



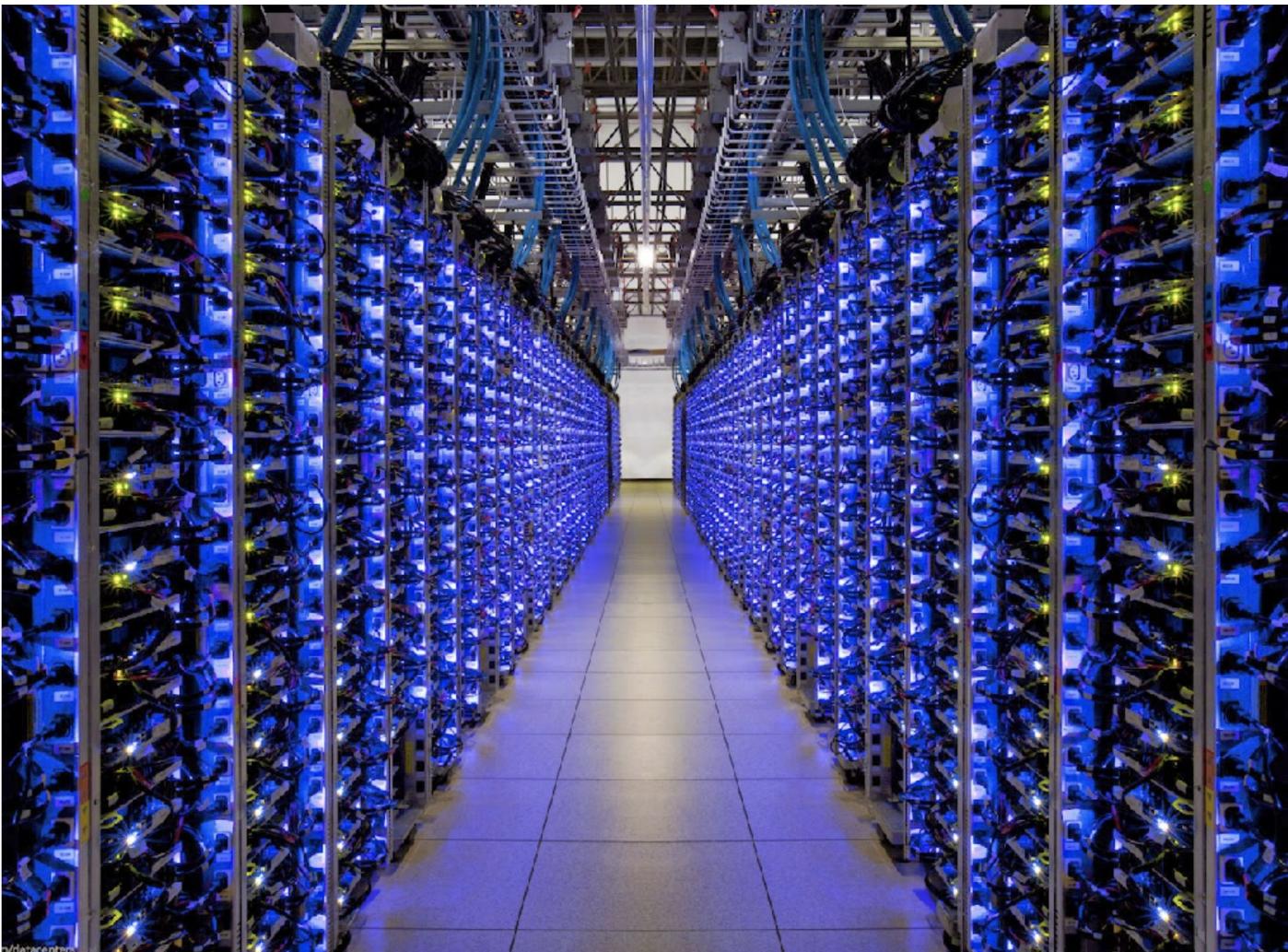
Introduction to DNA sequencing data analysis

Sina Majidian

Postdoctoral Fellow at the University of Lausanne.

SIB Member, Swiss Institute of Bioinformatics.

25 Nov. 2021

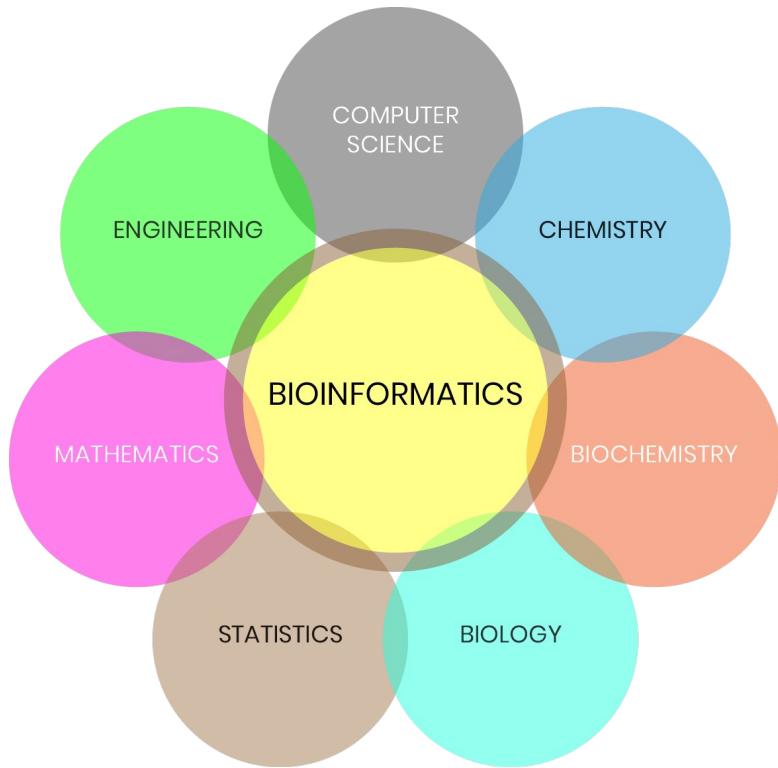




Learning objectives

- How machines sequence DNA?
- How to store data?
- How to download data?
- How to visualize and manipulate data?

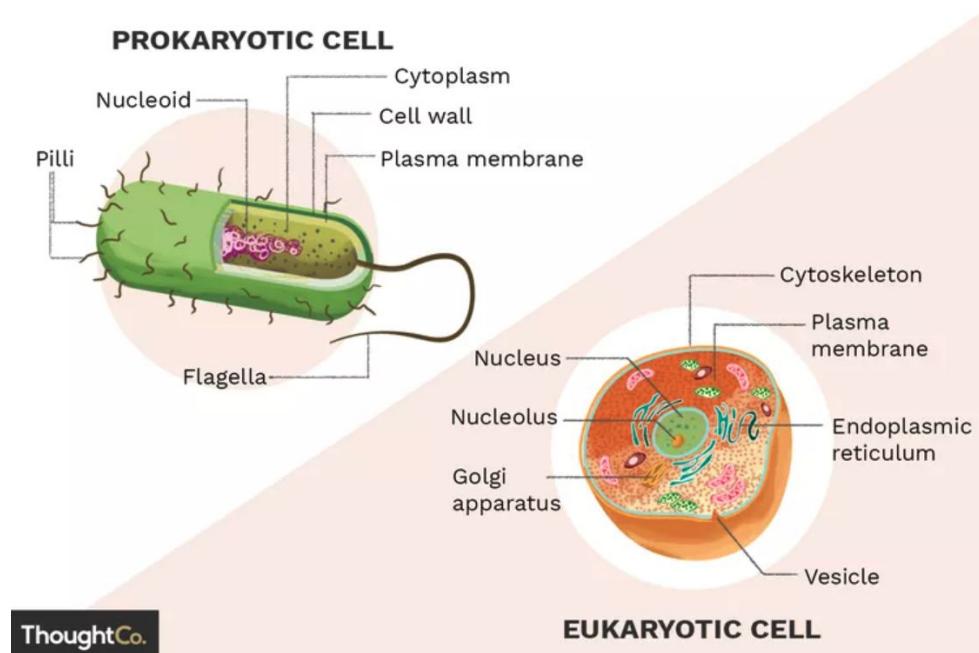
If you are participating with phone, we need computer for hands-on exercises!



Your background?

Cell

The basic unit from which a living organism is made.



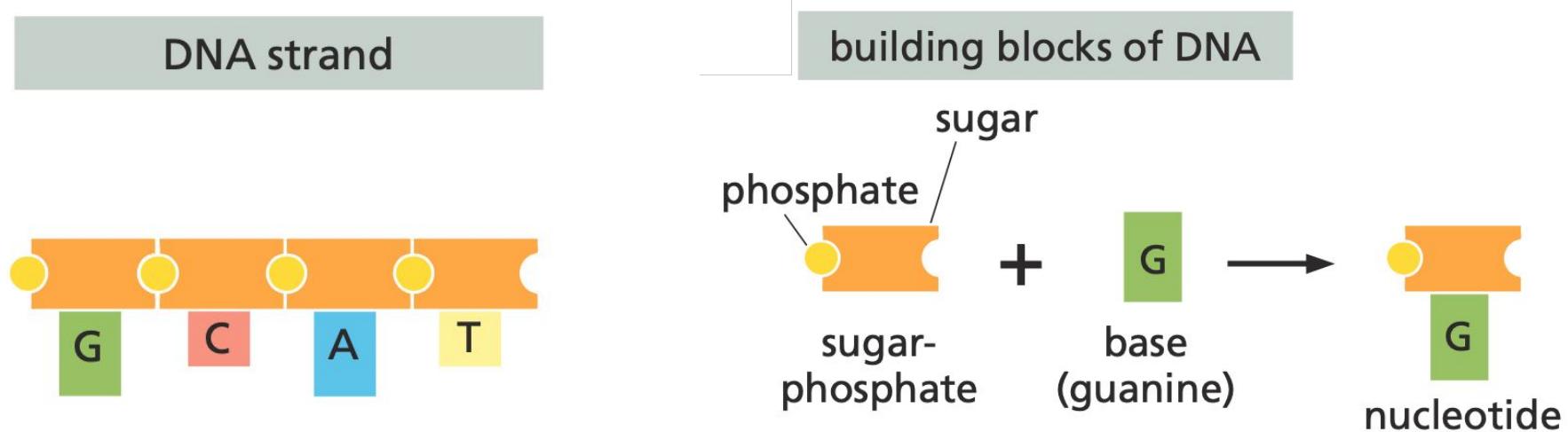
Genome

“The complete set of genes or genetic material present in a cell or organism.”

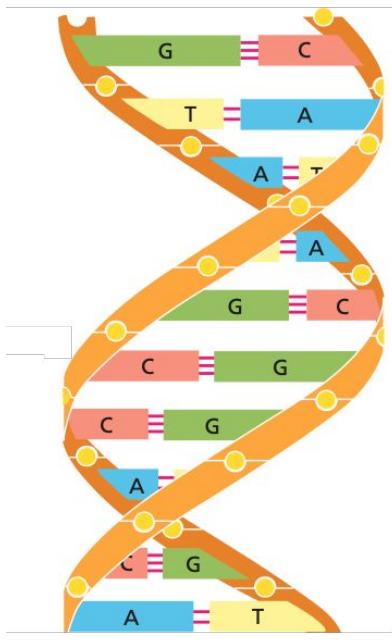
Oxford dictionaries

- containing the information needed to maintain organism's living.
- is made of DNA (or RNA in some viruses).

DNA: the code of life



DNA double helix

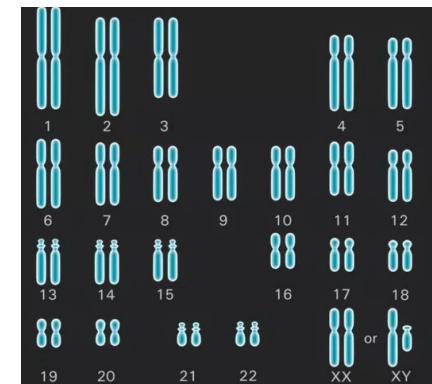
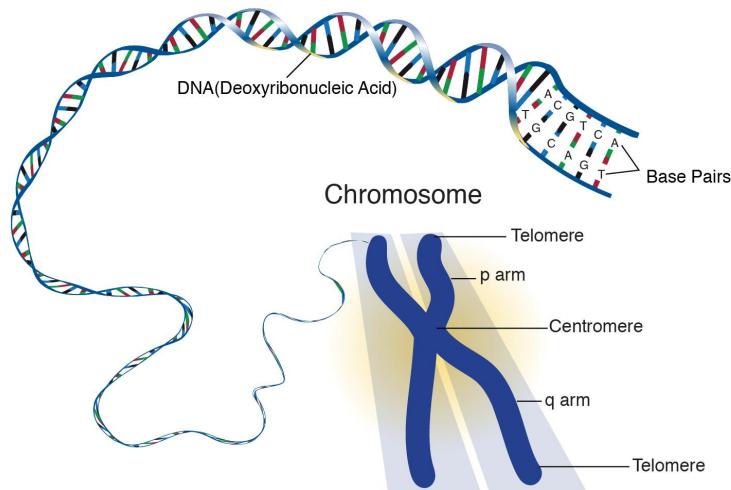


Genome size

Species	<i>T2 phage</i>	<i>Escherichia coli</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Paris japonica</i>
Genome Size	170,000 bp	4.6 million bp	130 million bp	3.2 billion bp	150 billion bp
Common Name	Virus	Bacteria	Fruit fly	Human	Canopy Plant

Chromosome

The DNA molecule is packaged into thread-like structures called chromosomes.





ucsc genome browser

[All](#)[Images](#)[Videos](#)[News](#)[Books](#)[More](#)[Tools](#)

About 2'920'000 results (1.08 seconds)

<https://genome.ucsc.edu>

⋮

UCSC Genome Browser Home

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the UCSC Genomics Institute.

Browser Gateway

As with the previous GRCh37 assembly, the Genome ...

Genome Browser

Use drop-down controls below and press refresh to alter tracks ...

In-Silico PCR

In-Silico PCR searches a sequence database with a pair ...

[More results from ucsc.edu »](#)

Human BLAT Search

Genome: Search all, Assembly:
Query type: Sort output: Output ...

hg19/GRCh37

Use drop-down controls below and press refresh to alter tracks ...

Table Browser

Use this tool to retrieve and export data from the Genome Browser ...



UCSC Genome Browser

The UCSC Genome Browser is an online and downloadable genome browser hosted by the University of California, Santa Cruz. [Wikipedia](#)



Our tools

- **Genome Browser**
interactively visualize genomic data
- **COVID-19 Research**
use the SARS-CoV-2 genome browser and explore coronavirus datasets
- **BLAT**
rapidly align sequences to the genome
- **Table Browser**
download data from the Genome Browser database
- **Variant Annotation Integrator**
get functional effect predictions for variant calls
- **Data Integrator**
combine data sources from the Genome Browser database
- **Genome Browser in a Box (GBiB)**
run the Genome Browser on your laptop or server
- **In-Silico PCR**
rapidly align PCR primer pairs to the genome
- **LiftOver**
convert genome coordinates between assemblies
- **Track Hubs**
import and view external data tracks
- **REST API**
returns data in JSON format

[More tools...](#)

Our story

On June 22, 2000, UCSC and the other members of the International Human Genome Project consortium completed the first working draft of the human genome assembly, forever ensuring free public access to the genome and the information it contains. A few weeks later, on July 7, 2000, the newly assembled genome was released on the web at <http://genome.ucsc.edu>, along with the initial prototype of a graphical viewing tool, the UCSC Genome Browser. In the ensuing years, the website has grown to include a broad collection of vertebrate and model organism assemblies and annotations, along with a large suite of tools for viewing, analyzing and downloading data. Learn more about our history on the UCSC Genome Browser Project History page and by watching this video.

What's new

- Nov. 23, 2021 - [New blog post about GenArk hubs](#)
- Nov. 18, 2021 - [New clinical rare disease track - Orphadata](#)
- Nov. 17, 2021 - [New Single-Cell track group and data for hg38](#)

[More news...](#)

[Subscribe](#)

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the UCSC Genomics Institute.

You might want to navigate to your nearest mirror - genome-euro.ucsc.edu

- User settings (sessions and custom tracks) will differ between sites.
- Take me to genome-euro.ucsc.edu
- Let me stay here genome.ucsc.edu

[Read more...](#)

Browse>Select Species

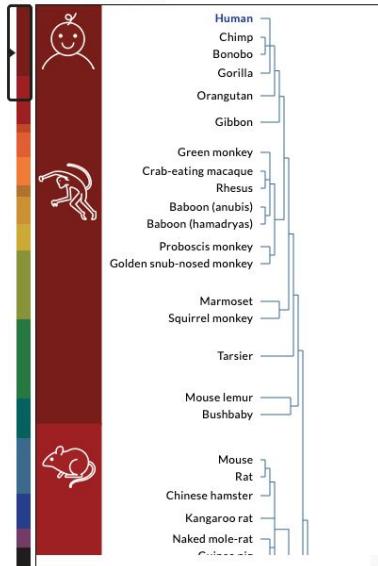
POPULAR SPECIES



Enter species, common name or assembly ID

Can't find a genome assembly?

REPRESENTED SPECIES



Find Position

Human Assembly

Dec. 2013 (GRCh38/hg38)

GO

Position/Search Term

Enter position, gene symbol or search terms

Current position: chrX:15,560,138-15,602,945

Human Genome Browser - hg38 assembly

view sequences

UCSC Genome Browser assembly ID: hg38

Sequencing/Assembly provider ID: Genome Reference Consortium Human GRCh38.p13 (GCA_000001405.28)

Assembly date: Dec. 2013 initial release; Dec. 2017 patch release 13

Assembly accession: GCA_000001405.28

NCBI Genome ID: 51 (Homo sapiens (human))

NCBI Assembly ID: GCF_000001405.39 (GRCh38.p13, GCA_000001405.28)

BioProject ID: PRJNA31257



Homo sapiens

(Graphic courtesy of CBSE)

Search the assembly:

- **By position or search term:** Use the "position or search term" box to find areas of the genome associated with many different attributes, such as a specific chromosomal coordinate range; mRNA, EST, or STS marker names; or keywords from the GenBank description of an mRNA. [More information](#), including sample queries.
- **By gene name:** Type a gene name into the "search term" box, choose your gene from the drop-down list, then press "submit" to go directly to the assembly location associated with that gene. [More information](#).
- **By track type:** Click the "track search" button to find Genome Browser tracks that match specific selection criteria. [More information](#).

Download sequence and annotation data:

- [Using rsync \(recommended\)](#)
- [Using HTTP](#)
- [Using FTP](#)
- [Data use conditions and restrictions](#)
- [Acknowledgments](#)

Assembly Details

The GRCh38 assembly is the first major revision of the human genome released in more than four years. As with the previous GRCh37 assembly, the Genome Reference Consortium (GRC) is now the primary source for human genome assembly data submitted to GenBank. Beginning with this release, the UCSC Genome Browser version numbers for the human assemblies now match those of the GRC to minimize version confusion. Hence, the GRCh38 assembly is referred to as "hg38" in the Genome Browser datasets and documentation. For a glossary of assembly-related terms, see the [GRC Assembly Terminology page](#).

Type: chr10

UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x 100x zoom out 1.5x 3x 10x 100x

multi-region chr10:1-133,797,422 133,797,422 bp. gene, chromosome range, or other position, see examples go examples

chr10 (p15.3-q26.3) 15.3 15.1 10p14 10p13 p12.31 10p12.1 11.21 q11.21 10q21.1 q21.2 10q21.3 q22.1 22.2 q22.3 10q23.1 23.31 24.2 10q25.1 25.2 q25.3 q26.13 26.2 10q26.3

Scale 50 Mb hg38

chr10: 10,000,000 20,000,000 30,000,000 40,000,000 50,000,000 60,000,000 70,000,000 80,000,000 90,000,000 100,000,000 110,000,000 120,000,000 130,000,000

Reference Assembly Fix Patch Sequence Alignments chr10_KN538367v1_fix

chr10_KN538367v1_fix

chr10_ML143354v1_fix

chr10_KN196480v1_fix

chr10_KQ090021v1_fix

Alt Haplotypes

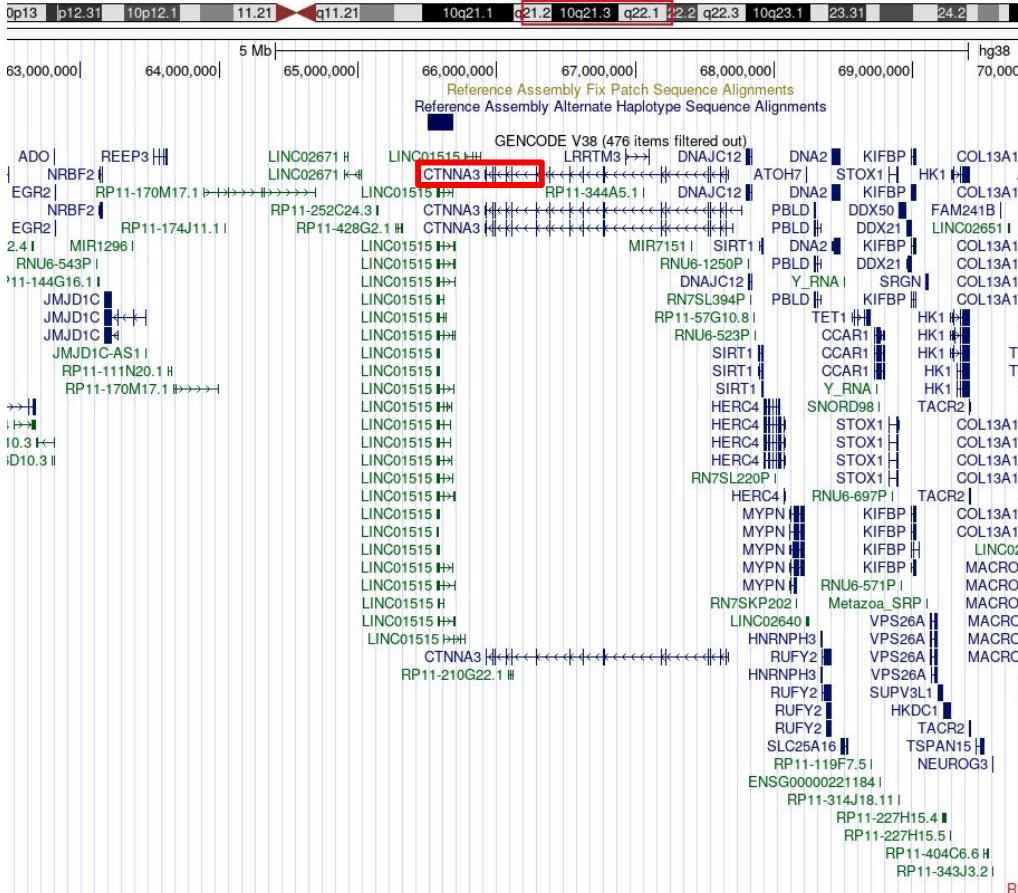
Reference Assembly Alternate Haplotype Sequence Alignments

