

# Medizinische Genomanalysen

LE 1 (06.04.2021 – 18:00)

Stephan Pabinger

[stephan.pabinger@gmail.com](mailto:stephan.pabinger@gmail.com)

# What to expect

- Finish with an understanding of major concepts and tools
- Know how to perform variant calling
- Ready to make informed choices about what kind of variant calling tools you may need
- Focus is on Variant Calling and Functional Annotation
  - Alignment refinement
  - Alignment metrics reports
  - Base quality score recalibration
  - Variant calling
  - Variant annotation and filtering

# Course Information

[https://github.com/spabinger/medizinische\\_genomanalysen\\_2021](https://github.com/spabinger/medizinische_genomanalysen_2021)

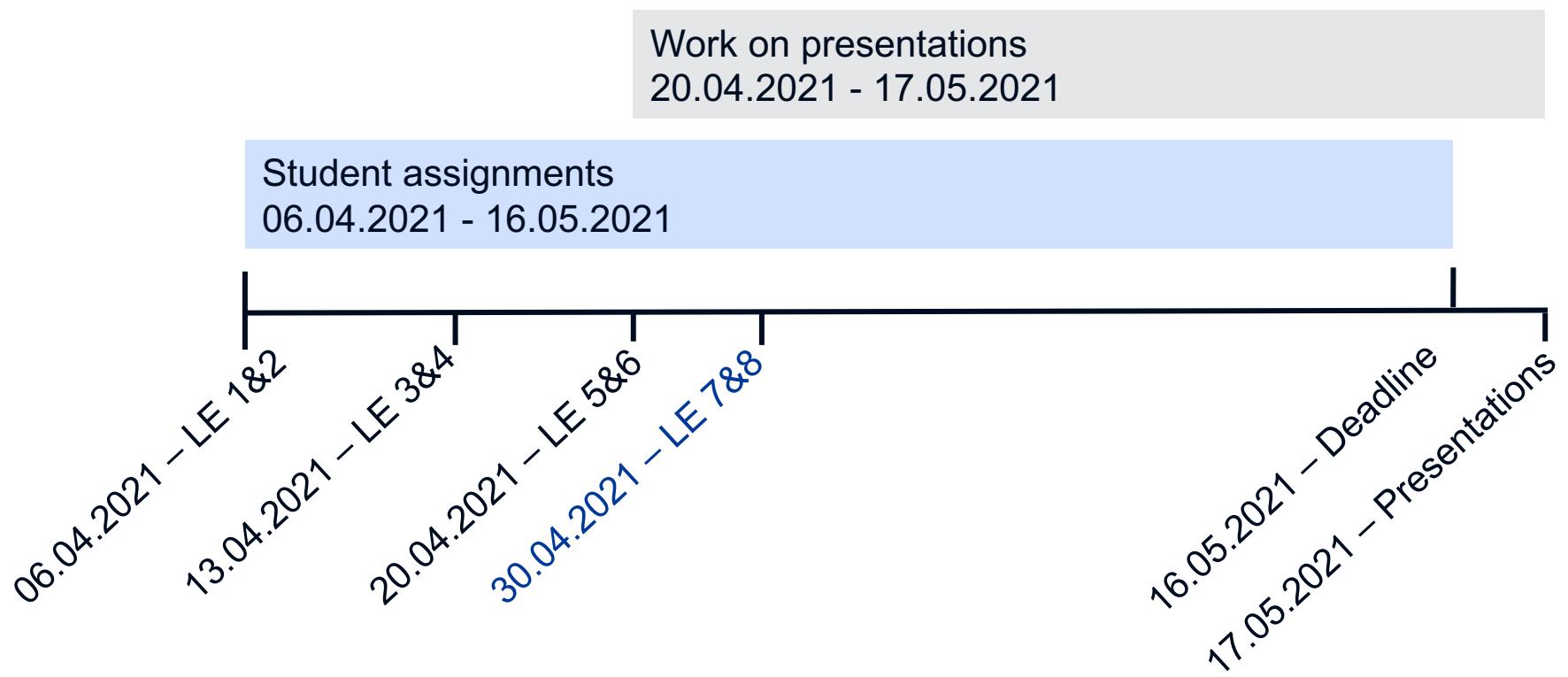
## Structure

- In addition to the lecture, we will work on 3 assignments (one per date)
- Assignment task: create a Python program
  - More information will follow
- Last lecture is reserved for student presentations – please see link above

## Grading

- 50% assignments
- 50% student projects & presentation

# Timeline



# Structure

## LE 1

- History, terms, sequencing, SAM/BAM

## LE 2

- Reference genome, FASTQ, cleaning

## LE 3

- SAMtools, genetic variation, variant calling

## LE 4

- Variant callers, structural variations & callers

## LE 5

- CNVs, somatic mutations, filtering, annotation

## LE 6

- Visualization, pipelines

## LE 7 & 8

- Albert Kriegner

## LE 9 & 10

- Presentations

# My background

## Bioinformatics

- Software development
- Web tools
- Pipeline design

## Working with sequencing data

- DNASeq
- RNASeq
- MethylationSeq

## Data analysis

- DNA data
- Protein data
- Peptide data
- Immunomics data

# Introduction

- Your background
- Python programming experience

**GIT**

# GIT

## Information

- Distributed revision control system
- Developed by Linus Torvalds (Linux developer)
- Local repositories & remote repositories

## Keep track of changes

- Code
- Manuscript
- Presentations
- Data analysis

## Master/PhD thesis

## Merging collaborators' changes



# "FINAL".doc



FINAL.doc!



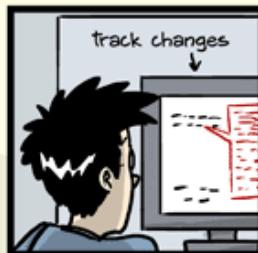
FINAL\_rev.2.doc



↑  
FINAL\_rev.6.COMMENTS.doc



FINAL\_rev.8.comments5.  
CORRECTIONS.doc



FINAL\_rev.18.comments7.  
corrections9.MORE.30.doc



FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRAD SCHOOL????.doc



JORGE CHAM © 2012

# Start with GIT

- Create a new directory
- Open it (cd into it)
- Perform  
`git init`
- Work ....
- Add files you want to store in the repository (locally)  
`git add XY.txt`  
`git add *` (to add everything)
- Commit files  
`git commit -m „Performed first analysis“`

# Start with GIT

```
stephan@shaq /tmp $ mkdir super_analysis
stephan@shaq /tmp $ cd super_analysis/
stephan@shaq /tmp/super_analysis $ git init
Initialized empty Git repository in /tmp/super_analysis/.git/
stephan@shaq /tmp/super_analysis $ touch XY.txt
stephan@shaq /tmp/super_analysis $ touch rawfile.txt
stephan@shaq /tmp/super_analysis $ git add XY.txt
stephan@shaq /tmp/super_analysis $ git add *
stephan@shaq /tmp/super_analysis $ git commit -m "Performed first analysis"
[master (root-commit) 915d159] Performed first analysis
2 files changed, 0 insertions(+), 0 deletions(-)
create mode 100644 XY.txt
create mode 100644 rawfile.txt
stephan@shaq /tmp/super_analysis $
```

# Features

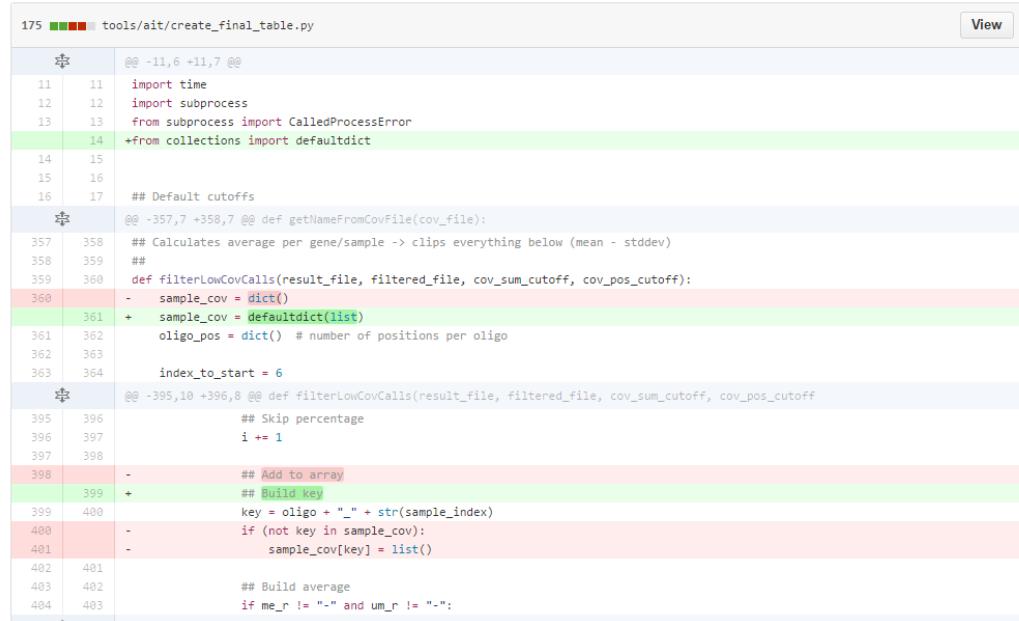
- Status of your project  
git status

```
stephan@shaq /tmp/super_analysis $ git status
On branch master
Untracked files:
  (use "git add <file>..." to include in what will be committed)

    nextanalysis.py

nothing added to commit but untracked files present (use "git add" to track)
stephan@shaq /tmp/super_analysis $
```

- History
- Go to a previous version
- Branching (implement a new feature without disrupting main code)
- Merging between different versions/branches



A screenshot of a GitHub interface showing a code diff. The top bar shows the file path 'tools/ait/create\_final\_table.py'. The diff shows several additions and deletions of code. The left column lists line numbers, and the right column shows the code changes with '+' for additions and '-' for deletions.

```
175 11,6 +11,7 @@  
11 import time  
12 import subprocess  
13 from subprocess import CalledProcessError  
14 +from collections import defaultdict  
15  
16 ## Default cutoffs  
17 @@ -357,7 +358,7 @@ def getNameFromCovFile(cov_file):  
357 358 ## Calculates average per gene/sample -> clips everything below (mean - stddev)  
358 359 ##  
359 360 def filterLowCovCalls(result_file, filtered_file, cov_sum_cutoff, cov_pos_cutoff):  
360 - sample_cov = dict()  
361 + sample_cov = defaultdict(list)  
361 362 oligo_pos = dict() # number of positions per oligo  
362 363  
363 364 index_to_start = 6  
364 @@ -395,10 +396,8 @@ def filterLowCovCalls(result_file, filtered_file, cov_sum_cutoff, cov_pos_cutoff  
395 396 ## Skip percentage  
396 397 i += 1  
397 398  
398 - # Add to array  
399 + ## Build key  
399 400 key = oligo + "_" + str(sample_index)  
400 - if (not key in sample_cov):  
401 + sample_cov[key] = list()  
401 402 ## Build average  
402 403 if me_r != "-" and um_r != "-":  
403 404
```

# Github - Gitlab

## Github

- Web-based Git repository hosting service
- Free to use
- Request for private repositories

## Gitlab

- Community version
- Open Source
- Share code, analyses, etc
- Easy to transfer to Github

The screenshot shows a GitHub repository page for 'tadKeys/tabsat'. At the top, there are tabs for Code, Issues (0), Pull requests (0), Wiki, Pulse, Graphs, and Settings. Below the tabs, it says 'Targeted Amplicon Bisulfite Sequencing Analysis Tool — Edit'. It shows 58 commits, 1 branch, 0 releases, and 1 contributor. The branch is set to 'master'. A 'New pull request' button is highlighted. Other buttons include New file, Upload files, Find file, HTTPS, and Download ZIP. The URL is https://github.com/tadKey/tabsat.

