



In God we trust; all
others bring data.

W. Edwards Deming

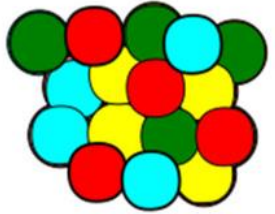


BULK RNASEQ **WORKSHOP**

What we want ?



heterogenous tissue

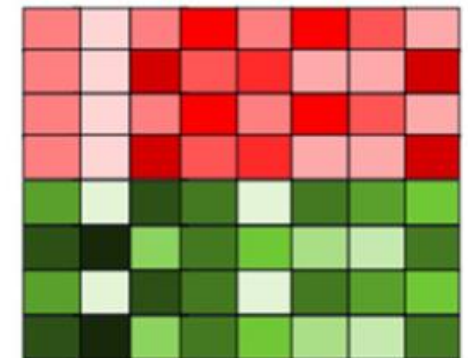


total RNA extraction



...CATCCTAGCTA...

sequencing



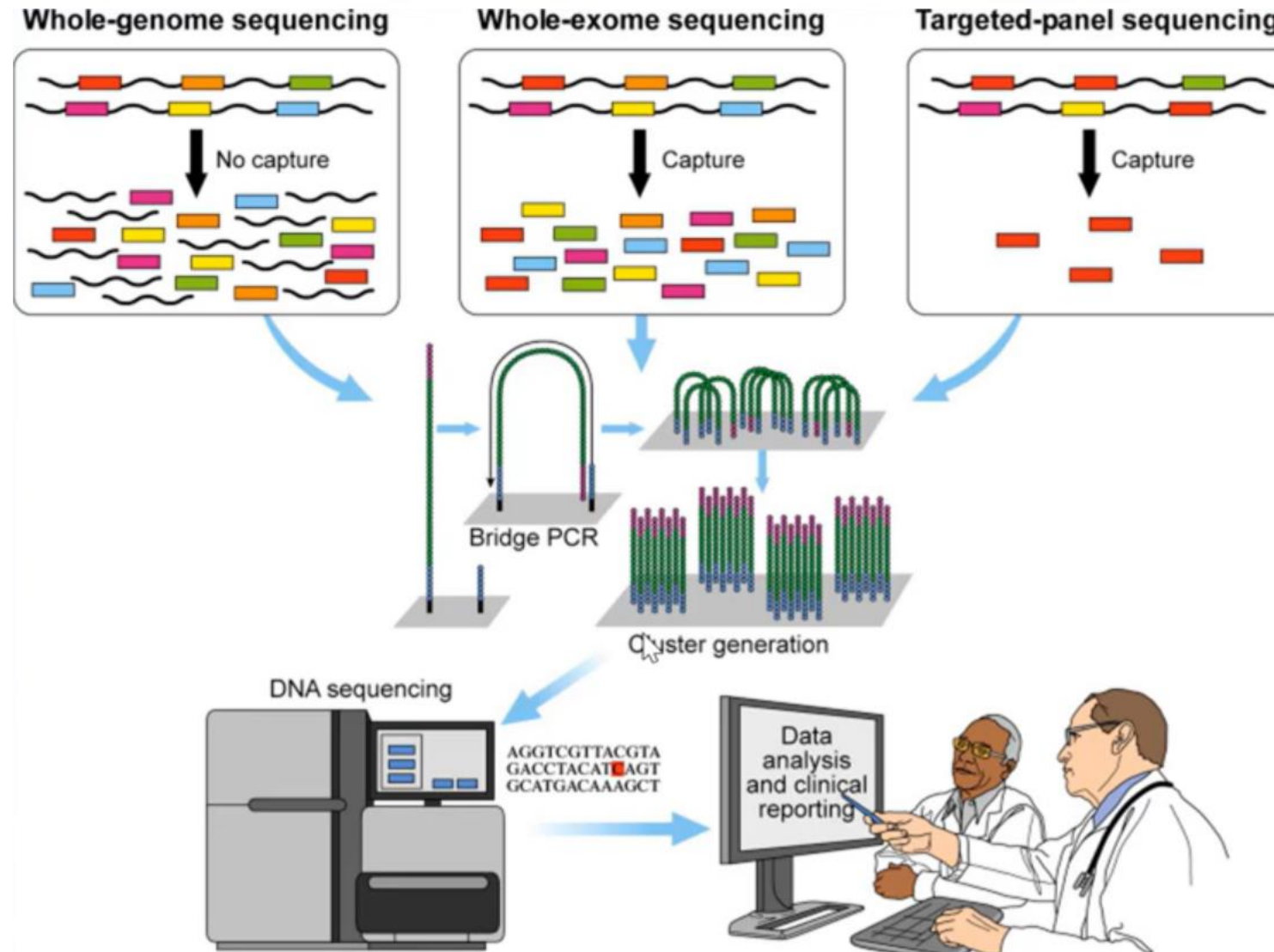
average expression data



What is Bulk RNA sequencing (Bulk RNA-seq)?

- Bulk RNA sequencing is the **method** of choice for transcriptomic analysis of tissue sections, or biopsies.
- It measures the **average expression level** of individual **genes** across hundreds to millions of input cells and is useful to get a global idea of gene expression differences between samples.

Bulk RNA-seq Pipeline

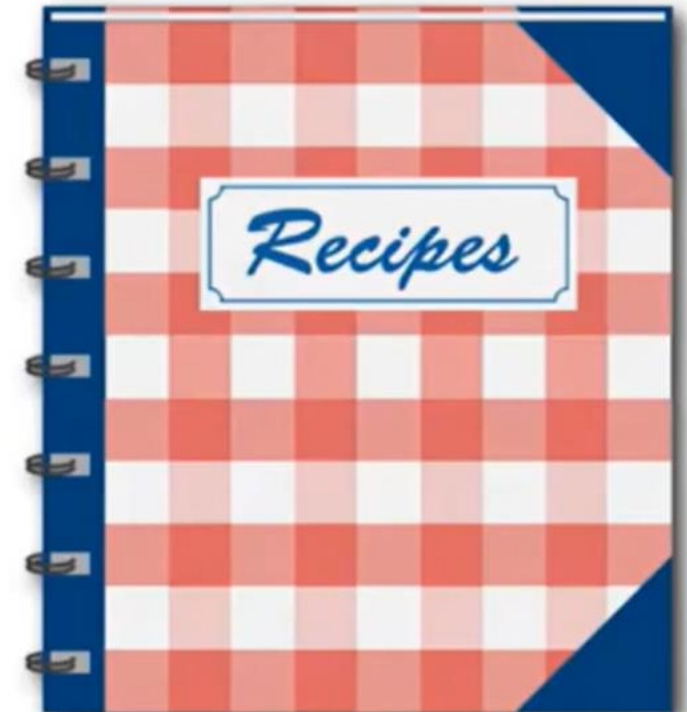
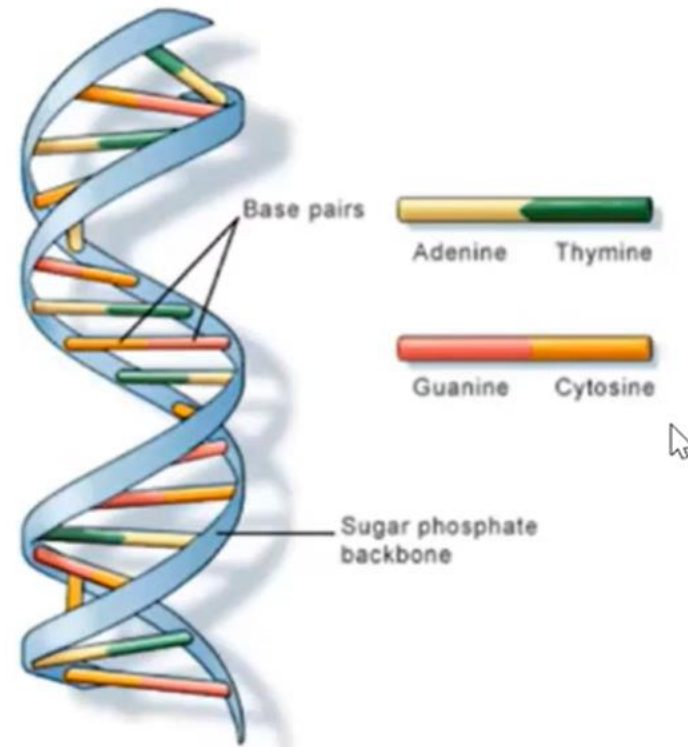




DNA

DNA

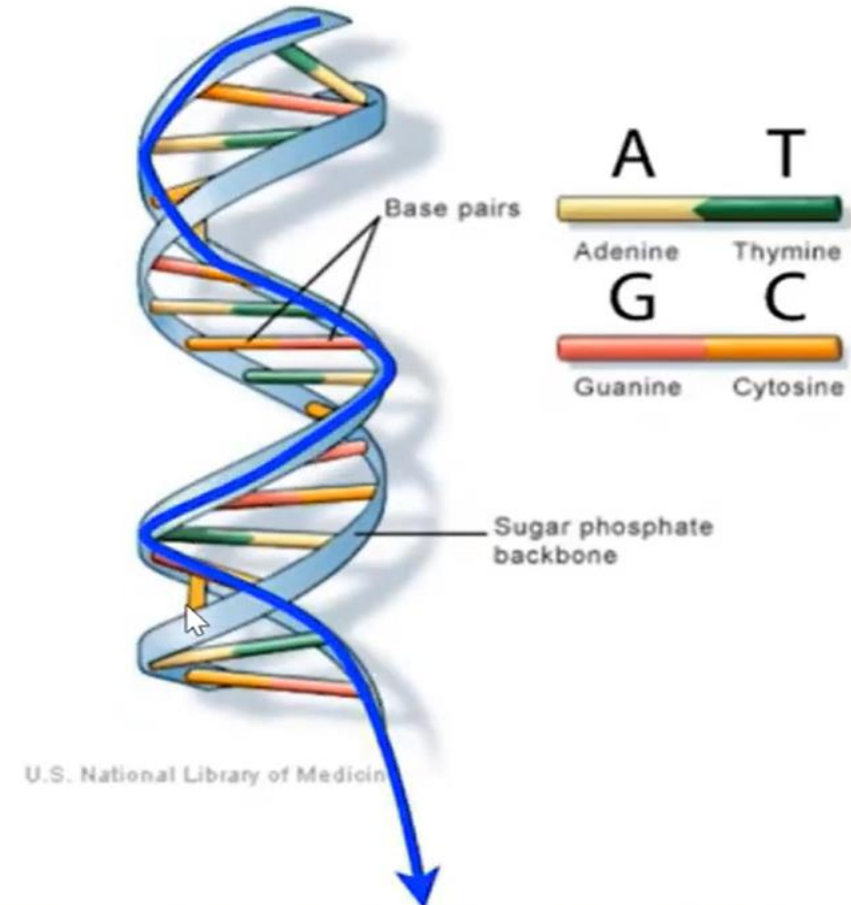
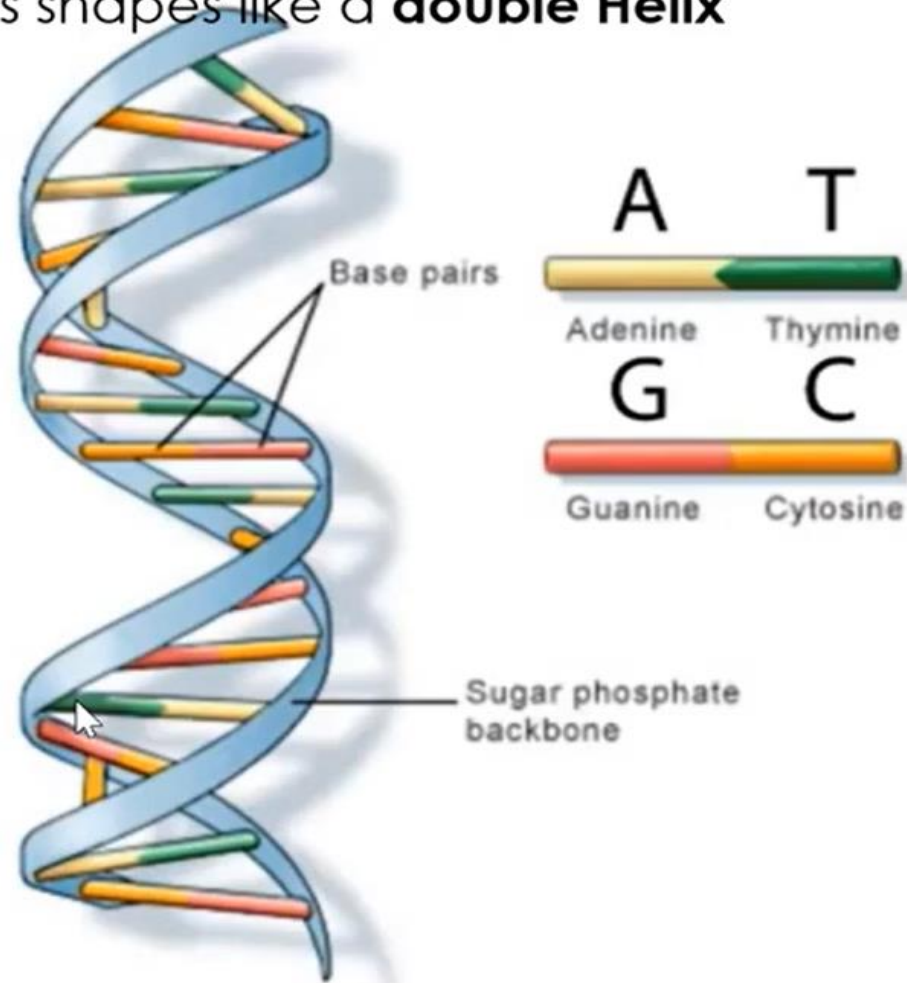
DNA is the kind of **Molecule** that **encodes your Genome** (all your **Genetic Information**). Separate **recipes** for each type of molecule (each little physical piece that make up brain cells, skin cells, heart cells building and maintaining you)