	GENCODE V38 (4700 items, filtered out)
RP11-631M21.1	U81 CELF2 PTER EBLN1 AB1 KIF6B Y_RNA1 BMS1 MAPK8 ZWINT RTRK2 SIRT1 MCU DLG5 NRG3 KLLN IDE AVP1 PSD XPNPEP1 GFR1 ATE1 NPS
TUBB8	U81 CELF2 PTER BM1 AB1 EPC1 ZNF25 RET VSTM4 ZWINT RTRK2 SIRT1 MCU ZM12 LRT2 LIPJ IDE LOXL4 U61 XPNPEP1 GFR1 ATE1 FOX12
TUBB8	U81 CELF2 PTER SPAG6 MKX ITGB1 LINC00839 GDF2 I PTKM1 ADO SIRT1 MCU ZM12 LRT2 LIPF IDE HPS1 INA RN4-5P1 GFR1 ATE1 CLRN3
TUBB8	UCN3 CELF2 SPAG6 MKX LINC00839 RBP2 MBL2 CISD1 LINC01515 PRF1 LRMDA NRG3 PTEN IDE HPS1 TAF5 ADD3 GFR1 ATE1 PTPRE
TUBB8	NET1 CELF2 C1QL3 SPAG6 MPP7 NRP1 LINC00839 MAPK8 ZWINT RTRK2 SIRT1 OIT3 PPIF1 RGR LIPF IDE HPS2 STN1 ADD3 GFR1 ATE1
TUBB8	NET1 CELF2 C1QL3 SPAG6 MPP7 NRP1 LINC00839 MAPK8 ZWINT RTRK2 SIRT1 OIT3 E1F5AL1 RGR LIPF IDE HPS2 STN1 ADD3 GFR1 ATE1
ZMYND11	ASB13 CELF2 RSU1 ARMC3 MPP7 NRP1 LINC01515 PRF1 KCNMA1 NRG3 RLNL5 IDE HPS2 SLK MX1 GFR1 TACC2 PTPRE
ZMYND11	GD12 CELF2 RSU1 ARMC3 MPP7 NRP1 LINC01515 PRF1 KCNMA1 NRG3 RLNL5 IDE HPS2 SLK MX1 GFR1 TACC2 MKI67
ZMYND11	GD12 CELF2 RSU1 ARMC3 MPP7 NRP1 LINC01515 PRF1 KCNMA1 NRG3 RLNL5 IDE HPS2 SLK MX1 GFR1 TACC2 MKI67
ZMYND11	FBH1 CELF2 CUBN ARMC3 WAC PARD3 LINC01518 MAPK8 RP11-179B15.6 LINC01515 PRF1 KCNMA1 NRG3 RLNL5 IDE HPS2 SLK MX1 C10orf82 TACC2 MGMT
ZMYND11	FBH1 CELF2 VIM MSR52 WAC PARD3 LINC01518 WDFY4 UBE2D1 LINC01515 UNC5B KCNMA1 C10orf99 LIPM LG1 GO1 SORCS1 GPAM VAX1 HTRA1 LINC0266
ZMYND11	FBH1 CELF2 Y_RNA1 PTFA1 WAC PARD3 LINC02632 WDHY4 TFAM LINC01515 CDH23 POLR3A CDHR1 FAS LGI1 NKKX2-3 SORCS1 GPAM VAX1 HTRA1 EBF
ZMYND11	IL15RA CELF2 HACD1 OTUD1 WAC PARD3 AL022344.71 LRCL18 TFAM LINC01515 CDH23 POLR3A CDHR1 FAS LGI1 NKKX2-3 SORCS1 ACSL5 EMX2 PSTK EBF
ZMYND11	IL15RA CELF2 HACD1 MIR603 WAC PARD3 RNU6-885P MIR2494 BICC1 CTNN3 P4HA1 SFTPA2 RGR FAS LGI1 COX15 SORCS1 TECTB EMX2 PSTK EBF
ZMYND11	IL15RA CELF2 STAM KIAA1217 SVIL CUL2 LINC02623 VSTM4 BICC1 CTNN3 P4HA1 SFTPA2 RGR FAS LGI1 COX15 SORCS1 ACSL5 CASC2 II HMX2 MIR429
ZMYND11	IL15RA UPF2 MRC1 PRTFDC1 SVIL CUL2 LINC01264 C10orf71 RN7SKP196 LINC01515 UNC5B KCNMA1 C10orf99 LIPM LG1 GO1 SORCS1 GPAM VAX1 HTRA1
ZMYND11	IL15RA UPF2 MIR511 KIAA1217 SVIL CUL2 MIR5100 DRGX LINC00844 CTNN3 P4HA1 SFTPA2 RGR FAS LGI1 CUTC LINC01435 HABP2 EIF3A OAT GLRX
ZMYND11	IL15RA UPF2 CACNB2 ENKUR LYZL2 GJD4 LINC02633 ERCC6 LINC00844 MIR2494 LRRTM3 VSIR1 RPS24I LIR11 LIPA LGI1 ABC2C LINC01435 HABP2 EIF3A OAT GLRX
ZMYND11	IL15RA DHTKD1 CACNB2 ENKUR JCAD CUL2 RET ERCC6 LINC00844 MIR1751 P4HA1 SFTPA1 RGR FAS LGI1 DNMBP LINC01435 NRAP SFNX1 LHPY Y_RI
DIP2C	IL15RA NUDT5 CACNB2 ENKUR LYZL2 GJD4 LINC02633 ERCC6 LINC00844 MIR1751 P4HA1 SFTPA1 RGR FAS LGI1 DNMBP LINC01435 NRAP GRK5 LHPH MIR3
DIP2C	IL2RA NUDT5 NSNU6 THNSL1 LYZL2 FZD8 FXDY4 ERCC6 LINC00844 MIR1751 P4HA1 SFTPA1 RGR FAS LGI1 DNMBP LINC01435 NRAP TIAL1 LHPH PPP
DIP2C	IL2RA NUDT5 ARLS5 THNSL1 ZNF438 Y_RNA1 FXDY4 ERCC6 PHYHIP1 RNU6-1250P Y_RNA1 NUMT2E WAPL1 BTAFL1 LCOR1 FGFB1 RNUS6-6P1 VWA2 TIAL1 CTBP2 E
RN7SL754P	IL15RA NUDT5 MALRD1 MYO3A ZEB1 ZNF248 ZNF32 CHAT PHYHIP1 RNU6-1250P Y_RNA1 NUMT2E WAPL1 BTAFL1 LCOR1 FGFB1 RNUS6-6P1 VWA2 TIAL1 CTBP2
RP11-490E15.2	IL15RA NUDT5 U3 MYO3A ZEB1 ZNF248 ZNF32 CHAT PHYHIP1 RNU6-1250P Y_RNA1 NUMT2E WAPL1 BTAFL1 LCOR1 FGFB1 RNUS6-6P1 VWA2 TIAL1 CTBP2
RP11-809C18.3	IL2RA CDC123 PLXDC2 GAD2 ZEB1 ZNF248 CXCL12 CHAT FAM13C RN7SL394P ANXA7 SFTP1 OPN4 Y_RNA1 SL1T1 FGFB1 ADD3-AS1 TACC2 C10orf14 LINC0264
MIR5699	RBM17 CDC123 PLXDC2 GAD2 ZEB1 ZNF248 CXCL12 CHAT FAM13C RN7SL394P ANXA7 SFTP1 OPN4 Y_RNA1 SL1T1 FGFB1 ADD3-AS1 TACC2 C10orf14 LINC0264
DIP2C-AS1	RBM17 MIR4480 MIR4675 GAD2 Y_RNA1 ZNF248 CXCL12 CHAT FAM13C HERC4 USP54 ANXA11 LDB3 HHEX ENTPD7 LINC02661 TRUB1 PLPP4 UROS DF
RP11-16A1C2.1	Y_RNA1 MIR4481 NEBL PDSS1 EPC1 Y_RNA1 CXCL12 CHAT FAM13C HERC4 USP54 ANXA11 LDB3 HHEX ENTPD7 LINC02661 TRUB1 PLPP4 UROS DF
LARP4B	LINP11 CCDC3 NEBL ABI1 EPC1 ZNF37A CXCL12 PARG FAM13C HERC4 USP54 ANXA11 LDB3 HHEX ENTPD7 LINC02661 TRUB1 PLPP4 UROS DF
RP11-36N22.2	LINP11 CCDC3 NEBL-AS1 ABI1 EPC1 ZNF37A CXCL12 PARG FAM13C HERC4 USP54 ANXA11 LDB3 HHEX ENTPD7 LINC02661 TRUB1 PLPP4 UROS DF
RP11-36N23.2	LINP11 OPTN RNU6-15P1 ABI1 EPC1 ZNF37A TMEM72 A1CF FAM13C MYPN FUT11 ANXA11 LDB3 EXOC6 CPN1 ADD3-AS1 HSPA12A TACC2 TCERI LINC0265
GTPBP4	SFMFT2 OPTN LINC02643 ABI1 CCDC7 AL022345.12I FAM170B RP11-443C13.31 MYPN FUT11 MAT1A BMPR1A EXOC6 ERLIN1 ADD3-AS1 HSPA12A TACC2 LINC0265
ID12	SFMFT2 OPTN Y_RNA1 ABI1 CCDC7 CGSALNACT2-DT SLC18A3 RP11-5C5.11 MYPN NDST2 DYDC1 MMRN2 EXOC6 ERLIN1 ADD3-AS1 HSPA12A TACC2 RP11-245I
ID12-AS1	SFMFT2 OPTN MIR1915HG1 ABI1 CCDC7 CGSALNACT2 C10orf53 RP11-135D11.21 MYPN MYOZ1 DYDC1 SNCG EXOC6 CHUK RP11-549L6.3 I ENO4 DMBT1 RP11-45
ID12-AS1	SFMFT2 OPTN MIR1915 ABI1 CCDC7 RASGEF1A C10orf53 RP11-135D11.21 MYPN MYOZ1 DYDC1 SNCG EXOC6 CWF19L1 MX1 ENO4 DMBT1 JAK
ID12-AS1	SFMFT2 OPTN SKIDA1 ABI1 CCDC7 RASGEF1A C10orf53 SLC16A9 ATOH7 NDST2 DYDC2 ADIRF1 EXOC6 SNORA12 MX1 ENO4 DMBT1 JAK
ID12-AS1	SFMFT2 MCM10 SKIDA1 ABI1 CCDC7 RASGEF1A C10orf53 SLC16A9 LINC02640 MYPF PLAU DYDC2 FAM25A MYOF BLOC1S2 MX1 SHTN1 DMBT1 JAK
ID11	SFMFT2 MCM10 MLLT10 ABI1 ITGB1 RP11-16BLL2.3I OGDHL MRLN PBLD PLAU DYDC2 GLUD1 MYOF BLOC1S2 MX1 SHTN1 DMBT1 JAK

UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

multi-region chr10:60,208,841-73,588,582 13,379,742 bp. gene, chromosome range, or other position, see examples go example





Genomes

Genome Browser

Tools

Mirrors

Downloads

My Data

Projects

Help

About Us

Human Gene CTNNA3 (ENST00000684154.1) from GENCODE V38

Description: catenin alpha 3 (from HGNC CTNNA3)

RefSeq Summary (NM_001127384): This gene encodes a protein that belongs to the vinculin/alpha-catenin family. The encoded protein plays a role in ventricular dysplasia, familial 13. Alternative splicing results in multiple transcript variants. [provided by RefSeq, Mar 2014].

Gencode Transcript: ENST00000684154.1

Gencode Gene: ENSG00000183230.18

Transcript (Including UTRs)

Position: hg38 chr10:65,912,457-67,763,637 **Size:** 1,851,181 **Total Exon Count:** 18 **Strand:** -

Coding Region

Position: hg38 chr10:65,920,330-67,647,513 **Size:** 1,727,184 **Coding Exon Count:** 17

Page Index	Sequence and Links	MalaCards	CTD	RNA-Seq Expression	Microarray Expression
RNA Structure	Other Species	mRNA Descriptions	Pathways	Other Names	Methods

Data last updated at UCSC: 2021-09-27 18:51:20

- Sequence and Links to Tools and Databases

Genomic Sequence (chr10:65,912,457-67,763,637)	mRNA (may differ from genome)	Protein (895 aa)
Gene Sorter	Genome Browser	Other Species FASTA
CGAP	Ensembl	ExonPrimer
PubMed		

- MalaCards Disease Associations

MalaCards Gene Search: [CTNNA3](#)

Diseases sorted by gene-association score: [arrhythmogenic right ventricular dysplasia, familial, 13*](#) (1242), [arrhythmogenic right ventricular dysplasia, ventricular dysplasia, right dominant form*](#) (101), [familial isolated arrhythmogenic ventricular dysplasia, biventricular form*](#) (101), fam

* = Manually curated disease association

- Comparative Toxicogenomics Database (CTD)

The following chemicals interact with this gene

- [C006253](#) pirinixic acid
- [D003485](#) Cyanates
- [D009151](#) Mustard Gas
- [D010634](#) Phenobarbital
- [D013749](#) Tetrachlorodibenzodioxin
- [C008261](#) lead acetate
- [C028007](#) nickel monoxide
- [C025643](#) vinclozolin

Phenotype and Literature

- OMIM Alleles**: dense ▾
- CADD**: hide ▾
- Cancer Gene Expr**: hide ▾
- ClinGen**: hide ▾
- Deprecated ClinGen CNVs**: hide ▾
- ClinVar Variants**: hide ▾
- Coriell CNVs**: hide ▾
- COSMIC Regions**: hide ▾
- Development Delay**: hide ▾
- Gene Interactions**: hide ▾
- GeneReviews**: hide ▾
- GWAS Catalog**: hide ▾
- HGMD Variants**: hide ▾
- LOVD Variants**: hide ▾
- OMIM Cyto Loci**: hide ▾
- OMIM Genes**: hide ▾
- New Orphanet**: hide ▾
- SNPedia**: hide ▾
- TCGA Pan-Cancer**: hide ▾
- UniProt Variants**: hide ▾
- Variants in Papers**: hide ▾

COVID-19

- COVID GWAS v4**: hide ▾
- COVID GWAS v3**: hide ▾
- Rare Harmful Vars**: hide ▾

Single Cell RNA-seq

- Colon Wang**: New hide ▾
- Ileum Wang**: New hide ▾
- Rectum Wang**: New hide ▾
- Blood (PBMC) Hao**: New hide ▾
- Cortex Velmeshev**: New hide ▾
- Fetal Gene Atlas**: 19 New hide ▾
- Heart Cell Atlas**: New hide ▾
- Kidney Stewart**: New hide ▾
- Liver MacParland**: New hide ▾
- Lung Travaglini**: New hide ▾
- Muscle De Michelis**: New hide ▾
- Pancreas Baron**: New hide ▾
- Placenta Vento-Tormo**: New hide ▾
- Skin Sole-Boldo**: 19 New hide ▾

mRNA and EST

- Human ESTs**: hide ▾
- Human mRNAs**: hide ▾
- Other ESTs**: hide ▾
- Other mRNAs**: hide ▾
- SIB Alt-Splicing**: hide ▾
- Spliced ESTs**: hide ▾

Expression

- GTEX Gene V8**: pack ▾
- GTEX RNA-Seq Coverage**: hide ▾
- Affy Archive**: New hide ▾
- EPDnew Promoters**: hide ▾
- GNF Atlas 2**: hide ▾
- GTEX Gene**: 19 hide ▾
- GTEX Transcript**: hide ▾
- GWIPS-viz Riboseq**: hide ▾
- miRNA Tissue Atlas**: hide ▾

Regulation

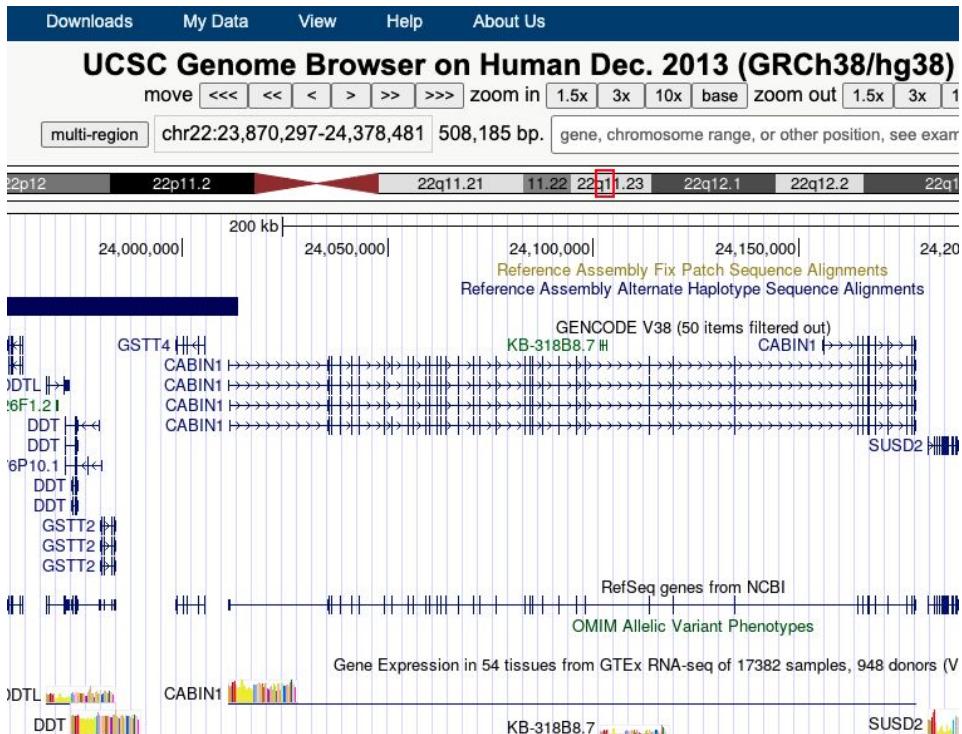
- ENCODE cCREs**: dense ▾
- ENCODE Regulation**: show ▾
- CpG Islands**: hide ▾
- GeneHancer**: hide ▾
- GTEX cis-eQTLs**: hide ▾
- Hi-C and Micro-C**: hide ▾
- JASPAR Transcription Factors**: New hide ▾
- OREGAnno**: hide ▾
- RefSeq Func Elems**: hide ▾

Comparative Genomics

- Conservation**: ..
- Cactus 241-way**: New hide ▾
- Cons 30 Primates**: hide ▾
- Primate Chain/Net**: hide ▾
- Placental Chain/Net**: hide ▾
- Vertebrate Chain/Net**: hide ▾

Exercise

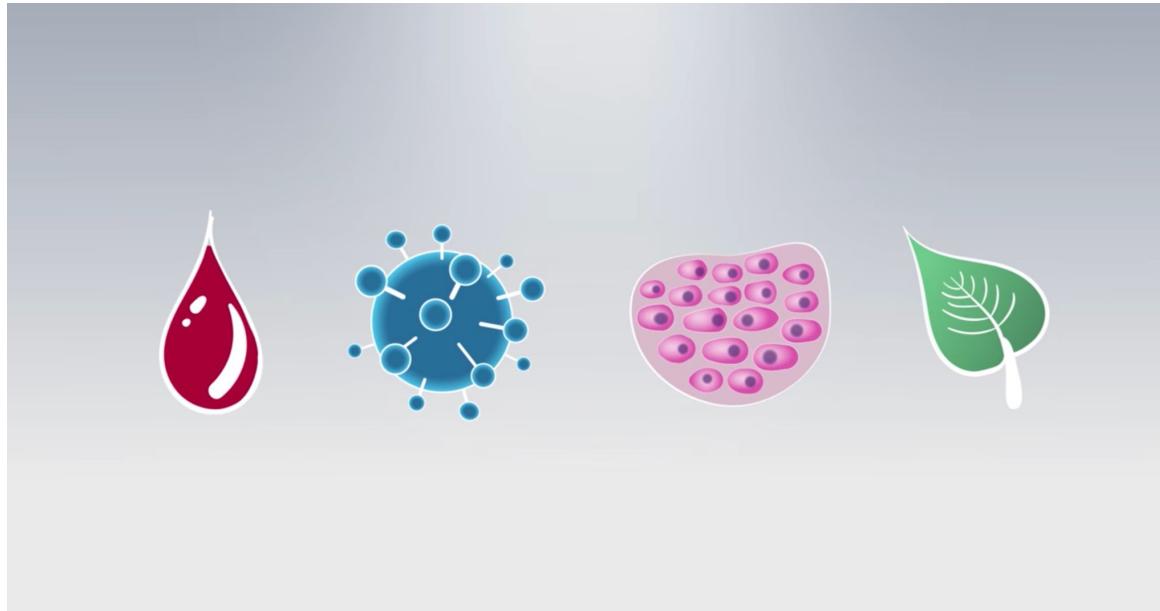
Mention a gene located at chr 22 position 24,100,000!

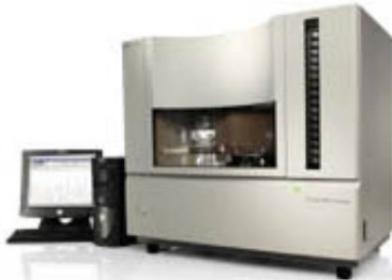


Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies

DNA sequencing





Sanger DNA sequencing

1977-1990s



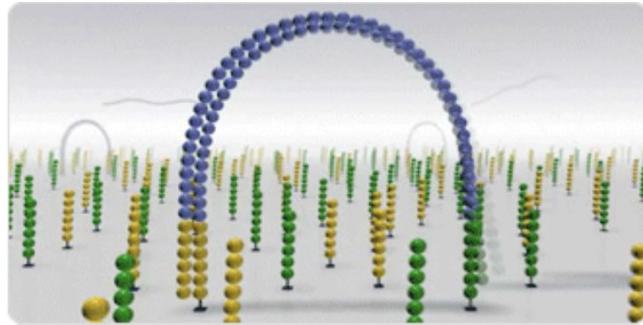
2nd-generation DNA sequencing

Since ~2007

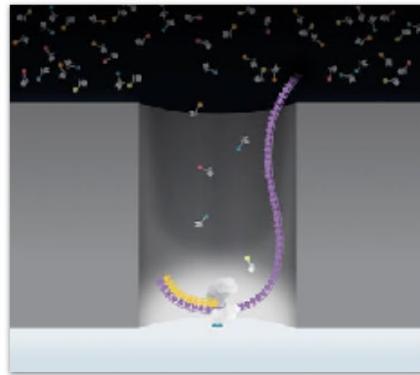


3rd-generation & single-molecule DNA sequencing

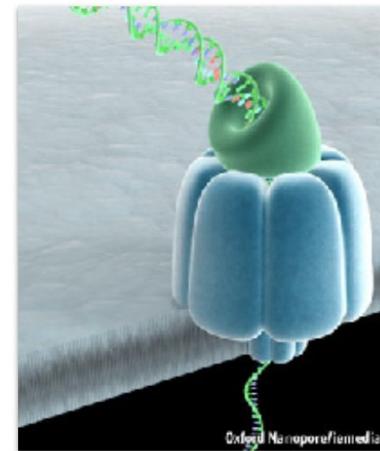
Since ~2010



Synthesis / ligation



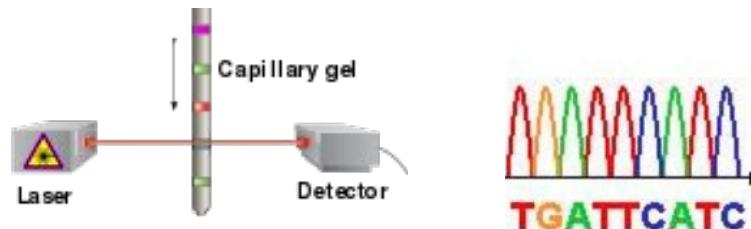
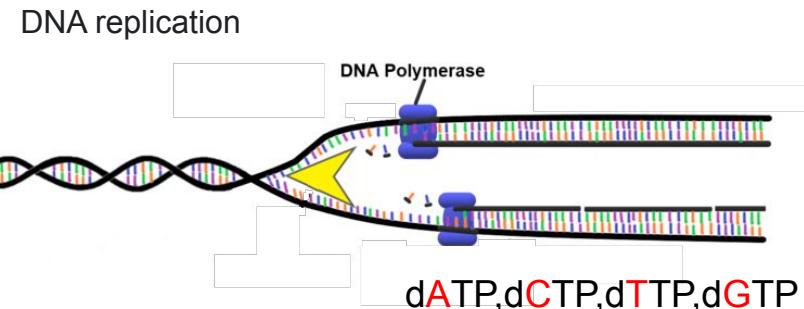
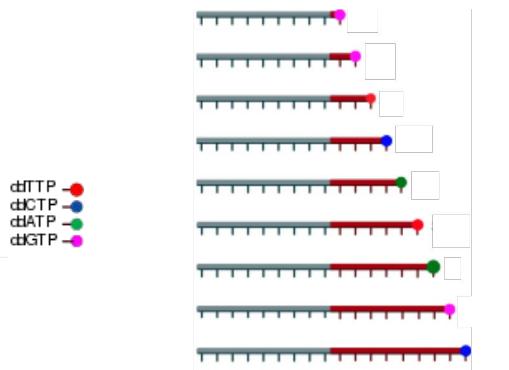
SMRT cell



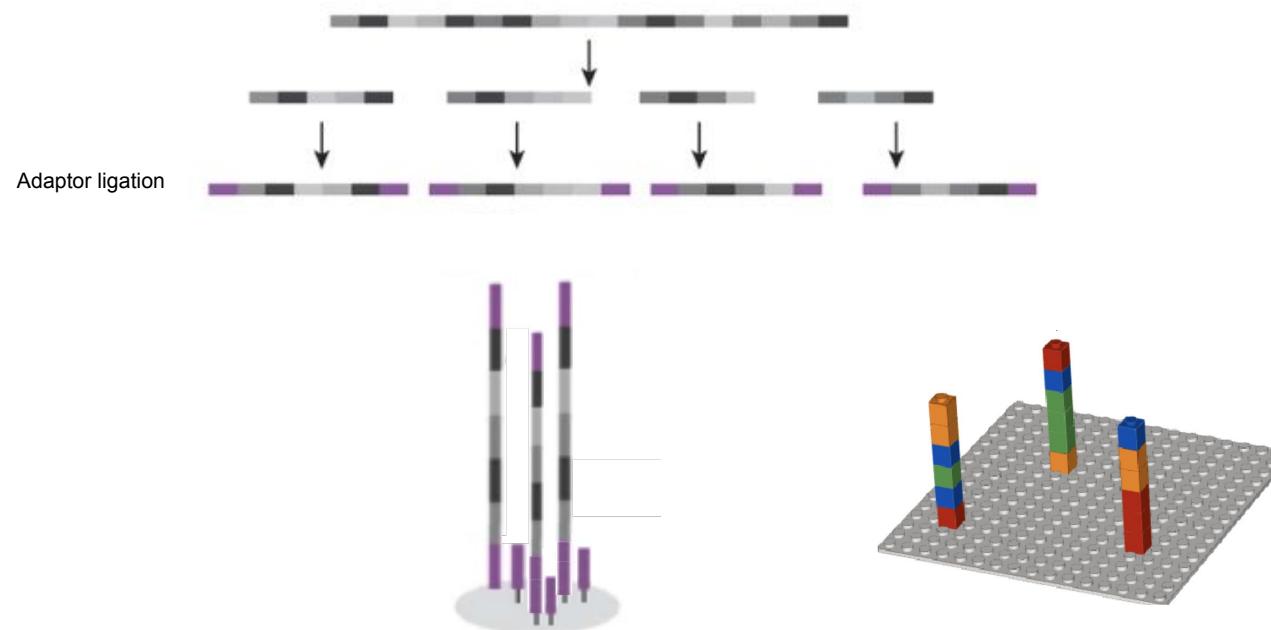
Nanopore

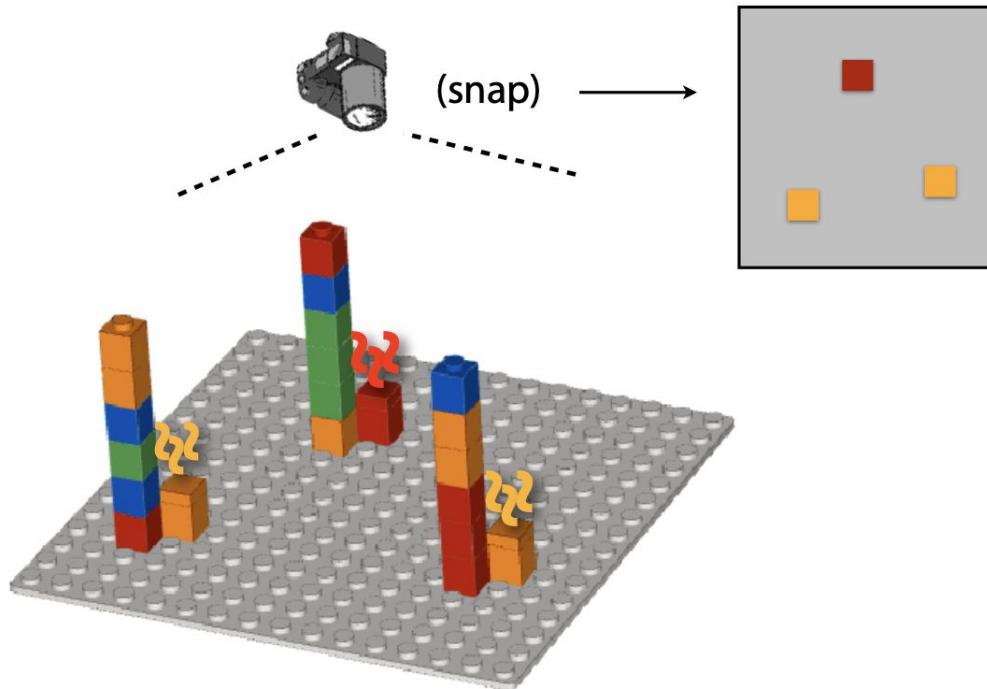
Sanger sequencing

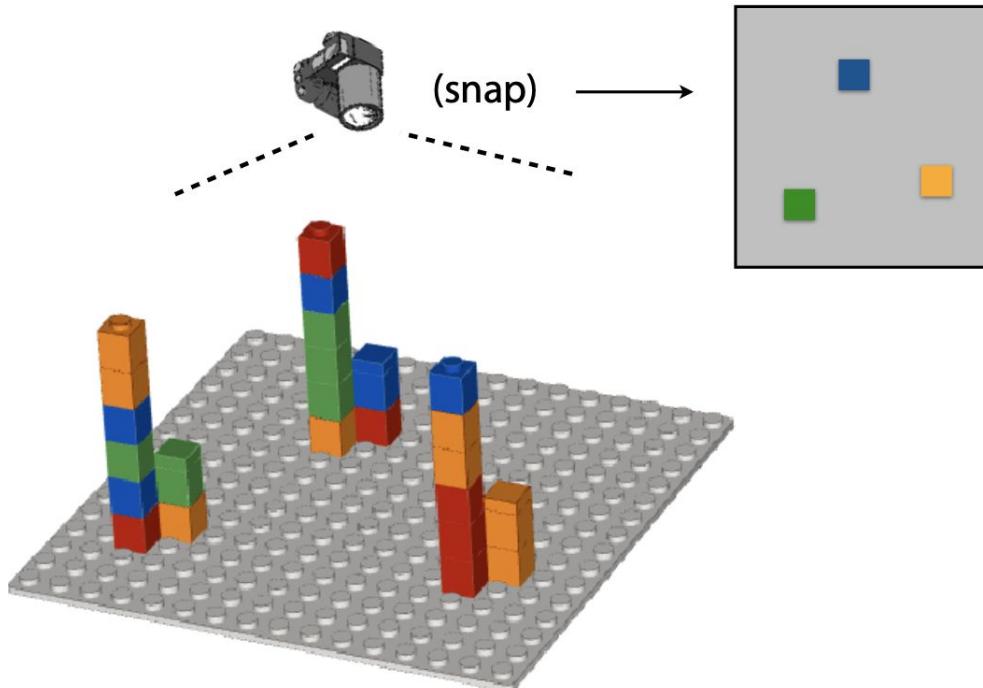
- sequencing by synthesis
- based on the chain termination procedure
- output sequences of different lengths

