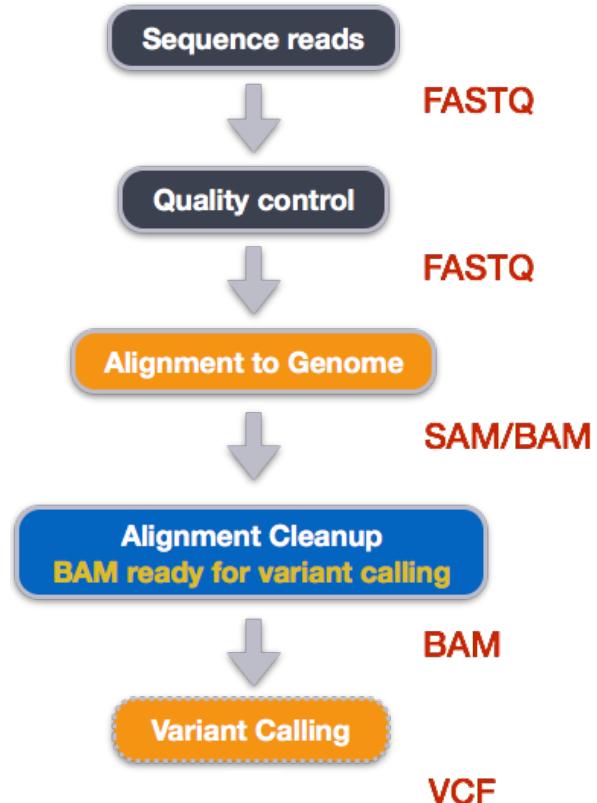
# Introduction to Bioinformatics

## Understanding the Digital Frontier in Biomedicine

---

### Part 3: Reference mapping

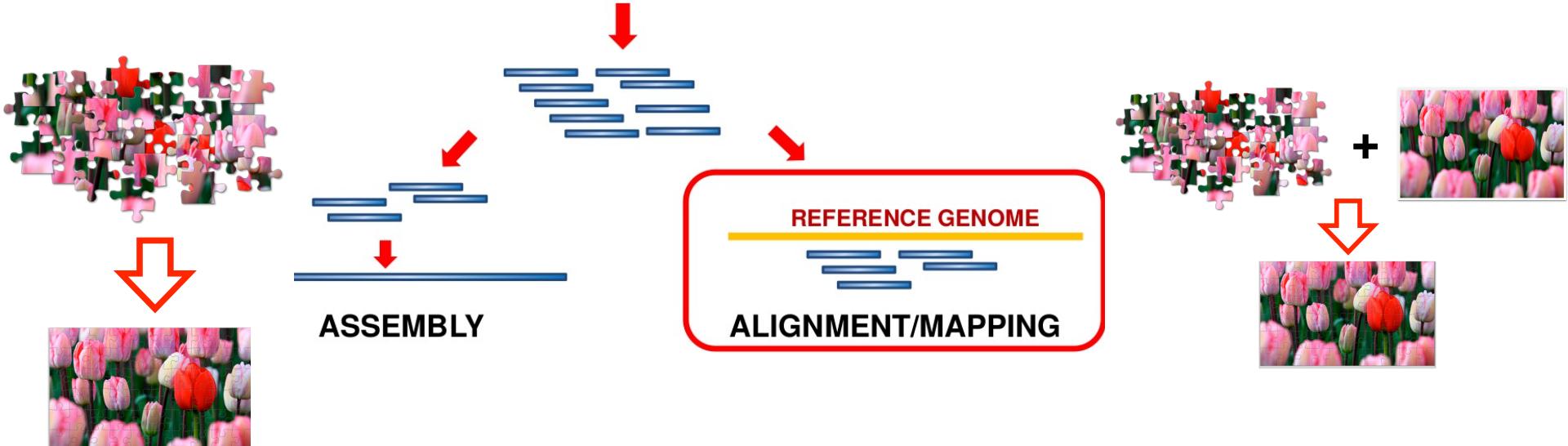
# Example of NGS Workflow



# What is Mapping or Alignment in NGS?



- Short reads must be combined into longer fragments
- **Mapping:** use a reference genome as a guide
- **De-novo assembly:** without reference genome