DNA

This recipe is not written in English.
DNA Molecule is shapes like a **double Helix**



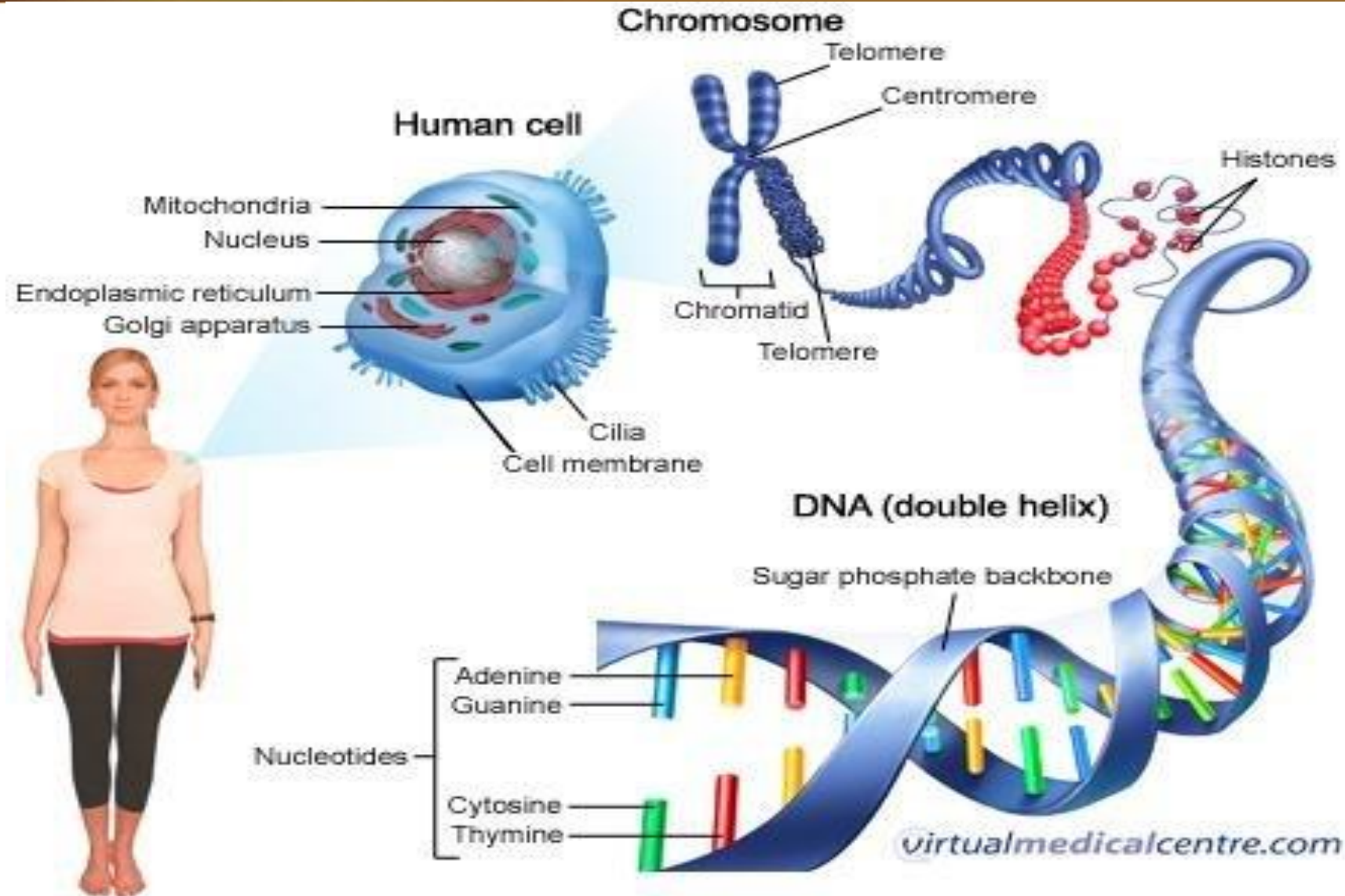
TCACACTGAGCGTGCTG

Miracle in compression

48 volumes of

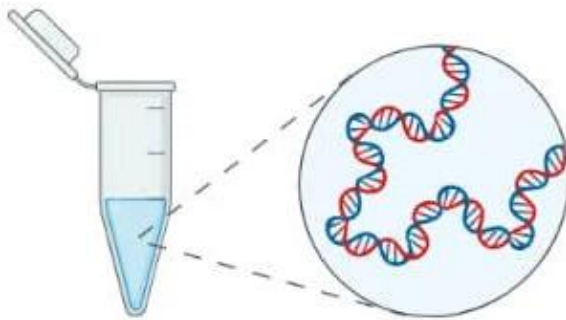
1000 pages

Every 1mm one letterA4

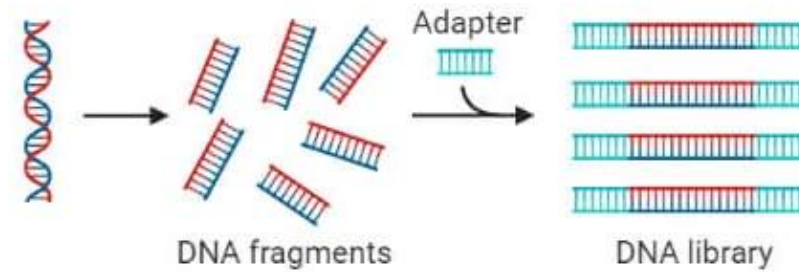


Pipeline

Step 1:
DNA extraction

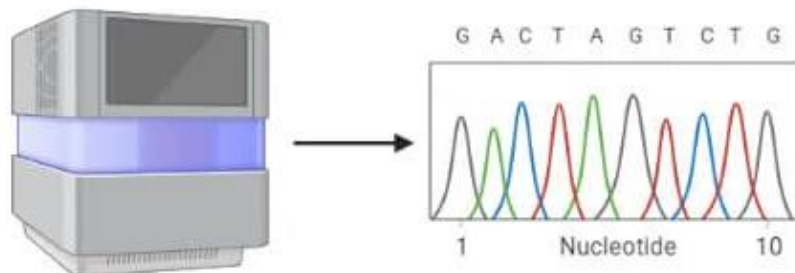


Step 2:
Library preparation

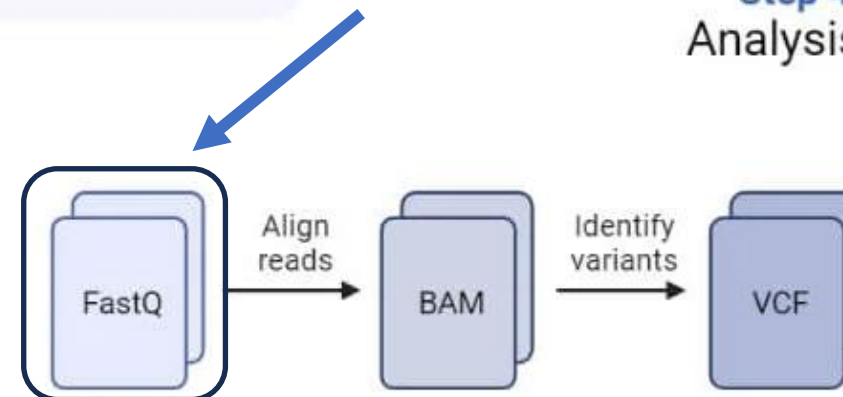


Next Generation Sequencing Workflow

Step 3:
Sequencing



Step 4:
Analysis





SAM File

chr11:5246500-5248500 (reverse strand):

```
ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTACGGCTGTC
ATCACTTAGACCTCACCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGGGAGGGCAGGAGCCAGG
GCTGGGCATAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACAACCTGTGTTCACTAGCAACCTCAAA
CAGACACCATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAGTT
GGTGGTGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAGACAGGTTTAAGGAGACCAATAGAAACTGGGCATGTGGAGA
CAGAGAAGACTCTTGGGTTTCTGATAGGCACTGACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAGGCTGCTGGTG
GTCTACCCTTGGACCCAGAGGTTCTTTGAGTCCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGT
GAAGGCTCATGGCAAGAAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCA
CACTGAGTGAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGGTACAGTTTAGAATGGGAAACAG
ACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTTGCTGTTTCATAACAATTGTTTTCTTTT
GTTTAATTCTTGCTTTCTTTTTTTTTCTTCTCCGCAATTTTTACTATTATACTTAATGCCTTAACATTGTGTATAACAAA
AGGAAATATCTCTGAGATACATTAAGTAACCTTAAAAAAAAAACTTTACACAGTCTGCCTAGTACATTACTATTTGGAATAT
ATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTTATTTCTTTTATTTTTAATTGATACATAATCATTATACATAT
TTATGGGTAAAGTGTAATGTTTTAATATGTGTACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAA
AATGCTTTCTTCTTTTAATATACTTTTTTGTATTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTCAGGGCAATAA
TGATACAATGTATCATGCCTCTTTGCACCATTCTAAAGAATAACAGTGATAATTTCTGGGTAAAGGCAATAGCAATATCT
```