The commit history lists several commits:

- Fixed bug in prepareReference.sh (6 months ago)
- Updated create\_final\_table - did some testing to ensure filtering of ... (18 days ago)
- New parameters for TABSAT: Improved configuration (8 months ago)
- CpG files are also copied to a summary directory. New CpG all file. I... (2 months ago)
- Added Dockerfile (3 months ago)
- Improved Readme. Updated test data (2 months ago)
- Update demo.md (10 days ago)
- Added version (4 days ago)

The README.md file contains the following text:

## TABSAT

TABSAT - Targeted Amplicon Bisulfite Sequencing Analysis Tool - is a tool for analyzing targeted bisulfite sequencing data generated on an Ion Torrent PGM / Illumina MiSeq. It performs

- Quality Assessment
- Alignment using BiMark

Below the README, there is a commit history for the 'master' branch, filtered by 'seqipenext'. It shows 6,885 commits in total. The commits listed are:

- \*) Use of BamUid to get the bam files ... (Stephan authored 26 days ago) (6fbbbe3c, Browse Files)
- \*) Sambamba program definition ... (Stephan authored 26 days ago) (8f61768a, Browse Files)
- \*) Cutadapt bam generates now an bam index for the trimmed bam file (Stephan authored 26 days ago) (bdb75c6, Browse Files)
- \*) Use of sambamba for mpileup file generation - dramatically increases processing speed (Stephan authored 26 days ago) (fee6b498, Browse Files)
- \*) Consolidated output directory creation (Stephan authored 26 days ago) (0d2068f3, Browse Files)
- \*) Include ExAC and Biotype output fields ... (Stephan authored 26 days ago) (Saaft540, Browse Files)
- \*) Fixed Refseq HGVS bug - assembly versions are not consistent throughout chrom... (Stephan authored 26 days ago) (c1b0a584, Browse Files)
- removed workspace.xml from git (Stephan authored 2 months ago) (08921a5b, Browse Files)
- \*) Added workspace.xml to .gitignore (Stephan authored 2 months ago) (76a55a8d, Browse Files)

# Remote repositories

- **Clone repository**

```
git clone username@host:/path/to/repository
```

- *Work and commit*

- **Push files to remote repository**

```
git push
```

# GIT - terms

- Committed means that the data is safely stored in your local database.
- Modified means that you have changed the file but have not committed it to your database yet.
- Staged means that you have marked a modified file in its current version to go into your next commit snapshot.

# GIT – Hints- Resources

- Make lots of commits
- Don't commit large files (they should be on the server or somewhere else)

## Information

<http://rogerdudler.github.io/git-guide/>

[http://kbroman.org/github\\_tutorial/](http://kbroman.org/github_tutorial/)

<http://nyucll.org/pages/GitTutorial/>

<https://swcarpentry.github.io/git-novice/>

## Windows

- <https://git-for-windows.github.io/>  
→ Graphical user-interface (for init, add, commit, push, compare)

## Linux & Mac

- Packages and GUIs available

# History

History of Molecular Diagnostics | Sequencing | and more

# History of Molecular Diagnostics

## The Molecular Biology Timeline

- 
- 1865      Gregor Mendel, Law of Heredity
  - 1866      Johann Miescher, Purification of DNA
  - 1949      Sickle Cell Anemia Mutation

# Sickle Cell Anemia

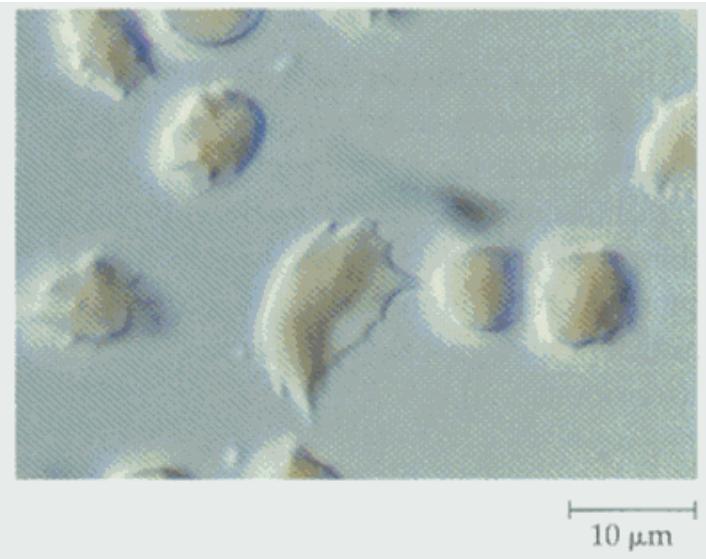
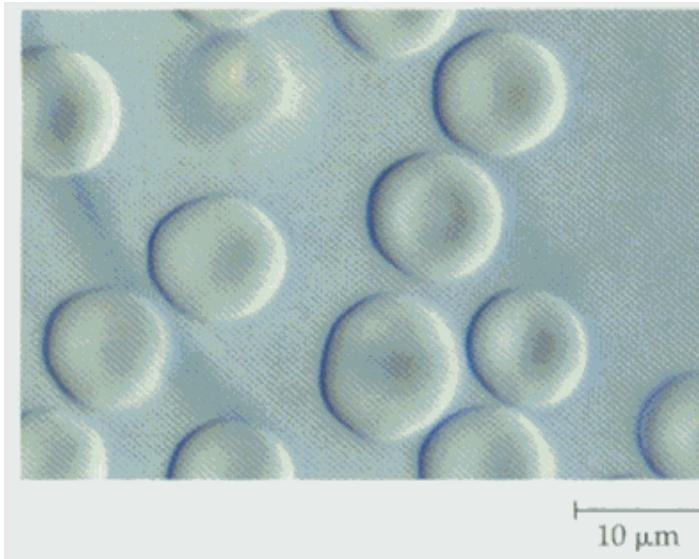
## Deformation of red blood cells

### Epidemiology:

1 out of 500 with African heritage; also common in Asia and Mediterranean region

### Genetics

Single point mutation (SNP) → 6. Codon of the  $\beta$ -globin Gens



# History of Molecular Diagnostics

## The Molecular Biology Timeline

1865	Gregor Mendel, Law of Heredity
1866	Johann Miescher, Purification of DNA
1949	Sickle Cell Anemia Mutation
1953	Watson and Crick, Structure of DNA



# Discovery of DNA Structure

James Watson & Francis Crick: **DNA structure - 1953**

- First correct double-helix model of DNA structure
- Based on one X-ray diffraction image taken by Rosalind Franklin and Raymond Gosling

NO. 4356 April 25, 1953 NATURE 737

equipment, and to Dr. G. E. R. Dickson, and the captain and officers of the R.R.S. *Discovery II* for their part in making the observations.

<sup>1</sup> Young, F. B., Gertrude, H., and Jeavons, W., *Phil. Mag.*, **40**, 149 (1925).

<sup>2</sup> Longworth-Higgins, M. S., *Mon. Not. Roy. Astr. Soc., Geophys. Suppl.*, **5**, 255 (1949).

<sup>3</sup> Von Hippel, P. S., Woods Hole Papers in Phys. Oceanogr. Kerec., **11** (33) (1960).

<sup>4</sup> Ekman, V. W., *Arkiv. Mat. Astron. Fysik* (Stockholm), **2** (11) (1960).

## MOLECULAR STRUCTURE OF NUCLEIC ACIDS

### A Structure for Deoxyribosyl Nucleic Acid

We wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey<sup>1</sup>. They kindly made their manuscript available to us in advance of publication. Their model consists of three interlocking chains, with phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray pattern is that salt, not the free acid. Without the acidic hydroxyl groups it is not clear how this would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain model has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for this reason we shall not comment on it.

We wish to put forward a really different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual assumption that each chain consists of phosphate diester groups joining  $\beta$ -D-deoxyribofuranose residues with 3',5'-linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the different directions of the atoms in the two chains, the helixes are in opposite directions. Each chain loosely resembles Furberg's<sup>2</sup> model No. 1. That is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furberg's "standard configuration", the sugar being roughly perpendicular to the attached base. There

This figure is purely schematic. The deoxyribose symbols show the two phosphate-sugar pental rod pairs of the structure joined together. The vertical line marks the fibre axis.

is a residue on each chain every 3.4 Å. in the direction. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 Å. The distance of a phosphorus atom from the fibre axis is 10 Å. As the phosphates are on the outside, esterions have ester bonds to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner

in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs in such a way that one hydrogen atom is added to a single base from the other chain, so that the two lie side by side with identical z-coordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bond is made as follows: purine position 1 to pyrimidine position 1; purine position 6 to pyrimidine position 6.

If it is assumed that the bases only occur in the structure in the most plausible standard forms (that is, with the keto rather than the enol configurations) it is found that only specific pairs of bases can bond together. These pairs are: adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In addition, if one member forms one member of a pair, on either chain, then on these assumptions the other member must be thymine; similarly for guanine and cytosine. The sequence of bases on a single chain does not appear to be restricted in any way; however, if only specific pairs of bases can be formed, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.

It has been found experimentally<sup>3,4</sup> that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribose nucleic acid.

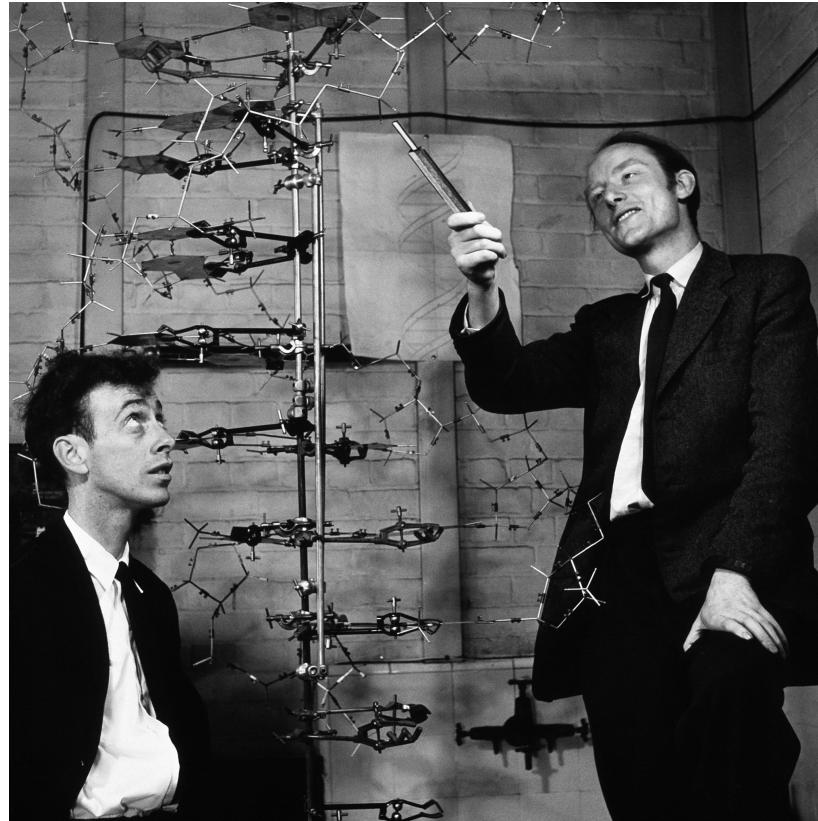
It is probably impossible to build this structure in the fibre form in place of the deoxyribose, as the extra oxygen atom would make too close a van der Waals contact.

The previously published X-ray data<sup>5,6</sup> on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. So far as we can see, it is roughly compatible with the experimental data, but it must be regarded as unproved until it has been checked against more exact results. Some of these are given in the following section. We were not able to get all the details of the results presented there, when we devised our structure, which rests mainly on stereochemical arguments.

We had not estimated the relative positions of pairing bases, but postulated immediately a possible copying mechanism for the genetic material.

Full details of the structure, including the conditions assumed in building it, together with a set of coordinates for the atoms, will be published elsewhere.