# Sequence Alignment

- Determine position of short read on the reference genome

**Reference:** . . . A A C G C C T T . . .

**Read:** A G G G G C C T T

# Sequence Alignment



أكاديمية كاوهست  
KAUST ACADEMY

- Determine position of short read on the reference

Reference: . . . A A - C G C C T T . . .  
| : - : | | | | |  
Read: A G G G G C C T T

| = match  
: = mismatch  
- = gap

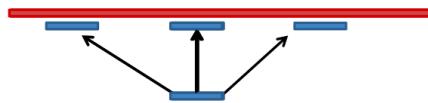


# Sequence Alignment

- Determine position of short read on the reference

Reference: . . . A A - C G G C C T T . . .	= match
. . .   : - :	: = mismatch
Read:            A G G G G C C T T	- = gap

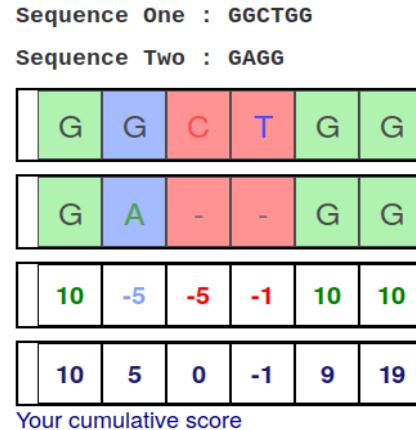
- Read could align to multiple places



- How to handle multi-mapped reads? Depends on tool:
  - Map to best region (but what is "best"? And what about ties?)
  - Map to all regions
  - Map to one region randomly
  - Discard read

# Alignment Scoring (basics)

- Reward for a match (e.g. +10), **penalty** for a mismatch (e.g. -5)
- **Penalty** for gaps
  - *Linear*: every gap same penalty (e.g. -5)
  - *Affine*: gap open vs gap extend (e.g. -5 and -1)
- Different tools use different scoring values (and give different results)



# Sequence Alignment



أكاديمية كاوهست

Reference: AAA CAGTGA GAA

Observed: AAA TCTCT GAA

## Alignment

AAA-CAGTGAGAA

||| - | -- | : : |||

AAATC--TCTGAA

Maybe like this?

AAACAGTGAGAA

||| - : | : : |||

AAA-TCTCTGAA

Or this?

AAACAGTGAGAA

||| : - | : : |||

AAAT-CTCTGAA

Or..?

AAACAGTCA----GAA

|||-----|||

AAA-----TCTCTGAA

What about this?