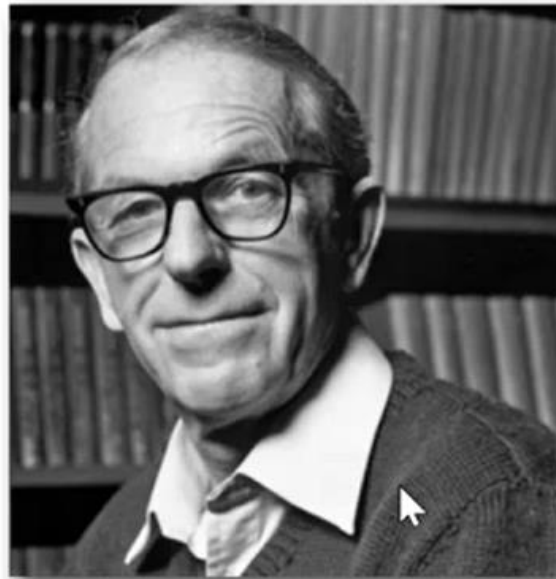

chr11:5246500-5248500 (reverse strand):

ATATCTTAGAGGGAGGGCTGAGGGTTTGAAGTCCAACCTCCTAAGCCAGTGCCAGAAGAGCCAAGGACAGGTACGGCTGTC
ATCACTTAGACCTCACCCTGTGGAGCCACACCCTAGGGTTGGCCAATCTACTCCCAGGAGCAGGGAGGGCAGGAGCCAGG
GCTGGGCATAAAAGTCAGGGCAGAGCCATCTATTGCTTACATTTGCTTCTGACACAACCTGTGTTCACTAGCAACCTCAA
CAGACACC**ATGGTGCATCTGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTGGGGCAAGGTGAACGTGGATGAAGTT**
GGTGGTGAGGCCCTGGGCAGGTTGGTATCAAGGTTACAAGACAGGTTTAAGGAGACCAATAGAACTGGGCATGTGGAGA
CAGAGAAGACTCTTGGGTTTCTGATAGGCACTGACTCTCTCTGCCTATTGGTCTATTTTCCCACCCTTAG**GCTGCTGGTG**
GTCTACCCTTGGACCCAGAGGTTCTTTGAGTCTTTGGGGATCTGTCCACTCCTGATGCTGTTATGGGCAACCCTAAGGT
GAAGGCTCATGGCAAGAAAGTGCTCGGTGCCTTTAGTGATGGCCTGGCTCACCTGGACAACCTCAAGGGCACCTTTGCCA
CACTGAGTGAGCTGCACTGTGACAAGCTGCACGTGGATCCTGAGAACTTCAGGGTGAGTCTATGGGACGCTTGATGTTTT
CTTTCCCCTTCTTTTCTATGGTTAAGTTCATGTCATAGGAAGGGGATAAGTAACAGGGTACAGTTTAGAATGGGAAACAG
ACGAATGATTGCATCAGTGTGGAAGTCTCAGGATCGTTTTAGTTTCTTTTATTGCTGTTTATAACAATTGTTTTCTTTT
GTTTAATTCTTGCTTTCTTAACATTGTGTATAACAAA
AGGAAATATCTCTGAGATACATTACTATTTGGAATAT
ATGTGTGCTTATTTGCATATTCATAATCTCCCTACTTTATTTCTTTTATTTTAAATTGATACATAATCATTATACATAT
TTATGGGTTAAAGTGTAATGTTTTAATATGTGTACACATATTGACCAAATCAGGGTAATTTTGCATTTGTAATTTTAAAA
AATGCTTTCTTCTTTTAATACTTTTTTTGTTTATCTTATTTCTAATACTTTCCCTAATCTCTTTCTTTCAGGGCAATAA
TGATACAATGTATCATGCCTCTTTGCACCATTTCTAAAGAATAACAGTGATAATTTCTGGGTAAAGGCAATAGCAATATCT
TTATGGGTTAAAGTTTCTGCATATAAATTGTAACCTGATGTAAGAGGTTTCATATTGCTAATAGCAGCTACAATCCAGCTA
AATGCTTTCTTCTTTTATGGTTGGGATAAGGCTGGATTATTCTGAGTCCAAGCTAGGCCCTTTTGCTAATCATGTTCA
TTATGGGTTAAAGTCTCCACAG**CTCCTGGGCAACGTGCTGGTCTGTGTGCTGGCCCATCACTTTGGCAAAGAATTCACC**
AATGCTTTCTTCTT**TGCCTATCAGAAAGTGGTGGCTGGTGTGGCTAATGCCCTGGCCACAAGTATCACTAAGCTCGCTT**

Homo sapiens hemoglobin, beta (HBB)



First generation DNA sequencing



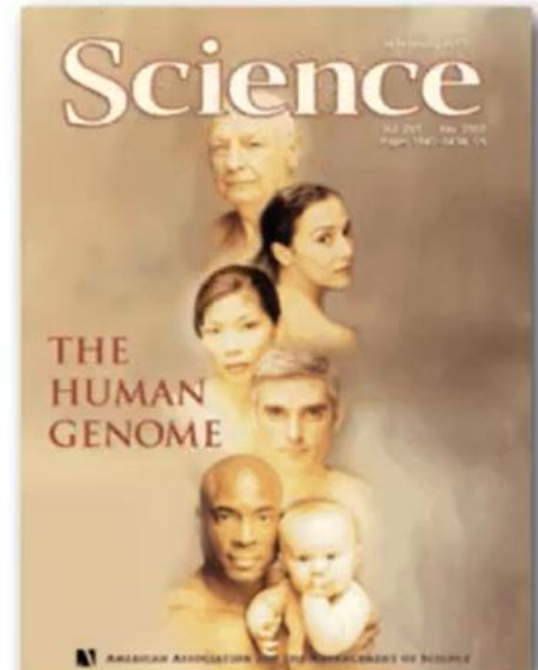
Fred Sanger

“Chain termination” sequencing



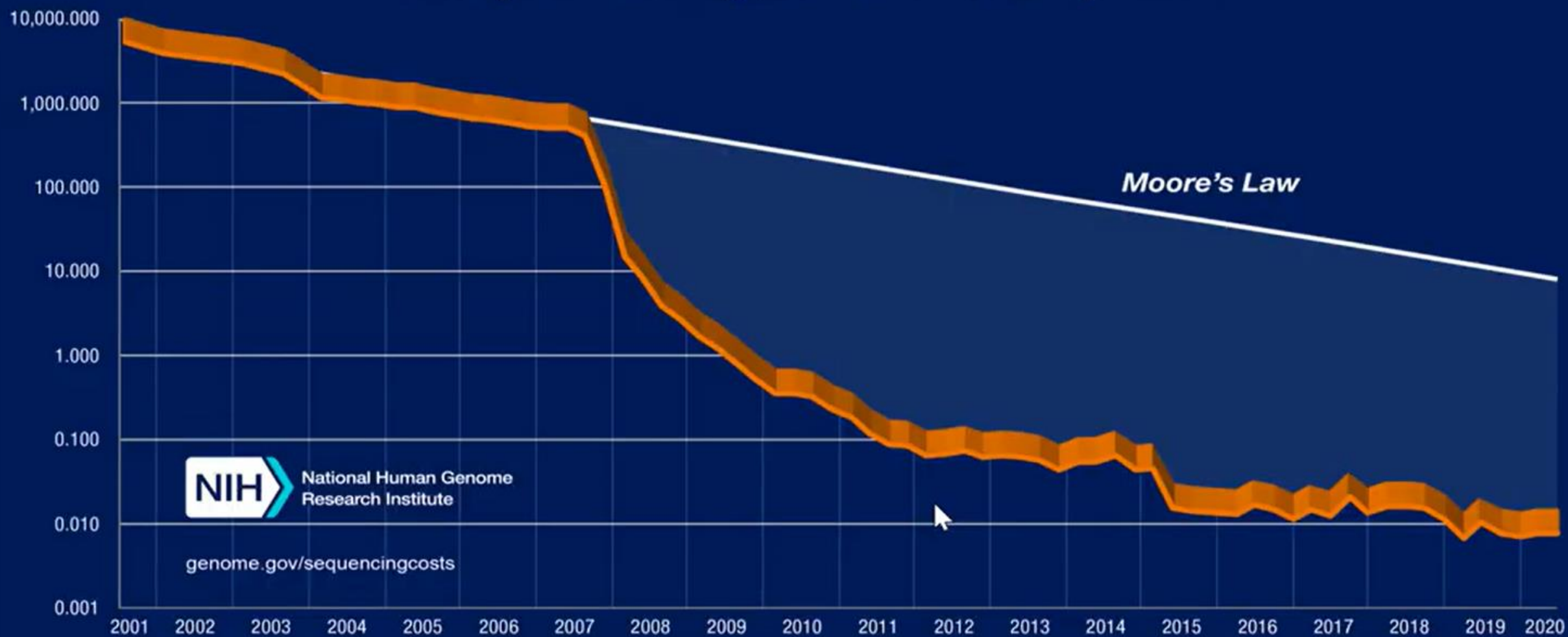


First generation DNA sequencing



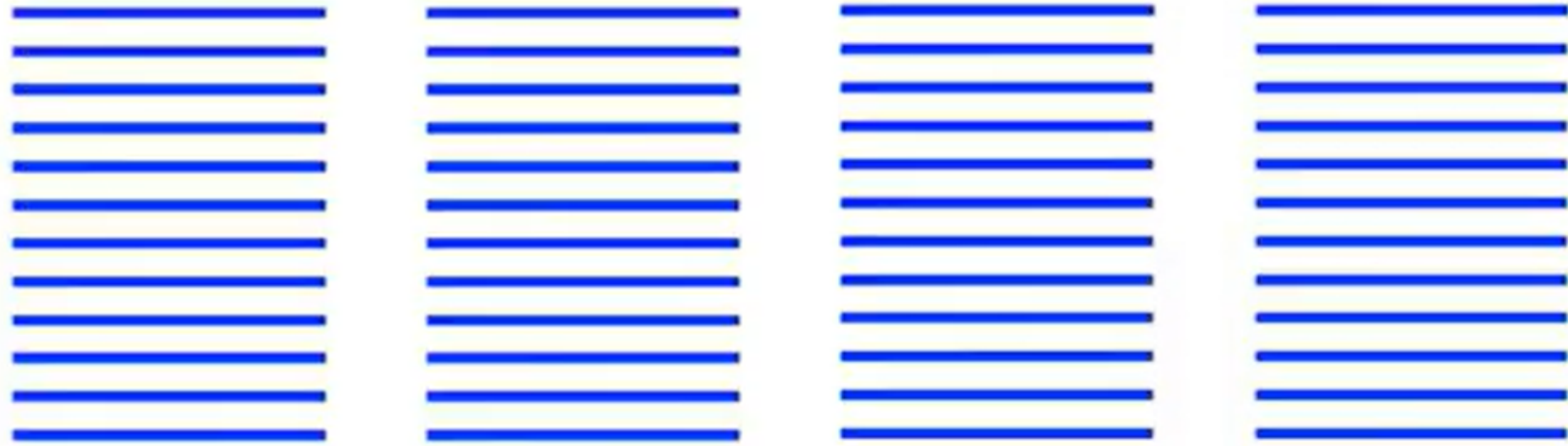


Cost per Raw Megabase of DNA Sequence





Reads



← 100 nt →

Your genome





Input DNA

CCATAGTATATCTCGGCTCTAGGCCCTCATTTTTT
CCATAGTATATCTCGGCTCTAGGCCCTCATTTTTT
CCATAGTATATCTCGGCTCTAGGCCCTCATTTTTT
CCATAGTATATCTCGGCTCTAGGCCCTCATTTTTT