We are much indebted to Dr. Jerry Donohue for constant advice and criticism, especially on interatomic distances. We have also been stimulated by a knowledge of the general nature of the work and experimental results and ideas of Dr. M. H. F. Wilkins, Dr. R. E. Franklin and their co-workers at



# History of Molecular Diagnostics

## The Molecular Biology Timeline

1865	Gregor Mendel, Law of Heredity
1866	Johann Miescher, Purification of DNA
1949	Sickle Cell Anemia Mutation
1953	Watson and Crick, Structure of DNA
1970	Recombinant DNA Technology
1977	DNA sequencing

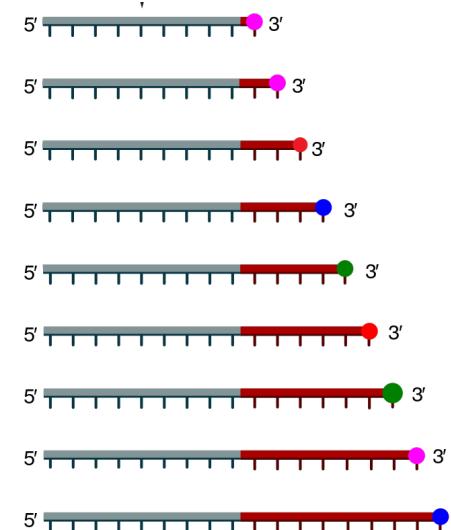
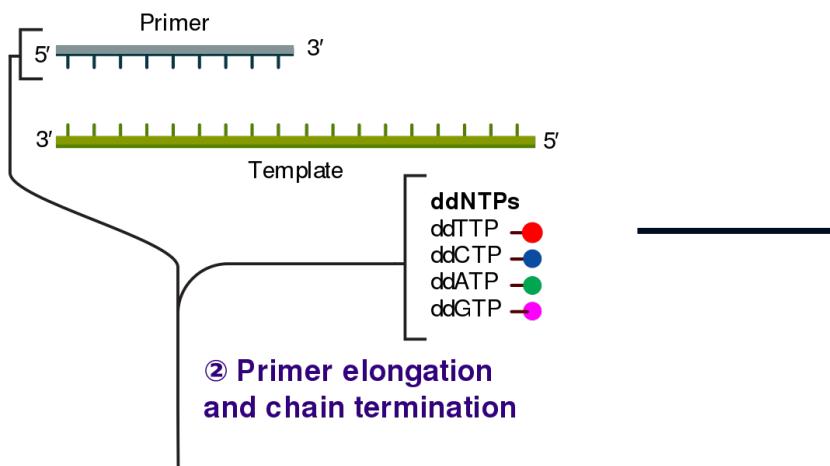
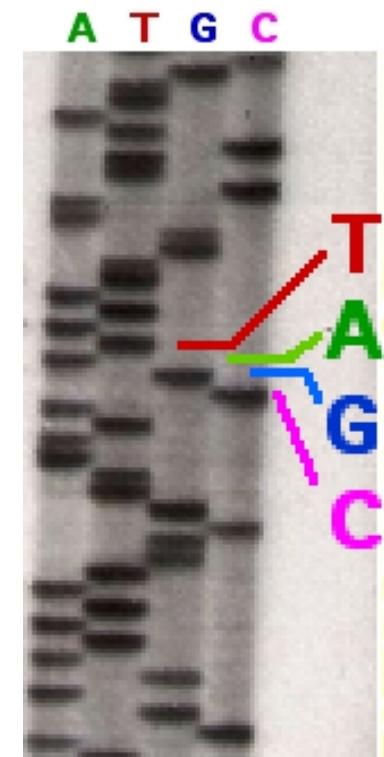


# DNA Sequencing - mid 1970s

Mutations/Variations can be detected using DNA sequencing

## Sanger Sequencing

- Dideoxy DNA sequencing paired with gel electrophoresis
- DNA is 5' labeled with radioactivity
- Small amount of Dideoxy base added to 4 separate primer extension reactions
- Run on a gel to determine bases at each position by size
- Still considered the gold standard for validating sequencing data



# History of Molecular Diagnostics

## The Molecular Biology Timeline

1865	Gregor Mendel, Law of Heredity
1866	Johann Miescher, Purification of DNA
1949	Sickle Cell Anemia Mutation
1953	Watson and Crick, Structure of DNA
1970	Recombinant DNA Technology
1977	DNA sequencing
1985	<i>In Vitro</i> Amplification of DNA (PCR)



# The PCR Revolution

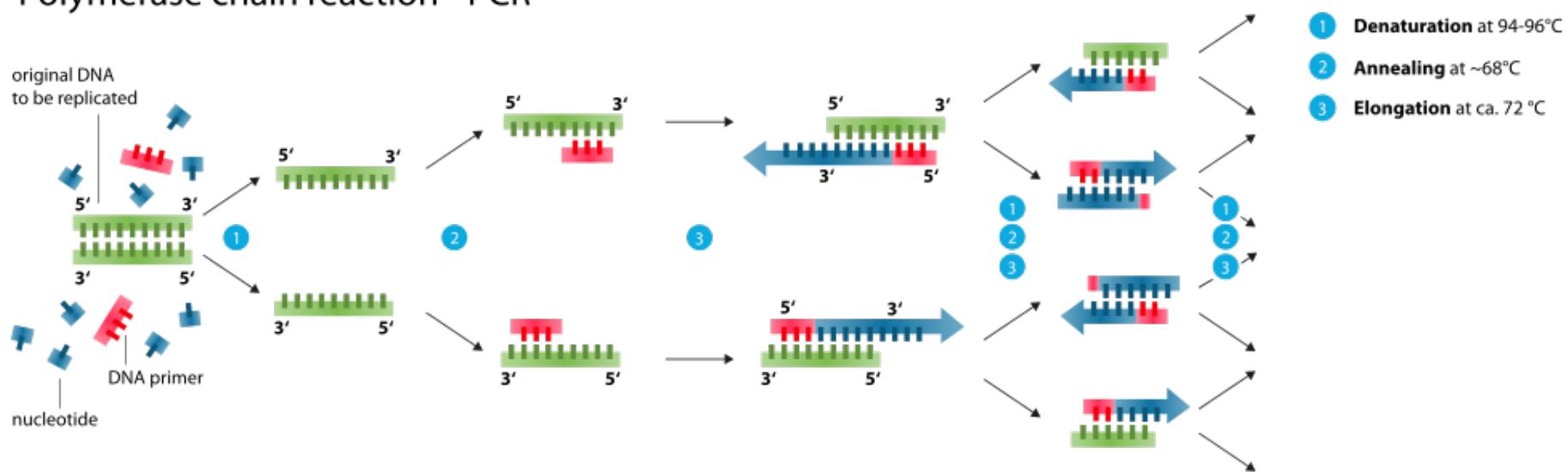
Kary Mullis

- 1983 - Invention of PCR
- 1993 - Received the Noble Prize



Driving late one night, when he had the idea to **use a pair of primers to bracket the desired DNA sequence** and to copy it using DNA polymerase

## Polymerase chain reaction - PCR



# History of Molecular Diagnostics

## The Molecular Biology Timeline

1865	Gregor Mendel, Law of Heredity
1866	Johann Miescher, Purification of DNA
1949	Sickle Cell Anemia Mutation
1953	Watson and Crick, Structure of DNA
1970	Recombinant DNA Technology
1977	DNA sequencing
1985	<i>In Vitro</i> Amplification of DNA (PCR)
2001	The Human Genome Project



# Human Genome Project

U.S. Government project coordinated by the Dept. of Energy and NIH

Total cost of over 3 billion \$

## Goals of the Human Genome Project (1990–2006)

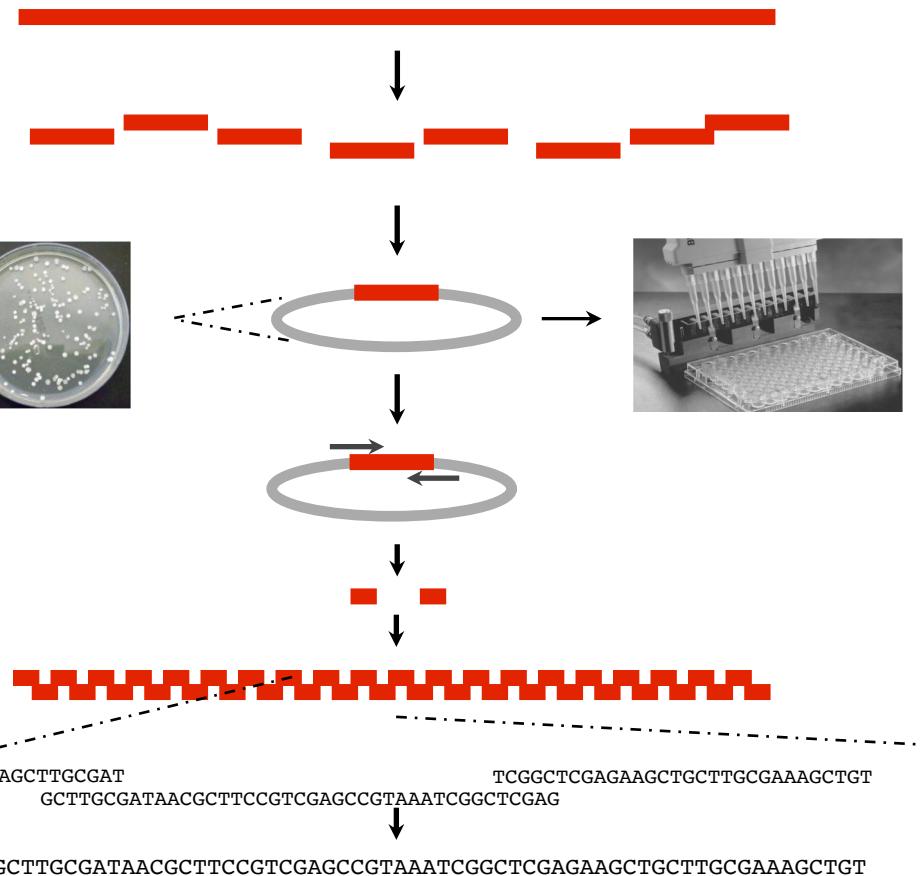
- Identify all of the genes in human DNA
- Determine the sequences of the 3 billion bases that make up human DNA
- Create databases
- Develop tools for data analysis
- Address the ethical, legal, and social issues that arise from genome research

## More information: Milestones in Genomic Sequencing

<https://www.nature.com/immersive/d42859-020-00099-0/index.html>

# Shotgun genome sequencing (Sanger, 1979)

1) Fragment the genome (or large BAC clones)



1977: Bacteriophage *fx* 174 (5kb)  
1995: *H. Influenza* (1Mb);  
1996: Yeast (12mb);  
2000: *Drosophila* (165Mb);  
2002: Human (3Gb)

# Technological advances of Sanger sequencing

**1977:** Fred Sanger



700 bases per day  
→ 118,000 years to  
sequence the human genome

**1985:** ABI 370 (first  
automated sequencer)



5000 bases per day  
→ 16,000 years

**1995:** ABI 377 (Bigger gels,  
better chemistry & optics, more  
sensitive dyes, faster computers)



19,000 bases per day  
→ 4,400 years

**1999:** ABI 3700 (96 capillaries,  
96 well plates, fluid handling  
robots)



400,000 bases per day  
→ 205 years

*To sequence with a coverage of 10X*

## Sequencing duration - calculation

<b>Genome size</b>	<b>Coverage</b>	<b>Throughput</b>	<b>Days</b>	<b>Years</b>
3.000.000.000	10	400.000	75.000	205
3.000.000.000	10	19000	1.578.947	4.326
3.000.000.000	10	5000	6.000.000	16.438
3.000.000.000	10	700	42.857.143	117.417

# Impact on Human Diseases

## Novelty

- Discovery of potential **novel molecular markers** of human diseases
- Utility of molecular markers to **develop useful molecular assays** for detection, diagnosis, and prediction of disease outcomes