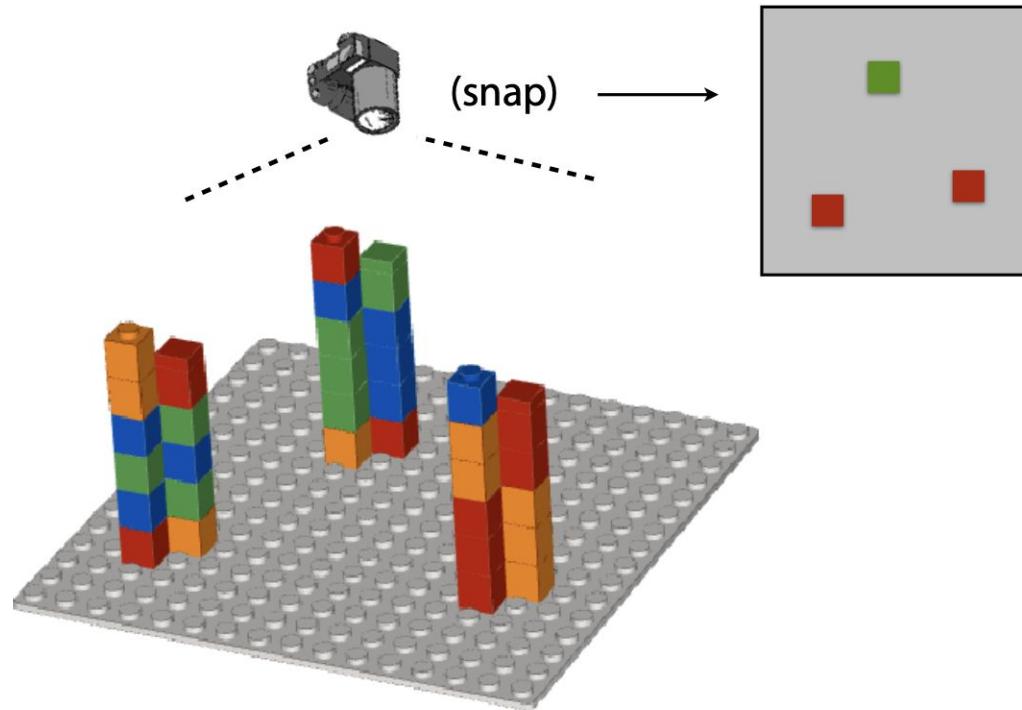


Illumina sequencing

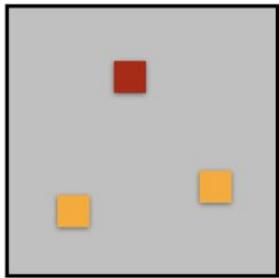




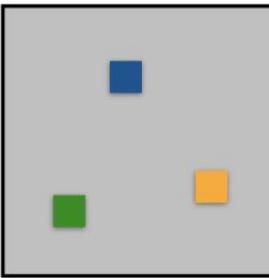




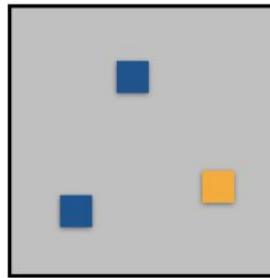
Cycle 1



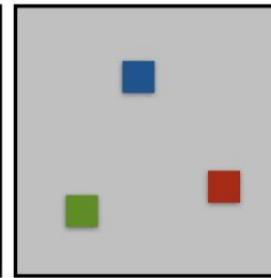
Cycle 2



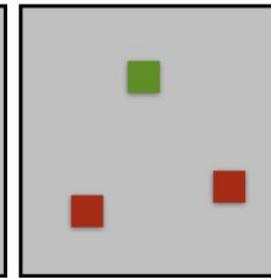
Cycle 3



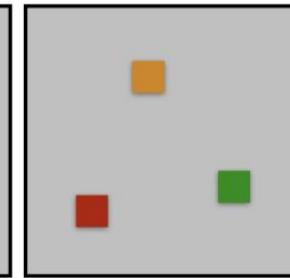
Cycle 4

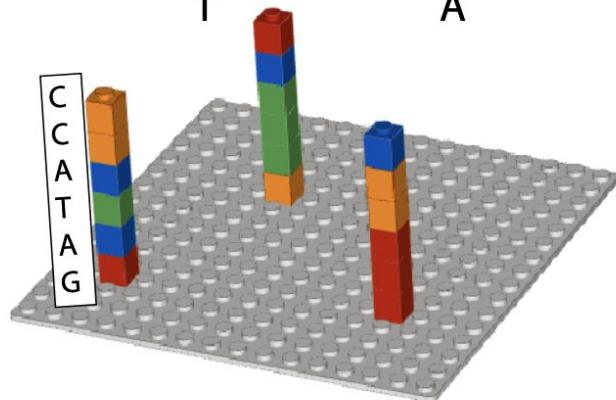
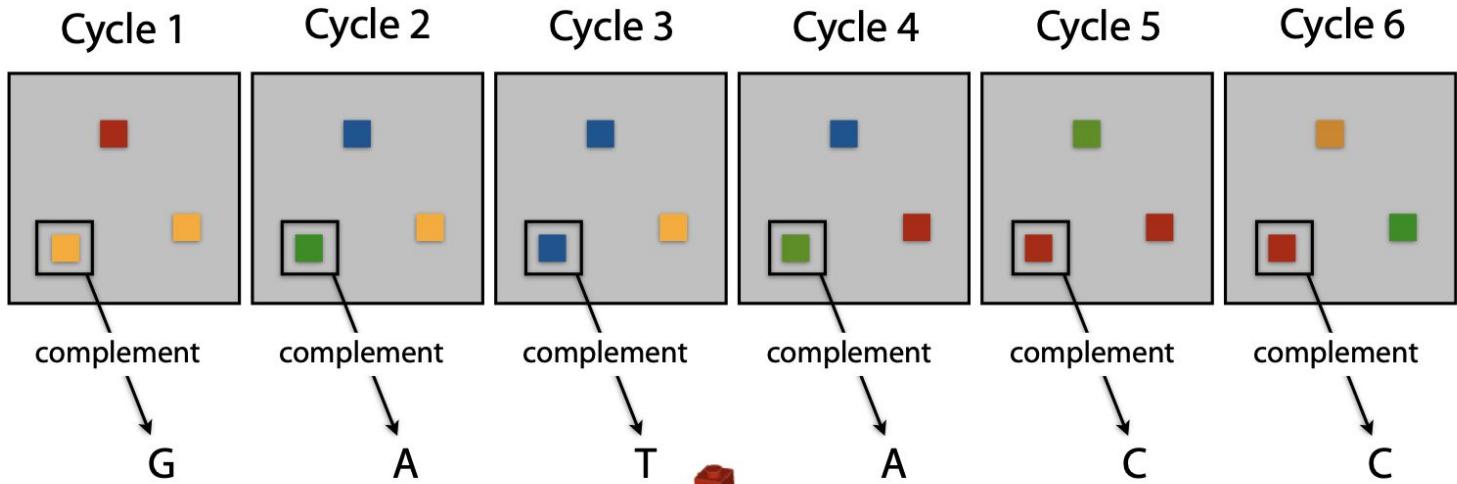


Cycle 5

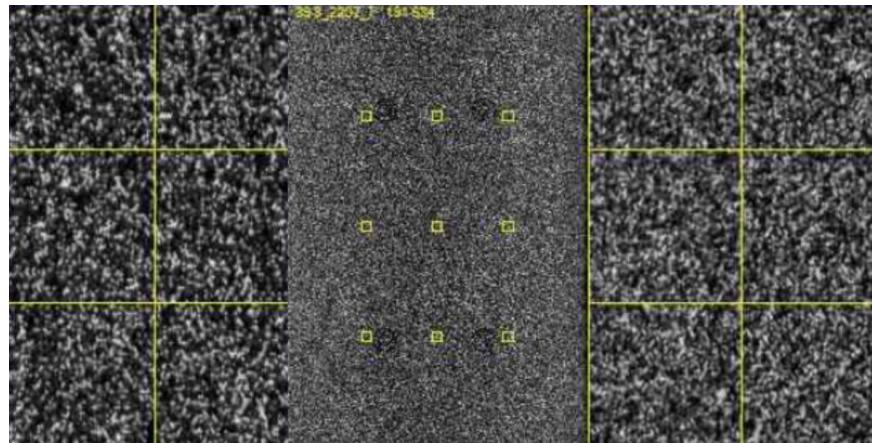


Cycle 6

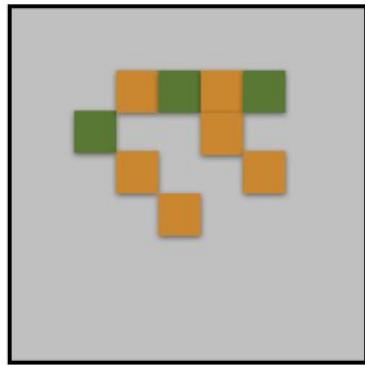




Actual Illumina HiSeq 3000



<https://www.youtube.com/watch?v=CZeN-lgjYCo>



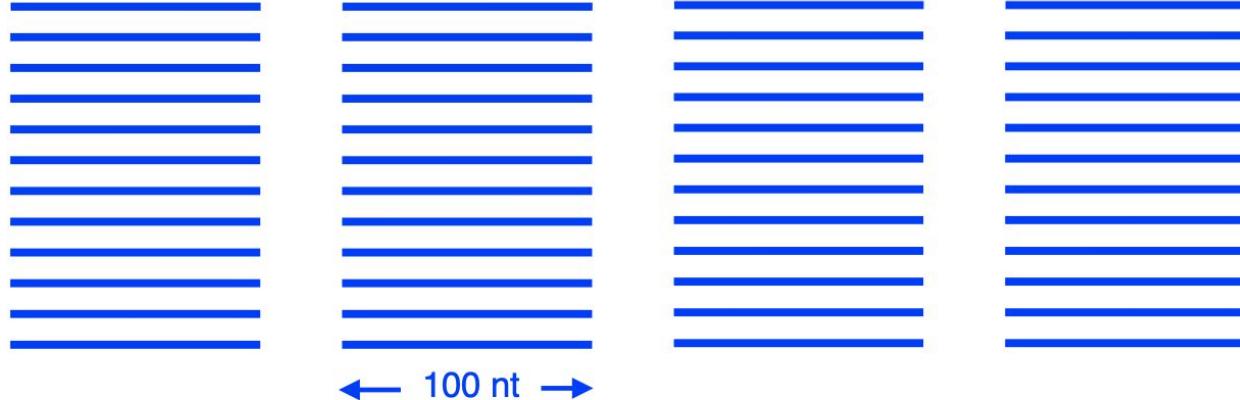
Call: orange (C)

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

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

Illumina error rate 0.01%

Reads



Your genome



Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell



PACBIO

Single Molecule, Real-Time (SMRT®) Sequencing



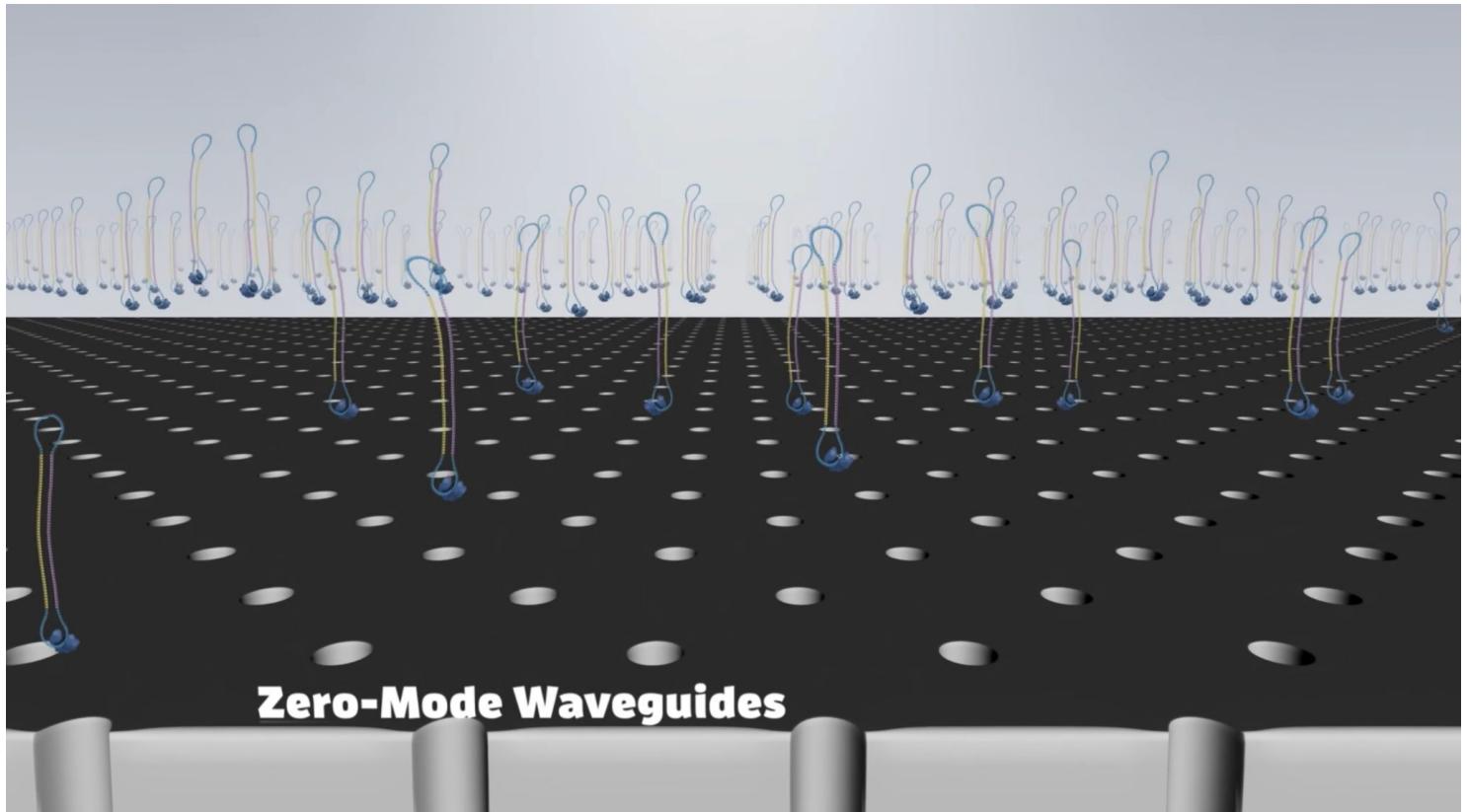


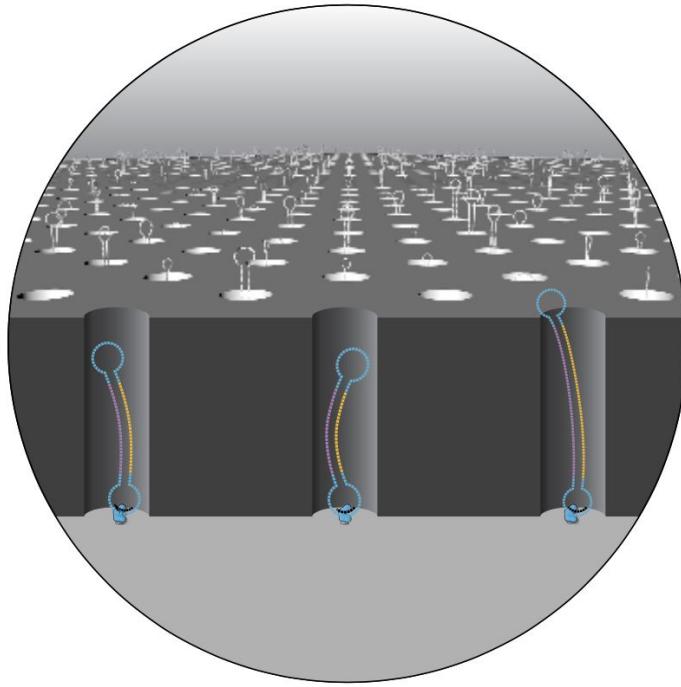
https://www.youtube.com/watch?v=_ID8JyAbwEo

SMRTbell® Library

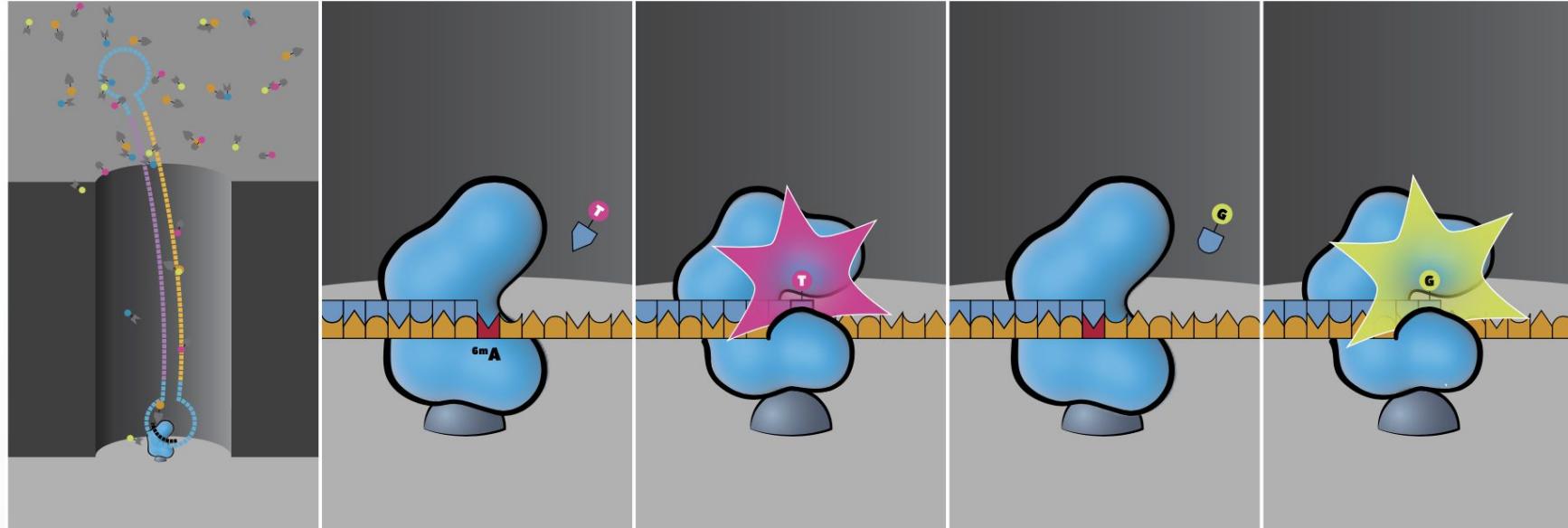


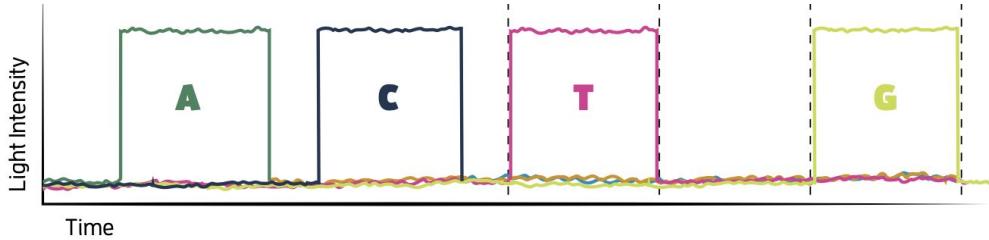
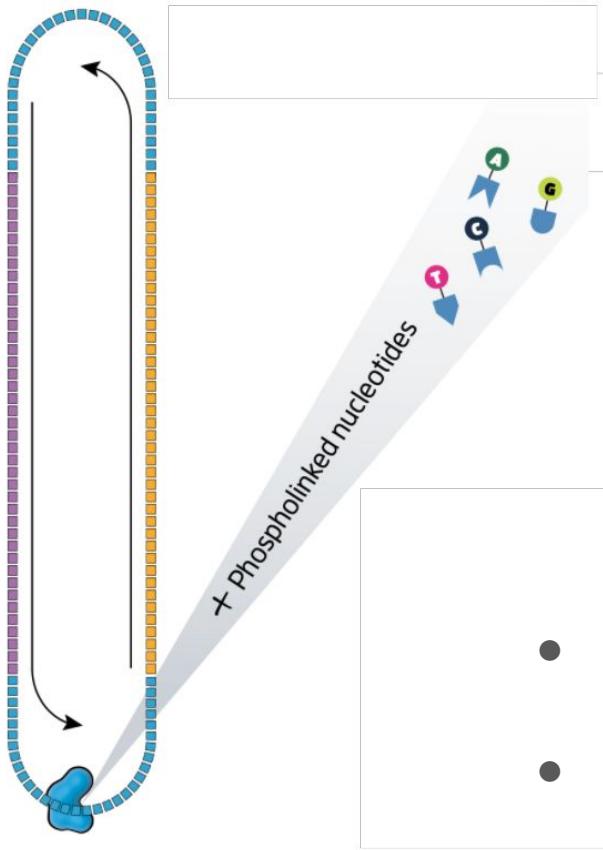
Primer + Polymerase





SMRT Cells contain millions of
zero-mode waveguides (ZMWs)





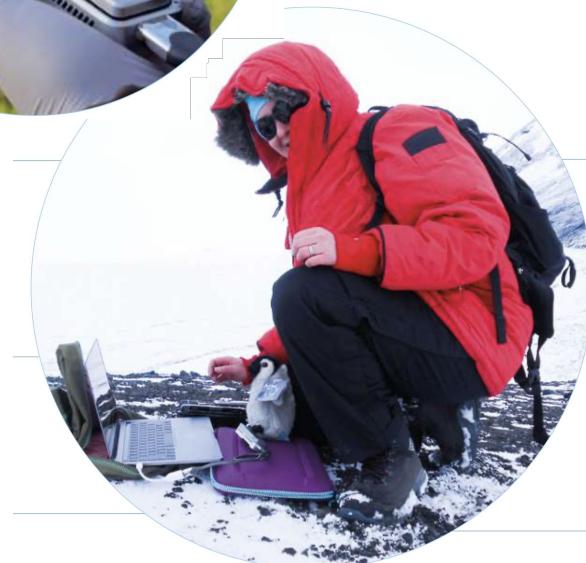
- Circular Consensus sequencing (CCS)
or High Fidelity (HiFi)
- ~12,000 bp