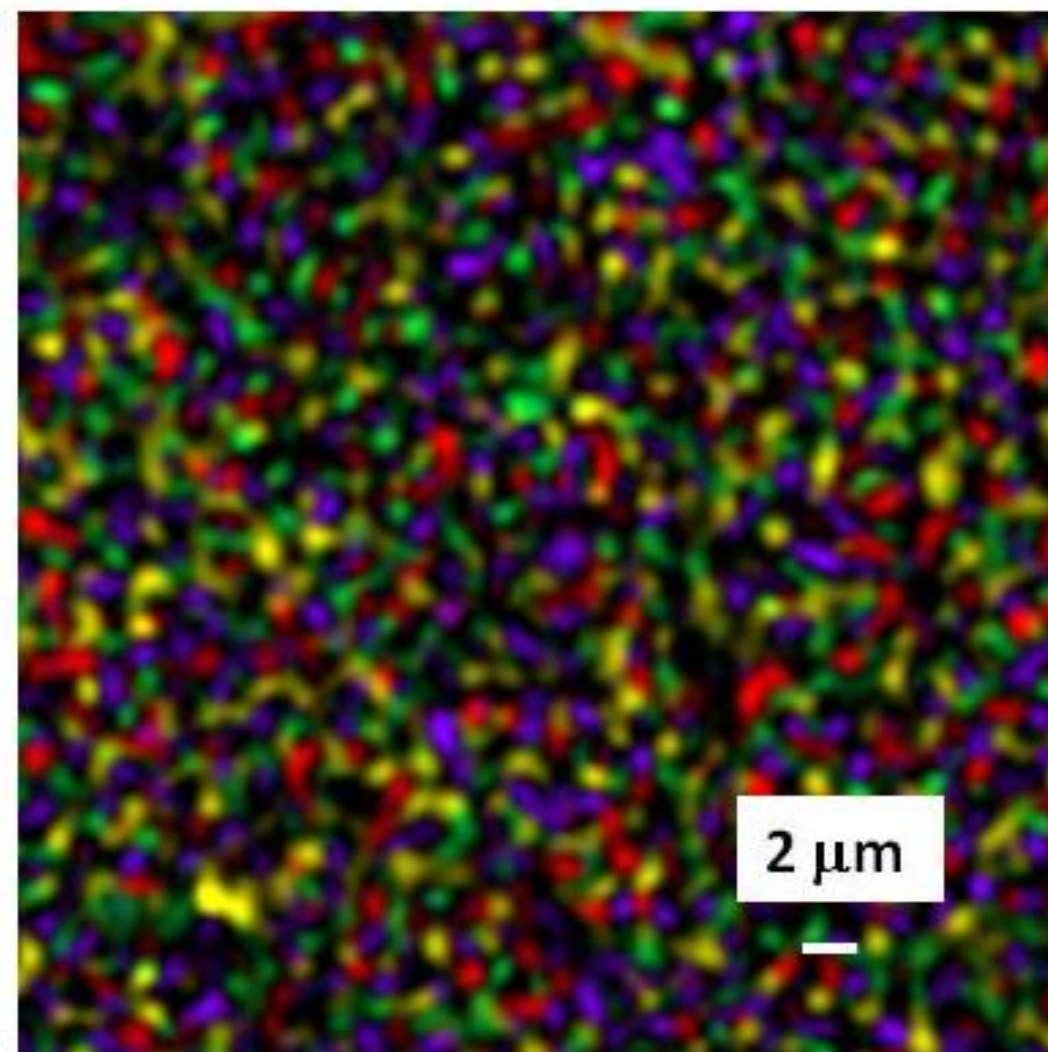
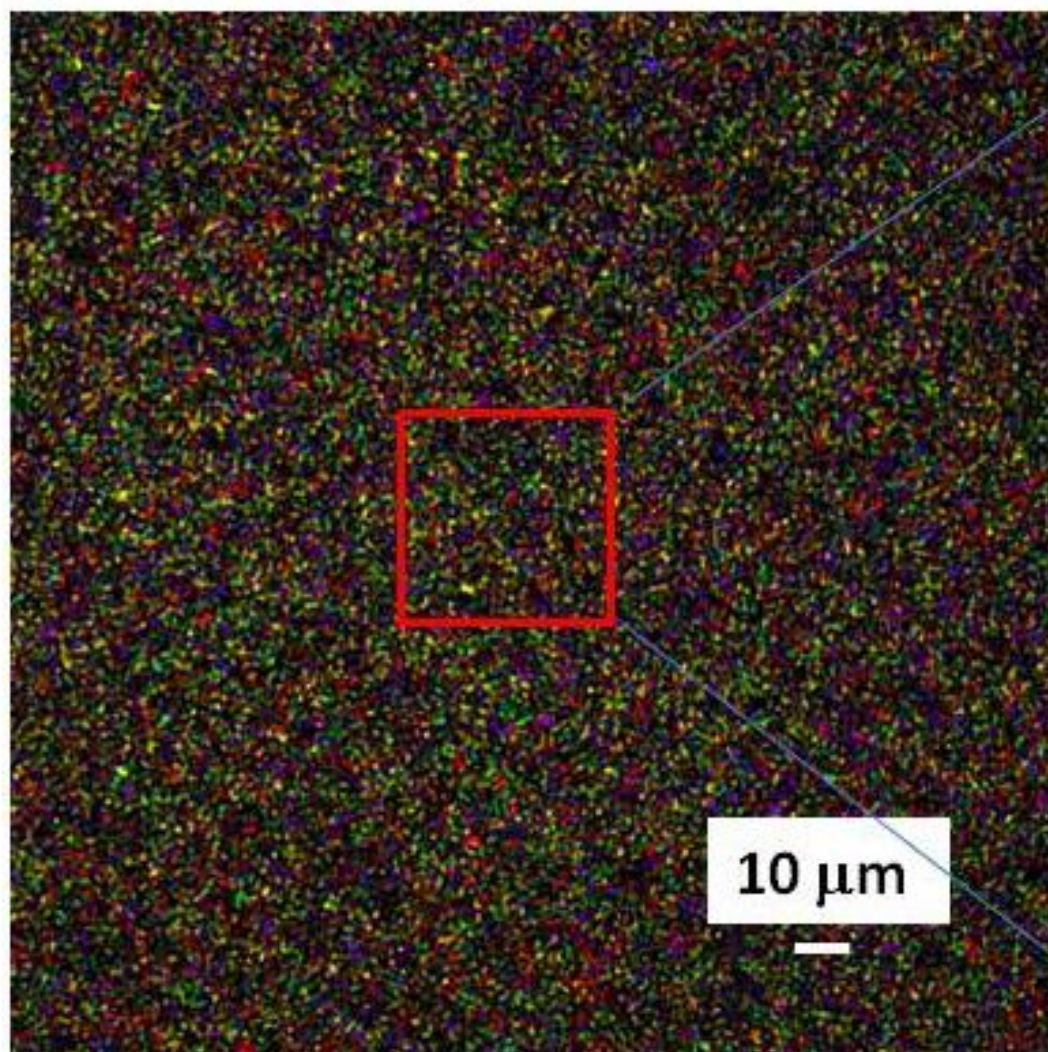
Cut into snippets

CCATAGTA TATCTCGG CTCTAGGCCCTC ATTTTTT
CCA TAGTATAT CTCGGCTCTAGGCCCTCA TTTTTT
CCATAGTAT ATCTCGGCTCTAG GCCCTCA TTTTTT
CCATAG TATATCT CGGCTCTAGGCCCT CATTTTTT

Deposit on slide

CCATAG







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

Base quality

Probability that
base call is incorrect

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

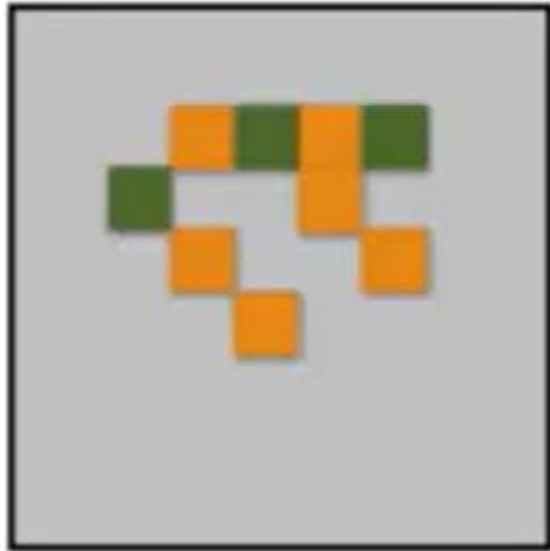
$Q = 20 \rightarrow 1$ in 100

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



The score values

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%



Call: orange (C)

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

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

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



A read in FASTQ format

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



Space	32	0	48	@	64	P	80
!	33	1	49	A	65	Q	81
"	34	2	50	B	66	R	82
#	35	3	51	C	67	S	83
\$	36	4	52	D	68	T	84
%	37	5	53	E	69	U	85
&	38	6	54	F	70	V	86
'	39	7	55	G	71	W	87
(40	8	56	H	72	X	88
)	41	9	57	I	73	Y	89
*	42	:	58	J	74	Z	90
+	43	;	59	K	75	[91
,	44	>	60	L	76	\	92
-	45	=	61	M	77]	93
.	46	<	62	N	78	^	94
/	47	?	63	O	79	_	95

Data gathering

An official website of the United States government. [Here's how you know](#)

NIH National Library of Medicine
National Center for Biotechnology Information

SRA [Advanced](#)

SRA - Now available on the

Sequence Read Archive (SRA) data, available through available repository of high throughput sequencing data, metagenomic and environmental surveys. SRA stores and facilitates new discoveries through data analysis.

Getting Started

- [Documentation](#)
- [How to submit](#)
- [How to search and download](#)
- [How to use SRA in the cloud](#)
- [Submit to SRA](#)

Tools and Software

- [Download SRA Toolkit](#)
- [SRA Toolkit Documentation](#)
- [SRA-BLAST](#)
- [SRA Run Browser](#)
- [SRA Run Selector](#)

FOLLOW NCBI

[X](#) [f](#) [in](#)

(14th February - mid-March) Suspension of the services due to the NIG Supercomputer replacement

DDBJ Bioinformatics and DDBJ Center provides sharing and analysis services for data from life science researches and advances science.

Search
Retrieve the data from the database

Submission
Navigation for how to submit your data

Services
Services available in DDBJ Center

Super Computer
NIG Supercomputer

Statistics
Statistics of DDBJ Center services

Activities
Training sessions and achievements of DDBJ Center

About us
About Bioinformatics and DDBJ Center

This website uses cookies. By continuing to browse this site, you are agreeing to the use of our site cookies. To find out more, see our [Terms of Use](#). [OK](#)

EMBL-EBI **ENA** European Nucleotide Archive

[Services](#) [Research](#) [Training](#) [About us](#)

[Search](#)
[Advanced Sequence](#)

[Home](#) [Search & Browse](#) [Submit & Update](#) [About ENA](#) [Support](#)

Genomes at EBI

- Complete genomes
- Archaea
- Archaeal virus
- Bacteria
- Eukaryota
- Organelle
- Phage
- Plasmid
- Viroid
- Virus
- Links
 - WGS info
 - EnsemblGenomes
 - Ensembl
 - Fasta33 Server

Genomes Pages - At the EBI

Please note that this page is no longer the primary point of access for assembly data. Up-to-date assembly data are searchable from <https://www.ebi.ac.uk/ena/data/warehouse/search?portal=assembly>. The last update to this page was in May 2015; the page has been retained as it includes organellar genomes that are not associated with nuclear genome assemblies and are not represented in the above search.

Access to Completed Genomes

The first completed genomes from viruses, phages and organelles were deposited into the EMBL Database in the early 1980's. Since then, molecular biology's shift to obtain the complete sequences of as many genomes as possible combined with major developments in sequencing technology resulted in hundreds of complete genome sequences being added to the database, including [Archaea](#), [Bacteria](#) and [Eukaryota](#). These web pages give access to a large number of complete genomes, [help](#) is available to describe the layout.