# Sequence Alignment



أكاديمية كاوفست  
KAUST ACADEMY

Reference: AAA CAGTGA GAA

Observed: AAA TCTCT GAA

## Alignment

## Tool

AAA-CAGTGAGAA  
| | | - | -- | : | : | |  
AAATC--TCTGAA

Novoalign

AAACAGTGAGAA  
| | | - : | : | : | |  
AAA-TCTCTGAA

Ssaha2

AAACAGTGAGAA  
| | | : - | : | : | |  
AAAT-CTCTGAA

BWA

AAACAGTCA----GAA  
| | | ----- | | |  
AAA-----TCTCTGAA

Complete Genomics



# Sequence Alignment



أكاديمية كاوهست  
KAUST ACADEMY

Reference: AAA CAGTGA GAA

Observed: AAA TCTCT GAA

## Alignment

AAA-CAGTGAGAA  
|||---|::|||  
AAATC--TCTGAA

AAACAGTGAGAA  
|||:-:|::|||  
AAA-TCTCTGAA

AAACAGTGAGAA  
|||:-:|::|||  
AAAT-CTCTGAA

AAACAGTG-----GAA  
||||-----|||  
AAA-----TCTCTGAA

## Variant calls

ins T  
del AG  
sub GA -> CT

del C  
sub AG -> TC  
sub GA -> CT

snp C -> T  
del A  
snp G -> C  
sub GA -> CT

del CAGTGA  
ins TCTCT



# Sequence Alignment



أكاديمية كاوهست  
KAUST ACADEMY

- Lego time! Who wants to volunteer?
- Or try this [online sequence alignment game](#):

Level 1   Level 2   Level 3   Customize Score Table   New Game

If you wish, you can change the default Match, Mismatch, and Gap scores in the below textboxes.

Match Score : 20   Mismatch Score : -10   Gap Score : -20   Max Score : 90   Your Score : -60

Sequence One : TATGATTACT

Sequence Two : TGTATACT

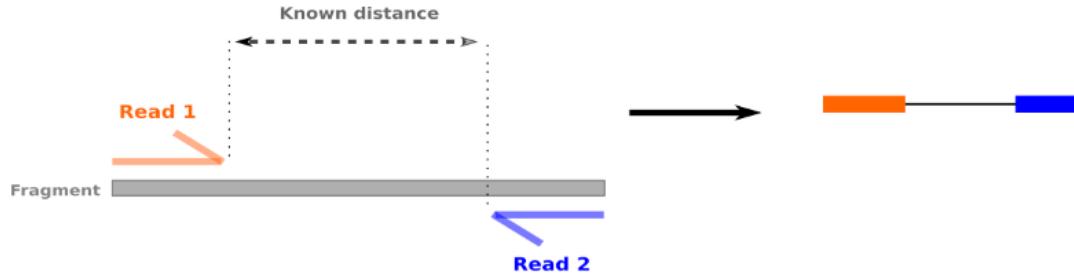
T	A	T	G	A	T	T	A	C	T
T	G	T	A	T	A	C	T	-	-
20	-10	20	-10	-10	-10	-10	-10	-20	-20
20	10	30	20	10	0	-10	-20	-40	-60