ONT

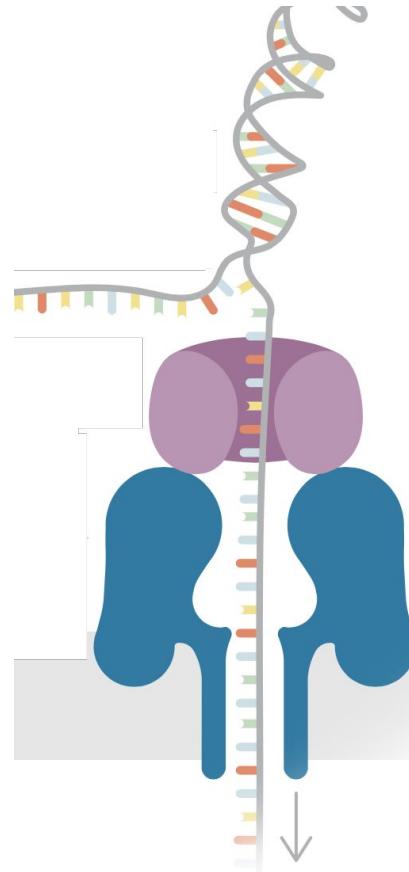


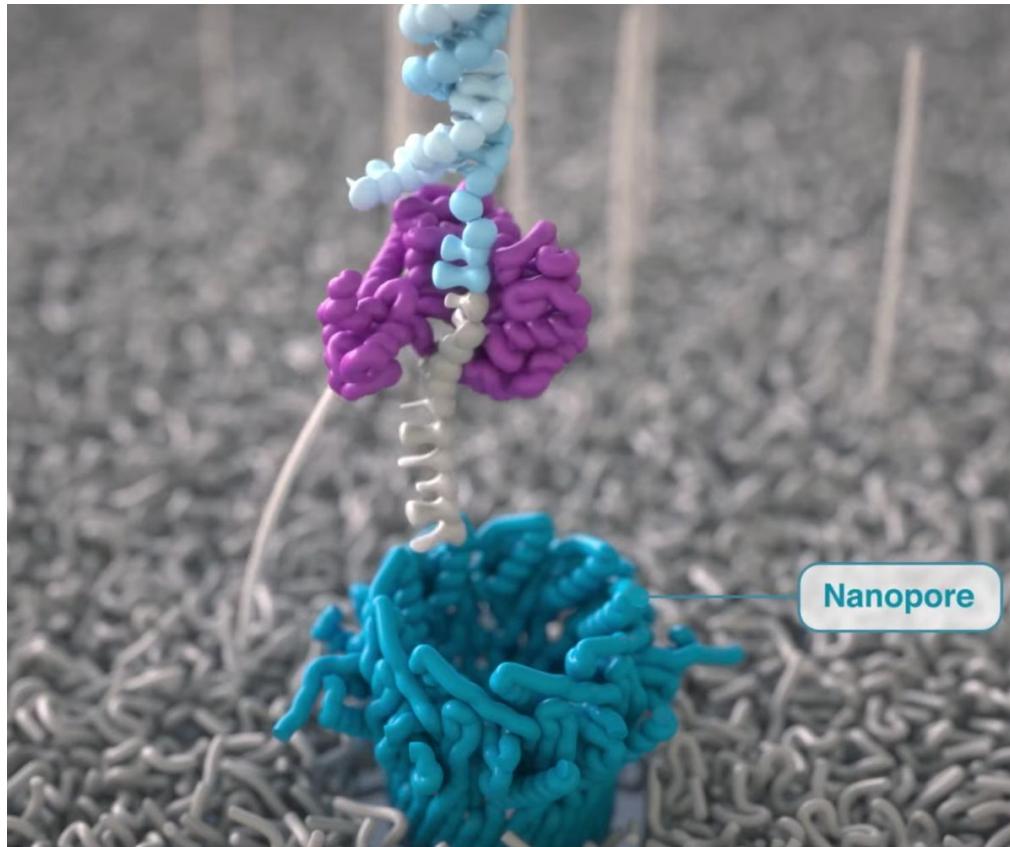
100,000- 4,000,000 bp

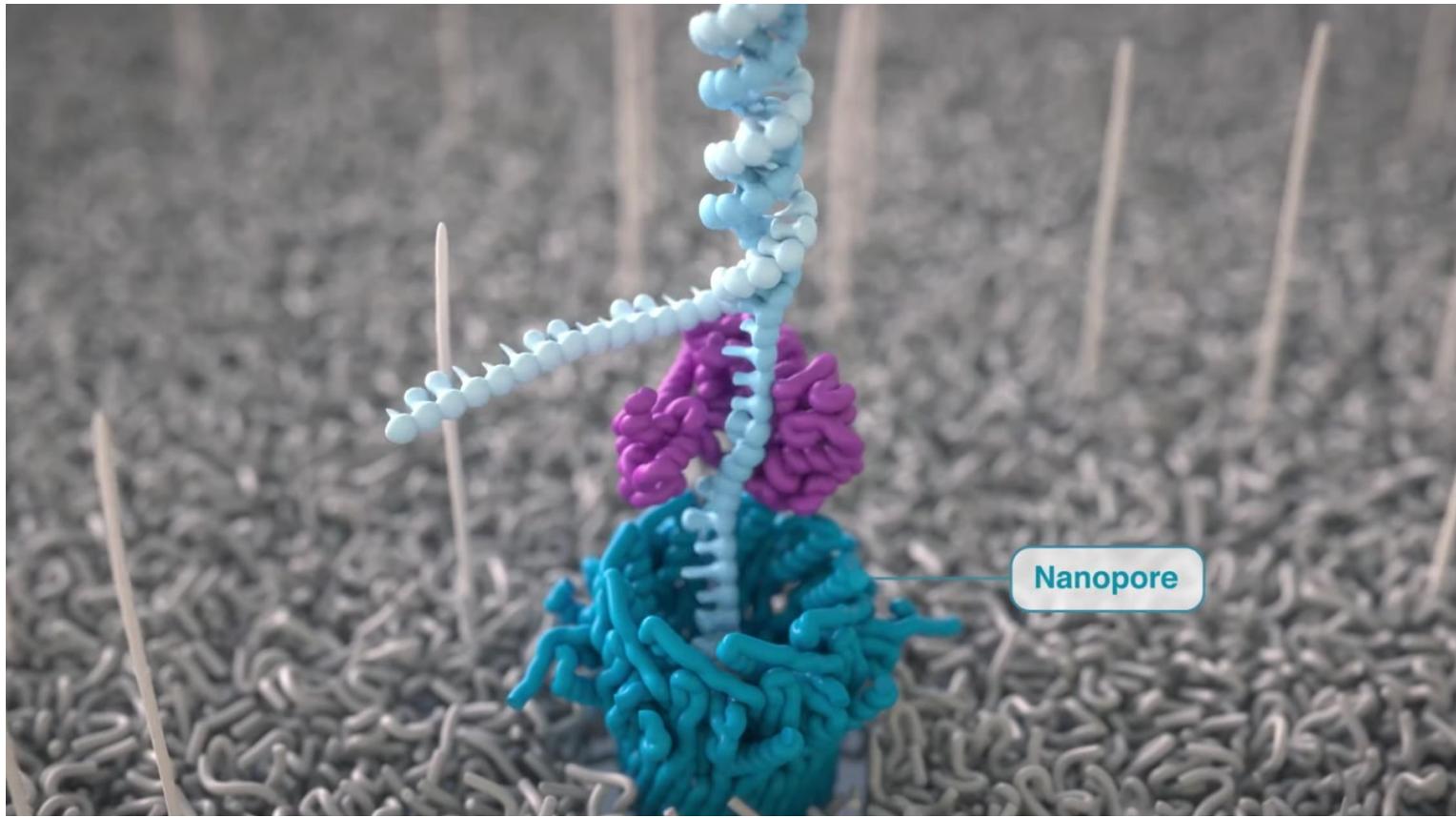
<https://www.youtube.com/watch?v=qzusVw4Dp8w>

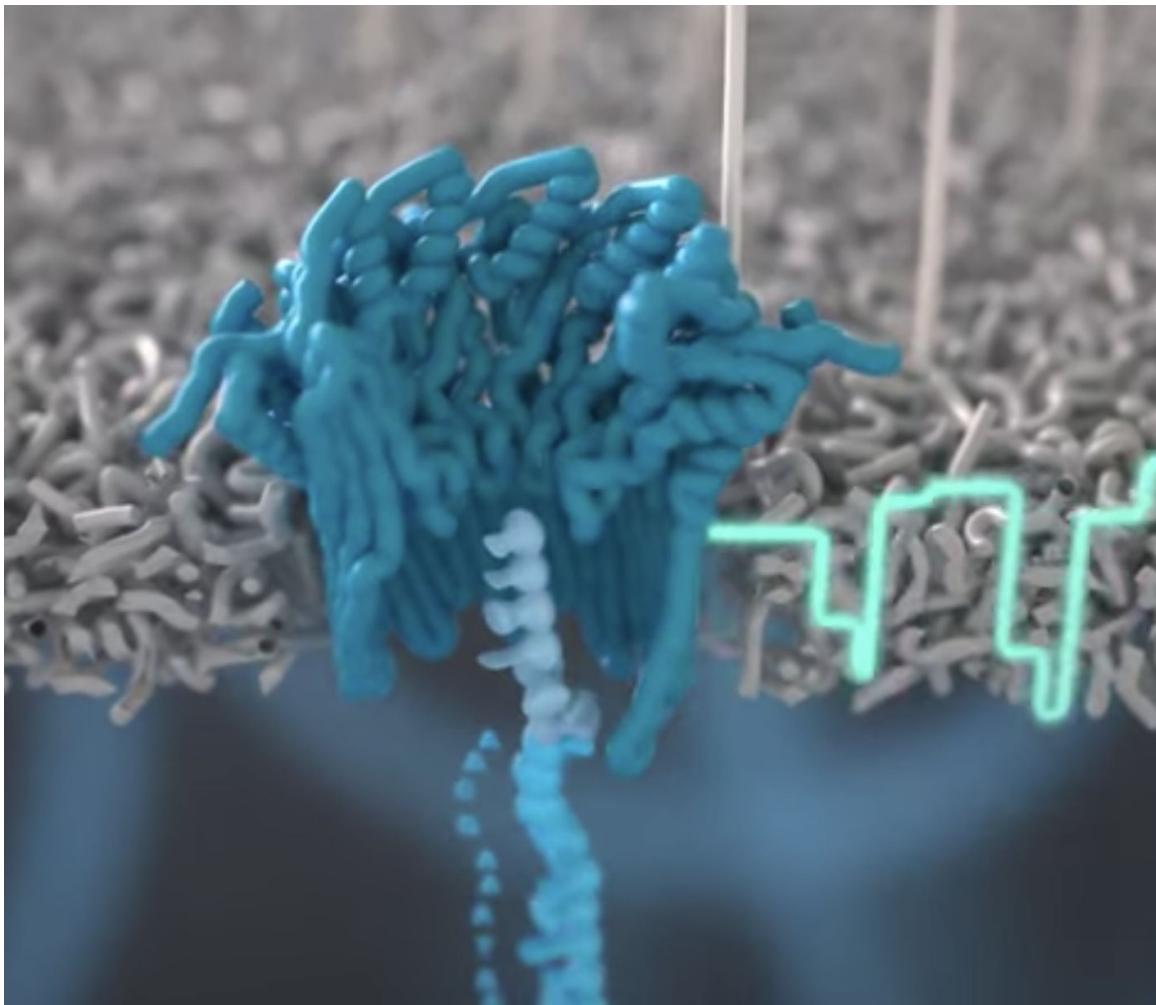


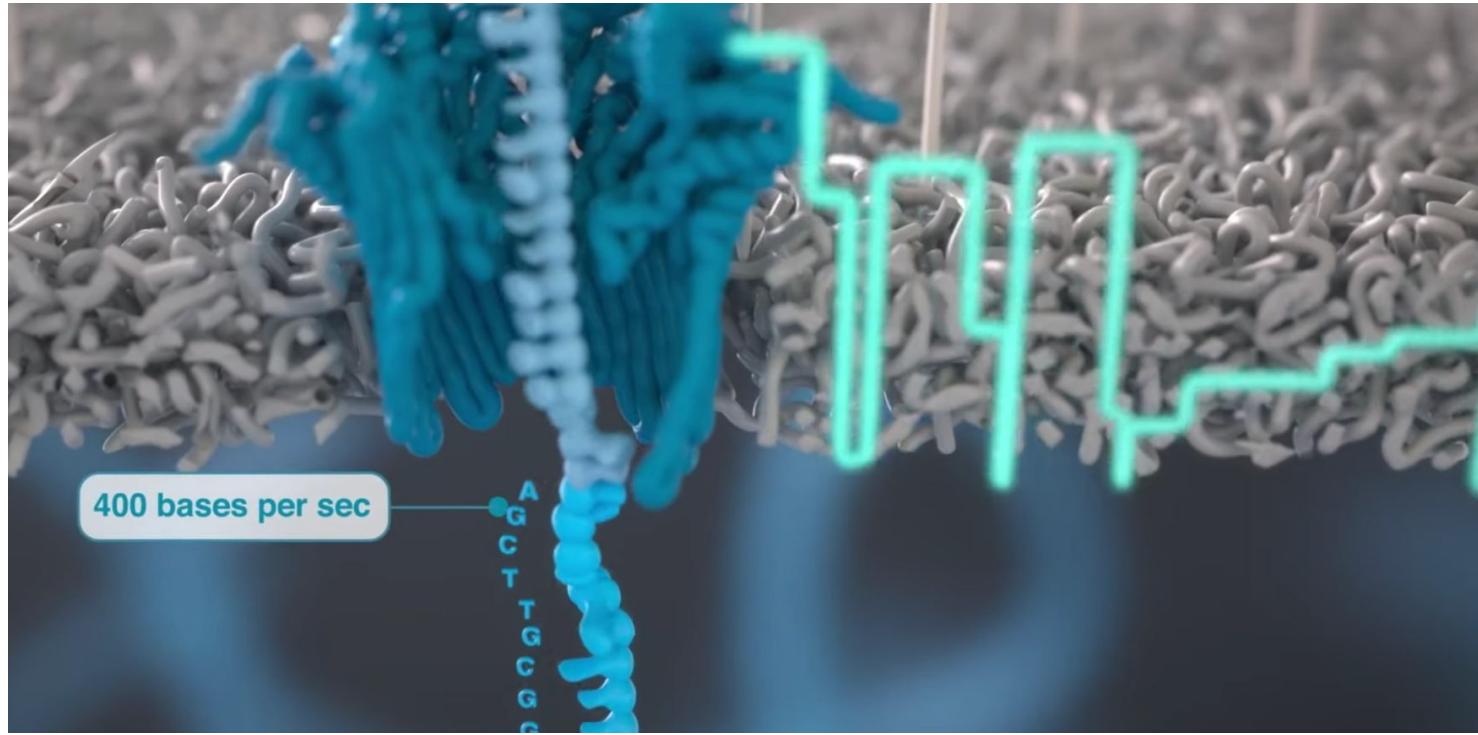
- DNA or RNA fragments pass through a nano-scale hole.
- The fluctuations in current during translocation are used.
- An enzyme motor controls the translocation of the DNA.

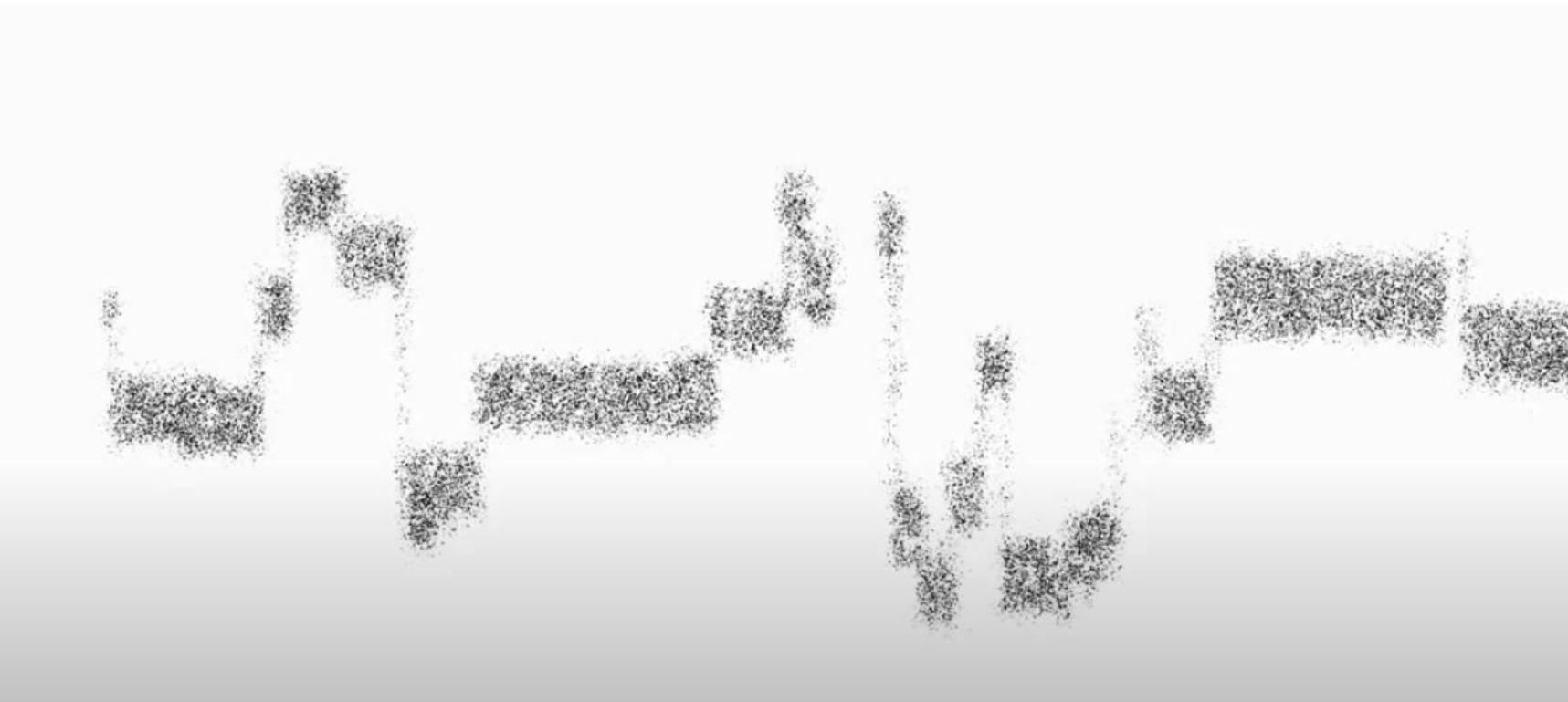




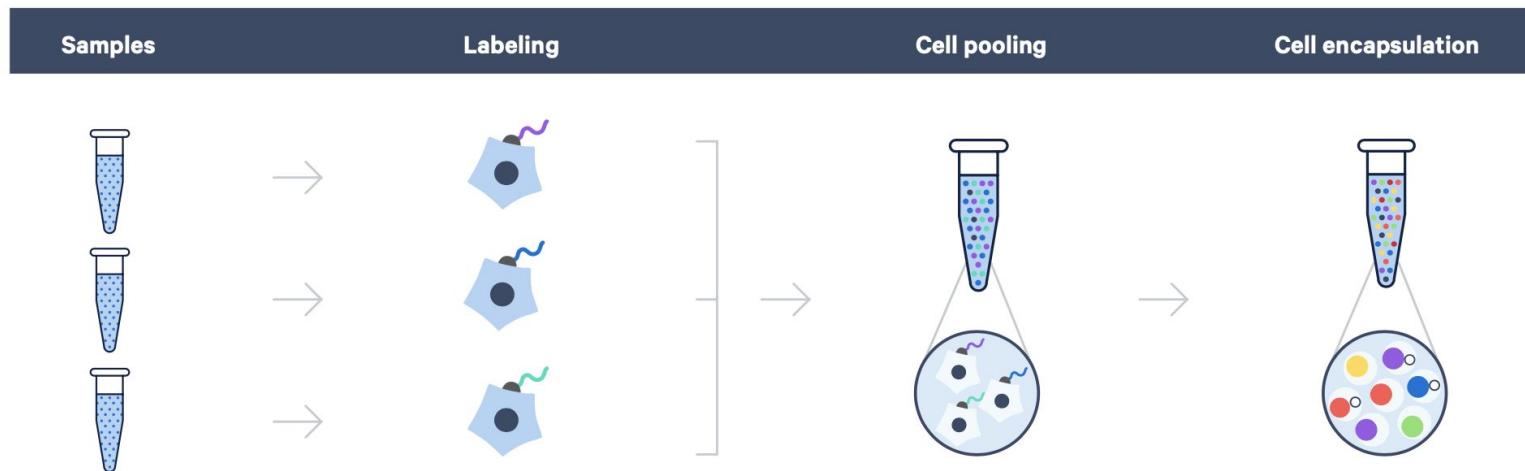




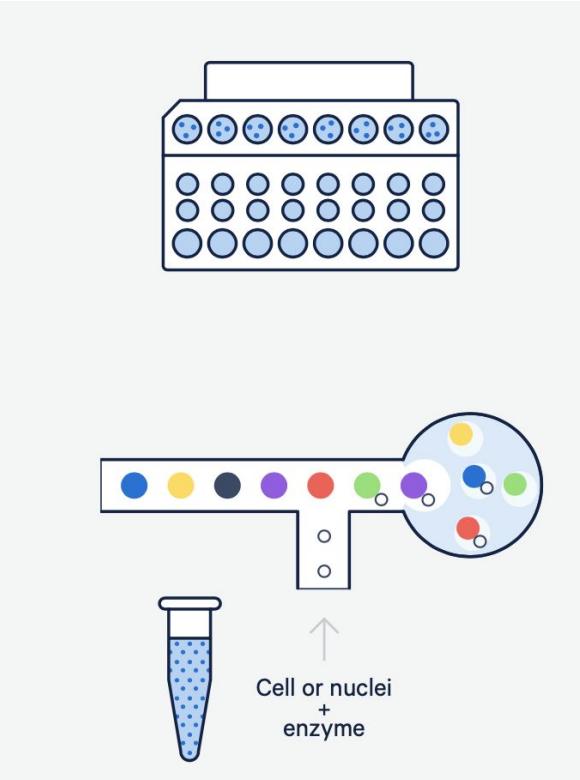


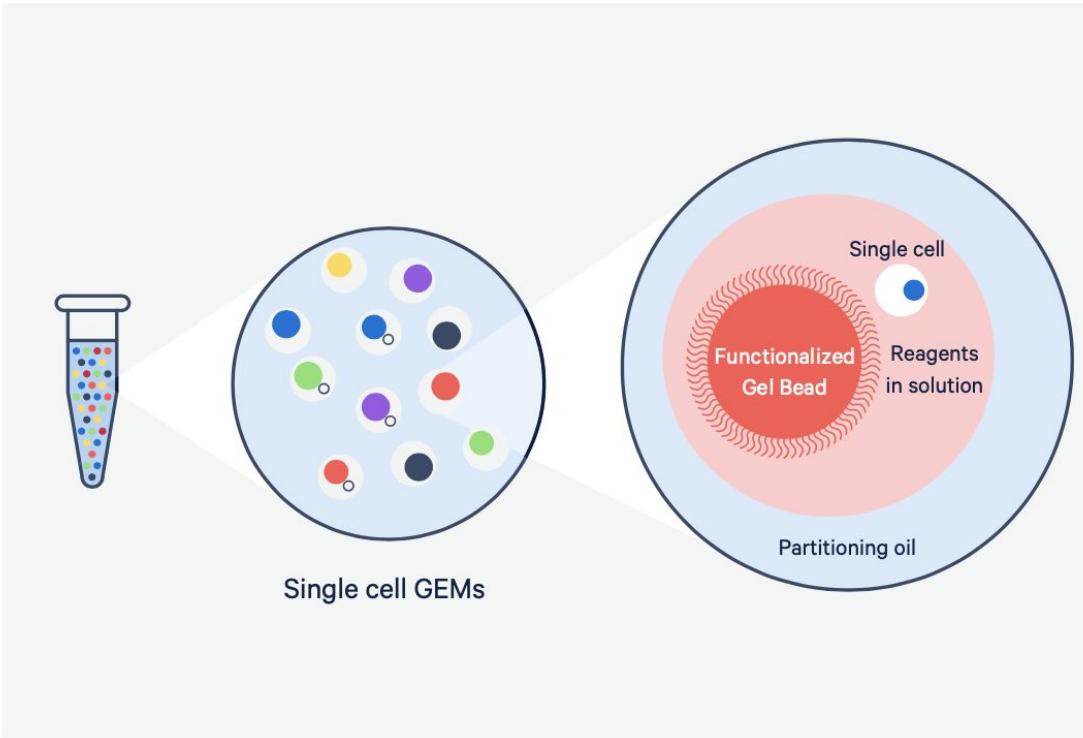


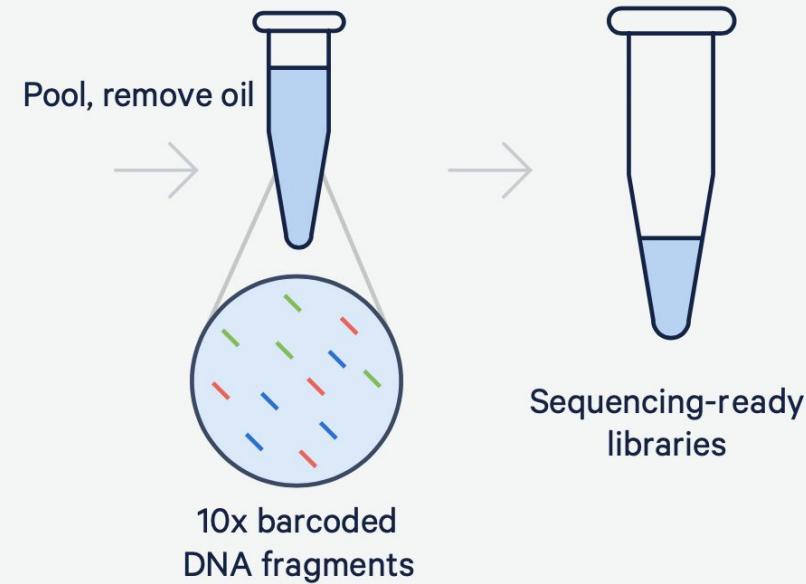
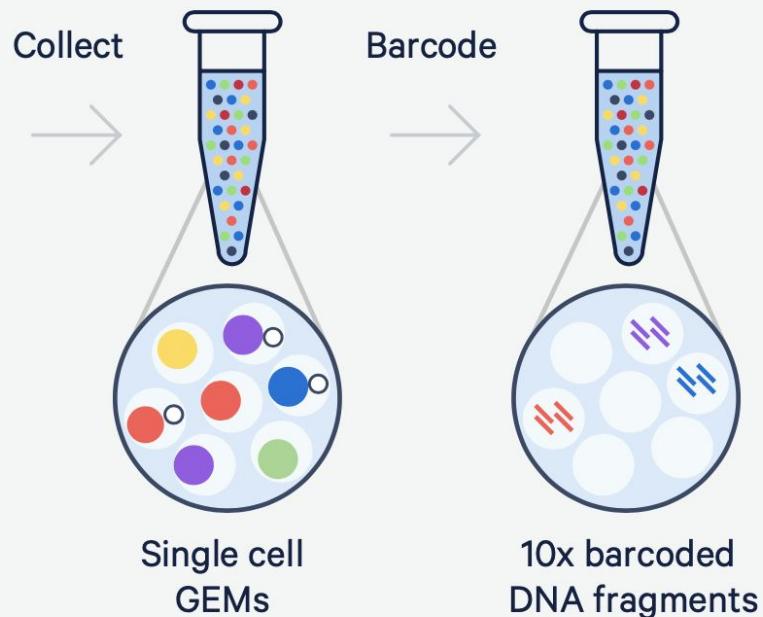
Single cell sequencing

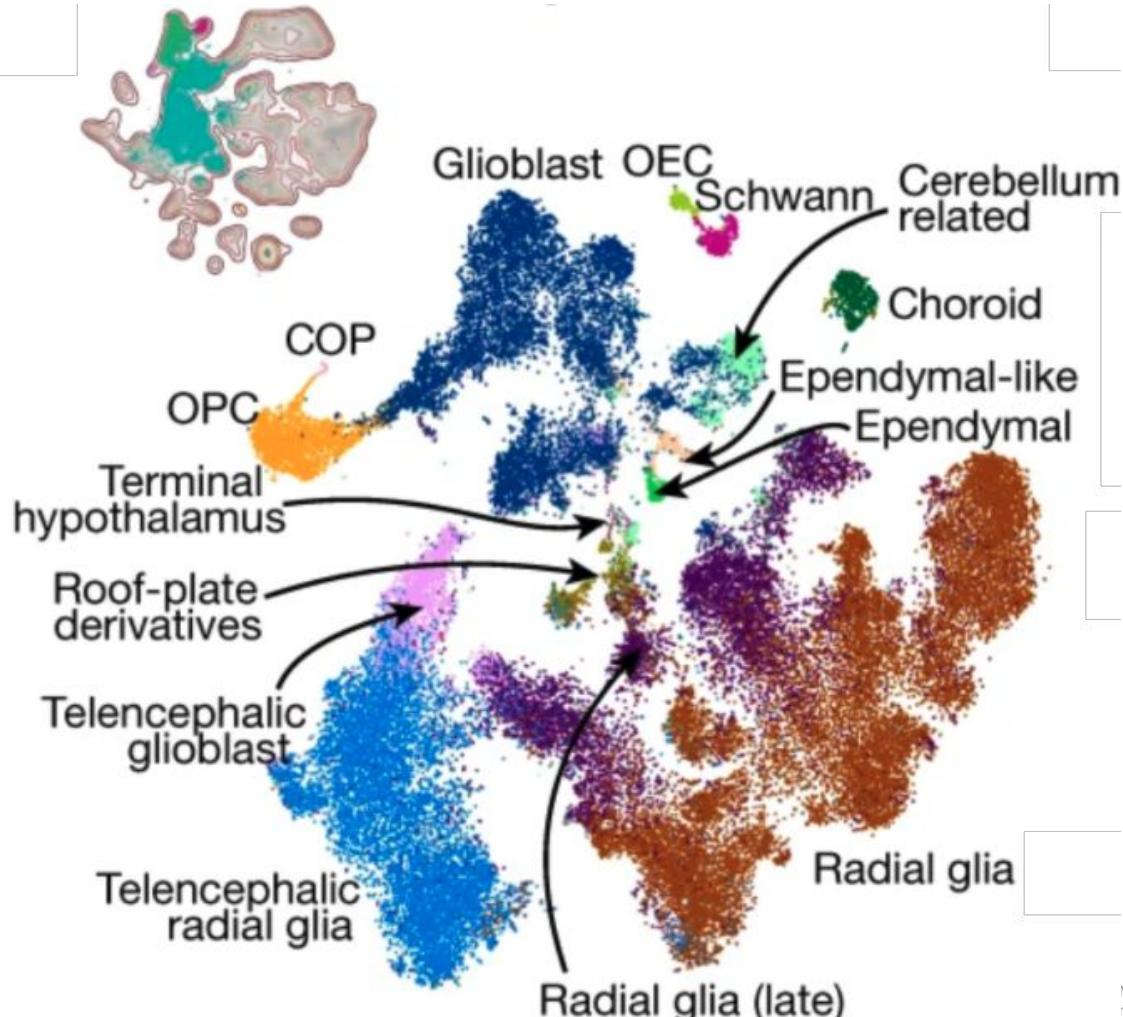


- Gel Beads, cells or nuclei, enzymes, and partitioning oil are loaded onto a chip.
- Barcoded gel beads are mixed with the cells or nuclei, enzymes, and partitioning oil.