Whole Genome Shotgun Sequences (WGS)

Methods using whole genome shotgun data are used to gain a large amount of genome coverage for an organism. WGS data for a growing number of organisms are being submitted to DDBJ/EMBL/GenBank. [More information about WGS projects...](#)

Last 40 Genome Entries

Date	Accession	Description
02-MAY-2015	CP011047.1	Cronobacter sakazakii strain ATCC 29544
02-MAY-2015	CP011330.1	Helicobacter pylori J99
02-MAY-2015	CP011331.1	Escherichia coli O104:H4 str. C227-11
02-MAY-2015	CP011341.1	Rhodococcus aetherivorans strain IcdP1
02-MAY-2015	KJ680300.1	Accipiter nisus mitochondrion
02-MAY-2015	KJ680301.1	Branta bernicla mitochondrion
02-MAY-2015	KJ680302.1	Pitta nympha mitochondrion
02-MAY-2015	K1701607.1	Fragaria chrysanthemum chloroplast

SRP : Study
SRX : Experiment
SRS : Sample
SRR: Run

DRP
DRX
DRS
DRR

ERP
ERX
ERS
ERR



SRP (Project/Study)

SRS (Sample)

SRX (Experiment)

SRR (Run)

A study is the overarching investigation, hypothesis and its associated tests

A sample refers to a biological sample (cell, mouse, human) on which an experiment is conducted

An experiment is a biological test/perturbation, conducted on one sample. eg: gene knockout, overexpression, or control

A run refers to a sequencing run, associated with one biological sample and experiment, replicated any number of times

study

sample

experiment

run

run

run

sample

experiment

run

run



SRR : SRR6468671.SRA (Too Dense) 326 M

SRAtoolkit : Fastq-dump

```
fasterq-dump --split-3 SRR11192680.SRA
```

FastQ : SRR6468671_1.fastq 738 M
SRR6468671_2.fastq 738 M
SRR6468671.fastq 1.5 K (Not pared)

1.6G	10:16	4	جولی	SRR6468540.sra
738M	18:00	29	جولی	SRR6468671_1.fastq
734M	18:00	29	جولی	SRR6468671_2.fastq
1.5K	18:00	29	جولی	SRR6468671.fastq
326M	15:01	26	مئی	SRR6468671.sra



Fastq File head !

```
SRR6468540.sra SRR6468671_1.fastq SRR6468671_2.fastq SRR6468671.fastq SRR6468671.sra
(base) my-pc :~/SRR$ head SRR6468671_1.fastq
@SRR6468671.1 1 length=101
TAGATCAATTCATTTATTTGGACTATGTTGTAATTTTATTCTTTGCAATGTTTGGAGATTCACTTTAGGAGATCGGAAGAGCACACGTCTGAACTCCA
+SRR6468671.1 1 length=101
CCCDFFDFHHDGDIIBGGGIGG?FEGIIFH<FHEHGGIIGGGEGHI@GGIIGHGGBDHG8DGG<FHEGIIIHGEEHIGEHFFFFCC@A???CD>AACDCC
@SRR6468671.2 2 length=101
GTACCGCCAATAAGCGTTTGAGAGATGGAGTGACACAGTAGGATAAGCTAACCGTACTGTTGGTTATGTACGGAGATCGGAAGAGCACACGTCTGAACTCC
+SRR6468671.2 2 length=101
@?@FFFDDFHHDHJIJFFGGHEGHHBHGHH0BDHIEIJBFEDEHIIJ>FGGGIH;=CAAEHF?;@CC>>CDDDD@>CB=8@BABCCBDDDDDDDB>ACCCD
@SRR6468671.3 3 length=101
CCGTCCCTTGGGTGCCGCTTTTTTGTTCATCAGATAAACAGGGTGGTACCGCGATGAGTCCCGTCGTCCTTGCAATTAGATCGGAAGAGCACACGT
(base) my-pc :~/SRR$ wc SRR6468671_1.fastq
 11191060  22382120 773622054 SRR6468671_1.fastq
(base) my-pc :~/SRR$ wc -l SRR6468671_1.fastq
11191060 SRR6468671_1.fastq
```

~ 12000000 line

~ 3000000 Read



FASTQC



Aims to provide a simple way to do some quality control checks on raw sequence data coming from high throughput sequencing pipelines.

<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

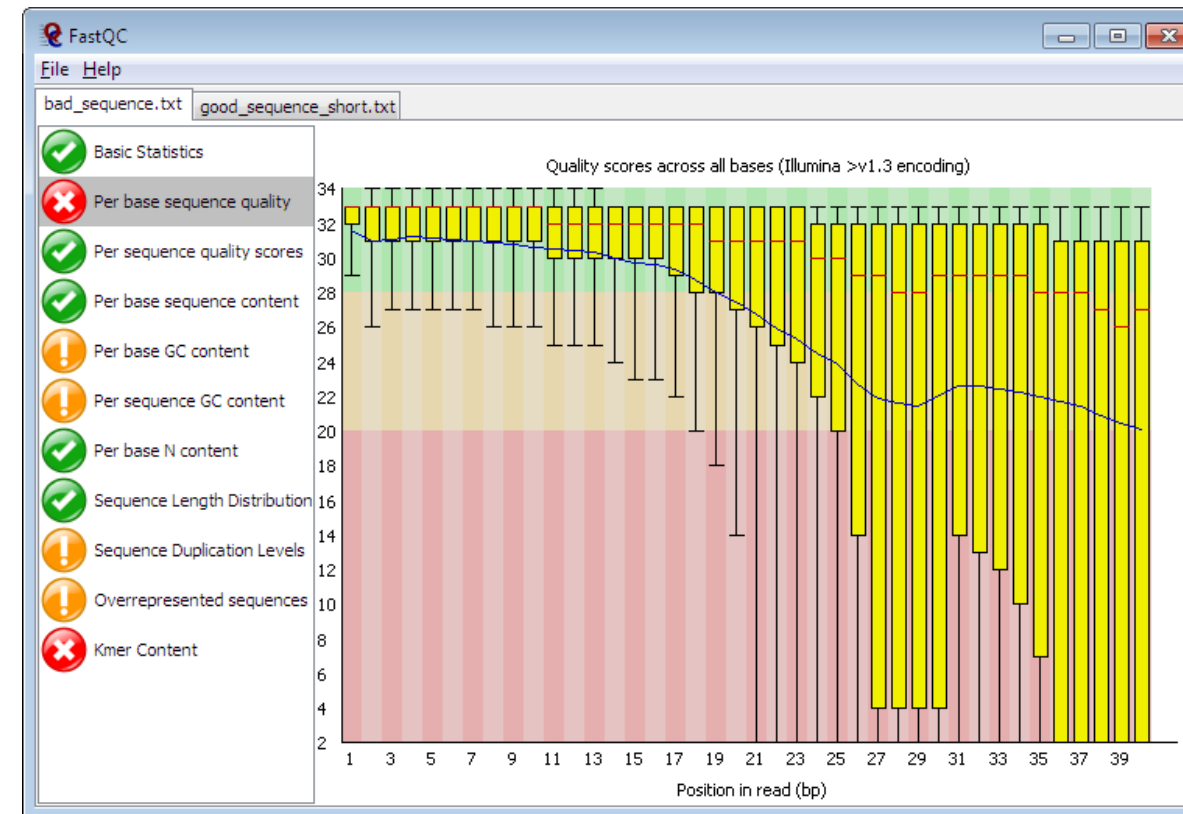
```
$ fastqc -o ~/results/ *.fq
```

Command

Output

Path

Filename





- You may have FastQC or Need to Download .
 - Then :

```
$ sudo apt-get update      ( Prepare to download )
```

```
$ sudo apt-get install fastqc
```

```
$ fastqc -f fastq SRR6468671_1.fastq SRR6468671_2.fastq
```



LET'S
CODE



FastQC Report

Summary

- ✓ Basic Statistics
- ✓ Per base sequence quality
- ✓ Per sequence quality scores
- ! Per base sequence content
- ! Per sequence GC content
- ✓ Per base N content
- ! Sequence Length Distribution
- ✓ Sequence Duplication Levels
- ✗ Overrepresented sequences
- ✗ Adapter Content

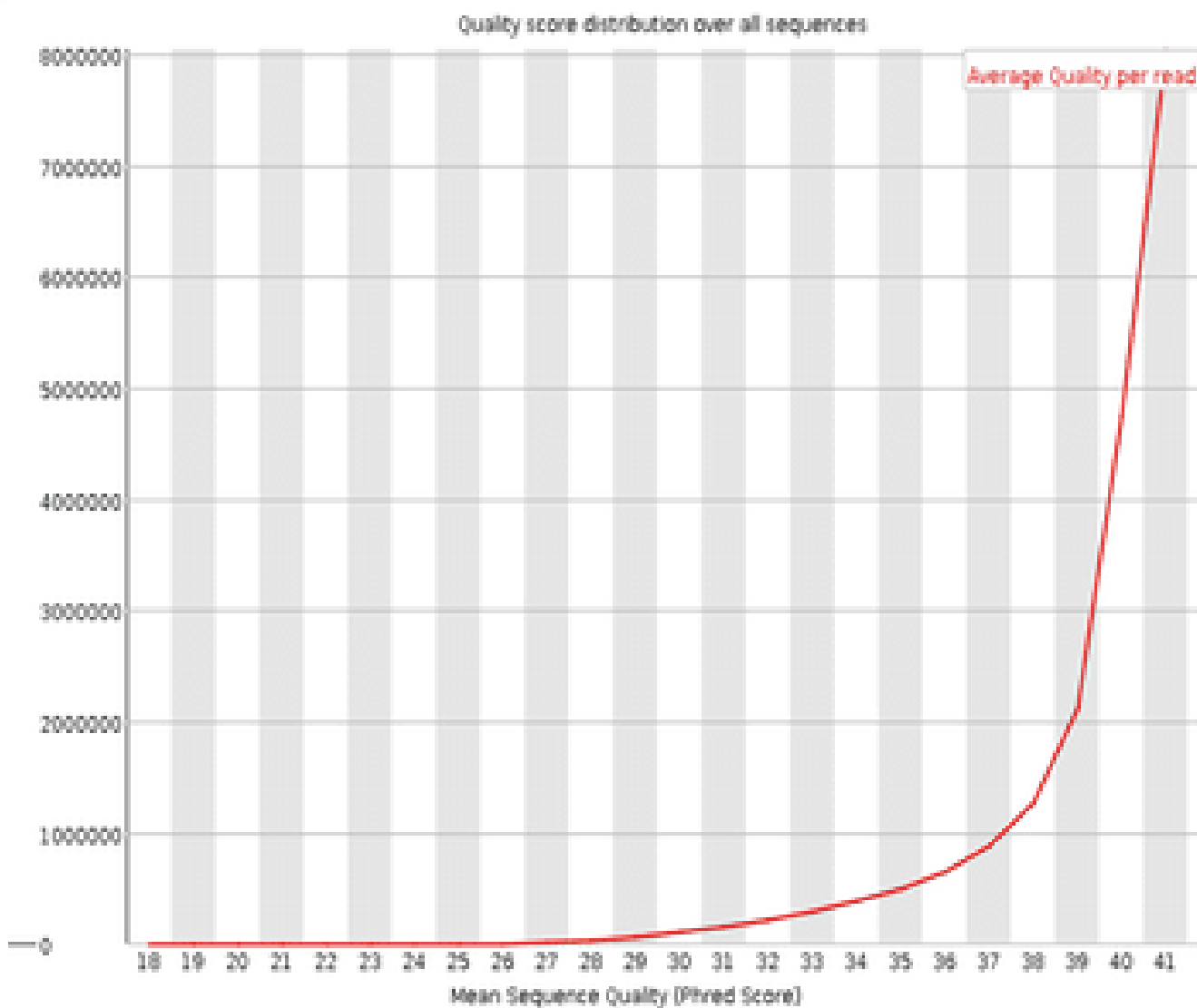
✓ Basic Statistics

Measure	Value
Filename	SRR6468671_1.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	2797765
Sequences flagged as poor quality	0
Sequence length	78-101
%GC	43