## Monitor diseases more accurately

- Allows for early treatment and better patient care

## Determine most appropriate treatment

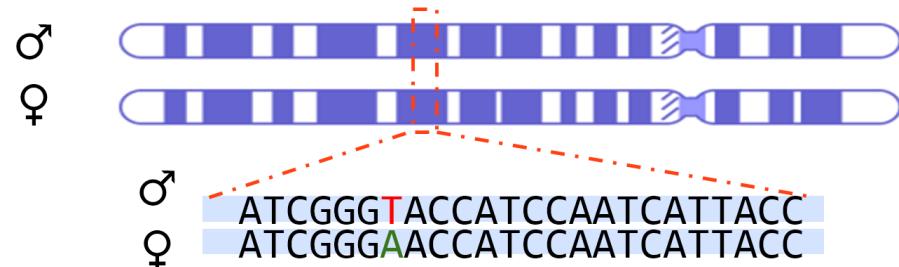
- Reduces or eliminates unnecessary treatment
- Reduces or eliminates inadequate treatment
- Yields greater cost effectiveness

## Reduce patient morbidity and mortality

# Terms

# Genetics Terms

**Humans are diploid:** Our genome is comprised of a paternal and a maternal "haplotype" → form our "genotype"



**Gene:** Any interval along the chromosomal DNA that is transcribed and then translated into a functional protein, or that is transcribed into a functional RNA molecule

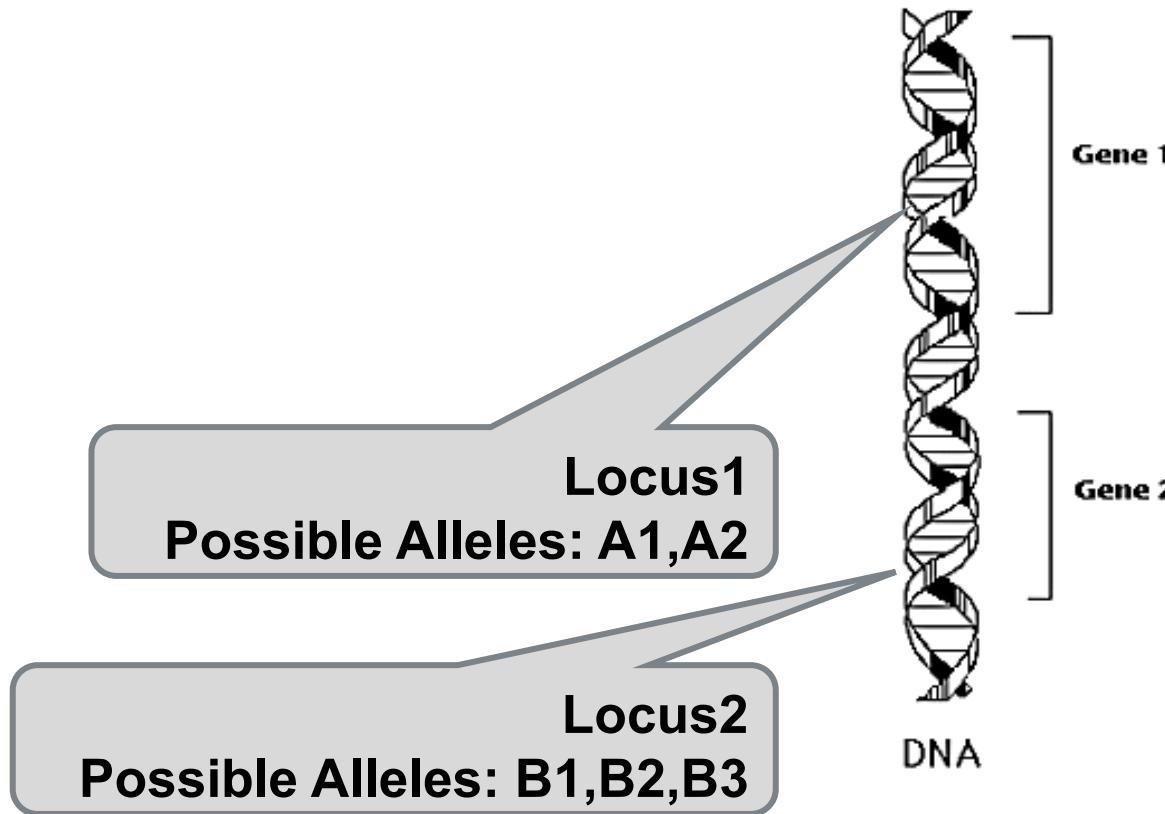
**Genotype:** Genetic makeup, **distinguished** from the physical appearance

**Phenotype:** The observable physical or biochemical characteristics as determined by both genetic makeup and environment

# Genetics terms

**Locus** location of a *gene/marker* on the chromosome

**Allele** one **variant form** of a gene/marker at a particular locus  
differ in their nucleotide sequence



# The human genome - basic stats

- ~ 3 billion base pairs (haploid)
- ~ 20,000 protein coding genes (RefSeq)
- ~ 200,000 coding transcripts (isoforms of a gene that each encode a distinct protein product)

**Table 1.** The number of human genes and transcripts in the RefSeq database and in the new CHESS (Comprehensive Human Expressed SequenceS) database built from 9,795 RNA-seq experiments. ncRNA: noncoding RNA; lncRNA: long noncoding RNA gene; miscRNA: miscellaneous RNA.

Type of gene	Number in RefSeq	Number in CHESS
Protein-coding genes	20,054	21,306
ncRNA genes		
- lncRNA	14,788	18,484
- antisense	23	2,144
- miscRNA	1,217	1,228
Total gene counts	<b>36,082</b>	<b>43,162</b>
Protein coding transcripts	127,718	267,476
ncRNA transcripts		
- lncRNA	28,015	49,307
- antisense	28	2,694
- miscRNA	2,005	4,347
Total transcripts	<b>157,766</b>	<b>323,824</b>

# Forms of genetic variations

## Single nucleotide substitution

Replacement of one nucleotide with another

(Recap) Tandem repeat  
ATTCG ATTCG ATTCG

## Microsatellites or mini-satellites

These tandem repeats often present high levels of inter- and intra-specific polymorphism

## Deletions or insertions

Loss or addition of one or more nucleotides

## Structural variations

Changes in chromosome number, segmental rearrangements and deletions

# Genetics Terms

## Polymorphism

- Variations in DNA sequence (substitutions, deletions, insertion, etc) that are present at a frequency greater than 1% in a population.
- Have a WEAK EFFECT or NO EFFECT at all.
- Ancient and COMMON

## Mutation

- Variations in DNA sequence (substitutions, deletions, etc) that are present at a frequency lower than 1% in a population.
- Can produce a gain of function and a loss of function.
- Recent and RARE.

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4502642/>

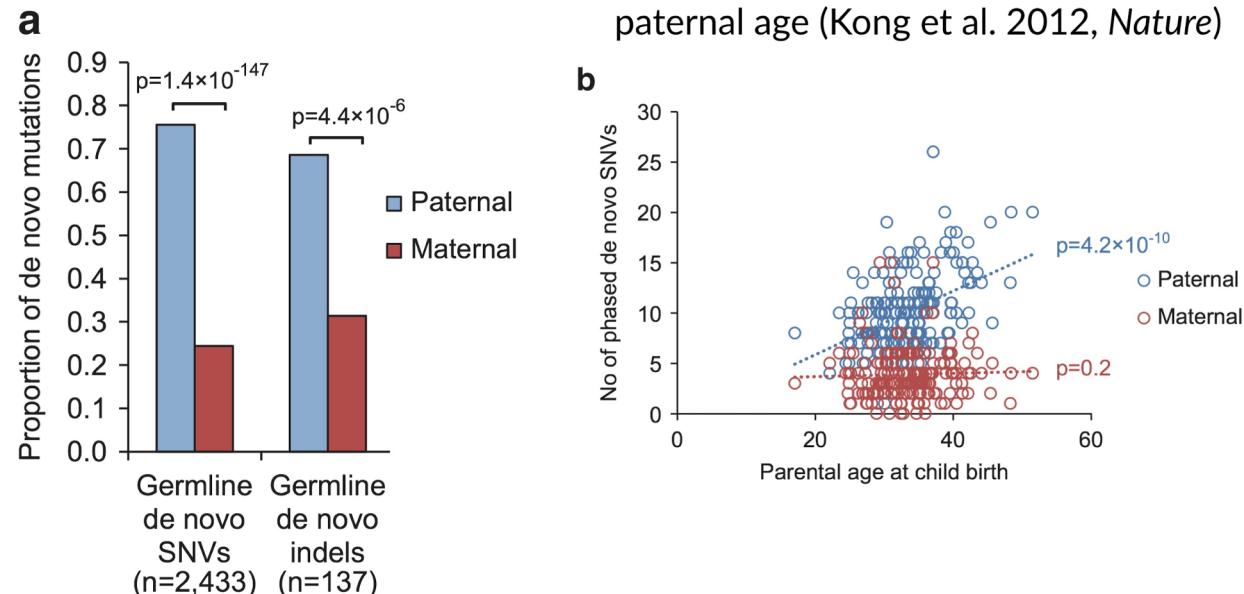
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4959039/>

# De-novo mutations (DNM)

## New mutation occurs in father's or mother's germ cell

- 28 new mutations (per haploid genome) -> 56 diploid genome  
*Roach et al. (2010) Science*
- 175 mutations -- *Nachman et al. (2000) Genetics*
- Gene4Denovo: an integrated database and analytic platform for de novo mutations in humans --- *Zhao et al (2020) NAR*

**DNMs are more likely to occur in the paternal germline & correlate with paternal age**



# Glossary and Definitions

## **Homozygote**

An organism with two identical alleles

## **Heterozygote**

An organism with two different alleles

## **Hemizygote**

Having only one copy of a gene

→ males are hemizygous for most genes on the sex chromosomes

## **Dominant trait**

A trait that shows in a heterozygote

## **Recessive trait**

A trait that is hidden in a heterozygote

# Mutations

# Humans

In human beings, **99.9% bases are same**

**Remaining 0.1% makes a person unique**

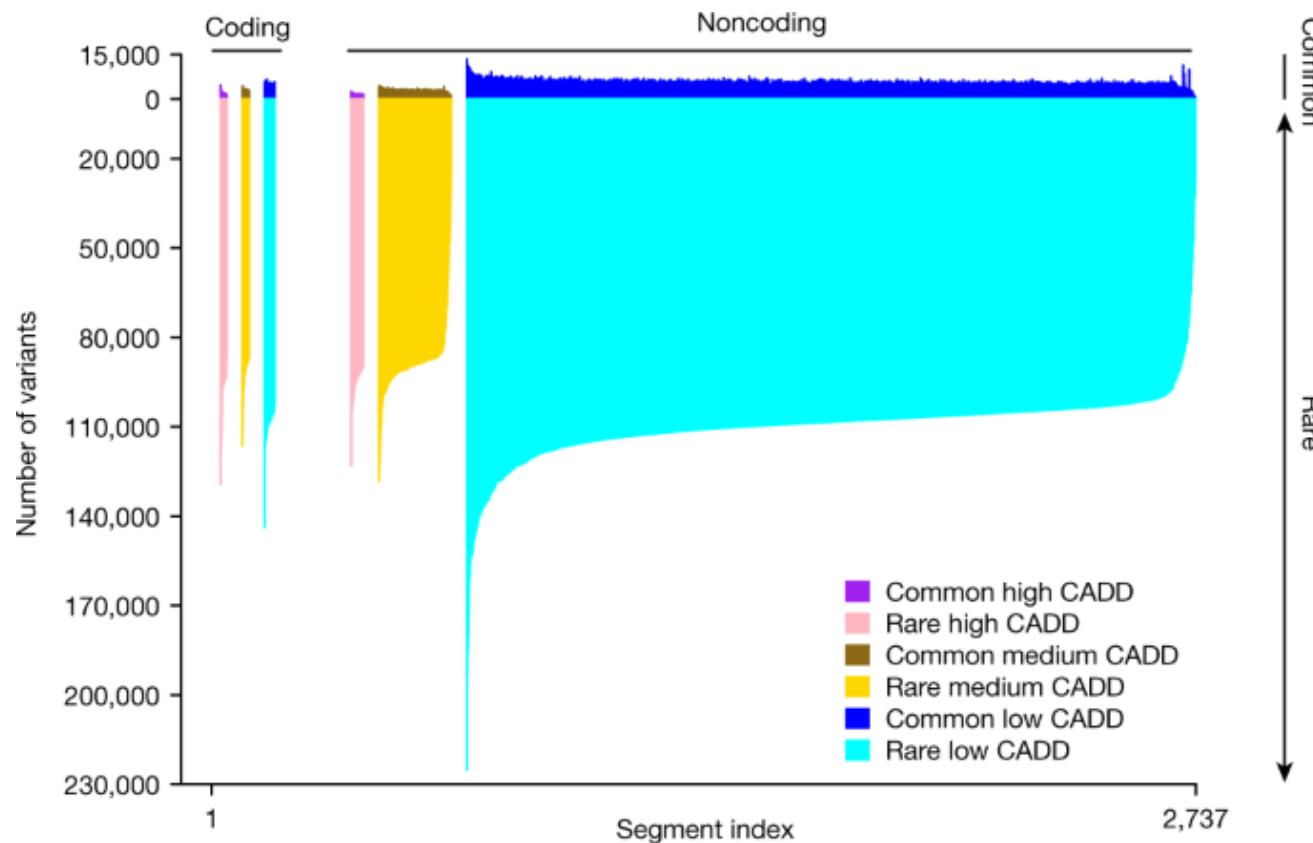
- Different attributes / characteristics / traits
- How a person looks
- Diseases