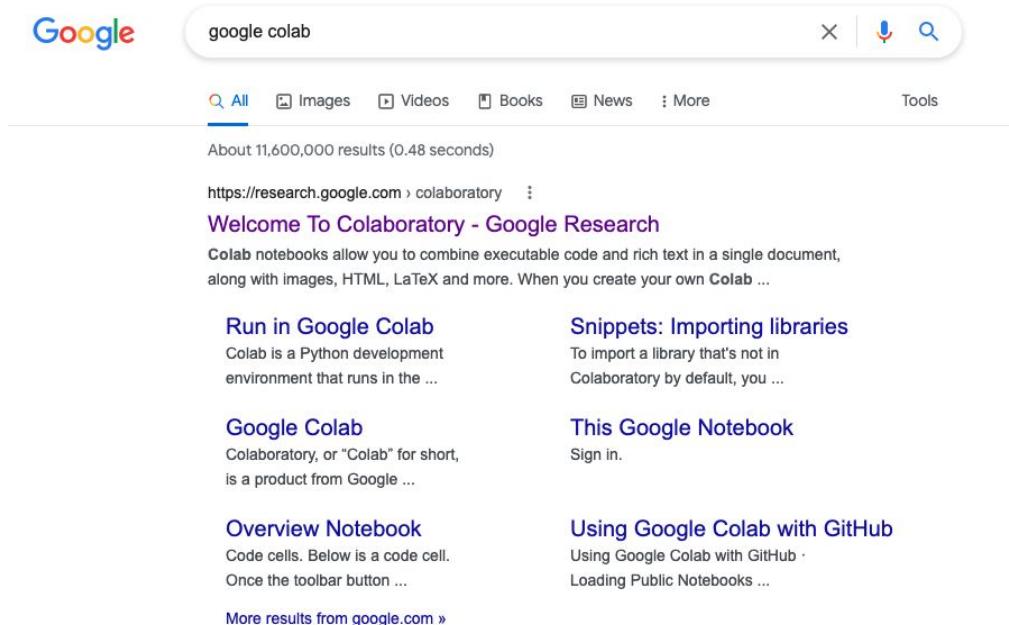
Other methods

- SOLiD
- Ion Torrent semiconductor sequencing
- 454 Life Sciences (pyrosequencing)

- Hi-C
- Optical mapping, BioNano
- Linked read, 10X genomics

To save time for the next section

Create a gmail account
(if you don't have)



A screenshot of a Google search results page. The search bar at the top contains the query "google colab". Below the search bar, there are several navigation links: "All" (which is underlined in blue), "Images", "Videos", "Books", "News", and "More". To the right of these links is a "Tools" button. The search results show approximately 11,600,000 results found in 0.48 seconds. The first result is a link to "https://research.google.com" titled "Welcome To Colaboratory - Google Research". The snippet for this result describes Colab notebooks as allowing executable code and rich text in a single document, along with images, HTML, LaTeX, and more. The second result is a link to "Run in Google Colab" with the snippet "Colab is a Python development environment that runs in the ...". The third result is a link to "Snippets: Importing libraries" with the snippet "To import a library that's not in Colaboratory by default, you ...". The fourth result is a link to "Google Colab" with the snippet "Colaboratory, or "Colab" for short, is a product from Google ...". The fifth result is a link to "This Google Notebook" with the snippet "Sign in." The sixth result is a link to "Overview Notebook" with the snippet "Code cells. Below is a code cell. Once the toolbar button ...". The seventh result is a link to "Using Google Colab with GitHub" with the snippet "Using Google Colab with GitHub · Loading Public Notebooks ...". At the bottom of the search results, there is a link "More results from google.com »".

The screenshot shows the Colab interface. At the top, there's a navigation bar with File, Edit, View, Insert, Runtime, Tools, and Help. Below it is a toolbar with Share, Connect, + Code, + Text, and Copy to Drive. A sidebar on the left titled 'Table of contents' lists sections like Getting started, Data science, Machine learning, More Resources, and Machine Learning Examples. The main area displays the 'What is Colaboratory?' notebook, which contains text about Colab and a code cell that prints '604800'.

Welcome To Colaboratory

File Edit View Insert Runtime Tools Help

Share Connect + Code + Text Copy to Drive

Table of contents

- Getting started
- Data science
- Machine learning
- More Resources
- Machine Learning Examples
- Section

What is Colaboratory?

Colaboratory, or "Colab" for short, allows you to write and execute Python in your browser, with zero configuration required. Free access to GPUs. Easy sharing.

Whether you're a student, a data scientist or an AI researcher, Colab can make your work easier. Watch [Introduction to Colab](#) to learn more, or just get started below!

Getting started

The document you are reading is not a static web page, but an interactive environment called a **Colab notebook** that lets you write and execute code.

For example, here is a **code cell** with a short Python script that computes a value, stores it in a variable, and prints the result:

```
[ ] seconds_in_a_day = 24 * 60 * 60  
seconds_in_a_day
```

604800

To execute the code in the above cell, select it with a click and then either press the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter". To edit the code, just click the cell and start editing.

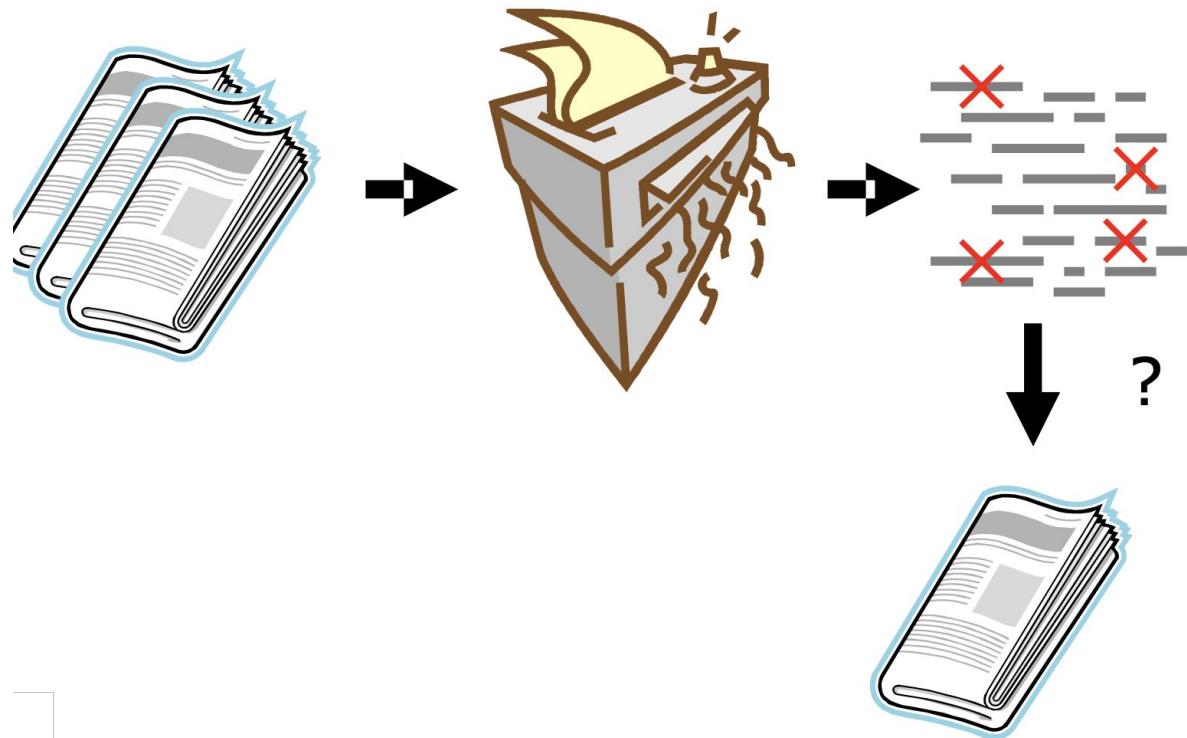
Variables that you define in one cell can later be used in other cells:

```
[ ] seconds_in_a_week = 7 * seconds_in_a_day  
seconds_in_a_week
```

Index

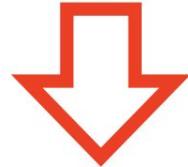
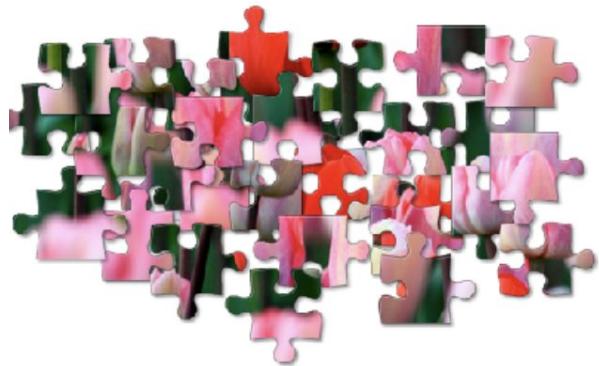
1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly

De novo Genome Assembly



Recovery of shredded newspaper

De novo= from the beginning



Overlap-layout-consensus approach

TACCTGGAACTGAAAGAAGGCTTACTGGAGCCGCTGGCAGTG

TGCGTGGGATCTCGCGAAATTCTTGCCGCA

GTGATGGTATGCGCACCTTGCCTGGGATCTC

CCGCTGGCAGTGACGGAACGGCTGCCATTATCTGGTGGTAGGTGATGGTATGCG

Overlap-layout-consensus approach

TACCTGGAACTGAAAGAAGGCTTACTGGAGCCGCTGGCAGTG

TGCGTGGGATCTCGGCGAAATTCTTGCCGCA

GTGATGGTATGCGCACCTTGCGTGGGATCTC

CCGCTGGCAGTGACGGAACGGCTGCCATTATCTGGTAGGTGATGGTATGCG

All versus all sequence similarity comparison

Overlap-layout-consensus approach

TACCTGGAACTGAAAGAAGGCTTACTGGAGCCGCTGGCAGTG

CCGCTGGCAGTGACGGAACGGCTGCCATTATCTCGTGGTAGGTGATGGTATGCG

GTGATGGTATGCGCACCTTGCGTGGGATCTC

TGCGTGGGATCTCGGCGAAATTCTTGCCGCA

Alignment of sequences with high sequence identity

Overlap-layout-consensus approach



Determination of consensus sequence

Overlap graph

Each node is a read

CTCGGCTCTAGCCCCTCATT

Draw edge A -> B when suffix of A overlaps prefix of B

CTCGGCTCTAGCCCCTCATT

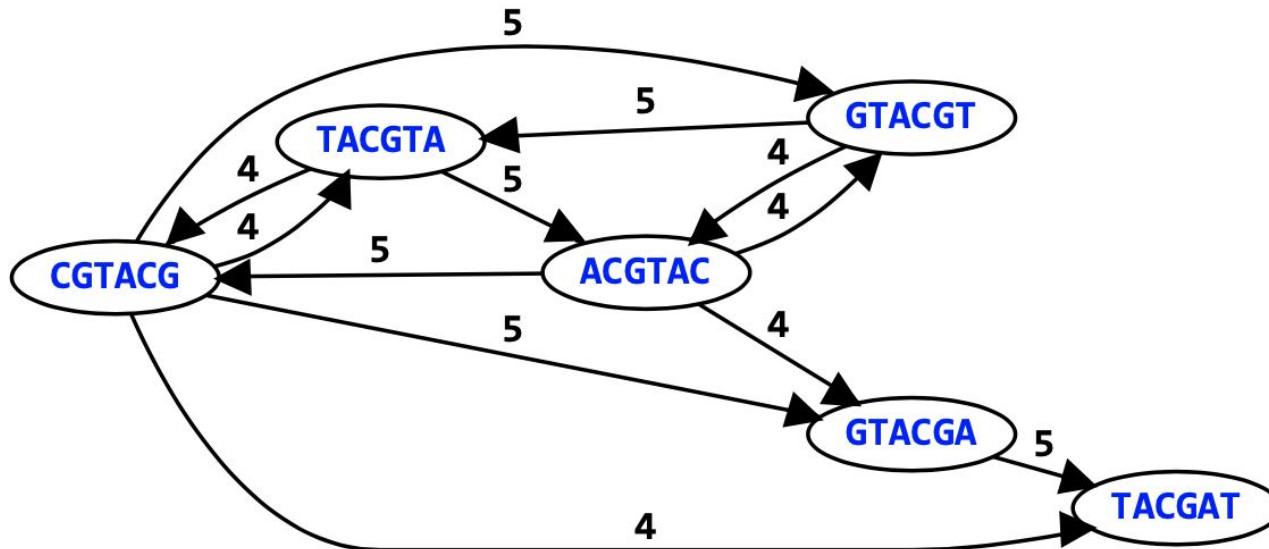


GGCTCTAGGCCCTCATT

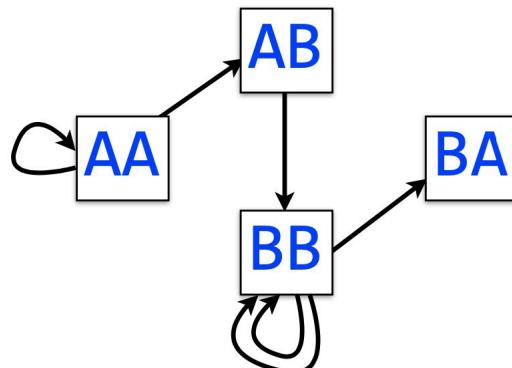
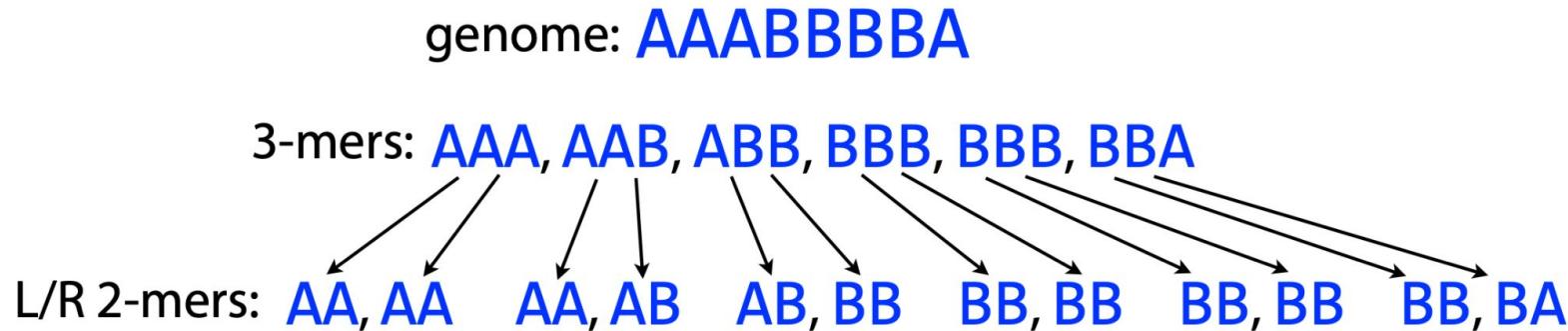
Overlap graph

Nodes: all 6-mers from **GTACGTACGAT**

Edges: overlaps of length ≥ 4



De Bruijn graph



One edge per k -mer
One node per distinct $k-1$ -mer

UNIT

Using the Velvet *de novo* Assembler for Short-Read Sequencing Technologies

Daniel R. Zerbino

First published: 01 September 2010 | <https://doi.org/10.1002/0471250953.bi1105s31> | Citations: 216

Genome Research

Method

Canu: scalable and accurate long-read assembly via adaptive *k*-mer weighting and repeat separation

Sergey Koren,^{1,5} Brian P. Walenz,^{1,5} Konstantin Berlin,² Jason R. Miller,³ Nicholas H. Bergman,⁴ and Adam M. Phillippy¹

¹Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland 20892, USA; ²Invincea Incorporated, Fairfax, Virginia 22030, USA; ³J. Craig Venter Institute, Rockville, Maryland 20850, USA; ⁴National Biodefense Analysis and Countermeasures Center, Frederick, Maryland 21702, USA

Long-read single-molecule sequencing has revolutionized *de novo* genome assembly and enabled the automated reconstruction of reference-quality genomes. However, given the relatively high error rates of such technologies, efficient and accurate assembly of large repeats and closely related haplotypes remains challenging. We address these issues with Canu, a successor of Celera Assembler that is specifically designed for noisy single-molecule sequences. Canu introduces support for nano-



brianwalenz Guard against OG == nullptr when called from layoutReads.

08b0347 on 18 Oct 10,284 commits



.github

Add an issue template.

5 years ago



documentation

Update quick-start.rst

4 months ago



scripts

Bump version number.

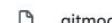
4 months ago



src

Guard against OG == nullptr when called from layoutReads.

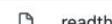
last month



.gitmodules

Use seqrequester submodule for making sqStoreDumpMetaData h...

2 years ago



.readthedocs.yml

Create .readthedocs.yml

2 years ago



README.license.GPL

Move old licensing/ files to root.

6 years ago



README.licenses

Add PacBio licenses.

5 years ago



README.md

Complain about .zip files missing submodules.

15 months ago

README.md

Canu

Canu is a fork of the [Celera Assembler](#), designed for high-noise single-molecule sequencing (such as the [PacBio RS II/Sequel](#) or [Oxford Nanopore MinION](#)).

Canu is a hierarchical assembly pipeline which runs in four steps:

About

A single molecule sequence assembler for genomes large and small.

[canu.readthedocs.io/](#)

[Readme](#)

Releases 15

Canu v2.2 Latest
on 27 Aug

+ 14 releases

Packages

No packages published

Contributors 20



+ 9 contributors

[Languages](#)

canu
latest

Search docs

Canu Quick Start
Canu FAQ
Canu Tutorial
Canu Pipeline
Canu Parameter Reference
Software Background



Read the Docs for Business: Private repos with Pull Request reviews. [Start your free trial today.](#)

Ad by EthicalAds · Monetize your site

Canu

Canu is a fork of the Celera Assembler designed for high-noise single-molecule sequencing (such as the PacBio RSII or Oxford Nanopore MinION).

Publications

Canu

Koren S, Walenz BP, Berlin K, Miller JR, Phillippy AM. [Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation](#). Genome Research. (2017).

TrioCanu

Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, Hiendleder S, Williams JL, Smith TPL, Phillippy AM. [De novo assembly of haplotype-resolved genomes with trio binning](#). Nature Biotechnology. (2018).

HiCanu

Nurk S, Walenz BP, Rhie A, Vollger MR, Logsdon GA, Grothe R, Miga KH, Eichler EE, Phillippy AM, Koren S. [HiCanu: accurate assembly of segmental duplications, satellites, and allelic variants from high-fidelity long reads](#). Genome Research. (2020).

Install

The easiest way to get started is to download a [release](#). If you encounter any issues, please report them using the [github issues](#) page.

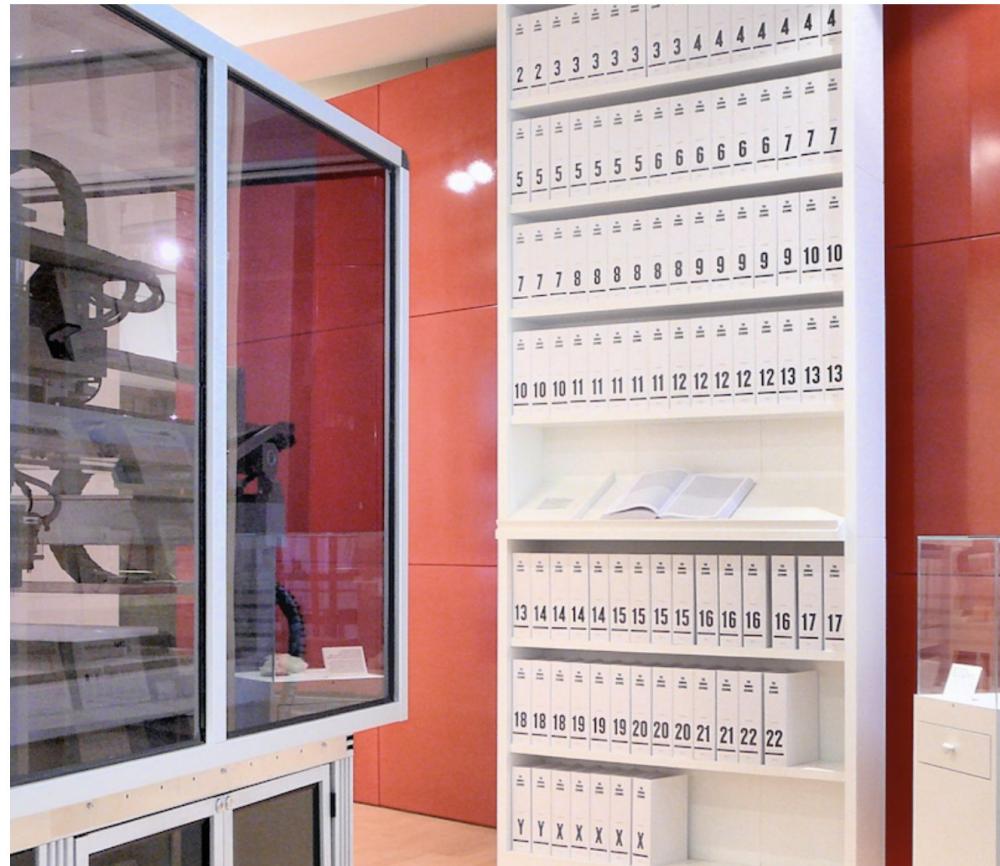
Alternatively, you can also build the latest unreleased from github:

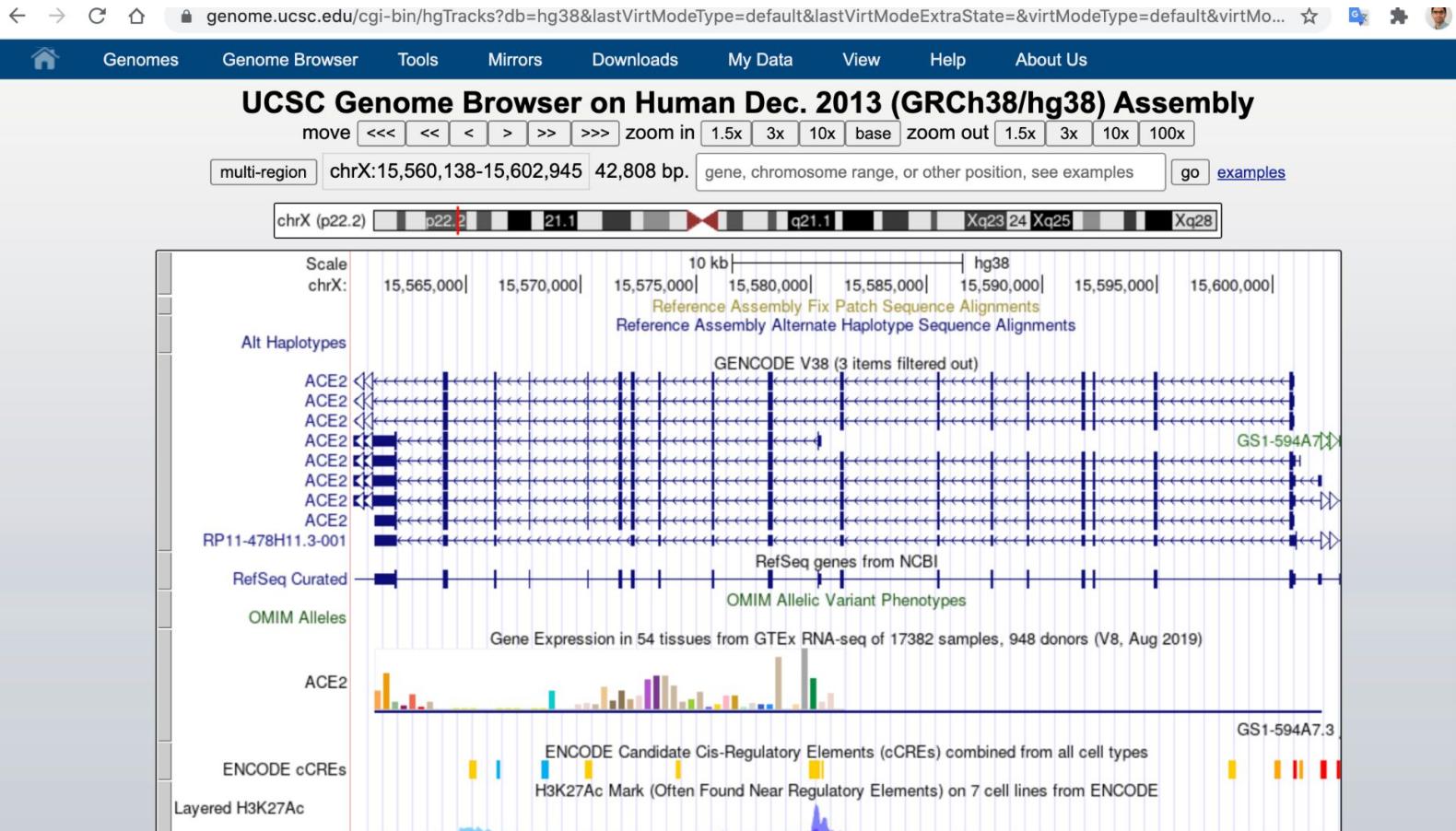
```
git clone https://github.com/marbl/canu.git
cd canu/src
make -j <number of threads>
```