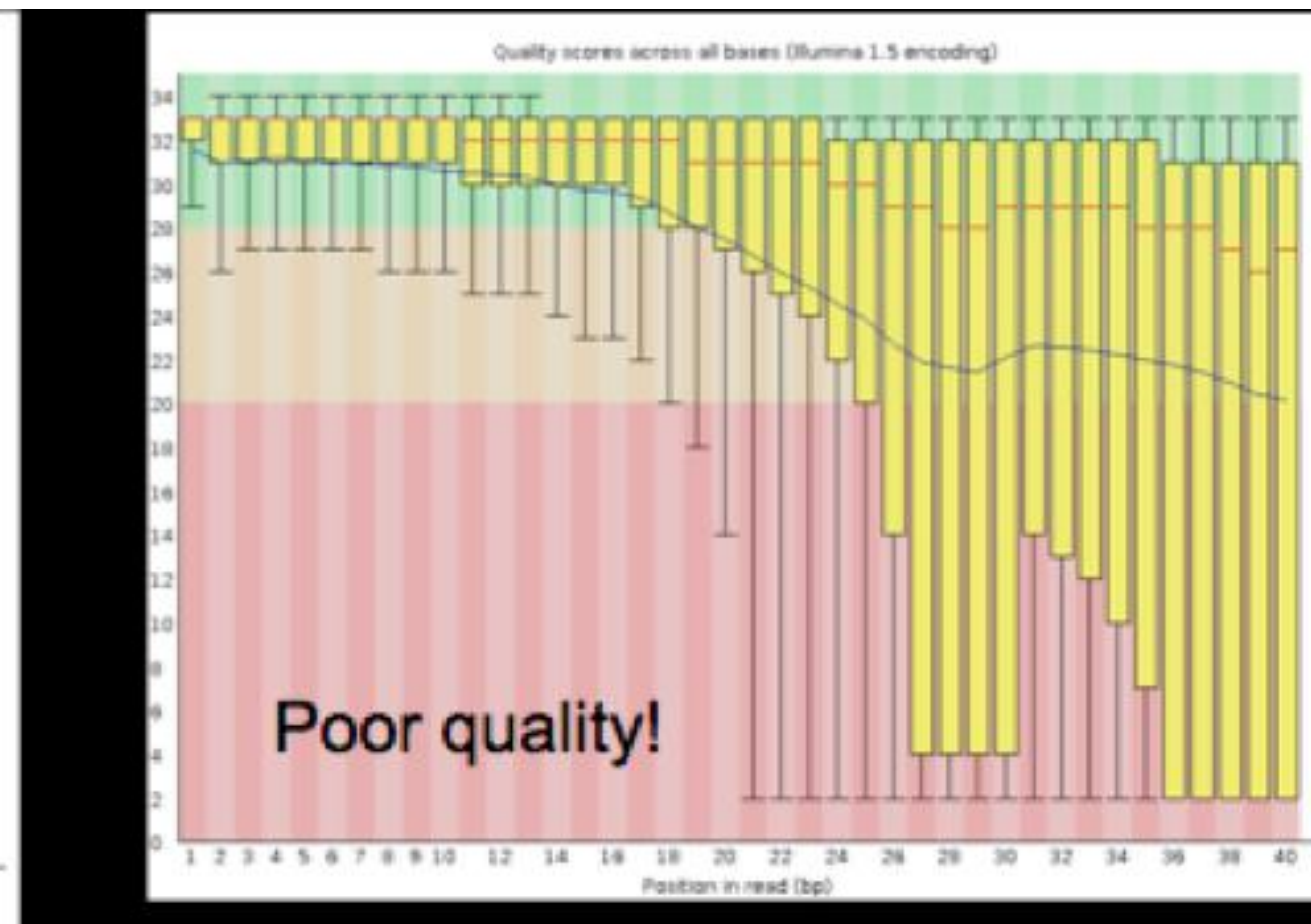
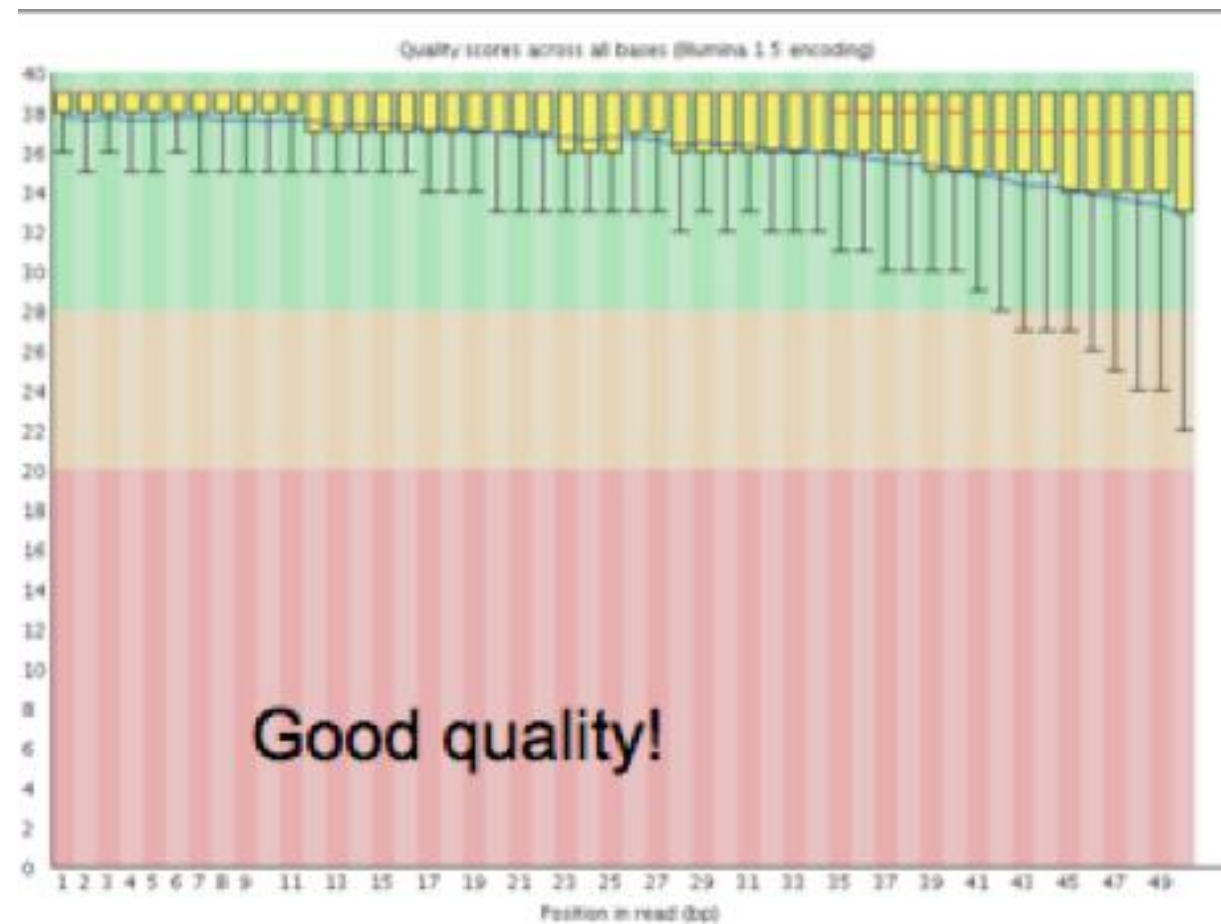
✓ Per base sequence quality



X : Position
Y : Quality



More than 8000000 base pairs Quality are 40





Trimming

Trimmomatic: A flexible read trimming tool for Illumina NGS data .

Download :

<http://www.usadellab.org/cms/index.php?page=trimmomatic>

Install :

```
$ sudo apt-get install trimmomatic
```

```
$ java -jar trimmomatic-0.35.jar SE/PE 2 input files and 4 outputs
```

```
LEADING:25 ( Trim from left side , quality < 25 )
```

```
TRAILING: 25 ( Trim from right side , quality < 25 )
```

```
SLIDINGWINDOW:4: 25 ( Slidin window size 4 pair, trim next by meanquality < 25 )
```

```
MINLEN:36 (Remove if Remaining read < 36 )
```

Required :

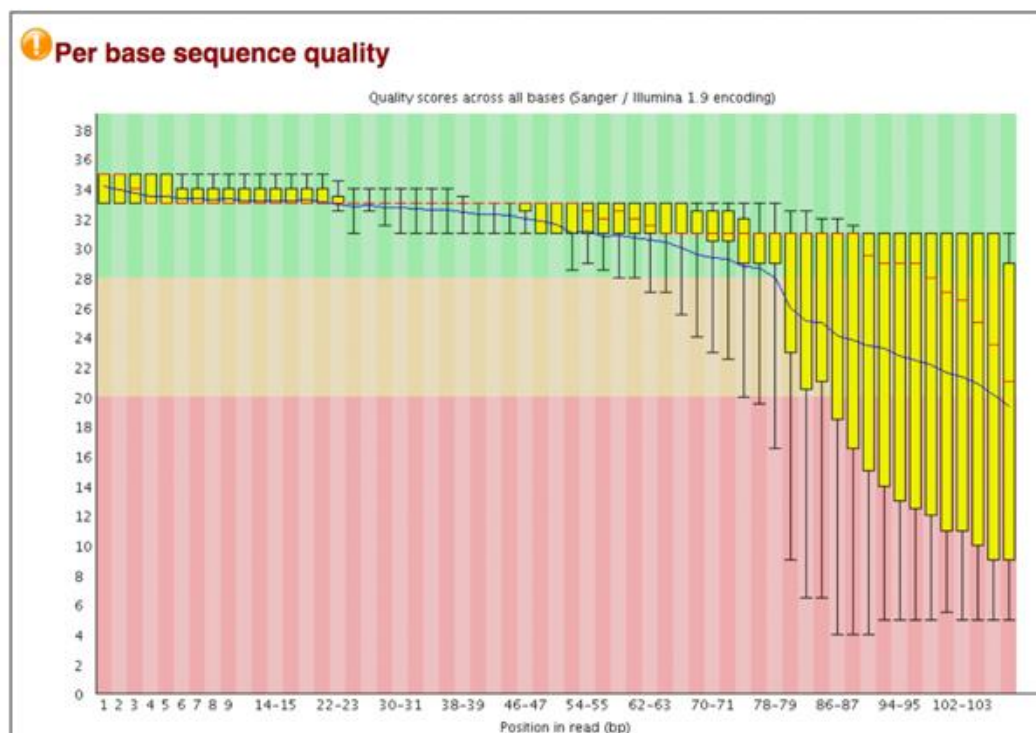
```
jre-8v261-linux-x64.tar.gz -----> tar xvzf jre-8v261-linux-x64.tar.gz
```