Your cumulative score

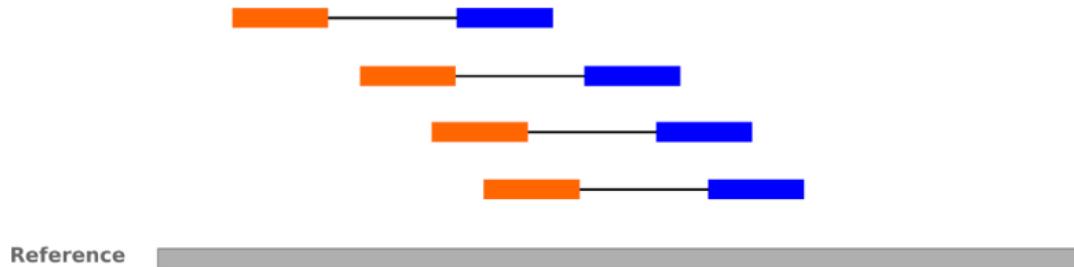


# Paired-end sequencing

- **Sequencing:** Cut longer fragments of DNA, sequence only the ends

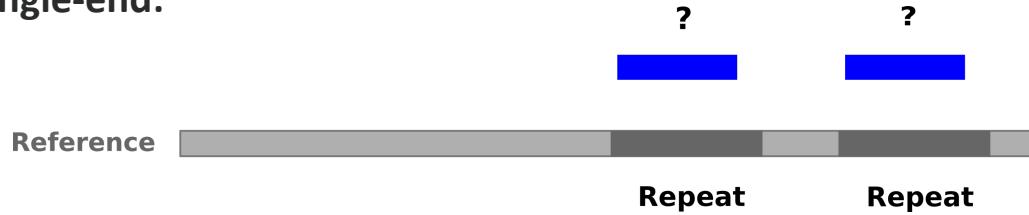


- **Mapping:** known distance between reads improves accuracy



# Repeats

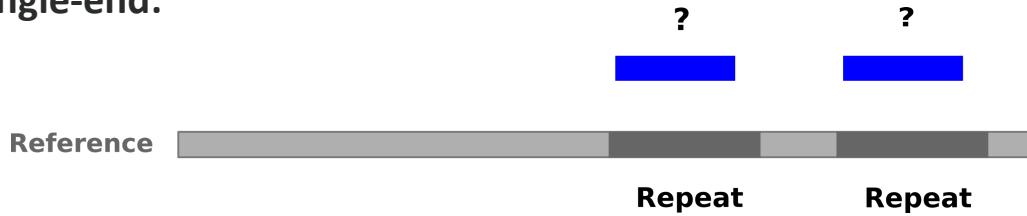
- Multi-mapped reads (e.g. because of repeats) may now be resolved
- **Single-end:**



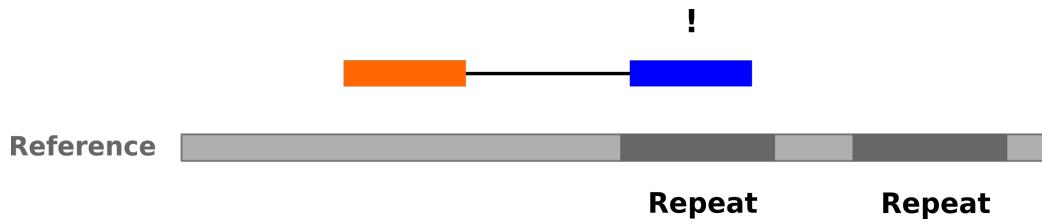
- **Paired-end:**

# Repeats

- Multi-mapped reads (e.g. because of repeats) may now be resolved
- **Single-end:**

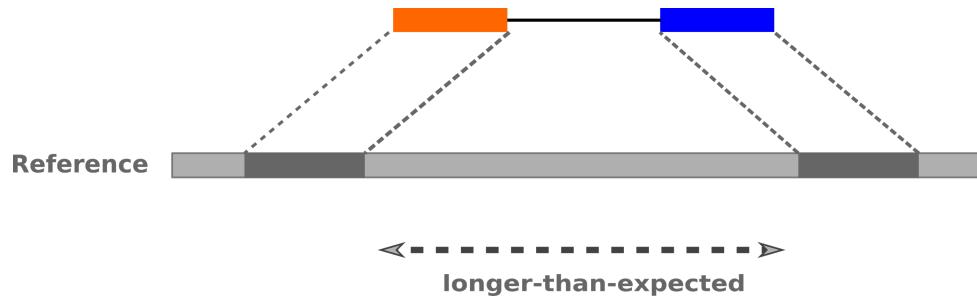


- **Paired-end:**

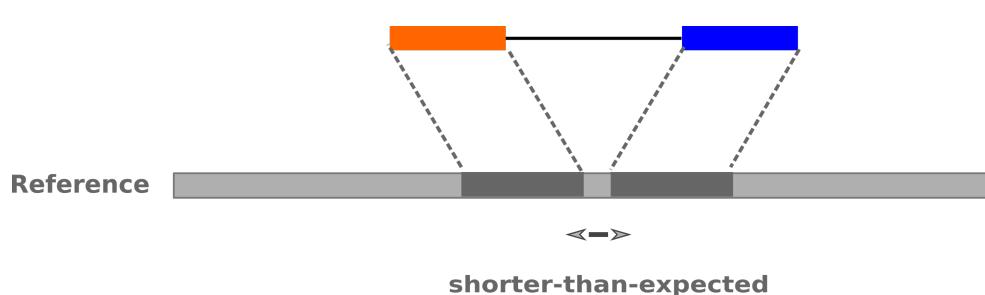


# InDels (Insertions / Deletions)

- Discordant insert size may indicate insertion or deletion between reads
- **Deletions:** Longer mapping distance than expected