<https://canu.readthedocs.io/en/latest/>

Reference Genome

Using illumina+pacbio+HiC

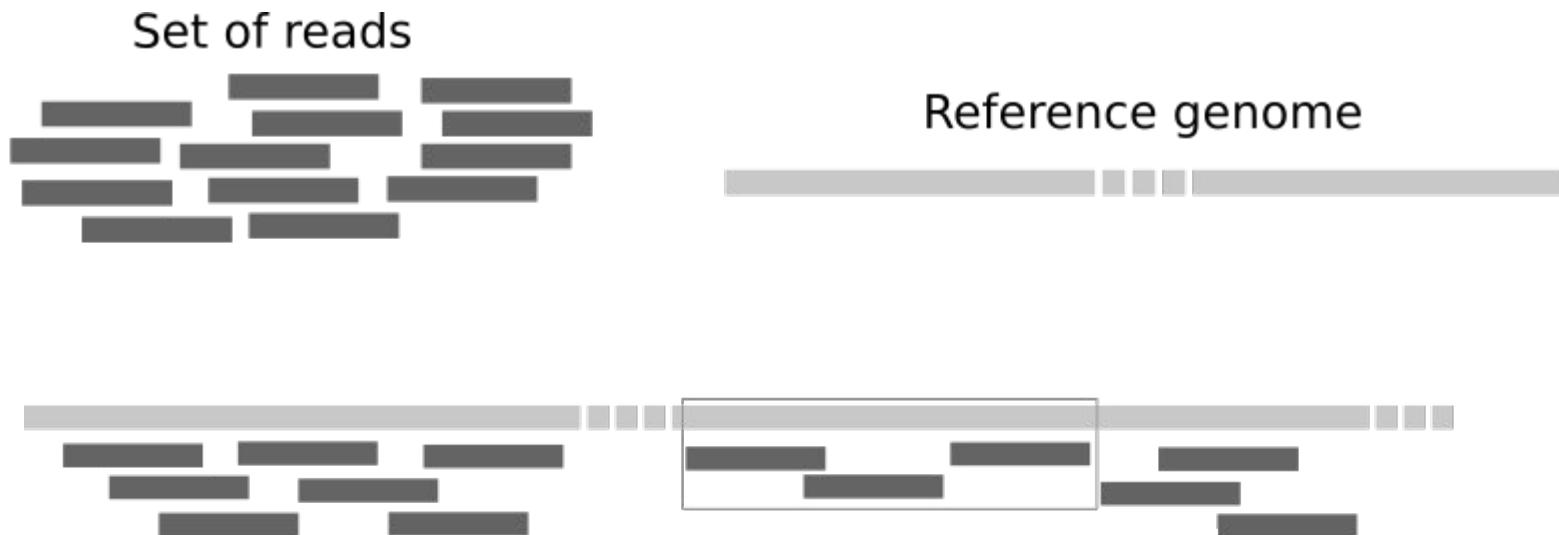


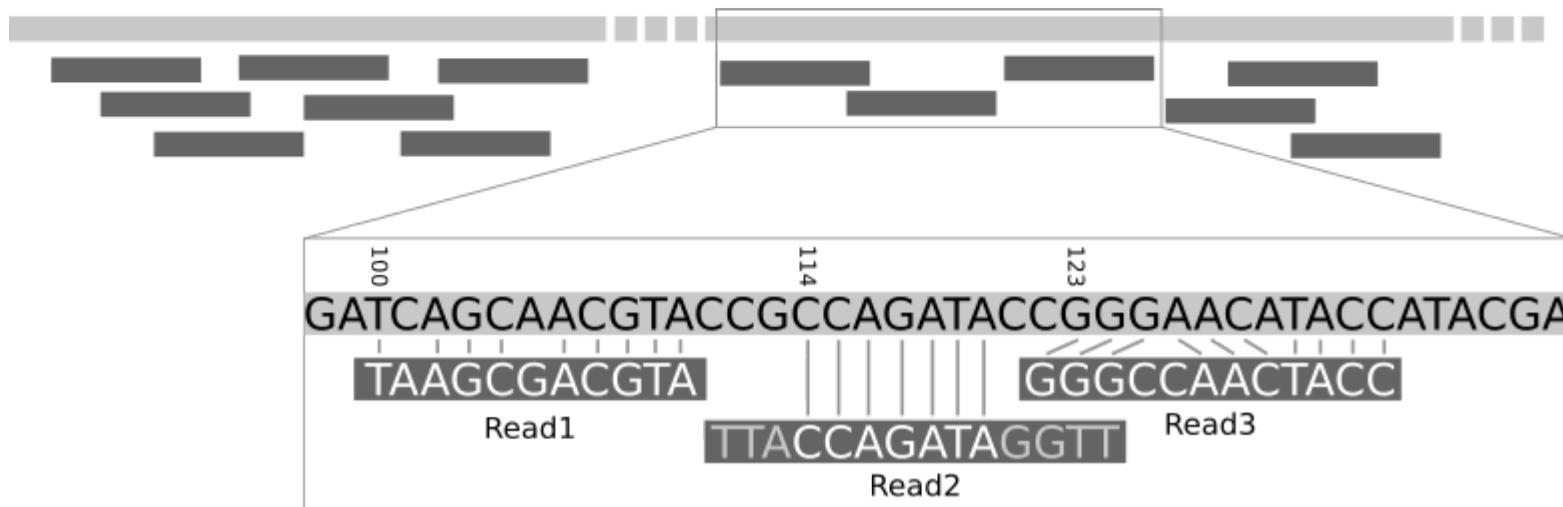


Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment

Read alignment/mapping





Sequence analysis

Fast and accurate short read alignment with Burrows–Wheeler transform

Heng Li and Richard Durbin*

Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK

Received on February 20, 2009; revised on May 6, 2009; accepted on May 12, 2009

Advance Access publication May 18, 2009

Associate Editor: John Quackenbush

ABSTRACT

Motivation: The enormous amount of short reads generated by the new DNA sequencing technologies call for the development of fast and accurate read alignment programs. A first generation of hash table-based methods has been developed, including MAQ, which is accurate, feature rich and fast enough to align short reads from a single individual. However, MAQ does not support gapped alignment for single-end reads, which makes it unsuitable for alignment of longer reads where indels may occur frequently. The speed of MAQ is also a concern when the alignment is scaled up to the resequencing of hundreds of individuals.

Results: We implemented Burrows–Wheeler Alignment tool (BWA), a new read alignment package that is based on backward search with Burrows–Wheeler Transform (BWT), to efficiently align short

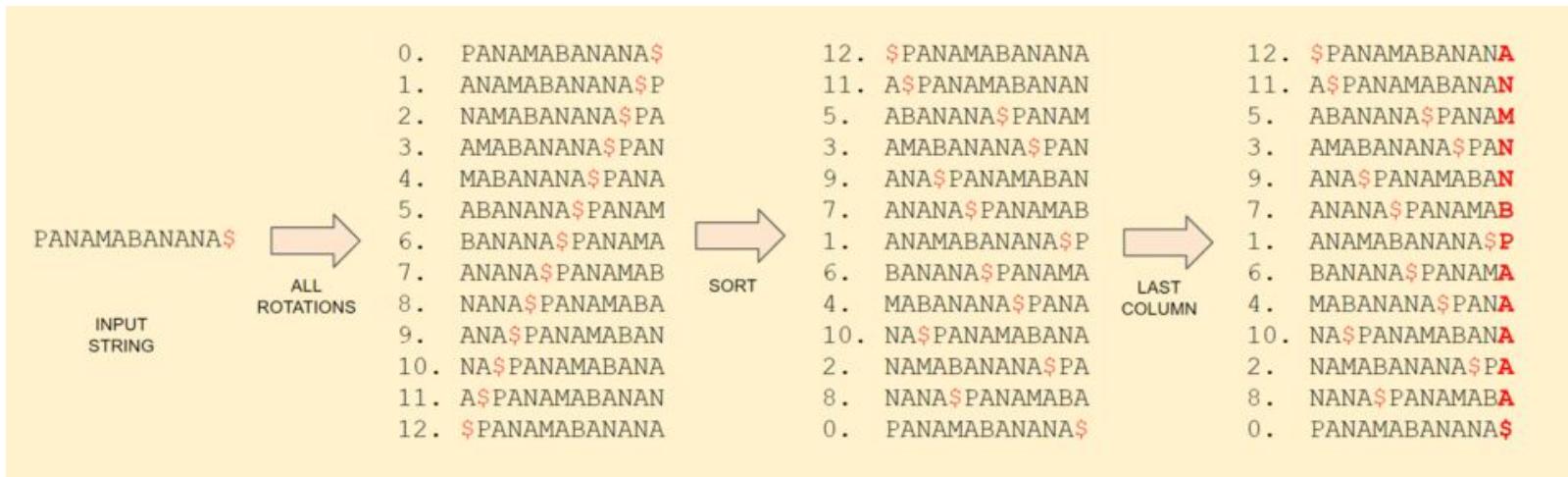
of scanning the whole genome when few reads are aligned. The second category of software, including SOAPv1 (Li *et al.*, 2008b), PASS (Campagna *et al.*, 2009), MOM (Eaves and Gao, 2009), ProbeMatch (Jung Kim *et al.*, 2009), NovoAlign (<http://www.novocraft.com>), ReSEQ (<http://code.google.com/p/re-seq>), Mosaik (<http://bioinformatics.bc.edu/marthlab/Mosaik>) and BFFAST (<http://genome.ucla.edu/bfast>), hash the genome. These programs can be easily parallelized with multi-threading, but they usually require large memory to build an index for the human genome. In addition, the iterative strategy frequently introduced by these software may make their speed sensitive to the sequencing error rate. The third category includes slider (Malhis *et al.*, 2009) which does alignment by merge-sorting the reference subsequences and read sequences.

BWA: <http://bio-bwa.sourceforge.net/>

minimap2: <https://github.com/lh3/minimap2>

Burrows–Wheeler transform

A Block-sorting Lossless Compression Algorithm. Used in bzip



PANAMABANANA" -> ANMNNBPAAAAA\$ -> ANM2NBP5A\$

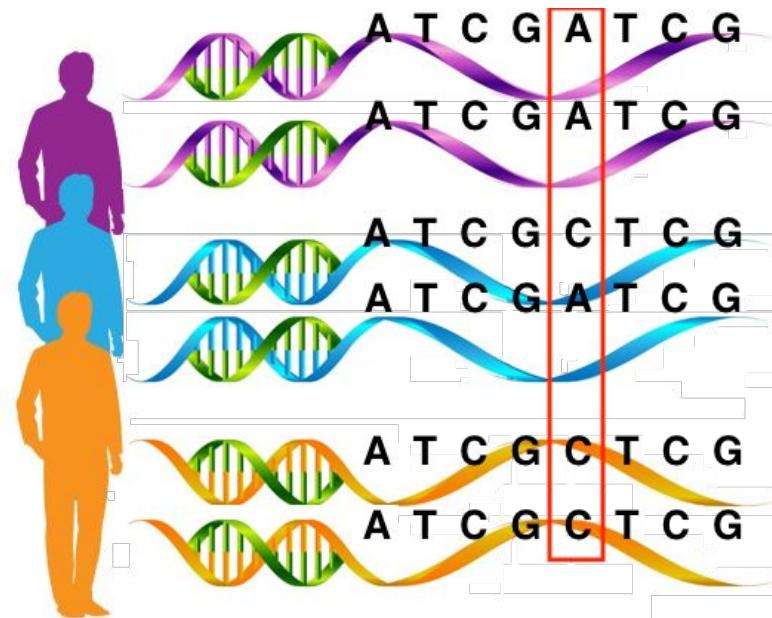
Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling

Genetic variation

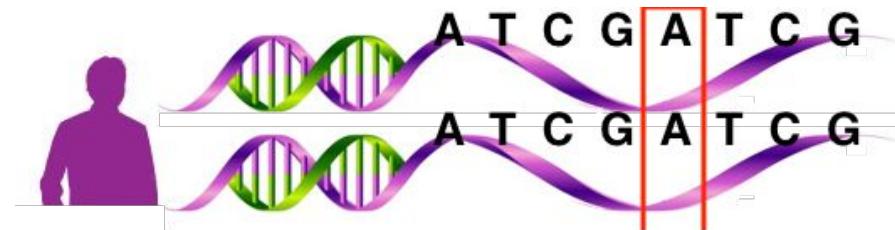
the difference in DNA among individuals.

polymorphism



Aligned reads person1

Ref	ATCG A TCG
Read1	ATCG A
Read2	ATCG A
Read3	G ATCG
Read4	G ATCG



Aligned reads of person2

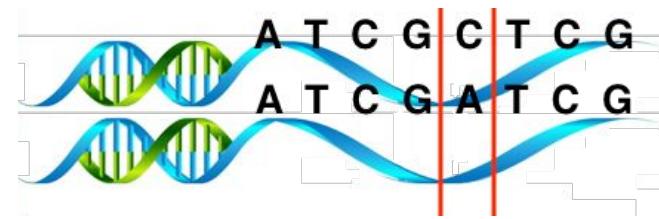
Ref ATCG**A**TCG

Read1 ATCG**C**

Read2 ATCG**A**

Read3 **G**CTCG

Read4 **G**ATCG



Variant calling & Genotyping

Individual 1

Ref ATCG**A**TCG

Read1 ATCG**A**

Read2 ATCG**A**

Read3 **G**ATCG

Read4 **G**ATCG

position 5: A/A

Individual 2

Ref ATCG**A**TCG

Read1 ATCG**C**

Read2 ATCG**A**

Read3 **G**CTCG

Read4 **G**ATCG

position 5: A/C

github.com/freebayes/freebayes

README.md

freebayes, a haplotype-based variant detector

user manual and guide

Overview

freebayes is a Bayesian genetic variant detector designed to find small polymorphisms, specifically SNPs (single-nucleotide polymorphisms), indels (insertions and deletions), MNPs (multi-nucleotide polymorphisms), and complex events (composite insertion and substitution events) smaller than the length of a short-read sequencing alignment.

freebayes is haplotype-based, in the sense that it calls variants based on the literal sequences of reads aligned to a particular target, not their precise alignment. This model is a straightforward generalization of previous ones (e.g. PolyBayes, samtools, GATK) which detect or report variants based on alignments. This method avoids one of the core problems with alignment-based variant detection--- that identical sequences may have multiple possible alignments:



GATK workflows

Official GATK workflows published by the Broad Institute's Data Sciences Platform

Cambridge, MA USA  <https://gatk.broadinstitute.org>  @gatk_dev

Overview  29   

Pinned

 [gatk4-germline-snps-indels](#)

Public

Workflows for germline short variant discovery with GATK4

 wdl  99  49

 [gatk4-data-processing](#)

Public

Workflows for processing high-throughput sequencing data for variant discovery with GATK4 and related tools

 wdl  100  64

 [gatk4-rnaseq-germline-snps-indels](#)

Public

Workflows for processing RNA data for germline short variant discovery with GATK v4 and related tools

 wdl  31  20

 [seq-format-conversion](#)

Public

Workflows for converting between sequence data formats

 wdl  22  24

 [seq-format-validation](#)

Public

Workflows for validating sequence data formats

 wdl  2  5

 [gatk4-basic-joint-genotyping](#)

Public

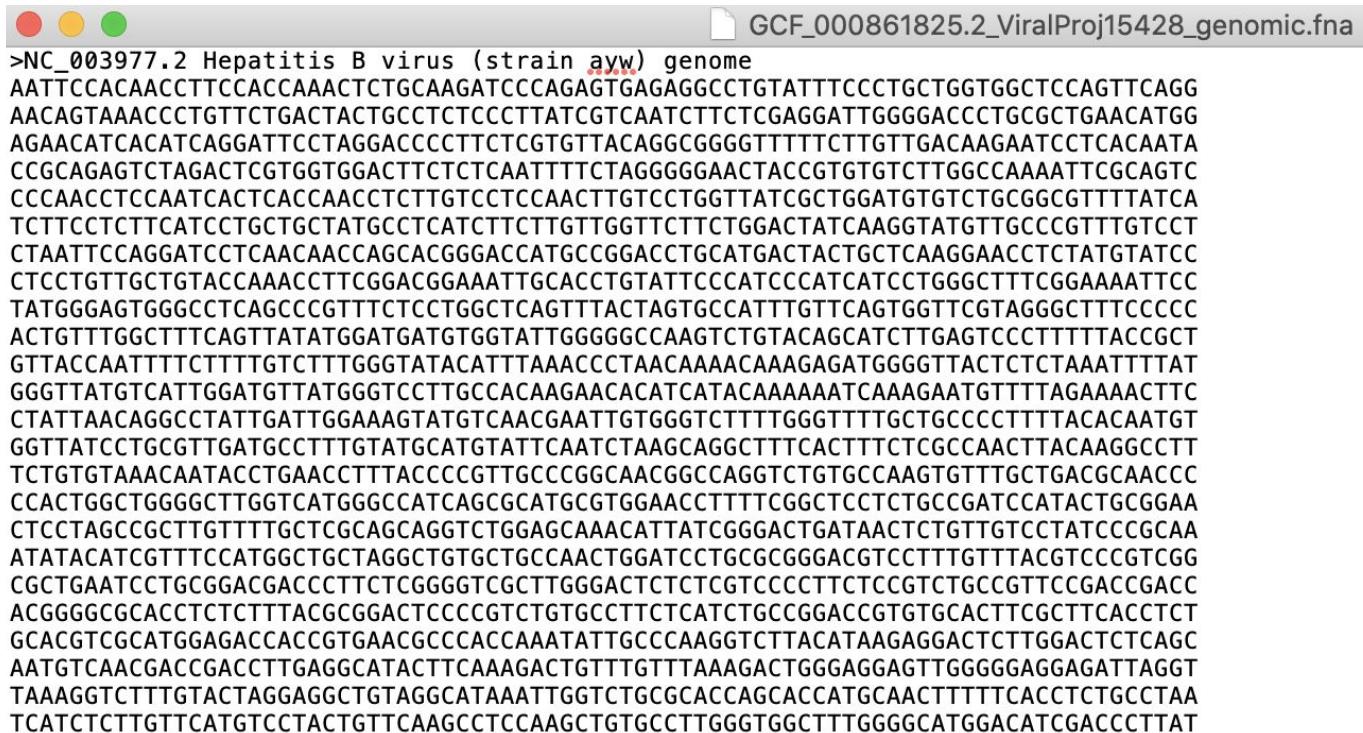
Basic joint genotyping with GATK4. NOT Best Practices, only for teaching/demo purposes.

 wdl  3

Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling
4. File formats
 - o fasta format

fasta Format



The screenshot shows a terminal window with a gray header bar containing three colored circles (red, yellow, green) and the text "GCF_000861825.2_ViralProj15428_genomic.fna". The main area of the terminal displays a FASTA sequence. The header line starts with ">NC_003977.2 Hepatitis B virus (strain ayw) genome". The sequence itself is a long string of DNA bases (A, T, C, G) representing the genome.

```
>NC_003977.2 Hepatitis B virus (strain ayw) genome
AATTCCACAAACCTTCCACCAAACCTGTCAAGATCCCAGAGTGAGAGGCCTGTATTTCCCTGCTGGTGGCTCCAGTTAGG
AACAGTAACACCTGTTCTGACTACTGCCTCTCCCTTATCGTAATCTTCGAGGATTGGGACCCCTGCGCTGAACATGG
AGAACATCACATCAGGATTCTAGGACCCCTTCTCGTTACAGGC GGTTAGGGGAACTACCGTGTCTGGCCAAAATTGCAGTC
CCGCAGAGTCTAGACTCGTGGGACTTCCTCTCAATTCTAGGGGAACTACCGTGTCTGGCCAAAATTGCAGTC
CCCAACCTCCAATCACTCACCAACCTTCTGCTCTCCAACTTGTCTGGTTATCGCTGGATGTGTCTGGCGTTTATCA
TCTTCCTCTTATCCTGCTCATGCCATCTTCTGGTTCTCTGGACTATCAAGGTATGTTGCCCTATGTATCC
CTAATTCAGGATCCTCAACAAACAGCACGGGACCATGCCGGA CTCATGACTACTGCTCAAGGAACCTCTATGTATCC
CTCCTGTTGCTGTACCAAACCTTCGGACGGAAATTGCCACCTGTATTCCCATCATCCTGGCCTTCGGAAAATTCC
TATGGGAGTGGGCTTCAGCCCCGTTCTCTGGCTCAGTTACTAGTGCCATTGTTCA GTGGTCTGGTAGGGCTTCCCCC
ACTGTTGGCTTCAGTTATGGATGATGTGGTATTGGGGCAAGTCTGTACAGCATCTTGAGTCCCTTTTACCGCT
GTTACCAATTCTTTGTCTGGGTATACATTAAACCTAACAAAACAAAGAGATGGGGTTACTCTCAAATTTTAT
GGGTTATGTCATTGGATGTTATGGGTCTTGCCACAAGAACATCATACAAAAAAATCAAAGAATGTTTAGAAAACCTC
CTATTACAGGCCATTGATTGGAAAGTATGTCACAGAATTGTTGGCTTTGGCTGGCTGCCAACCTTACACAATGT
GGTTATCCTGCGTTGATGCCCTTGATGCTATGCTATTCAATCTAAGCAGGCTTCACTTCTGCCAACCTACAAGGCCCT
TCTGTGTAACAAATACCTGAACCTTACCCGTTGCCGGCAACGGCCAGGTCTGTGCCAACGTGTTGCTGACGCAACCC
CCACTGGCTGGGCTTGGTCATGGGCATCAGCGCATGCGTGGAACCTTCCGCTCTGCCGATCCACTGCGGAA
CTCCTAGCCGCTTGTGCTCGCAGCAGGTCTGGAGCAAACATTATGGGACTGATAACTCTGTTGCTCATCCGCAA
ATATAACATCGTTCCATGGCTGCTAGGCTGTGCTGCCAACACTGGATCCTGCGGGGACGTCTTGTGCTGCCGACCGACC
CGCTGAATCCTGCGGACGACCCCTCTGGGGTCGCTGGGACTCTCTGTCCTCTCATCTGCCGGACCGTGTGCACTCGCTTACCTCT
ACGGGGCGCACCTCTTACCGGGACTCCCGTCTGCGCTTCTCATCTGCCGGACCGTGTGCACTCGCTTACCTCT
GCACGTCGATGGAGACCAACCGTGAACGCCACCAAATATTGCCAAGGTCTTACATAAGAGGACTCTGGACTCTCAGC
AATGTCAACGACCGACCTTGAGGCATACTTCAAGACTGTTGTTAAAGACTGGGAGGAGTTGGGGAGGAGATTAGGT
TAAAGGTCTTGACTAGGAGGCTGTAGGCATAAATTGGCTGCGCACCAGCACCAGTGAACCTTTCACCTCTGCCTAA
TCATCTTGTTCATGTCCTACTGTTCAAGCCTCAAGCTGTGCCCTGGTGGCTTGGGATGGACATGACCCCTAT
```

fasta Format

- DNA/RNA/Protein sequence
- Header line, begins with >
- Each line 50/60 characters
- Extension: .fa .fna .fasta

fastq format

● ● ● sra_data.fastq.txt

```
@SRR16356243.194720.1 194720 length=151
GGCAGCAAAGTTTATTGTAAAATAAGAGATCGATATAAAATGGGATATAAAAGGGAGAAGGGAGGGGAAGGGTGGGTGAAAATGCAG
ATGTGCTTGCAGAATGTAAGATGTTGACCTCCAGCTGGACGTGGCTCAATTGT
+SRR16356243.194720.1 194720 length=151
AAFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF/
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFAFFFFFFFFFFFFFFF
@SRR16356243.194720.2 194720 length=151
CAATTGGAGGCCACCACGTCCAGCTGAAAGGTCAACATCTTACATTCTGCAAGCACATCTGCATTTCACCCACCCCTCCCTCCTT
CTCCCTTTATATCCATTTCATCGATCTCTTACAAATAAAACTTGCTGGCA
+SRR16356243.194720.2 194720 length=151
FFFFFFFAFFFFF/FFFFFFFAFFFFFFFFFFFAFFFFFFFFFFFFFFF
FFFFFFFAFFFAAAFFFFFFFFFF=AFF/FFFFFFFFF/FFF=FFF=F/FAFFFFFF6FF/FF
@SRR16356243.247542.1 247542 length=151
GTTGACTTAGTATGAATGTGGTACGTTGGAAGCAAATGTGCTTCACTTATCATGAAAAAGTCTGCAAGTGCTCTGCACGTCCAG
GGAAATGATCCTACCCCTAACATCTCAGCTAAAGGGACCTTGCTTTCAAGTGA
+SRR16356243.247542.1 247542 length=151
FAFF/FA=FAFF/FFFF/FFFA//F/A/FF/FFFFAAA///FF/FFFF//A///FFFFFA==//FF/F6FFA//A/FFF//
=F//F/AFA//F//F=/FAA/AF//A/F/FF/AFFF//FFF/F//FFFF/F/////
@SRR16356243.247542.2 247542 length=151
ATTGGACACTGAAAAGAGAAAAGGTTCCCTGTGAGCTGAAGAGTGAGTTAGGATCATTCCAAGGGACGTGCAAGGCACTTGCAGAC
TTTATCATGAAATAGTTAACGACACATTGATTCAACGAAACCACATTCAAACAAAGTCAA
+SRR16356243.247542.2 247542 length=151
//F/F/FF/=FAFFF/FAFF/F/FF//F//FF/FFFA//FA/F//A/F//FF//AFAF//F//AF6//FAFFFFAA/F//FF//F
FF//F/F/FF/A/F/FF//AA/A/FFA/FA//FFFFF/A/AA///F/A/FFFF/F/FF///=A
```

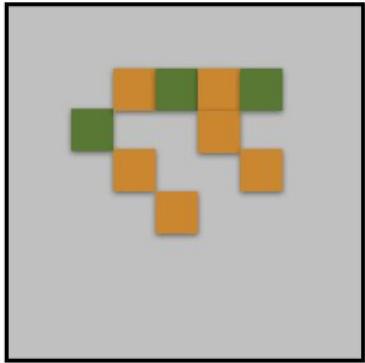
Storing sequencing reads

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

Bases and qualities line up:

AGCTCTGGTGACCCATGGGCAGCTGCTAGGGA
||||| | | | | | | | | | | | | | | | | |
HHHHHHHHHHHHHHHHHGCGC5FEFFF GHHHHHH

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



Call: orange (C)

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

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

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

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

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

ASCII table

phred33ToQ(qual): $\text{ord}(\text{qual}) - 33$

QtoPhred33(Q): $\text{chr}(\text{int}(\text{round}(Q)) + 33)$

43	+	75	K	107	k
44	,	76	L	108	l
45	-	77	M	109	m
46	.	78	N	110	n
47	/	79	O	111	o
48	0	80	P	112	p
49	1	81	Q	113	q
50	2	82	R	114	r
51	3	83	S	115	s
52	4	84	T	116	t
53	5	85	U	117	u
54	6	86	V	118	v
55	7	87	W	119	w
56	8	88	X	120	x
57	9	89	Y	121	y
58	:	90	Z	122	z

Sequence Alignment Map (SAM)

- for storing reads aligned to reference
- Tab-delimited text
- header section: lines start with @

Read1 163 chr20 7 20 8M2I4M1D3M ..

- read name
- Flag
- chr
- Position of start alignment on the reference sequence
- Mapping quality
- CIGAR

BAM: Binary representation of SAM; compressed version.

VCF (Variant Call Format)

```
##fileformat=VCFv4.3  
##fileDate=20090805
```

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001
20	14370	rs6054257	G	A	29	PASS	AF=0.5	GT	0 0
20	17330	.	T	A	3	q10	AF=0.017	GT	0/1

Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling
4. File formats
 - o fasta format
 - o fastq format
 - o sam format
 - o vcf format
5. Basic linux

OS: operating system

A system software that manages computer hardware, software resources, and provides common services for computer programs.