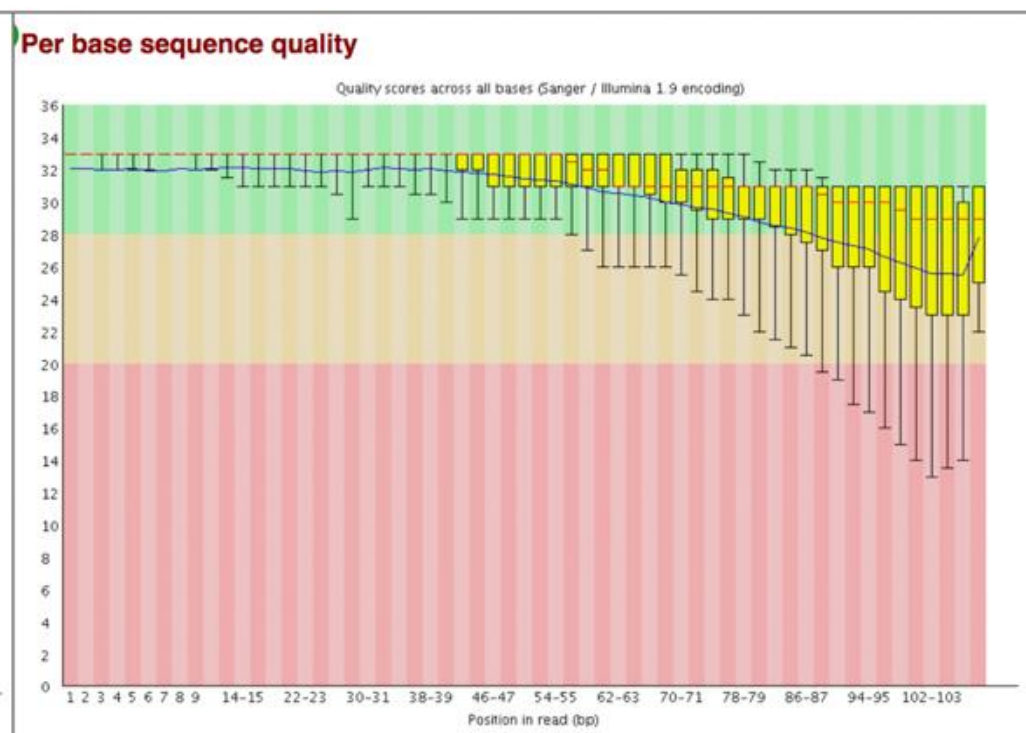

Input : SRR6468671_1.fastq SRR6468671_2.fastq

Output : SRR6468671_1_trimmed.fastq SRR6468671_1_untrimmed.fastq
SRR6468671_2_trimmed.fastq SRR6468671_2_untrimmed.fastq

Before trimming

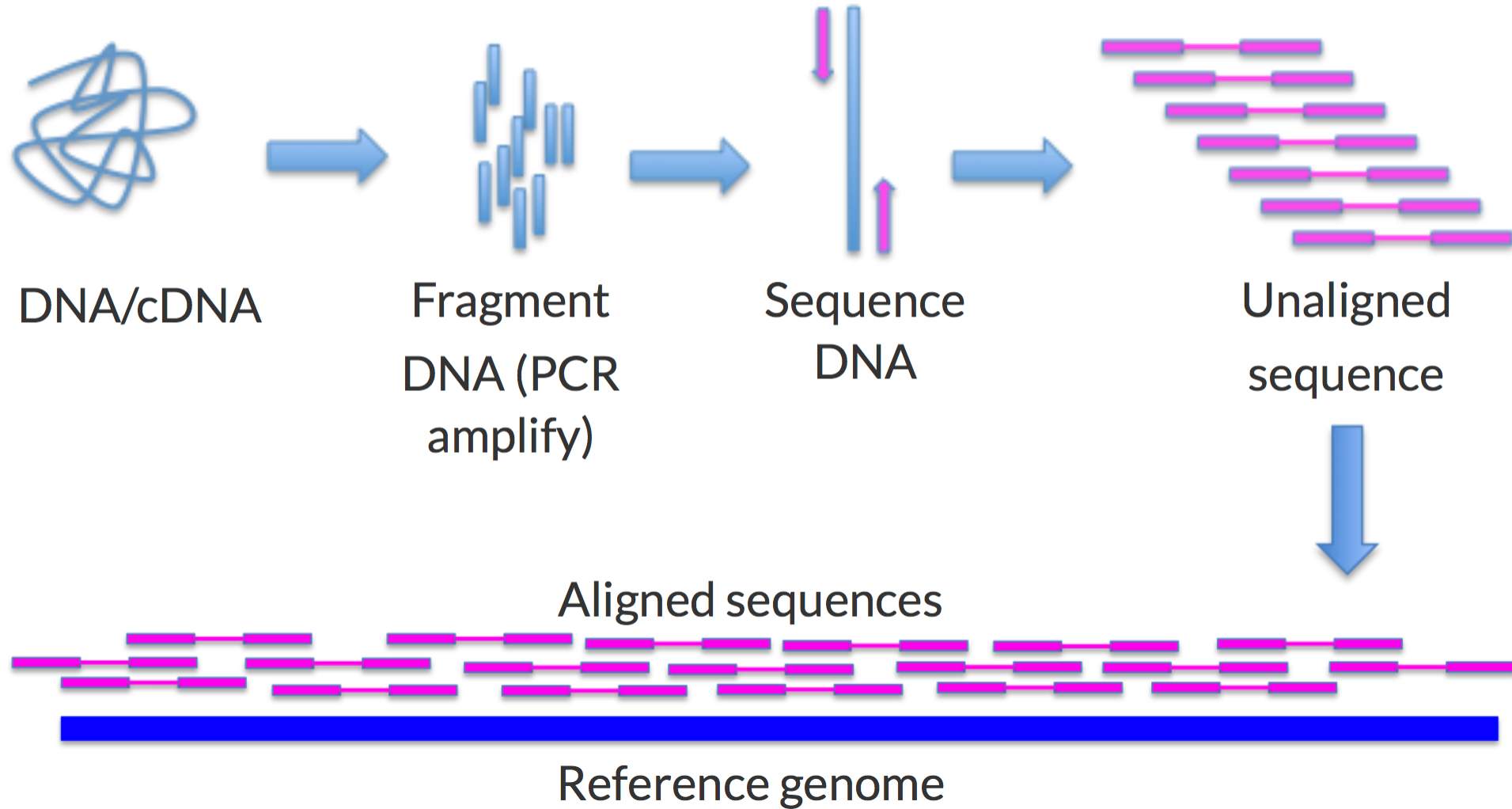


After trimming





Aligning





Hisat2

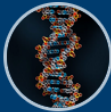
HISAT2 is a fast and sensitive alignment program

for mapping next-generation sequencing reads

(both DNA and RNA) to a population of human

genomes as well as to a single reference genome.

<http://daehwankimlab.github.io/hisat2/>



HISAT2

graph-based alignment of next generation sequencing reads to a population of genomes

HISAT2 is a fast and sensitive alignment program for mapping next-generation sequencing reads (both DNA and RNA) to a population of human genomes as well as to a single reference genome. Based on an extension of BWT for graphs (Sirén et al. 2014), we designed and implemented a graph FM index (GFM), an original approach and its first implementation. In addition to using one global GFM index that represents a population of human genomes, **HISAT2** uses a large set of small GFM indexes that collectively cover the whole genome. These small indexes (called local indexes), combined with several alignment strategies, enable rapid and accurate alignment of sequencing reads. This new indexing scheme is called a Hierarchical Graph FM index (HGFM).

The HISAT-3N paper published at *Genome Research*. 7/1/2021

HISAT-3N beta release 12/14/2020

HISAT-3N is a software system for analyzing nucleotide conversion sequencing reads. See the [HISAT-3N](#) for more details.

Index files are moved to the AWS Public Dataset Program. 9/3/2020

We have moved HISAT2 index files to the AWS Public Dataset Program. See the [link](#) for more details.

HISAT 2.2.1 release 7/24/2020

This patch version includes the following changes.

- Python3 support
- Remove the HISAT-genotype related scripts. HISAT-genotype moved to <http://daehwankimlab.github.io/hisat-genotype/>
- Fixed bugs related to `--read-lengths` option

HISAT 2.2.0 release 2/6/2020

[Main](#)

[About](#)

[Manual](#)

[HISAT-3N](#)

[Download](#)

[HowTo](#)

[Links](#)

Funding

This work was supported in part by the National Human Genome Research Institute under grants R01-HG006102 and R01-HG006677, and NIH grants R01-LM06845 and R01-GM083873 and NSF grant CCF-0347992 to Steven L. Salzberg and by the Cancer Prevention Research Institute of Texas under grant RR170068 and NIH grant R01-GM135341 to Daehwan Kim

Getting Help

Please use hisat2.genomics@gmail.com for private communications only. Please do not email technical questions to HISAT2 contributors directly.



Binaries

Version: HISAT2 2.2.1

Release Date: 7/24/2020

Source	https://cloud.biohpc.swmed.edu/index.php/s/fE9QCsX3NH4QwBi/download
OSX_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/zMgEtnF6LjnJFrr/download
Linux_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/oTtGWbWjaxesQ2Ho/download

Version: HISAT2 2.2.0

Release Date: 2/6/2020

Source	https://cloud.biohpc.swmed.edu/index.php/s/hisat2-220-source/download
OSX_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/hisat2-220-OSX_x86_64/download
Linux_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/hisat2-220-Linux_x86_64/download

Version: HISAT2 2.1.0

Release Date: 6/8/2017

Source	https://cloud.biohpc.swmed.edu/index.php/s/hisat2-210-source/download
OSX_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/hisat2-210-OSX_x86_64/download
Linux_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/hisat2-210-Linux_x86_64/download
Windows	http://www.di.fc.ul.pt/~afalcao/hisat2_windows.html

Index

HISAT2 indexes are hosted on AWS (Amazon Web Services), thanks to the AWS Public Datasets program. Click this [link](#) for more details.

H. sapiens

- GRCh38

genome	https://genome-idx.s3.amazonaws.com/hisat/grch38_genome.tar.gz
genome_snp	https://genome-idx.s3.amazonaws.com/hisat/grch38_snp.tar.gz
genome_tran	https://genome-idx.s3.amazonaws.com/hisat/grch38_tran.tar.gz
genome_snp_tran	https://genome-idx.s3.amazonaws.com/hisat/grch38_snptran.tar.gz
genome_rep(above 2.2.0)	https://genome-idx.s3.amazonaws.com/hisat/grch38_rep.tar.gz
genome_snp_rep(above 2.2.0)	https://genome-idx.s3.amazonaws.com/hisat/grch38_snprep.tar.gz

- UCSC hg38

genome	https://genome-idx.s3.amazonaws.com/hisat/hg38_genome.tar.gz
genome_tran	https://genome-idx.s3.amazonaws.com/hisat/hg38_tran.tar.gz

- GRCh37

genome	https://genome-idx.s3.amazonaws.com/hisat/grch37_genome.tar.gz
genome_snp	https://genome-idx.s3.amazonaws.com/hisat/grch37_snp.tar.gz
genome_tran	https://genome-idx.s3.amazonaws.com/hisat/grch37_tran.tar.gz
genome_snp_tran	https://genome-idx.s3.amazonaws.com/hisat/grch37_snptran.tar.gz



If the Index file is not available , we must build one , Then Download The reference Genom as a Fasta format and use given command :

```
$ hisat-build reference.fasta
```

Here we have it !

The index files are given for ChrX as 8 files . (chrX_tran_1.ht2 , ..)

```
$ hisat2 -q -x < Index.ht2 folder > { -1 m1 -2 m2 } --add-chrname S name.sam
```

To investigate samfile we must use samview but “head” command could be beneficial !