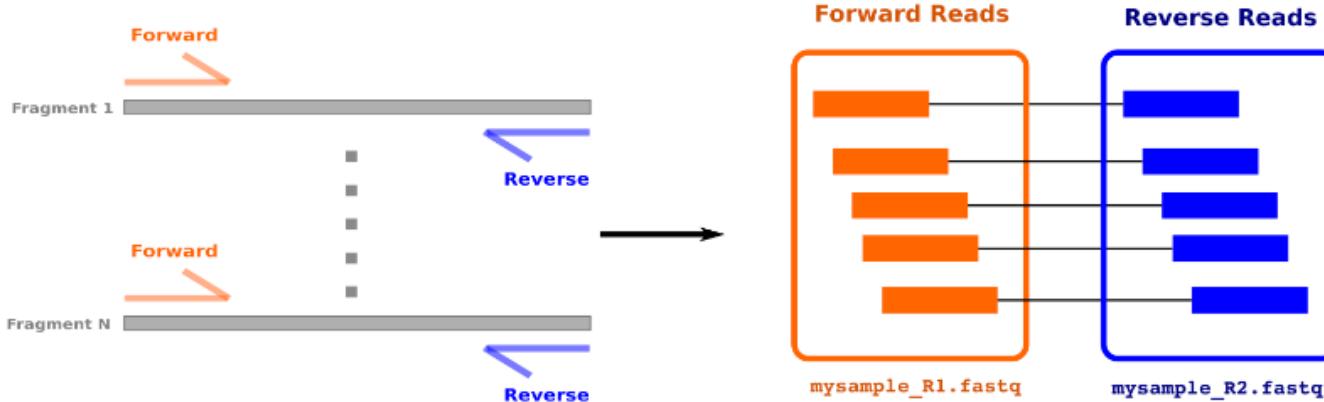


- **Insertions:** Shorter mapping distance than expected



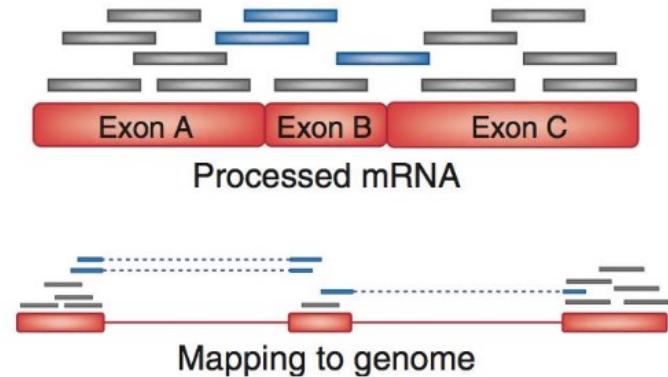
# Paired-end FASTQ Files

- Sequencer produces two FASTQ files:
  - **Forward** reads (usually **\_1** or **\_R1** in file name)
  - **Reverse** reads (usually **\_2** or **\_R2** in file name)



# Choosing an Aligner

- Each tool makes **different choices** during alignment
  - Choice of aligner may **affect downstream results**
  - Default options may not be best for your data
- Best tool for your data **depends on many factors**
  - Type of experiment (e.g. DNA, RNA, Bisulphite)
  - Sequencing platform
  - Compute resources vs sensitivity
  - Read characteristics (paired/single end, read length)