**These variations can be:**

- Harmless (change in phenotype)
- Harmful (diabetes, cancer, heart disease, Huntington's disease, hemophilia)
- Latent (e.g. susceptibility to heart attack; **variations** found in coding and regulatory regions, are **not harmful on their own**, and the change in each gene only **becomes apparent under certain conditions**)

# TOPMed Program (2021)

- Sequencing of 53,831 diverse genomes
- Identified more than 400 million genetic variations
  - Very rare --> less than 1 percent of the population
  - 78% of which had not been described before



# Consequences of mutations

## Most mutations are neutral

- 97% DNA neither codes for protein or RNA, nor indirectly affects gene function
- A new variant in the 1.5% coding regions may not result in a change in amino acid
- Variants that change amino acid may not affect function

## Certain mutations have functional effect and even cause disease

- Gain-of-function mutations often produce dominant disorders
- Loss-of-function mutations result in recessive disease

# Sequencing

# Second Generation Sequencing

- Developed to **increase throughput of Sanger sequencing**
- Can sequence **many molecules in parallel**
  - Does not require homogenous input
  - DNA sequenced as clusters or in nanowells
  - Single machine can sequence 3-10 Billion independent DNA fragments at once
  - Single Sanger Sequencer maxes out at 1152 reactions per machine
- Time from DNA to genome reduced from 10 years to 1 day!



# Disadvantages of 2nd Generation Tech

## Rely on amplification to create libraries and clusters

- All polymerases have an inherent error rate ( $10^{-6}$ - $10^{-7}$ )
- Errors introduced every 10 million to 100 million bases
- Secondary validation of variants is key

## GC bias

- PCR bias against GC rich sequences
- Exome capture bias against GC rich sequences

## Short reads can miss large structural variations

- Genome Translocations and inversions likely will be missed
- Require significant read depth at break points for these variations to be detected

## Trouble detecting small insertions and deletions

- Short reads computationally hard to align and call

# Third Generation Sequencing

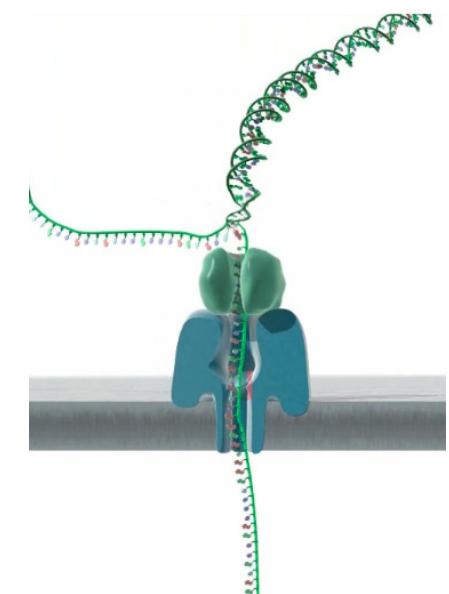
## Single molecule sequencing

- Less complex sample prep
- Much longer read length (1-100kb)
- Many technical hurdles with very high error rates
- Expensive



## Sequencing by synthesis - Pacific Biosciences

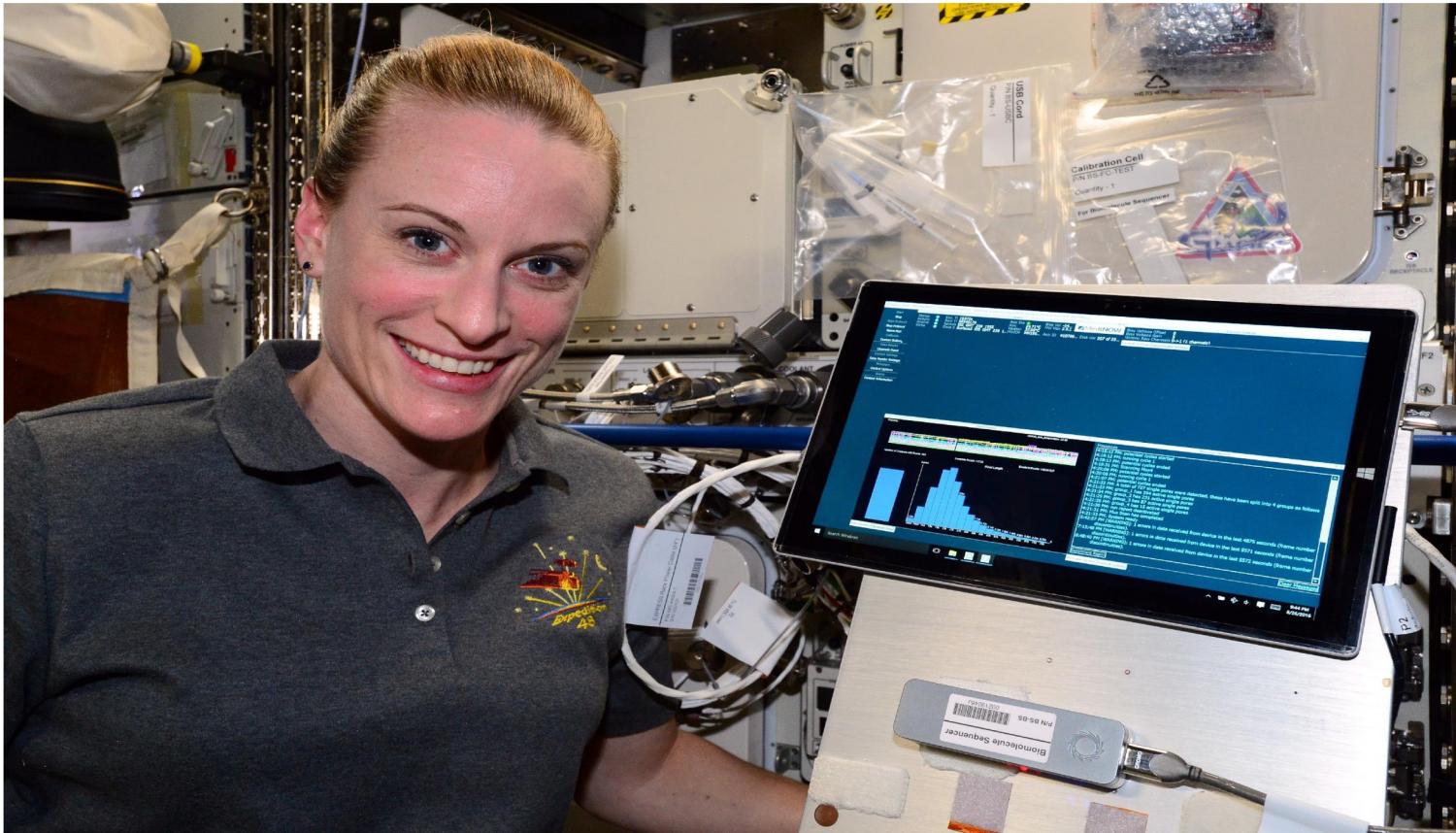
- 1 DNA molecule and 1 polymerase in each well (zero-mode waveguide)
- 4 different marker on the phosphate of the nucleotide  
→ Polymerase interacts; marked freed
- No “theoretical” limit to DNA fragment length



## Direct sequencing by passing DNA through a nanopore

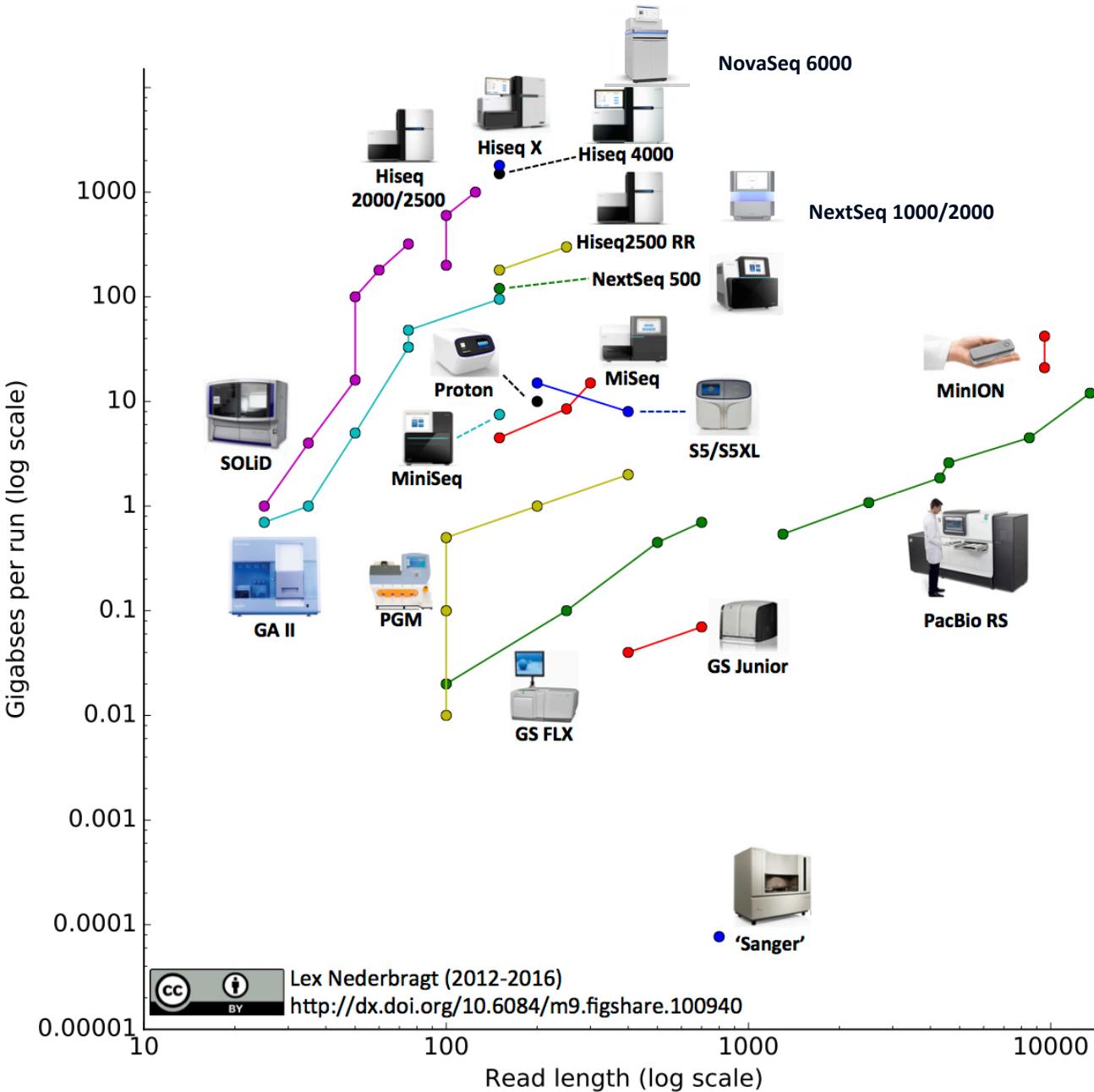
- Bases fed through a membrane bound nanopore
- Ionic difference between both sides of the membrane