For computers:

- Windows
- Linux
- macOS

Why linux?

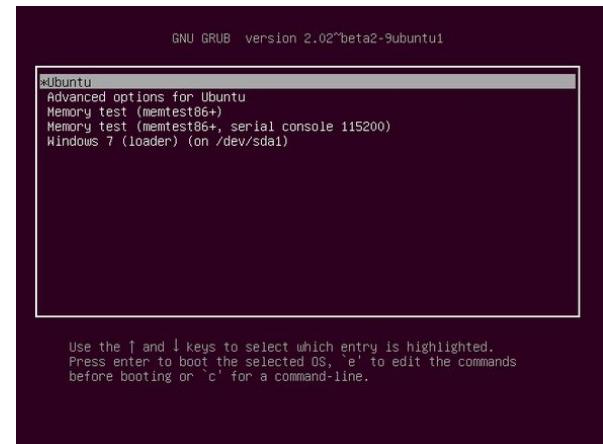
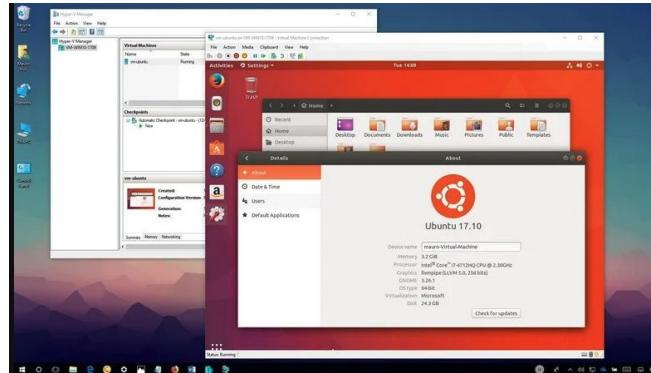
- many bioinformatics software developed only for Linux.
 - Graphical or web user interfaces exists but often struggle to provide flexibility
- free (!)
- performance
- security (virus)

Linux on a computer with windows

1- connect to a Linux machine (server) using PuTTy

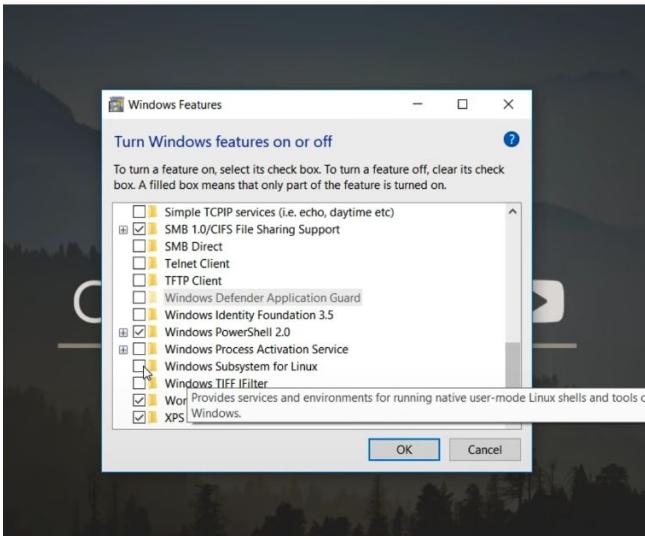
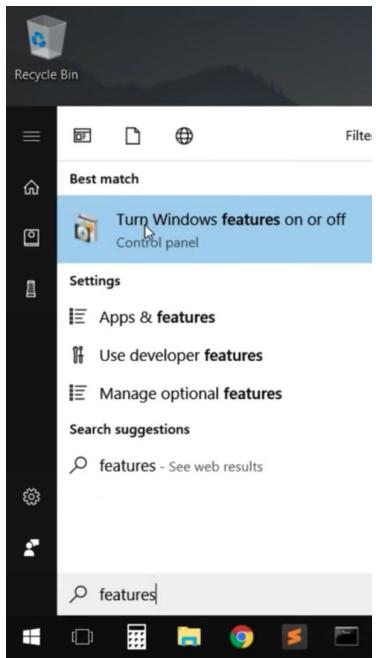
2- install on a drive (D/E not C) [danger]

3- use VirtualBox or VMware



4-install bash on windows

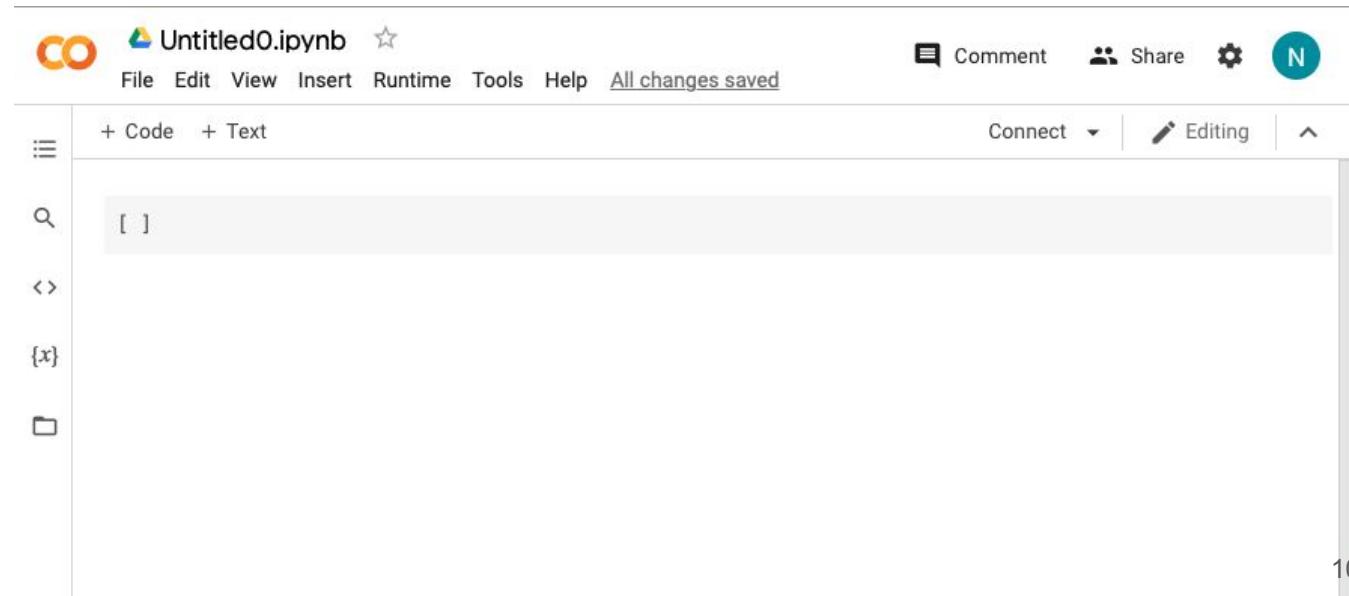
<https://www.laptopmag.com/articles/use-bash-shell-windows-10>



A screenshot of the Microsoft Store page titled "Run Linux on Windows". The page features a large blue header with the text "Run Linux on Windows" and "Install and run Ubuntu, openSUSE, SLES, and Fedora side-by-side with the Windows Subsystem for Linux (WSL)". Below the header are three colored cards: red for Ubuntu, green for openSUSE Leap 42, and dark blue for SUSE Linux Enterprise 99. Each card displays its respective logo.

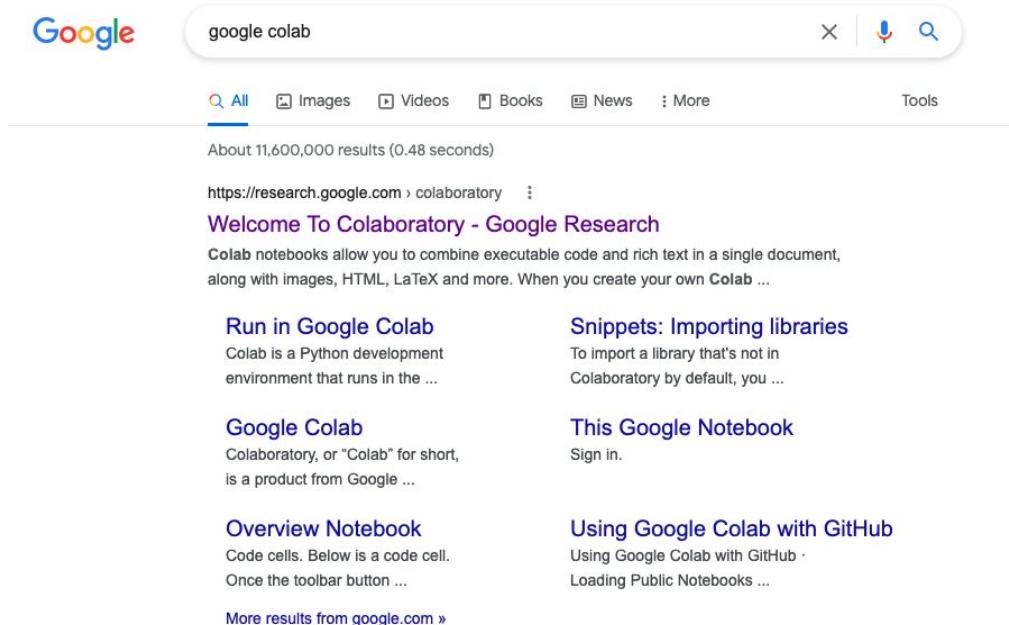
5- use online Linux terminal

- best is <https://research.google.com/colaboratory/>
- <https://copy.sh/v86/?profile=buildroot>
- <https://bellard.org/jslinux/vm.html?url=alpine-x86.cfg&mem=192>
- <http://cb.vu/>



To save time for the next section

Create a gmail account
(if you don't have)



A screenshot of a Google search results page. The search bar at the top contains the query "google colab". Below the search bar, there are several navigation links: All (highlighted), Images, Videos, Books, News, More, and Tools. A status message indicates "About 11,600,000 results (0.48 seconds)". The first result is a link to "https://research.google.com" titled "Welcome To Colaboratory - Google Research". The snippet for this result describes Colab notebooks as allowing executable code and rich text in a single document, along with images, HTML, LaTeX, and more. The second result is a link to "Run in Google Colab" with the snippet "Colab is a Python development environment that runs in the ...". The third result is a link to "Snippets: Importing libraries" with the snippet "To import a library that's not in Colaboratory by default, you ...". The fourth result is a link to "Google Colab" with the snippet "Colaboratory, or "Colab" for short, is a product from Google ...". The fifth result is a link to "This Google Notebook" with the snippet "Sign in." The sixth result is a link to "Overview Notebook" with the snippet "Code cells. Below is a code cell. Once the toolbar button ...". The seventh result is a link to "Using Google Colab with GitHub" with the snippet "Using Google Colab with GitHub · Loading Public Notebooks ...". At the bottom of the search results, there is a link "More results from google.com »".

 Welcome To Colaboratory

File Edit View Insert Runtime Tools Help

Share Sign In

+ Code + Text Copy to Drive Connect Editing

Table of contents

- Getting started
- Data science
- Machine learning
- More Resources
- Machine Learning Examples

What is Colaboratory?

Colaboratory, or "Colab" for short, allows you to write and execute Python in your browser, with

- Zero configuration required
- Free access to GPUs
- Easy sharing

Whether you're a **student**, a **data scientist** or an **AI researcher**, Colab can make your work easier.

Watch [Introduction to Colab](#) to learn more, or just get started below!

Getting started

The document you are reading is not a static web page, but an interactive environment called a **Colab notebook** that lets you write and execute code.

For example, here is a **code cell** with a short Python script that computes a value, stores it in a variable, and prints the result:

```
[ ] seconds_in_a_day = 24 * 60 * 60
seconds_in_a_day
86400
```

To execute the code in the above cell, select it with a click and then either press the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter". To edit the code, just click the cell and start editing.

Variables that you define in one cell can later be used in other cells:

```
[ ] seconds_in_a_week = 7 * seconds_in_a_day
seconds_in_a_week
```

 Welcome To Colaboratory

File Edit View Insert Runtime Tools Help

Share Sign In

+ Code + Text Copy to Drive Connect Editing

Table of contents

- Getting started
- Data science
- Machine learning
- More Resources
- Machine Learning Examples

What is Colaboratory?

Colaboratory, or "Colab" for short, allows you to write and execute Python in your browser, with

Examples Recent Google Drive GitHub Upload

Filter notebooks

Title	Last opened	First opened
Welcome To Colaboratory	9:45 PM	9:45 PM

New notebook Cancel

604989

Colab notebooks allow you to combine executable code and rich text in a single document, along with images, HTML, LaTeX and more. When you create your own Colab notebooks, they are stored in your Google Drive account. You can easily share your Colab notebooks with co-workers or friends, allowing them to comment on your notebooks or even edit them. To learn more, see [Overview of Colab](#). To create a new Colab notebook you can use the File menu above, or use the following link: [create a new Colab notebook](#)

```
▶ mkdir myfolder
```

FILE EDIT VIEW NOTE

+ Code + Text

Insert code cell below
⌘/Ctrl+M B

```
Q mkdir myfolde
```

Os

{}
x}

```
✓ [1] mkdir myfolde
```

```
✓ [1] mkdir myfolder  
0s
```

```
✓ [2] ls  
0s
```

myfolder/ sample_data/

```
✓ [3] cd myfolder  
0s
```

/content/myfolder

```
✓ [4] pwd  
0s
```

'/content/myfolder'

```
✓ [5] cd a  
0s
```

↳ [Errno 2] No such file or directory: 'a'
/content/myfolder

Basic linux

- pwd

Print Working Directory

- ls

List files and directories

- man ls
- man pwd



Colab limitation

use ! before each command

pwd

mkdir yzd

ls yzd

cd yzd

mkdir medicine

mkdir math

ls -alht

cd ..

df -h

```
echo "yazd"
```

```
echo "yazd" > file.txt
```

```
cat file.txt
```

```
cp file file2
```

```
cp yazd/* folder1
```

```
rm mv zip unzip
```

mkdir

head -n 2

tail -n 2

grep

sed

awk

```
echo ">chr1" > file
```

```
echo "aaaa tttt gggg" >> file
```

```
head -n 2 chr.fa
```

```
tail -n 5 chr.fa
```

```
wc -l chr.fa
```

```
sed -n 2,3p chr.fa
```

```
sed s/a/mm/ chr.fa
```

tiny dataset

copy the link to fasta file from github

<https://github.com/brainstorm/tiny-test-data>

```
!fgrep -o G hg19.fa | wc -l
```

GC content=
$$\frac{G + C}{A + T + G + C} \times 100\%$$

Exercise: Calculate number of C

Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling
4. File formats
 - o fasta format
 - o fastq format
 - o sam format
 - o vcf format
5. Basic linux
6. Databases
 - o Genome datasets (NCBI, UCSC)
 - o



UNITE

A new NIH initiative to end structural racism and achieve racial equity in the biomedical research enterprise.

[LEARN MORE](#)

NCBI Home

Resource List (A-Z)

All Resources

Chemicals & Bioassays

Data & Software

DNA & RNA

Domains & Structures

Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [NCBI News & Blog](#)

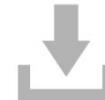
Submit

Deposit data or manuscripts into NCBI databases



Download

Transfer NCBI data to your computer



Learn

Find help documents, attend a class or watch a tutorial



Popular Resources

[PubMed](#)

[Bookshelf](#)

[PubMed Central](#)

[BLAST](#)

[Nucleotide](#)

[Genome](#)

[SNP](#)

[Gene](#)

[Protein](#)

[PubChem](#)

Literature	Genes	Proteins
Bookshelf	Gene 825	Conserved Domains 30
MeSH	GEO DataSets 1,081	Identical Protein Groups 98,143
NLM Catalog	GEO Profiles 54,675	Protein 195,051
PubMed	HomoloGene 3	Protein Family Models 104
PubMed Central	PopSet 1,254	Structure 151

Genomes	Clinical	PubChem
Assembly 63	ClinicalTrials.gov 1,277	BioAssays 3,958
BioCollections 0	ClinVar 27	Compounds 5
BioProject 238	dbGaP 40	Pathways 0
BioSample 2,438	dbSNP 0	Substances 29
Genome 1	dbVar 0	
Nucleotide 179,903	GTR 44	

Assembly

Assembly ▾

Hepatitis B virus

[Create alert](#) [Advanced](#) [Browse by organism](#)

Organism group

Viruses (63)

[Customize ...](#)

Status

[clear](#) [Latest](#) (63)

Latest GenBank (63)

Latest RefSeq (1)

[Download Assemblies](#)[Send to:](#) ▾

Assembly level

Complete genome (53)

Chromosome (9)

Scaffold (1)

Contig (0)

RefSeq category

Reference (0)

Representative (0)

Exclude

[clear](#)

Exclude partial (0)

Exclude from large multi-isolate project (0)

 [Exclude anomalous](#) (0)[Customize ...](#)

Annotation status

Has annotation (21)

GenBank has annotation (20)

RefSeq has annotation (1)

Taxonomy check

status

TAXONOMY

Hepatitis B virus

Hepatitis B virus is a species of dsDNA-RT virus in the family Hepadnaviridae (hepatitis B-type viruses).

Taxonomy ID: [10407](#)

I want to ...

- [search the Assembly database for Hepatitis B virus as an organism](#)
- [see assembly details for the Hepatitis B virus reference genome](#)
- [browse and download genomes for Hepatitis B virus](#)

Search results

Items: 1 to 20 of 63

<< First < Prev Page of 4 Next > Last >>

Assembly

Assembly



Hepatitis B virus

Create alert

Advanced

Browse by organism

Organism group

Viruses (63)

Customize ...

Status

clear

Latest (63)

Latest GenBank

Latest RefSeq (

Assembly level

Complete genome

Chromosome (9)

Scaffold (1)

Contig (0)

RefSeq category

Reference (0)

Representative

Exclude

Exclude partial (0)

Exclude from large multi-isolate project (0)

Exclude anomalous (0)

Customize ...

Annotation status

Has annotation (21)

GenBank has annotation (20)

RefSeq has annotation (1)

Taxonomy check
status

Download Assemblies

Source database (GenBank or RefSeq) ?

RefSeq

File type ?

Genomic FASTA (.fna)

Estimated size is 1.2 KB

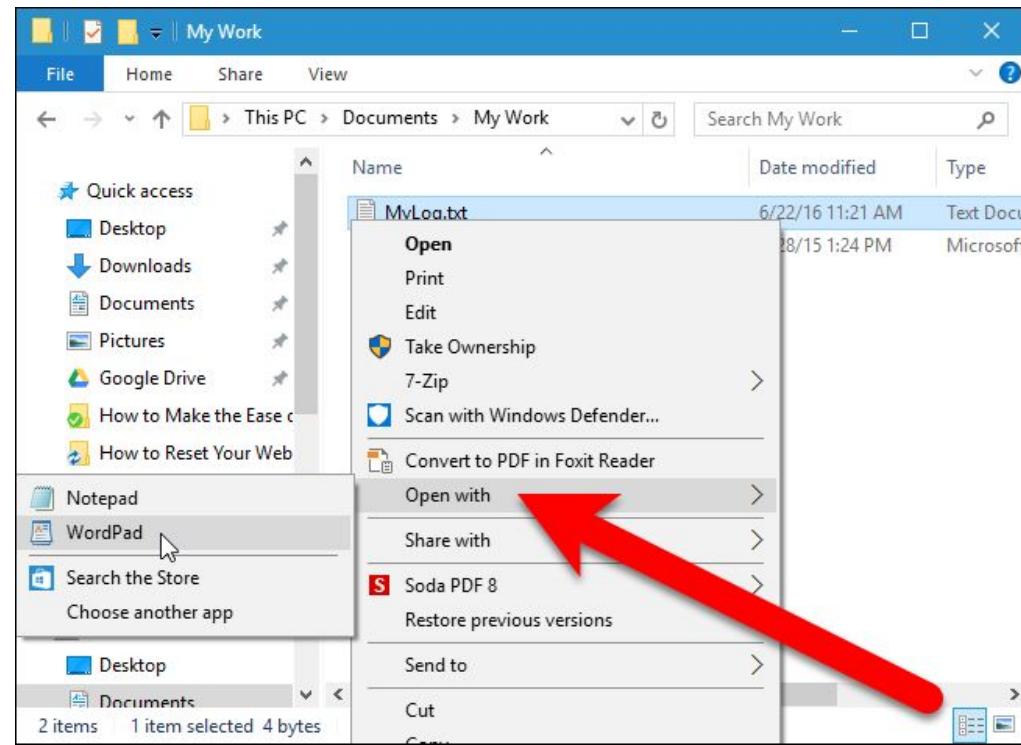
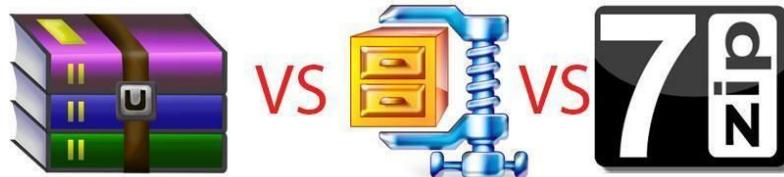
Hepatitis B virus in the family H

Download

I want to ...

- [search the Assembly database for Hepatitis B virus as an organism](#)
- [see assembly details for the Hepatitis B virus reference genome](#)
- [browse and download genomes for Hepatitis B virus](#)

Unzip





Genome Data

[Source Code](#)

[Genome Browser Store](#)

[Utilities](#)

[FTP](#)

[MySQL Access](#)

[REST API](#)

Table Browser

download data from the Genome Browser database

Variant Annotation Integrator

get functional effect predictions for variant calls

Data Integrator

combine data sources from the Genome Browser database

Genome Browser in a Box (GBiB)

run the Genome Browser on your laptop or server

In-Silico PCR

rapidly align PCR primer pairs to the genome

LiftOver

convert genome coordinates between assemblies

Track Hubs

import and view external data tracks

REST API

returns data in JSON format

Human genome

Dec. 2013 (GRCh38/hg38)

- Genome sequence files and select annotations (2bit, GTF, GC-content, etc) ▶
- Sequence data by chromosome
- Annotations ▶
- SNP-masked fasta files ▶
- LiftOver files
- Pairwise alignments ▶
- Multiple alignments ▶
- Patches ▶
- Data archive

Feb. 2009 (GRCh37/hg19)

- Genome sequence files and select annotations (2bit, GTF, GC-content, etc)
- Sequence data by chromosome
- Annotations ▶
- GC percent data
- Protein database for hg19
- SNP-masked fasta files ▶
- LiftOver files
- Pairwise alignments (primates) ▶

		2014-01-23 16:39	1.7M
<u>chr2v1.fa.gz</u>			
<u>chr20_GL383577v2_alt.fa.gz</u>	2014-01-23 16:40	42K	
<u>chr20_KI270869v1_alt.fa.gz</u>	2014-01-23 16:40	37K	
<u>chr20_KI270870v1_alt.fa.gz</u>	2014-01-23 16:40	55K	
<u>chr20_KI270871v1_alt.fa.gz</u>	2014-01-23 16:40	18K	
<u>chr21.fa.gz</u>	2014-01-23 16:39	12M	
<u>chr21_GL383578v2_alt.fa.gz</u>	2014-01-23 16:40	21K	
<u>chr21_GL383579v2_alt.fa.gz</u>	2014-01-23 16:40	65K	
<u>chr21_GL383580v2_alt.fa.gz</u>	2014-01-23 16:40	25K	
<u>chr21_GL383581v2_alt.fa.gz</u>	2014-01-23 16:40	36K	
<u>chr21_KI270872v1_alt.fa.gz</u>	2014-01-23 16:40	26K	
<u>chr21_KI270873v1_alt.fa.gz</u>	2014-01-23 16:  47K		
<u>chr21_KI270874v1_alt.fa.gz</u>	2014-01-23 16:  53K		
<u>chr22.fa.gz</u>	2014-01-23 16:39	12M	
<u>chr22_GL383582v2_alt.fa.gz</u>	2014-01-23 16:40	52K	
<u>chr22_GL383583v2_alt.fa.gz</u>	2014-01-23 16:40	31K	
<u>chr22_KB663609v1_alt.fa.gz</u>	2014-01-23 16:40	21K	
<u>chr22_KI270731v1_random.fa.gz</u>	2014-01-23 16:39	48K	
<u>chr22_KI270732v1_random.fa.gz</u>	2014-01-23 16:39	14K	
<u>chr22_KI270733v1_random.fa.gz</u>	2014-01-23 16:39	57K	
<u>chr22_KI270734v1_random.fa.gz</u>	2014-01-23 16:39	53K	
<u>chr22_KI270735v1_random.fa.gz</u>	2014-01-23 16:39	14K	
<u>chr22_KI270736v1_random.fa.gz</u>	2014-01-23 16:39	42K	
<u>chr22_KI270737v1_random.fa.gz</u>	2014-01-23 16:39	25K	
<u>chr22_KI270738v1_random.fa.gz</u>	2014-01-23 16:39	24K	
<u>chr22_KI270739v1_random.fa.gz</u>	2014-01-23 16:39	16K	
<u>chr22_KI270875v1_alt.fa.gz</u>	2014-01-23 16:40	83K	
<u>chr22_KI270876v1_alt.fa.gz</u>	2014-01-23 16:40	84K	
<u>chr22_KI270877v1_alt.fa.gz</u>	2014-01-23 16:40	33K	
<u>chr22_KI270878v1_alt.fa.gz</u>	2014-01-23 16:40	60K	
<u>chr22_KI270879v1_alt.fa.gz</u>	2014-01-23 16:40	97K	
<u>chr22_KI270928v1_alt.fa.gz</u>	2014-01-23 16:40	55K	
<u>chrM.fa.gz</u>	2014-01-23 16:40	5.4K	
<u>chrUn_GL000195v1.fa.gz</u>	2014-01-23 16:40	60K	
<u>chrUn_GL000213v1.fa.gz</u>	2014-01-23 16:40	53K	
<u>chrUn_GL000214v1.fa.gz</u>	2014-01-23 16:40	42K	
<u>chrUn_GL000216v2.fa.gz</u>	2014-01-23 16:40	43K	
<u>chrUn_GL000218v1.fa.gz</u>	2014-01-23 16:40	52K	
<u>chrUn_GL000219v1.fa.gz</u>	2014-01-23 16:40	57K	
<u>chrUn_GL000220v1.fa.gz</u>	2014-01-23 16:40	51K	
<u>chrUn_GL000224v1.fa.gz</u>	2014-01-23 16:40	49K	
<u>chrUn_GL000226v1.fa.gz</u>	2014-01-23 16:40	1.9K	
<u>chrUn_KI270302v1.fa.gz</u>	2014-01-23 16:40	735	

Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling
4. File formats
 - o fasta format
 - o fastq format
 - o sam format
 - o vcf format
5. Basic linux
6. Databases
 - o Genome datasets (NCBI, UCSC)
 - o Sequence Read Archive (SRA)
 - o

SRA database

NCBI Resources How To Sign in to NCBI

NCBI National Center for Biotechnology Information

Ending Structural Racism NIH nih.gov/ending-structural-racism

NCBI Home Resource List (A-Z) All Resources Chemicals & Bioassays Data & Software DNA & RNA Domains & Structures Genes & Expression Genetics & Medicine Genomes & Maps Homology Literature Proteins Sequence Analysis

✓ All Databases Assembly Biocollections BioProject BioSample BioSystems Books ClinVar Conserved Domains dbGaP dbVar Gene Genome GEO DataSets GEO Profiles GTR HomoloGene Identical Protein Groups MedGen MeSH NCBI Web Site NLM Catalog Nucleotide OMIM PMC PopSet Protein Protein Clusters Protein Family Models PubChem BioAssay

Search

and structural racism and achieve racial equity in the biomedical research enterprise.