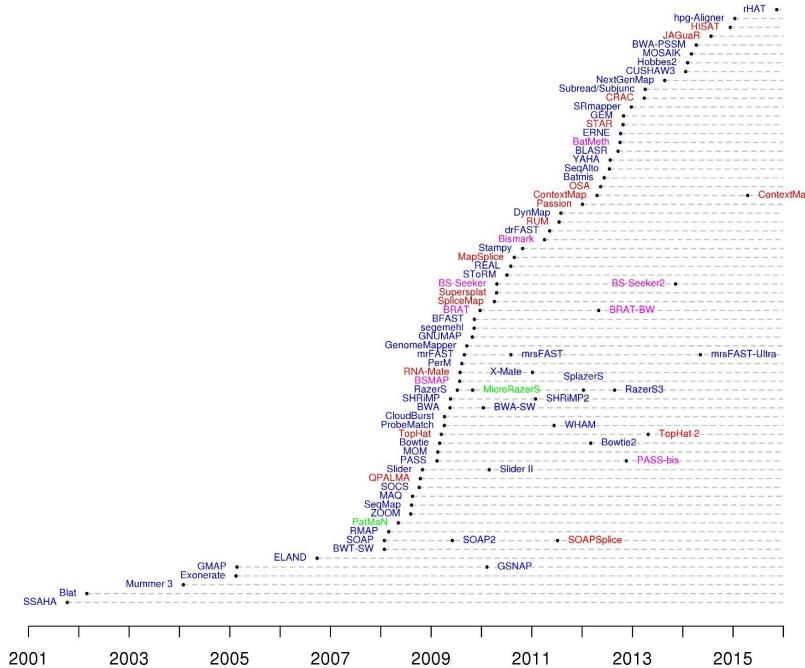


- Know your data!
- “... there is no tool that outperforms all of the others in all the tests. Therefore, the end user should clearly specify [their] needs in order to choose the tool that provides the best results.” - Hatem et al BMC Bioinformatics 2013, 14:184
  - [DOI: 10.1186/1471-2105-14-184](https://doi.org/10.1186/1471-2105-14-184)

# Mapping Tools



أكاديمية كاوهست  
KAUST ACADEMY



60+ different mappers, many comparison papers. Figure from [10.1109/bioinformatics.bts605](https://doi.org/10.1109/bioinformatics.bts605)

Mapping tool	Uses	Characteristics	KAUST ACADEMY
HISAT2	DNA/RNA	Short reads. Based on <a href="#">GCSA</a> . <a href="#">Reference</a> .	
RNASTAR	RNA	Short reads. Extremely fast. High sensitive and accuracy. Based on Maximal Mappable Prefixes (MMPs). <a href="#">Reference</a> .	
BWA-MEM2	DNA	Short reads. Twice as faster as BWA-MEM. Memory efficient. Based on <a href="#">Burrows-Wheeler</a> . <a href="#">Reference</a> .	
Minimap2	DNA/RNA	Long reads (PacBio and ONT). Extremely fast. Based on <a href="#">DALIGN</a> and <a href="#">MHAP</a> . <a href="#">Reference</a> .	
Bismark	DNA/RNA	Short reads. Bisulfite treated sequencing. Based on <a href="#">GCSA</a> <a href="#">Reference</a> .	
BBMap	DNA/RNA	Short and long reads (PacBio and ONT). Memory demanding. <a href="#">Reference</a> .	
Whisper 2	DNA	Short reads. Indel sensitive. Variant-calling oriented. <a href="#">Reference</a> .	
S-conLSH	DNA	Long reads (ONT). High sensitivity and accuracy. <a href="#">Reference</a> .	