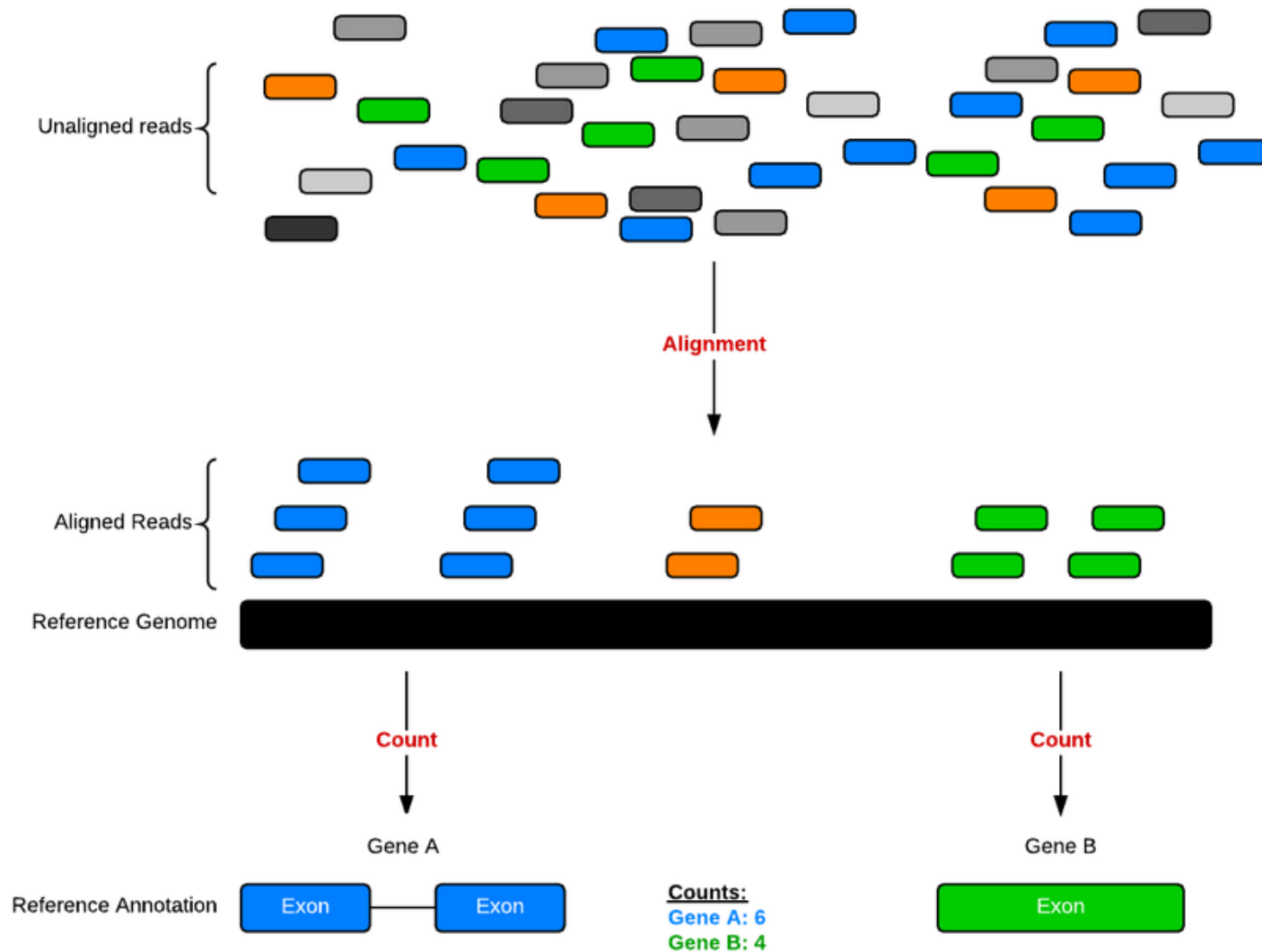
Sam tools

- To identify read count first we must have BAM file .
- Then , converting sam file to bam file by using samtools !
- \$ samtools view -F 4 -b filename.sam output.bam

173M	17:46	15	آقۇست	ERR188044.bam
641M	11:19	10	آقۇست	ERR188044.sam
646M	11:22	10	آقۇست	ERR188104.sam



Gene Counting



Requirements :

- *.sam
- *.bam
- *.gtf




The Gene Transfer Format (GTF) is a widely used format for storing gene annotations.

You can obtain GTF files easily from the
UCSC table browser and Ensembl.

Gtf file

ucsc genome browser



UNIVERSITY OF CALIFORNIA
SANTA CRUZ Genomics
Institute **UCSC** **Genome Browser**

Genomes Genome Browser Tools Mirrors Downloads My Data Projects Help About Us

Try our new clinical tutorial!

Search genes, data, help docs and more... Search

Tools

-  **Genome Browser** - Interactively visualize genomic data
- BLAT** - Rapidly align sequences to the genome
- In-Silico PCR** - Rapidly align PCR primer pairs to the genome
- Table Browser** - Download and filter data from the Genome Browser
- LiftOver** - Convert genome coordinates between assemblies
- REST API** - Returns data requested in JSON format
- Variant Annotation Integrator** - Annotate genomic variants
- More tools...**

News

- Feb. 14, 2025 - **CIVIC track for hg19 and hg38 is now available**
- Jan. 31, 2025 - **New COSMIC track for hg19 and update for hg38 COSMIC track**
- Dec. 20, 2024 - **NCBI Gene Orthologs track available for hg38, mm39, danRer11, ca...**
- Dec. 10, 2024 - **DECIPHER Population CNVs track for Human (hg19/hg38)**
- Nov. 8, 2024 - **New GENCODE gene tracks: Human V47 (hg19/hg38) - Mouse M36**
- Nov. 4, 2024 - **GIAB Problematic Regions tracks for human (hg38 and hs1)**

[More news...](#) [Subscribe](#)

Meetings and Workshops: Come see us in person!

- **CSHL Biology of Genomes** -- Cold Spring Harbor, NY. May 6-10, 2025
- **ESHG: European Human Genetics** -- Milan, Italy. May 24-27, 2025. **Visit us at booth 2020!**

Feel free to [contact us](#) if you are interested in attending a workshop, or meeting someone from the team to collaborate, get help, or ask any questions at the meetings.

To obtain “gtf” file, identifying the situation of each gene on each chromosome .



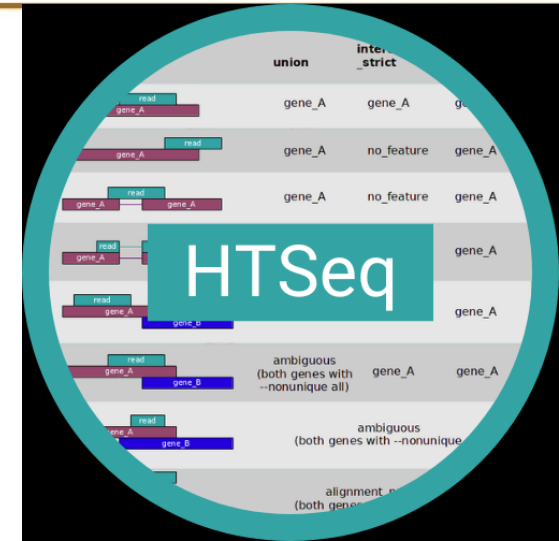
For a specific chromosome .

chrUn_KI270751v1.fa.gz	2014-01-23	16:40	48K
chrUn_KI270752v1.fa.gz	2014-01-23	16:40	9.2K
chrUn_KI270753v1.fa.gz	2014-01-23	16:40	20K
chrUn_KI270754v1.fa.gz	2014-01-23	16:40	8.7K
chrUn_KI270755v1.fa.gz	2014-01-23	16:40	11K
chrUn_KI270756v1.fa.gz	2014-01-23	16:40	16K
chrUn_KI270757v1.fa.gz	2014-01-23	16:40	14K
chrX.fa.gz	2014-01-23	16:40	47M
chrX_KI270880v1_alt.fa.gz	2014-01-23	16:40	80K
chrX_KI270881v1_alt.fa.gz	2014-01-23	16:40	46K
chrX_KI270913v1_alt.fa.gz	2014-01-23	16:40	77K
chrY.fa.gz	2014-01-23	16:40	8.0M
chrY_KI270740v1_random.fa.gz	2014-01-23	16:40	5.4K
md5sum.txt	2014-08-21	14:18	26K



Htseq-count

A tool developed with HTSeq that preprocesses RNA-Seq data for differential expression analysis by **counting** the **overlap** of **reads** with **genes**.



Install :

```
$ sudo apt-get install build-essential python2.7-dev python-numpy python-matplotlib python-pysam python-htseq
```

Usage :

```
$ htseq-count [options] <alignment_files(*.sam)> <gff_file> > output.count
```



- Count to csv







For each sample , All
The pipeline must be
done separately .

Genes

Samples

GENE ID	KD.2	KD.3	OE.1	OE.2	OE.3	IR.1	IR.2	IR.3
1/2-SBSRNA4	57	41	64	55	38	45	31	39
A1BG	71	40	100	81	41	77	58	40
A1BG-AS1	256	177	220	189	107	213	172	126
A1CF	0	1	1	0	0	0	0	0
A2LD1	146	81	138	125	52	91	80	50
A2M	10	9	2	5	2	9	8	4
A2ML1	3	2	6	5	2	2	1	0
A2MP1	0	0	2	1	3	0	2	1
A4GALT	56	37	107	118	65	49	52	37
A4GNT	0	0	0	0	1	0	0	0
AA06	0	0	0	0	0	0	0	0
AAA1	0	0	1	0	0	0	0	0
AAAS	2288	1363	1753	1727	835	1672	1389	1121
AACS	1586	923	951	967	484	938	771	635
AACSP1	1	1	3	0	1	1	1	3
AADAC	0	0	0	0	0	0	0	0
AADACL2	0	0	0	0	0	0	0	0
AADACL3	0	0	0	0	0	0	0	0
AADACL4	0	0	1	1	0	0	0	0
AADAT	856	539	593	576	359	567	521	416
AAGAB	4648	2550	2648	2356	1481	3265	2790	2118
AAK1	2310	1384	1869	1602	980	1675	1614	1108
AAMP	5198	3081	3179	3137	1721	4061	3304	2623
AANAT	7	7	12	12	4	6	2	7
AARS	5570	3323	4782	4580	2473	3953	3339	2666



Name	
 FastQC	
 hisat2-2.2.1	
 NGS	
 samtools-1.17	
 sratoolkit.3.0.0-ubuntu64	
 Trimmomatic-0.39	



- `fastqc -f fastq ../NGS/samples/ERR188044_chrX_1.fastq ERR188044_chrX_2.fastq`
- `java -jar trimmomatic-0.39.jar PE ../NGS/samples/ERR188044_chrX_1.fastq ../NGS/samples/ERR188044_chrX_2.fastq
ER1trim.fastq ER1untrim.fastq ER2trim.fastq ER2untrim.fastq
LEADING:25 TRAILING:25
SLIDINGWINDOW:4:20
hisat2 -q -x ../bulk/NGS/indexes/chrX_tran -1 ../bulk/NGS/samples/ER1trim.fastq -2 ../bulk/NGS/samples/ER2trim.fastq
--add-chrname -S mapped.sam`
- `htseq-count mapped.sam ../genes/chrX.gtf > ERR18.count`



```
sudo -i
```

```
wget http://archive.ubuntu.com/ubuntu/pool/main/o/openssl/libssl1.1_1.1.1f-  
1ubuntu2_amd64.deb
```

```
sudo dpkg -i libssl1.1_1.1.1f-1ubuntu2_amd64.deb
```

```
sudo -i
```

```
root@s:~# apt --fix-broken install
```

```
sudo apt-get install build-essential python3-numpy python3-matplotlib python3-  
pysam python3-htseq
```



Thanks