# Nanopore – very portable



Kate Rubins sequencing DNA on the ISS

# High throughput sequencing



# Properties of different technologies

Instrument	Amplification	Run time	Bases / read	Gbp/run
Illumina MiSeq	BridgePCR	5-55h	50-600	0.3-13.2
Illumina NextSeq 500	BridgePCR	11-30h	75-300	19.5-120
Illumina HiSeq 2500	BridgePCR	10h - 11days	50-300	15-500
Illumina NovaSeq 6000	BridgePCR	44h	150	Up to 3TB
Ion Torrent - PGM	emPCR	2-7h	200-400	0.095-1.9
Ion Torrent - Proton	emPCR	4-6h	175	12.25-87.5
Pacific Biosciences RS II	None	2 hrs.	3000	0,09
Oxford Nanopore MinION (forecast)	None	≤6 hrs.	9000	0,9

<http://www.molecularecologist.com/next-gen-fieldguide-2016/>

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7649788/>

# Error Rates

<b>Instrument</b>	<b>Primary Errors</b>	<b>Single-pass Error Rate (%)</b>	<b>Final Error Rate (%)</b>
Illumina	Substitutions	~0.1	~0.1
Ion Torrent	INDELs	~1	~1
Oxford Nanopore	Deletions	≥4	4
PacBio RS	INDELs	~13	≤1

# Flavors of Sequencing

## Whole Genome Sequencing

- Obtain whole blood or tissue sample
- Create sequencing libraries of all DNA fragments

## Whole Exome Sequencing

- Utilizes a selection protocol
- Attach complimentary RNA or DNA strands to beads
- Fish out **only** coding DNA sequences
- Create sequencing libraries from enriched DNA
- Reduces cost and analysis time

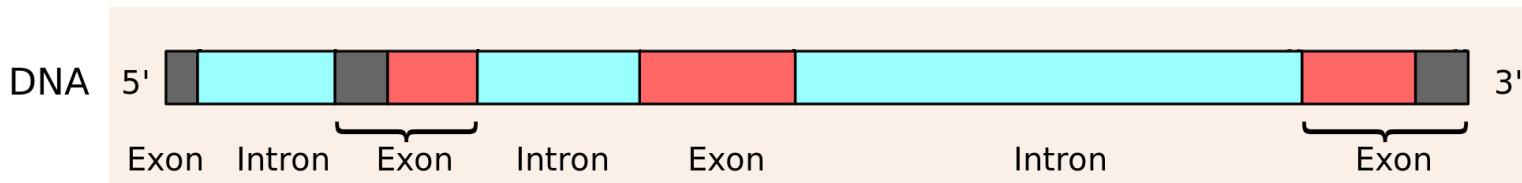
## Custom Capture

- Same protocol as Exome sequencing
- Only target desired DNA sequences

## Amplicon Sequencing

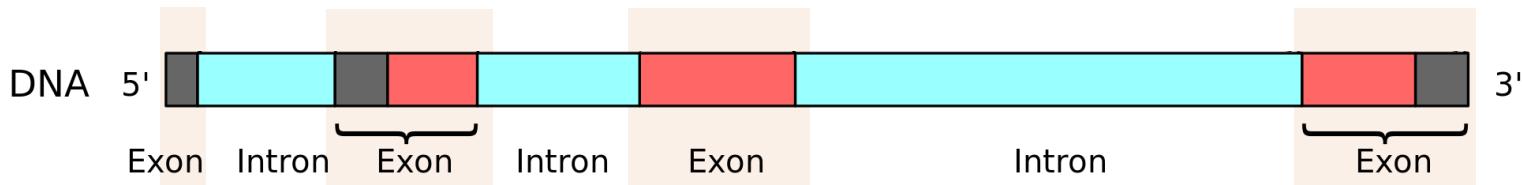
- Use PCR to amplify target DNA
- Sequence amplified DNA (Amplicon)

# Sequencing Techniques



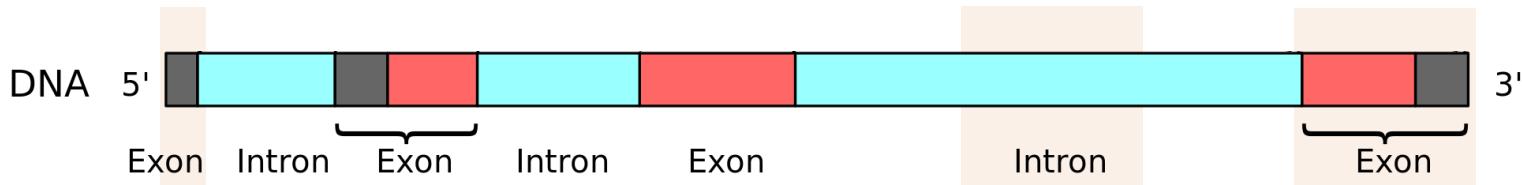
- Whole genome sequencing

# Sequencing Techniques



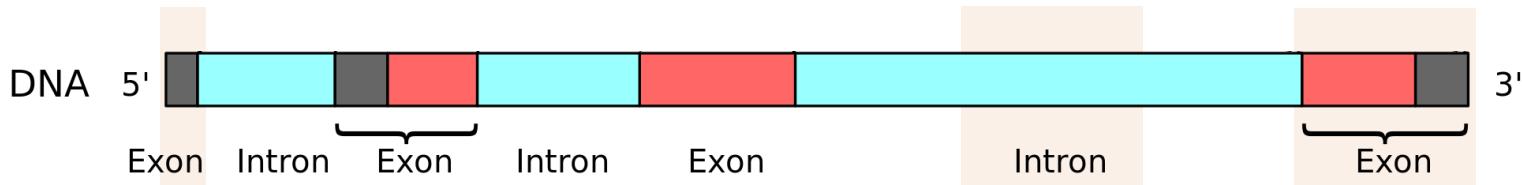
- Whole genome sequencing
- Whole exome sequencing
- Custom capture

# Sequencing Techniques



- Whole genome sequencing
- Whole exome sequencing
- Custom capture
- Amplicon sequencing

# Sequencing Techniques



- Whole genome sequencing
- Whole exome sequencing
- Custom capture
- Amplicon sequencing

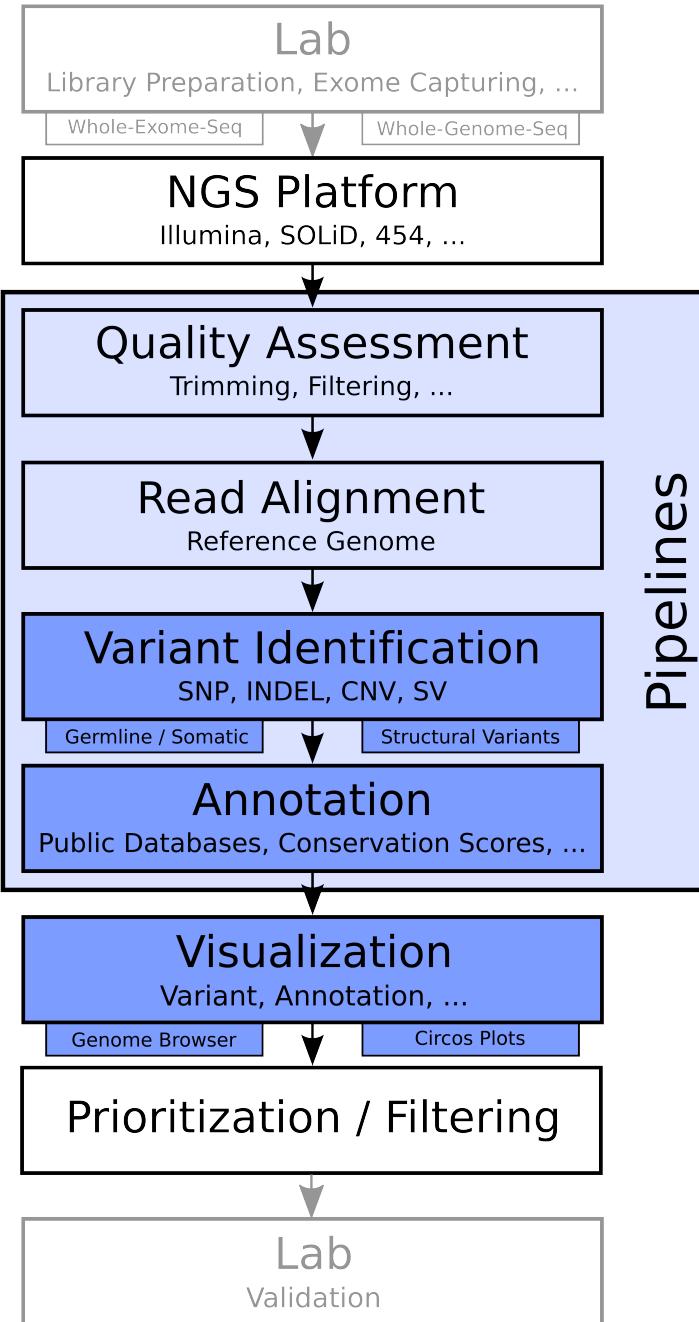
What is the best technology for my use-case?

- Clinical question?
- Number of samples?
- Cost?
- Future strategies?

# NGS Data analysis

- Absolutely the largest roadblock for next generation sequencing
- Terabytes of data are useless if we can't efficiently analyze the data
- How long should data be kept?
  - Depends on application
  - Human Diagnostic sequencing?
  - Research sequencing?
- Where should data be kept and processed?
  - Local or Cloud (Amazon, etc)?
  - Cost of infrastructure vs cost of cloud service
  - Security issues
- Future
  - Cloud based solutions will become more attractive

## Workflow overview



# Genome reference

# Genome reference

Reference genome is a “consensus” across all chromosomes of DNA pooled from multiple individuals

UCSC and Genome Reference Consortium (GRCh)

hg18, hg19, hg38 ↔ GRCh36, GRCh37, GRCh38

Newest version (hg38) release on Dez 24<sup>th</sup> 2013



	HG38 (UCSC)	GRCh38
Prefix	Chr	-
Mitochondrial	chrM	MT
Order	chrM, chr1, chr2, ...chrX, chrY	1,2, ..., X, Y, MT

# Genome reference

## Indexing

- Fai file (created by samtools faidx)  
contig, size, location, bases-per-line and for efficient random
- Dict file (created by Picard CreateSequenceDictionary)  
SAM style header describing the contents of the FASTA file
- Different mapping programs

## Important

- Choose one reference genome (well sorted, indexed) and stick to it
- Be sure that previous variant calls use same reference - otherwise convert coordinates (lift-over)

<http://www.broadinstitute.org/gatk/guide/best-practices?bpm=DNaseq#data-processing-ovw>

FASTQ

# FASTQ

Storing and defining sequences from next-generation sequencing technologies

Sequence ID

@SEQ\_ID

Sequence

GATTGGGTTCTATCGGATCTCCAAAGCCTAGATGCCCATCG

Separator

+

Quality Score

! ’ ’ \* ( ( ( \*\*\*+ ) ) % % + + ( ) % % % ) + ‘ ’ + \* \* >> CCCCCCCC65

## Old format

@HWUSI-EAS100R:6:73:941:1973#0/1

HWUSI-EAS100R	the unique instrument name
6	flowcell lane
73	tile number within the flowcell lane
941	'x'-coordinate of the cluster within the tile
1973	'y'-coordinate of the cluster within the tile
#0	index number for a multiplexed sample (0 for no indexing)
/1	the member of a pair, /1 or /2 ( <i>paired-end or mate-pair reads only</i> )