# SAM/BAM File Format

```
1:497:R:-272+13M17D24M 113    chr1 497  37      37M      15      100338662 0
CGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAG 0;=====9;>>>>=>>>>>>>=>>>>>>>> XT:A:U  NM:i:0
SM:i:37  AM:i:0  X0:i:1  X1:i:0  XM:i:0  X0:i:0  XG:i:0  MD:Z:37
19:20389:F:275+18M2D19M 99      chr1    17644   0       37M      =      17919
314     TATGACTGCTAATAATACCTACACATGTTAGAACCAT >>>>>>>>>>>>>><>>><>>4:>>><9
RG:Z:UM0098:1  XT:A:R  NM:i:0  SM:i:0  AM:i:0  X0:i:4  X1:i:0  XM:i:0  X0:i:0
XG:i:0  MD:Z:37
19:20389:F:275+18M2D19M 147     chr1    17919   0       18M2D19M =
17644   -314     GTAGTACCAACTGTAAGTCCTTATCTTCATACTTG<:44999;499<8<8<<<8<<><<<><7<,<<<><<
XT:A:R  NM:i:2  SM:i:0  AM:i:0  X0:i:4  X1:i:0  XM:i:0  X0:i:1  XG:i:2  MD:Z:
18^CA19
9:21597+10M2I25M:R:-209  83      chr1    21678   0       8M2I27M =
21469   -244     CACCACATCACATATACCAAGCCTGGCTGTGTTCT<;9<<5><<<><<<><<><><><9>><>>9>>><>
XT:A:R  NM:i:2  SM:i:0  AM:i:0  X0:i:5  X1:i:0  XM:i:0  X0:i:1  XG:i:2  MD:Z:35
```

**SAM:** Sequence Alignment Map

**BAM:** Binary (compressed) SAM; not human-readable Formats

## SAM/BAM File Format

Read ID	Read sequence	Read position	Alignment (37 Matches)			
1_497:R:-272+13M17D24M-113	chr1 497 37	37M	15	100338662	0	XT:A:U NM:i:0
CGGGCTGACCCGAGGGAGAACGTGCTCCGCCCTTCAG	0;=====9;>>>>=>>>>>>=>>>>>>					
SM:i:37 AM:i:0 X0:i:1 X1:i:0 XM:i:0	X0:i:0 XC:i:0 MD:Z:37					
19:20389:F:275+18M2D19M 99	chr1 17644 0 37M =					17919
314 TATGACTGCTAATAATACCTACACATGTTAGAACCAT	>>>>>>>>>>>>><>><>>4:>><9					
RG:Z:UM0098:1 XT:A:R NM:i:0 SM:i:0	AM:i:0 X0:i:4 X1:i:0 XM:i:0 X0:i:0					
XG:i:0 MD:Z:37						
19:20389:F:275+18M2D19M 147 chr1 17919 0 18M2D19M =						
17644 -314 GTAGTACCAACTGTAAGTCCTTATCTTCATACTTTGT	:44999:499<8<8<<8<<<7;<<><<					
XT:A:R NM:i:2 SM:i:0 AM:i:0 X0:i:4 X1:i:0	XM:i:0 X0:i:1 XG:i:2 MD:Z:					
18^CA19						
9:21597+10M2I25M:R:-209 83 chr1 21678 0 8M2I27M =						
21469 -244 CACCACATCACATATACCAAGCCTGGCTGTCTTCT	<;9<<5><<<><<>><>><9>><>>9>><>					
XT:A:R NM:i:2 SM:i:0 AM:i:0 X0:i:5 X1:i:0	XM:i:0 X0:i:1 XG:i:2 MD:Z:35					

- Original read information (from FASTQ) plus mapping information
    - Position on reference, alignment, quality score, uniqueness, ...

# Genome Browsers

- Visualise aligned reads (BAM files)



This is [IGV \(Integrative Genome Browser\)](#) DOI: [10.1038/nbt.1754](https://doi.org/10.1038/nbt.1754)

# Key points

- Mapping is not trivial
- There are many mapping tools, best choice depends on your data
- Choice of mapper can affect downstream results
- Know your data!
- Genome browsers can be used to view aligned reads

# Hands-on and Practical Part



## Part 3: Sequence alignments

- Map reads on a reference genome
- Inspect a BAM/SAM file



# **Done with Day 1, Heyyyyy!**

**Thank You !**