NCBI
National Center for Biotechnology Information advances science and health by providing access to genomic information.

[Mission](#) | [Organization](#) | [NCBI News & Blog](#)

Submit **Download** **Learn**

manuscripts
uses

Transfer NCBI data to your computer

Find help documents, attend a class or watch a tutorial

Analyze Research

Popular Resources

PubMed Bookshelf PubMed Central BLAST Nucleotide Genome SNP Gene Protein PubChem

NCBI News & Blog

NCBI on YouTube: Customize MSA

123

SRA

SRA



lung tumor

[Create alert](#) [Advanced](#)

COVID-19 Information

[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#) |**Access**Controlled (7,943)
Public (3,691)**Source**DNA (11,120)
RNA (113)**Type**

exome (8,054)

Library Layoutpaired (10,342)
single (1,308)**Platform**BGISEQ (6)
Illumina (11,408)
Gn Torrent (236)**Strategy**[clear](#) [Exome \(11,650\)](#)**Data in Cloud**GS (11,613)
S3 (11,624)**File Type**

Summary ▾ 20 per page ▾

Send to: ▾

F**R****T**[View results as an expanded interactive table using the RunSelector.](#) [Send results to Run selector](#)**N**

Search results

Items: 1 to 20 of 11650[<< First](#) [< Prev](#) [Page](#) of 583 [Next >](#) [Last >>](#)**M****i** Filters activated: Exome. [Clear all](#) to show 51747 items.

- [DNA-targeted sequencing of lung adenocarcinoma cell lines](#)
 - 1. 1 ILLUMINA (Illumina MiSeq) run: 2.5M spots, 745.9M bases, 259.1Mb downloads
Accession: SRX12633731
 - [DNA-targeted sequencing of lung adenocarcinoma cell lines](#)
 - 2. 1 ILLUMINA (Illumina MiSeq) run: 2.3M spots, 696M bases, 289.3Mb downloads
Accession: SRX12633730
 - [DNA-targeted sequencing of lung adenocarcinoma cell lines](#)
 - 3. 1 ILLUMINA (Illumina MiSeq) run: 2.3M spots, 680.2M bases, 282.6Mb downloads
Accession: SRX12633729

T**P****S****E****E****C****C**

Main **Browse** Search Download Submit Software Trace Archive Trace BLAST

Studies Samples Analyses Run Browser Run Selector Provisional SRA

COVID-19 Information ×

[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#) | [Español](#)

DNA-targeted sequencing of lung adenocarcinoma cell lines (SRR16356243) Change accession

Metadata Analysis Reads Data access

Run	Spots	Bases	Size	GC content	Published	Access Type
SRR16356243	2.5M	745.9Mbp	271.7M	50%	2021-10-15	public

Quality graph ([bigger](#))

This run has 2 reads per spot:

[Legend](#)

Experiment	Library Name	Platform	Strategy	Source	Selection	Layout	Action
SRX12633731	Sam11	Illumina	Targeted-Capture	GENOMIC	PCR	PAIRED	BLAST

Show design

Biosample	Sample Description	Organism
SAMN2247257 (SRS10587109)		Homo sapiens

DNA-targeted sequencing of lung adenocarcinoma cell lines (SRR16356243)

Metadata Analysis Reads Data access

Filter: Find Filtered Download [What does it do?](#)

[What can the filter be applied to?](#)

< 1 246983 >

View: biological reads technical reads quality scores [advanced options](#)

Reads (separated)

1. SRR16356243.1 SRS10587109

name: 1, member: 9

2. SRR16356243.2 SRS10587109

name: 2, member: 9

3. SRR16356243.3 SRS10587109

name: 3, member: 9

4. SRR16356243.4 SRS10587109

name: 4, member: 9

5. SRR16356243.5 SRS10587109

name: 5, member: 9

6. SRR16356243.6 SRS10587109

name: 6, member: 9

7. SRR16356243.7 SRS10587109

name: 7, member: 9

8. SRR16356243.8 SRS10587109

name: 8, member: 9

9. SRR16356243.9 SRS10587109

name: 9, member: 9

10. SRR16356243.10 SRS10587109

name: 10, member: 9

>gnl|SRA|SRR16356243.1.1 1 (Biological)

CCAAAGTTCCAAAANAAAAGAAATGCAGGGGATACGGCCAGGCATTGAAGTCTCATGGAA
GCCAGCCCCCTCAGGGCAACTGACCGTGCAAGTCACAGACTTNGCTGTCCCAGAAATGCAAG
AAGCCCAGACGGAAACCGNAGCTGCCCCGGT

>gnl|SRA|SRR16356243.1.2 1 (Biological)

TGTCCTTCCCAGAAAAACCTACCAAGGCAGCTACGGTTCCGTCTGGGCTTCTTGCATTC
TGGGACACCCAAGTCCTGTGACTTGCACGGTCAGTTGCCCTGAGGGGNTGGCTTCCATGAG
ACTTCAATGCCCTGGCCGTATCCCCCTGCAATT

[Write to the Help Desk](#) | [Privacy Notice](#) | [Disclaimer](#) | [Accessibility](#)

Last update: Fri May 8 15:05:51 EDT 2020

National Center for Biotechnology Information | U.S. National Library of Medicine

Click to go forward, hold to see history

Sequence Read Archive

Main   Download Submit Software Trace Archive Trace BLAST

SRA Objects BLAST Entrez

Download for Experiment SRX12633731

Accession	# of bases	# of spots
	total	filtered
<input checked="" type="checkbox"/> SRR16356243	745.9M	2.5M

Filter

Search:

 What can the filter be applied to?

Download Format

filtered clipped FASTA FASTQ

[Download](#)

^ATTGGA

NCBI Site map All databases Search

Sequence Read Archive NCBI Home (access key '1')

Main Browse Search Download Submit Software Trace Archive Trace BLAST

SRA Objects BLAST Entrez

Download for Experiment SRX12633731

Accession	# of bases	# of spots
		total filtered
<input checked="" type="checkbox"/> SRR16356243	745.9M	2.5M 

Filter

Search: ^ATTGGA

[What can the filter be applied to?](#)

Download Format

filtered clipped FASTA FASTQ

[Write to the Help Desk](#) | [Privacy Notice](#) | [Disclaimer](#) | [Accessibility](#)
[National Center for Biotechnology Information](#) | [U.S. National Library of Medicine](#)

Last update: Fri May 8 15:05:51 EDT 2020

NIH  FIRSTGOV.gov

<https://ncbi.nlm.nih.gov>



sra_data.fastq.txt

@SRR16356243.194720.1 194720 length=151
GGCAGCAAAGTTTATTGTAAAATAAGAGATCGATATAAAATGGGATATAAAAGGGAGAAGGAGGGGAAGGGTGGGTGAAAATGCAG
ATGTGCTTCAGAACATGTAAAAGATGTTGACCTTCAGCTGGACGTGGCTCAATTGT
+SRR16356243.194720.1 194720 length=151
AAFFFFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/
FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/
@SRR16356243.194720.2 194720 length=151
CAATTGGAGGCCACCACGTCCAGCTGGAAGGGTCAACATCTTACATTCTGCAAGCACATCTGCATTTCACCCCACCCCTCCCTCCTT
CTCCCTTTATATCCCATTATCGATCTTATTTACAATAAAACTTGCTGGCA
+SRR16356243.194720.2 194720 length=151
FFFFFFFAFFFFF/FFFFFFFAFFFFF/FFFFFFFAFFFFF/FFFFFFFAFFFFF/FFFFFFFAFFFFF/FFFFFFFAFFFFF/
FFFFFFFAFFFAFAAAFFFFFFFFFF=AFF/FFFFFFFFF/FFFRAFF=FFF/=F/FAFFFFFFF6FF/FF
@SRR16356243.247542.1 247542 length=151
GTTTGACTTAGTATGAATGTGGTTACGTGGAAGCAAATGTGTCTTCACTTATCATGAAAAAGTCTGCAAGTGCTCTGCGACGTCCAG
GGAAATGATCCTACCCCTCAATCTCAGCTCAAAGGGAACCTTGCTTTTCAGTGACCA
+SRR16356243.247542.1 247542 length=151
FAFF/FA=FAFF/FFFF/FFAFA//F/A/FF/FFFFAAA//FF/FFFF//A///FFFFA==//FF/F6FFA//A/FFF//
=F//F/AFA//F//=//F//=FAA/AF//A/F/FF/AFFF//FFF/F//FFFFF/F//
@SRR16356243.247542.2 247542 length=151
ATTGGACACTGAAAAGAGAAAAGGTTCCCTGTGAGCTGAAGAGTGAGTTAGGATCATTCAAGGGACGTGCGAGAGCACTGCAGAC
TTTATCATGAAATAGTTAACGACACATTGATTCAACGAAACCACTAAACTAAGTCAA
+SRR16356243.247542.2 247542 length=151
//F/F/FF//=FAFFFAFF/F/FF//F//FF/FFFA//FA/F//A/F//FF//AFAF//F//AF6//FAFFFFAA/F//FF//
FF//F/F/FF/A/F/FF//AA/A/FFA/FA//FFFFF/A/AA//F/A/FFFF/F/FF//=A
@SRR16356243.315284.1 315284 length=151
AAAAGATCTTATTGGCTTGGTCTTCAAGTAGCCAAAGGCATGAAATATCTTGCAGCAAGCAAAAGTTGTCCACAGAGACTGGCTGCAAG
AAACTGTATGTAAGTATCAGAACATCTGTGCCACAATCAAATTAAAGTGACAAGGAGGAAT
+SRR16356243.315284.1 315284 length=151
FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/
FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/FFFFFFF/

1000 Genome project

VCF files

<http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>

Genetic map

https://mathgen.stats.ox.ac.uk/impute/1000GP_Phase3.html

GIAB

https://github.com/genome-in-a-bottle/giab_data_indexes

AshkenazimTrio

Son:HG002 https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG002_NA24385.son/

Father:HG003 https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG003_NA24149_father/

Mother:HG004 https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG004_NA24143_mother/

Sequencing Platform	Sequence	Alignment
Illumina WGS 2x150bp 300X per individual	All HG002 HG003 HG004	novoalign: All HG002 HG003 HG004
Illumina 6KB Matepair	All HG002 HG003 HG004	bwamem:hg19 All HG002 HG003 HG004
Illumina WGS 2X250bp	All HG002 HG003 HG004	isaac:hg19 All HG002 HG003 HG004 novoalign: All HG002 HG003 HG004
Moleculo	All HG002 HG003 HG004	
Illumina Whole Exome	-	bwamem:hg19 All HG002 HG003 HG004
SOLiD 60x for son	All HG002	LifeScope:hg19 All HG002
CompleteGenomics	-	CGAtools:hg19 All HG002 HG003 HG004

Another linux website

← → C ⌂ copy.sh



You can email me at copy@copy.sh. Use my [GnuPG key](#).

Projects

Virtual x86
Run KolibriOS, Linux or Windows 98 in your browser.

Disable mouse.

No internet

Select profile

Arch Linux 12 MB	A complete Arch Linux
Damn Small Linux 50 MB	Graphical Linux with 2
Buildroot Linux 5.0 MB	Minimal Linux with bus

copy.sh/v86/profile=buildroot

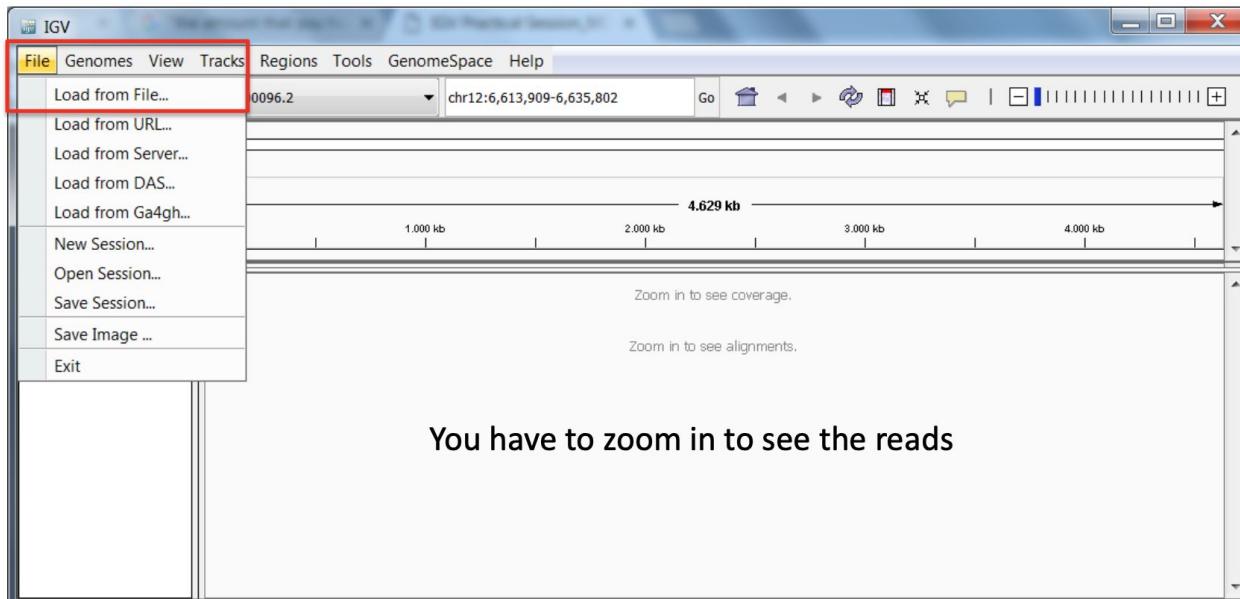
serial18250: ttyS0 at I/O 0x3f8 (irq = 4, base_baud = 115200) is a 16550A
ne2k-pci 0000:00:05.0 found PCI INT A -> IRQ 10
ne2k-pci 0000:00:05.0 eth0: Realtek RTL-8029 found at 0xc140, IRQ 10, 00:22:15:00:00:00
0:0:0:0:0:0
18042: PNP: No PS/2 controller found.
18042: Probing ports directly.
serio: i8042 KBD port at 0x60, irq 1
serio: i8042 KBD port at 0x64, irq 12
NET: Registered protocol family 10
input: AT Translated Set 2 keyboard as /devices/platform/i8042/serio0/input/input0
0:0
Segment Routing with IPv6
NET: Registered protocol family 17
ppp: module loaded
sched_clock: Maxxing stable (59550599999, 85360000) ->(6130590000, -90170000)
Freeing unused kernel image (initmem) memory: 2576K
RAM: 128M total, 128M free, 128M used, 0M available, 0M reserved, 0M read-only data: 4920K
Run /init as init process
Enable networking using 'udhcpc'
Files send via emulator appear in /mnt/
-k wget "C
-k s"C
-k -

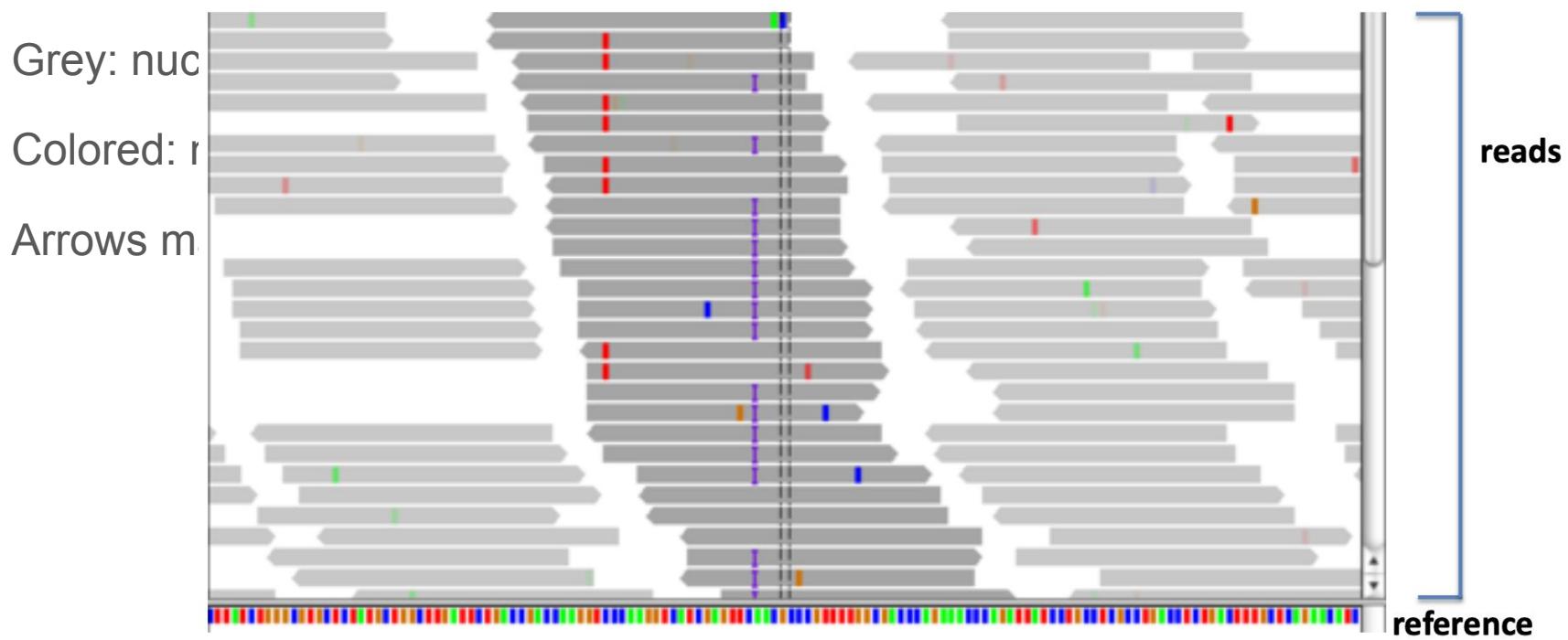
This is the serial console. Whatever you type or paste here will be sent to COM1
mount: mounting hostip on /mnt failed: Device or resource busy
Enable networking using 'udhcpc'
Files send via emulator appear in /mnt/
-k

Index

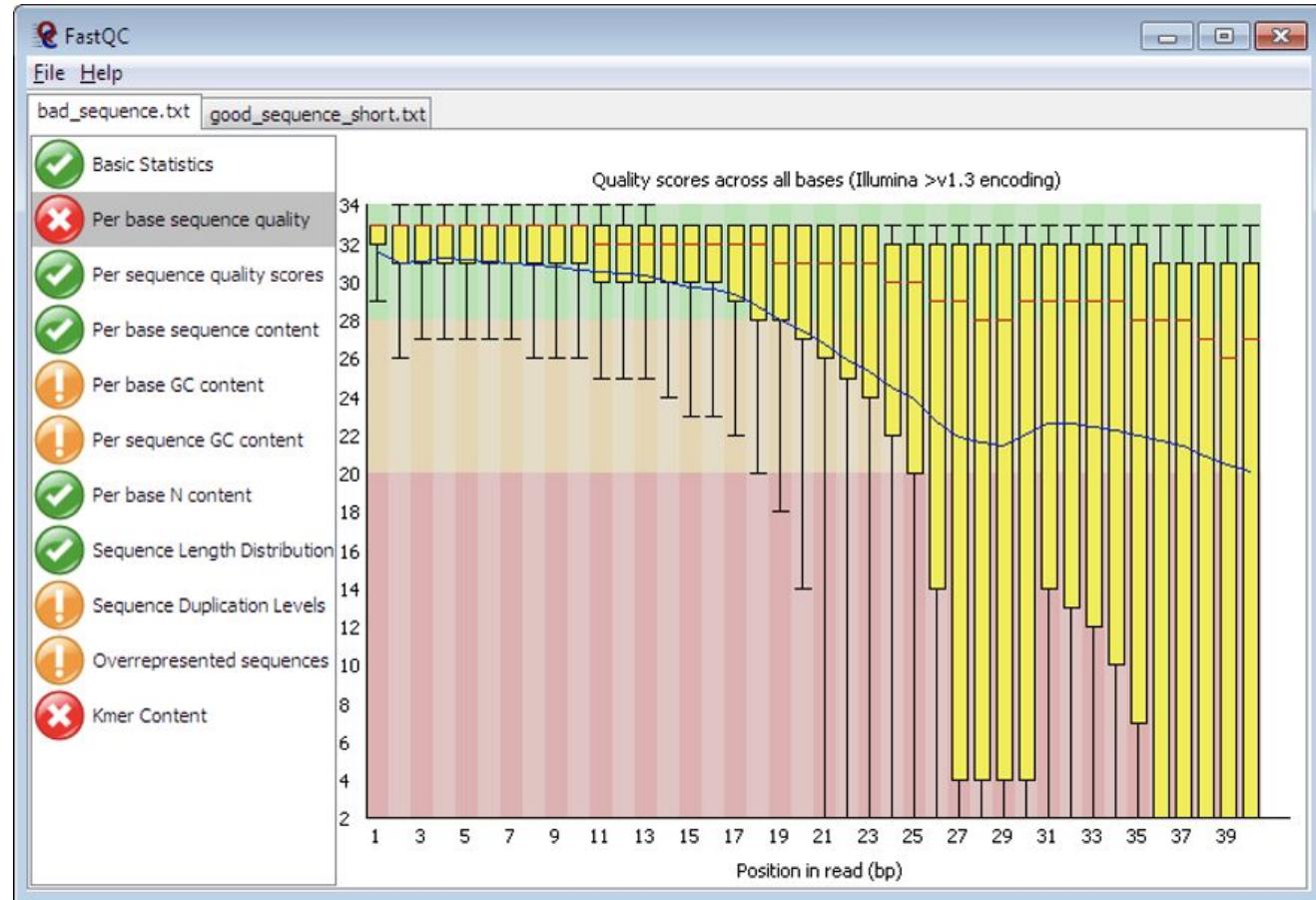
1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling
4. File formats
 - o fasta format
 - o fastq format
 - o sam format
 - o vcf format
5. Basic linux
6. Databases
 - o Genome datasets (NCBI, UCSC)
 - o Sequence Read Archive (SRA)
 - o 1000 Genome project
 - o GIAB
7. Integrative Genomics Viewer (IGV)
8. conda FastQC

Visualization of mapping results: Loading data into IGV





FASTQC



Trimming low quality bases from the ends of the reads

G	A	T	T	T	C	T	T	A	G	T	T	G	T	A	T	A	A	T
12	19	17	40	23	30	36	37	12	33	34	26	21	28	34	24	26	12	14

5' trimming bases with quality below 20

X	X	X	T	T	C	T	T	A	G	T	T	G	T	A	T	A	A	T
X	X	X	40	23	30	36	37	12	33	34	26	21	28	34	24	26	12	14

3' trimming bases with quality below 20

G	A	T	T	T	C	T	T	A	G	T	T	G	T	A	T	A	A	X
12	19	17	40	23	30	36	37	12	33	34	26	21	28	34	24	26	12	14

REVIEW

Open Access

Best practices for variant calling in clinical sequencing



10 recommendations 8 PC

Daniel C. Koboldt^{1,2}

Strategy

Alignment and pre-processing

Read alignment	BWA-MEM [25], Bowtie 2 [26], minimap2 [27], Novoalign
Marking duplicates	Picard tools [28], Sambamba [29], SAMBLASTER [30]
BAM file creation	Samtools [31], GATK [19]
Sequencing metrics	BEDTools [32], Picard tools [28], QualiMap 2 [33]
Sample quality control	KING [34], VerifyBamID [35]

Variant calling

Inherited SNVs/indels	FreeBayes [36], GATK HaplotypeCaller [19], Platypus [20], Samtools/BCFtools [37]
Somatic mutations	deepSNV [38], MuSE [39], MuTect2 [40], SomaticSniper [41], Strelka2 [42], VarDict [43], VarScan2 [44]
Copy number variants	cn.MOPS [45], CONTRA [46], CoNVEX [47], ExomeCNV [48], ExomeDepth [49], XHMM [50]
Structural variants	DELLY [51], Lumpy [52], Manta [53], Pindel [54], SVMerge [55]
Gene fusions (RNA-seq)	fusionCatcher [56], fusionMap [57], mapSplice [58], SOAPfuse [59], STAR-Fusion [60], TopHat-Fusion [61]

Variant review/storage

Visualization and review	Artemis [62], Integrative Genomics Viewer [63]
VCF/BCF file manipulation	BCFtools [37]

conda

conda install -c conda-forge jupyterlab

conda install samtools

conda install fastqc

conda install jupyter

Google search results for "install miniconda". The top result is the official Miniconda documentation page on conda.io/miniconda.

Miniconda — Conda documentation

There are two variants of the installer: **Miniconda** is Python 2 based and **Miniconda3** is Python 3 based. Note that the choice of which **Miniconda** is installed only ...

Miniconda hash information
Built with Sphinx using a theme provided by Read the Docs ...

Help and support
Courses are available to individuals online, at numerous ...

[More results from conda.io »](#)

bioconda / packages / samtools 1.14

Tools for dealing with SAM, BAM and CRAM files

Conda Files Labels Badges

License: MIT Home: <https://github.com/samtools/samtools> 2350081 total downloads Last upload: 27 days and 22 hours ago

Install > linux

Install74+ce35d868 documentation

Python installations or packages. The fastest way to i version of ...

Installers

Info: This package contains files in non-standard labels.

conda install

linux-64 v1.14
osx-64 v1.14

To install this package with conda run one of the following:

```
conda install -c bioconda samtools
conda install -c bioconda/label/cf201901 samtools
```

Resources

- Coursera
- <https://maktabkhooneh.org/> (دکتر شریفی زارچی)
- Galaxy: Work on public server: <https://usegalaxy.org/>
- Slides on github
- Algorithms: <https://langmead-lab.org/teaching-materials/>

Biology & Genetics

Targeted studies of one or a few genes

Targeted,
low-throughput experiments

Clever experimental design, painstaking experimentation



Genomics

Studies considering all genes in a genome

Global,
high-throughput experiments

Tons of data,
uncertainty,
computation



Index

1. Basic biology
 - o DNA
 - o UCSC genome browser
2. DNA sequencing technologies
 - o Sanger
 - o Second generation (Illumina)
 - o Third (Pacbio & ONT)
 - o Single Cell
3. Bioinformatics workflow
 - o Genome assembly
 - o Read alignment
 - o Variant calling
4. File formats
 - o fasta format
 - o fastq format
 - o sam format
 - o vcf format
5. Basic linux
6. Databases
 - o Genome datasets (NCBI, UCSC)
 - o Sequence Read Archive (SRA)
 - o 1000 Genome project
 - o GIAB
7. Integrative Genomics Viewer (IGV)
8. Conda & FastQC

Thanks for your attention!



Sina.Majidian
@gmail.com