## New format

@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	'x'-coordinate of the cluster within the tile
197393	'y'-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 ( <i>paired-end or mate-pair reads only</i> )
Y	Y if the read is filtered, N otherwise
18	0 when none of the control bits are on, otherwise it is an even number
ATCACG	index sequence

# Phred quality score

Characterize the quality of DNA sequences

$$q = -10 \log_{10}(p)$$

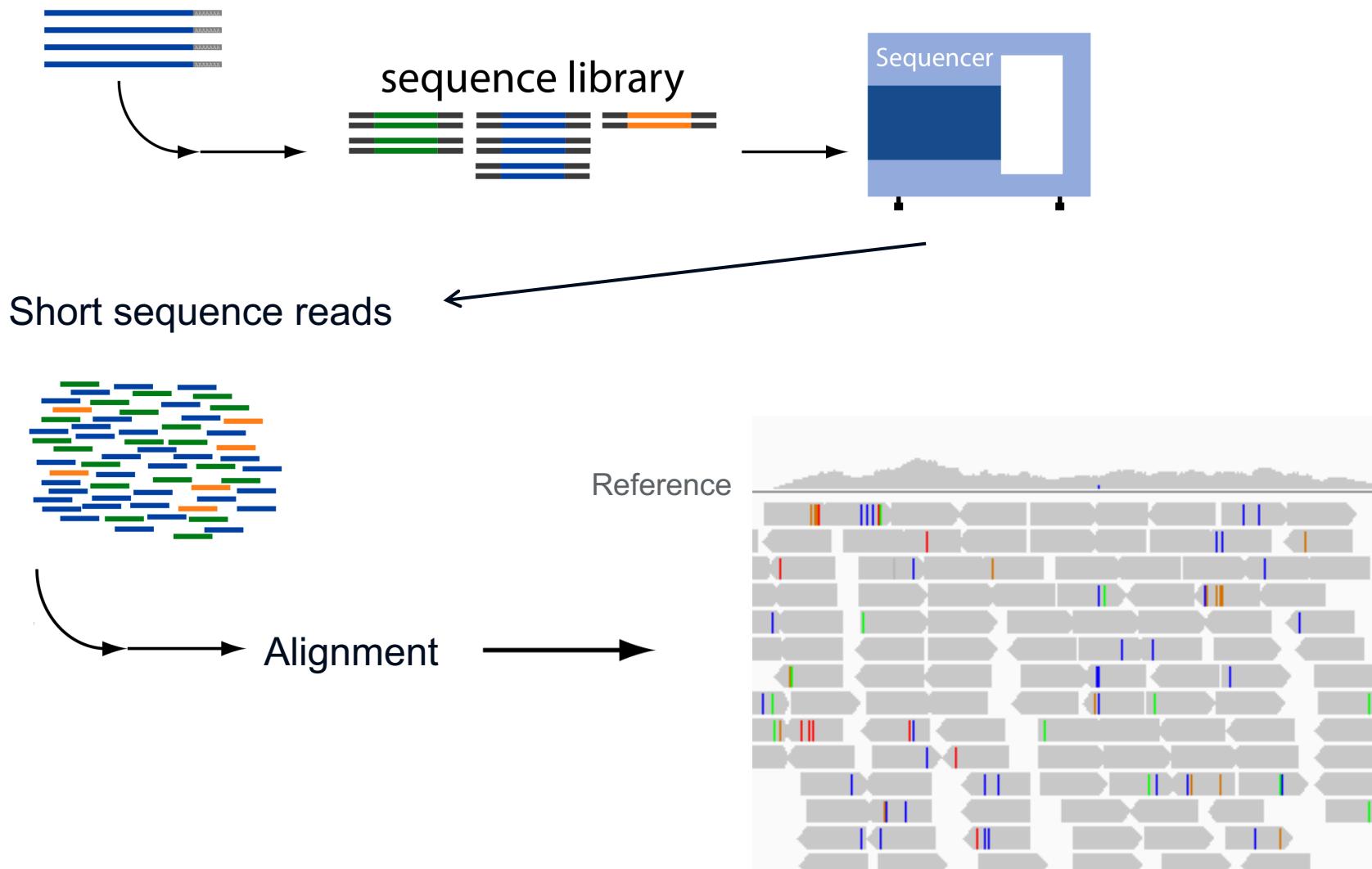
p = error probability for the base

Phred quality score	Probability	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

## Mapping and QC

SAM - Sequence Alignment/Map Format

# Mapping - Principle



Adapted from <http://raetschlab.org//members/research/transcriptomics/images/RNA-Sequencing.png>

Sequencing Reads



Alignment



SAM file

# Sequence **mapping** versus **alignment**

**Mapping:** (quickly) find the best possible loci to which a sequence could be aligned

**Alignment:** for each locus to which a **sequence can be mapped**, determine the **optimal base by base alignment** of the query sequence to the reference sequence

## Tools:

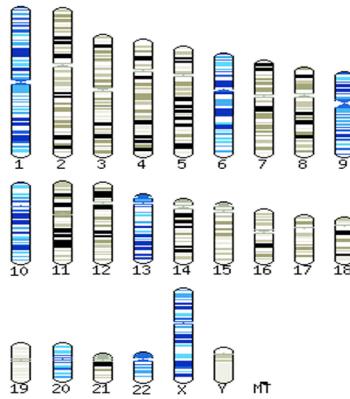
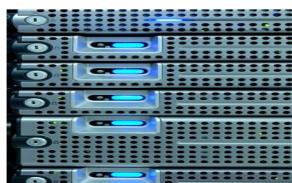
- BWA
- Bowtie2

# Paired-end sequencing

## 1) A read-pair

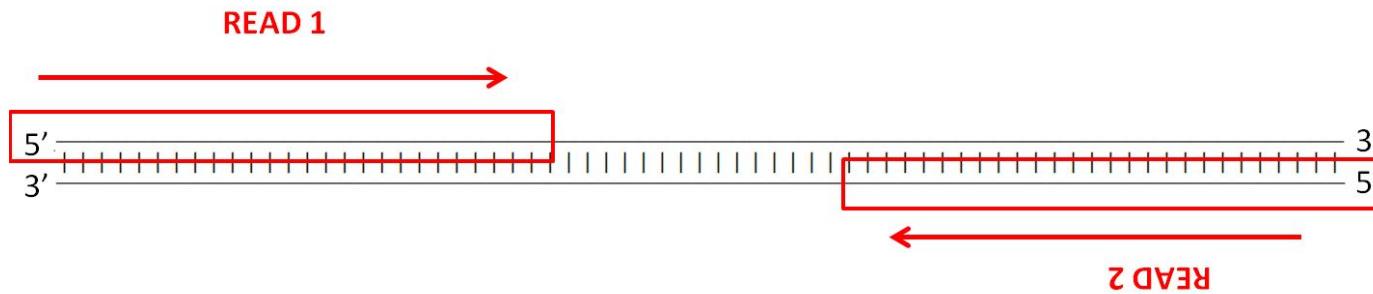
5' GGTGTACGAATAGTTCTTTACCTCCTGACCATCCTAGC ----- //----- GGACTGAAACTCATCTGTCTTATAGATATGCGTCAGCAGC 5'

## 2) A reference genome on a computer



TATGGCAATTAAAATTGGTATCAAATGGTTGGCTGTATGCCGTATCGTATTCCGTGAGCACACAC  
CGTGATGACATITGAAGTTGAGGTATAACGACTTAATCGACGTTGAATACTGGCTTATATGTTGAAT  
ATGATCAACTCAGCTACGGCTGTTGAGGACTGTTGAGGATGGTAAAGATGGTAAATGGTAA  
AACTATCGGTGTAACTCAGAACCGTGATCCAGCAAACCTTAAACTGGGTGCAATCGGTGTTGATATCGCT  
GTTGAAGCCACTGGTTTATCTTAACTGATGAAACTCTCGTAAACATATACCTGAGGGCCAAAAAAG  
TTGTATTAACCTGCCCATCTAAAGATGCAACCCCTATGTTGGTTGTAACCTCAACGCATAGCC  
AGGTCAGATATCGTTCTAACGCACTTGTACAACAAACTGTTTACGCTCTTACGCACTGTTGTCAT  
GAAACTTTCGGTATCAAAGATGGTTAATGACCACTGTCAAGCAACGACTGCAACTCAAACACTGTTG  
ATGGTCCATGCTAAAGACTGGCGGGCGCGCGTGCATCACAAACATCATTCCATCTTCAACAGG  
TATGGCAATTAAAATTGGTATCAAATGGTTGGCTGATCGGCGTATCGTATTCCGTGAGCACACAC  
CGTGTGACATTTGAGTTGAGGTATAACGACTTAACTGACGTTGAAATACCTGCGTTATATGTTGAAAT  
ATGATCAACTCAGCTGCTTGGACGGCAGTGGAGGATGGTAACTTGTGTTAATGGTAA  
AACTATCGGTAAATGCGAACCGTGATCCAGCAAACCTTAAACTGGGTGCAATCGGTGTTGATACTCGCT  
GTTGAAGCCACTGGTTTATCTTAACTGATGAAACTCTGCGTAAACATATACCTGAGGGCCAAAAAAG  
TTGTATTAACCTGCCCATCTAAAGATGCAACCCCTATGTTGGTTGTAACCTCAACGCATAGCC  
AGGTCAGATATCGTTCTAACGCACTTGTACAACAAACTGTTGAGCTGAGCAGACTGCAACTCAAACACTGTTG  
GAAACTTTCGGTATCAAAGATGGTTAATGACCACTGTCAAGCAACGACTGCAACTCAAACACTGTTG  
ATGGTCCATGCTAAAGACTGGCGGGCGCGCGTGCATCACAAACATCATTCCATCTTCAACAGG

## 3) Alignment of the read-pair to the reference genome gives coordinates describing where in the human genome the read-pair came from



# SAM file format

The Sequence Alignment Map (SAM) Format and SAMtools

Heng Li et al. Bioinformatics, 2009

Tab-delimited text file

Output of most alignment programs

- Header section
- Alignment section
- 11 Required columns
- Optional fields

<https://samtools.github.io>

# SAM Header

- Header lines start with @ symbol
- Always at top of file
- Contain lots of information about what was mapped, what it was mapped to, and how (metadata)
  - The version information for the SAM/BAM file
  - Whether or not and how the file is sorted
  - Information about the reference sequences
  - Any processing that was used to generate the various reads in the file
  - Software version

# SAM

## 1.4 The alignment section: mandatory fields

In the SAM format, each alignment line typically represents the linear alignment of a segment. Each line has 11 mandatory fields. These fields always appear in the same order and must be present, but their values can be ‘0’ or ‘\*’ (depending on the field) if the corresponding information is unavailable. The following table gives an overview of the mandatory fields in the SAM format:

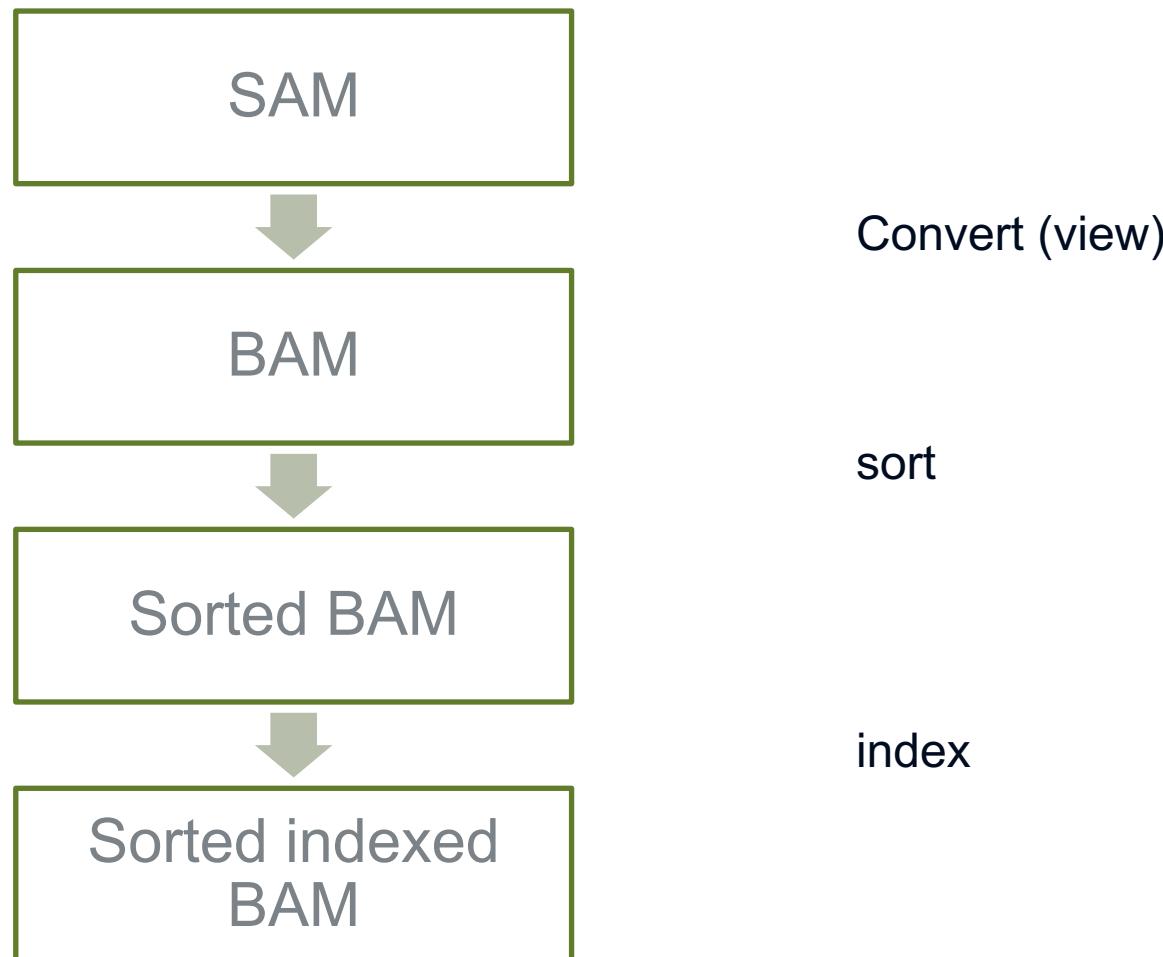
Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,255}	Query template NAME
2	FLAG	Int	[0,2 <sup>16</sup> -1]	bitwise FLAG
3	RNAME	String	\*  [!-()+-<>-~] [!-~]*	Reference sequence NAME
4	POS	Int	[0,2 <sup>31</sup> -1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0,2 <sup>8</sup> -1]	MAPping Quality
6	CIGAR	String	\*  ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	\* =  [!-()+-<>-~] [!-~]*	Ref. name of the mate/next read
8	PNEXT	Int	[0,2 <sup>31</sup> -1]	Position of the mate/next read
9	TLEN	Int	[-2 <sup>31</sup> +1,2 <sup>31</sup> -1]	observed Template LENgth
10	SEQ	String	\*  [A-Za-z.=.]+	segment SEQuence
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33

# Example SAM file

```
@PG      ID:bwa_sam      PN:bwa    PP:bwa_aln_fastq      VN:0.5.9-r16    CL:bwa sampe -a 1050 -r $rg_line -f $sam_file $reference fasta $sai_file(s) $fastq_file(s)
@PG      ID:sam_to_fixed_bam    PN:samtools    PP:bwa_sam      VN:0.1.17 (r973:277)    CL:samtools view -bSu $sam_file | samtools sort -n -o - samtools_nsort_tmp | :
ort_tmp | samtools fillmd -u - $reference fasta > $fixed_bam_file
@PG      ID:gatk_target_interval_creator PN:GenomeAnalysisTK    PP:sam_to_fixed_bam      VN:1.2-29-g0acaf2d    CL:java $jvm_args -jar GenomeAnalysisTK.jar -T Realig:
dels_file(s)
@PG      ID:bam_realignment_around_known_indels PN:GenomeAnalysisTK    PP:gatk_target_interval_creator VN:1.2-29-g0acaf2d    CL:java $jvm_args -jar GenomeAnalysis:
bam_file -targetIntervals $intervals_file -known $known_indels_file(s) -LOD 0.4 -model KNOWNS_ONLY -compress 0 --disable_bam_indexing
@PG      ID:bam_count_covariates PN:GenomeAnalysisTK    PP:bam_realignment_around_known_indels VN:1.2-29-g0acaf2d    CL:java $jvm_args -jar GenomeAnalysisTK.jar -:
recal_data.csv -knownSites $known_sites_file(s) -l INFO -L '1;2;3;4;5;6;7;8;9;10;11;12;13;14;15;16;17;18;19;20;21;22;X;Y;MT' -cov ReadGroupCovariate -cov QualityScore:
@PG      ID:bam_recalibrate_quality_scores      PN:GenomeAnalysisTK    PP:bam_count_covariates VN:1.2-29-g0acaf2d    CL:java $jvm_args -jar GenomeAnalysisTK.jar -:
.csv -I $bam_file -o $recalibrated_bam_file -l INFO -compress 0 --disable_bam_indexing
@PG      ID:bam_calculate_bq      PN:samtools    PP:bam_recalibrate_quality_scores      VN:0.1.17 (r973:277)    CL:samtools calmd -Erb $bam_file $reference fasta > $:
@PG      ID:bam_merge      PN:picard      PP:bam_calculate_bq      VN:1.53 CL:java $jvm_args -jar MergeSamFiles.jar INPUT=$bam_file(s) OUTPUT=$merged_bam VALIDATION_ST:
@PG      ID:bam_mark_duplicates PN:picard      PP:bam_merge      VN:1.53 CL:java $jvm_args -jar MarkDuplicates.jar INPUT=$bam_file OUTPUT=$markdup_bam file ASSUME_SOR:
@PG      ID:bam_merge.1  PN:picard      PP:bam_mark_duplicates VN:1.53 CL:java $jvm_args -jar MergeSamFiles.jar INPUT=$bam_file(s) OUTPUT=$merged_bam VALIDATION_ST:
@CO      $known_indels_file(s) = ftp://ftp.1000genomes.ebi.ac.uk/voll/ftp/technical/reference/phase2_mapping_resources/ALL.wgs.indels.mills_devine hg19_leftAligned.co:
@CO      $known_indels_file(s) . = ftp://ftp.1000genomes.ebi.ac.uk/voll/ftp/technical/reference/phase2_mapping_resources/ALL.wgs_low_coverage_vqsr.2010123.indels.site:
@CO      $known_sites_file(s) = ftp://ftp.1000genomes.ebi.ac.uk/voll/ftp/technical/reference/phase2_mapping_resources/ALL.wgs.dbsnp.build135.snp.sites.vcf.gz
SRR070823.24480225 163 11 60464 0 100M = 60720 355 CCATGGTTAACATAAATGCAAATATGTTACTGAATAACTTATCTGTGCCAGTGGTGTATTAAATGATTG:
MNKJKNKNJNNNMGNNDNGNGKKGNKNKKJNMILLLCFLKBLHJGFHED X0:i:3 X1:i:6 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@Ioooooooooooo
@oooooooooooo@GI
SRR070823.24480518 1187 11 60464 0 100M = 60720 355 CCATGGTTAACATAAATGCAAATATGTTACTGAATAACTTATCTGTGCCAGTGGTGTATTAAATGATTG:
DJKBEIJBEFJFICJKGKGKHIG@IHHDHNLCJLIAKLLIKMIJLFBEE?EB X0:i:3 X1:i:6 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@Eoooooooooooo
@oooooooooooo@GG
SRR070823.24480225 83 11 60720 0 100M = 60464 -355 TATGGAGTTTGATGTTATGTCAGGGTAATTACATGATTATAATTAAACAGGTTCTTTAAATCAGCTATATCAA:
GHHKKJJJJGJJJJGHGGJJJJGGGGJGGGJIHGGHGGGGHIEJIIEDCF6 X0:i:9 X1:i:0 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@oooooooooooo
@oooooooooooo@OO
SRR070823.24480518 1107 11 60720 0 100M = 60464 -355 TATGGAGTTTGATGTTATGTCAGGGTAATTACATGATTATAATTAAACAGGTTCTTTAAATCAGCTATATCAA:
=EGJJF->?@DBADD>BBB<AD:@2B@BABBC>B:::E@DBDG;6E5=A??@6 X0:i:9 X1:i:0 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@oooooooooooo
@oooooooooooo@OO
SRR070531.23281260 99 11 61942 0 100M = 62035 192 TCCATGCCCTTGATGACAGGATAATATGTAAGCTTTCTATATTTCAGAAACTATATGACATGACGAAAAGTAA:
JJMLKHLIJMJJHKNHFGKNHMGGMJHJFGIGIINOLNIKMENNIFILKGGHEA X0:i:8 X1:i:1 MD:Z:100 RG:Z:SRR070531 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@oooooooooooo
@oooooooooooo@OO
SRR070531.23281260 147 11 62035 0 100M = 61942 -192 GGAGGTATCCTGAATTGACTGAGAAATAAGGAGGTATCCACAGAGAAATATAAAACATATACTTAGTGTGTTCAAG:
MKKJKJKCCKJNKJJJJMLJMLLLLIMJMMLJJIIJIIIFIKKKKHGJECG6 X0:i:8 X1:i:2 MD:Z:100 RG:Z:SRR070531 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@oooooooooooo
@oooooooooooo@OO
SRR070823.17243685 99 11 62388 0 100M = 62452 163 CCTTGCCAATTGTGTTCTCTTATTCTCTGCTGGATATGACCAGTGTGCTTCCATTGCATTGTGTGTTAA:
MGHMMMKIMLKHHJEKKFDHJEKDEIDCCDF?IEGJHIDEHJIGE?GAFDCCC X0:i:10 X1:i:0 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@oooooooooooo
@oooooooooooo@OO
SRR070823.17243685 147 11 62452 0 100M = 62388 -163 ATGTGTTTTAAATAGACTTAATGGTCTCAAGTGTGCAATTAGTTGGTTCTTGGAAACTTATATAATGAA:
MJNMMMLJMLJNLIMKKJFLJGGJJJJLLJJGLIJKKKJKHLIKHKJGJHDA6 X0:i:10 X1:i:0 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@oooooooooooo
@oooooooooooo@OO
SRR070823.1968642 163 11 62725 0 100M = 62918 292 TTTTTAACCATACAAACATGCTGCTATGAACATTCTTTGTAATCACCTGGTTCATATGTGCAAGATATCCTCTI:
LKMLMNNNNHKKMMKJHNNHNMNNHHKNNKNLNNNMGHIFIEIIIGJDHHGCCBD X0:i:9 X1:i:1 MD:Z:100 RG:Z:SRR070823 AM:i:0 NM:i:0 SM:i:0 MQ:i:0 XT:A:R BQ:Z:@EDDB@oooooooo
@oooooooooooo@AF
```

BAM – Binary SAM

# SAM / BAM



# BAM

## SAM

- Information on the alignment of each read
- Optimized for readability and sequential access

## BAM (Binary SAM)

- Compressed -> saves disk space; with BGZF (Blocked GNU Zip Format) - a variant of GZIP
- Can be sorted & indexed - quick viewing/searching (bigger than GZIP files)
- Cannot be read without a tool (samtools)

## uBAM

- unmapped BAM → compress FASTQ files

## CRAM

- Better lossless compression than BAM
- Cramtools for conversion from/to BAM
- [http://www.ebi.ac.uk/ena/about/cram\\_toolkit](http://www.ebi.ac.uk/ena/about/cram_toolkit)

# Assignment 1

General information:

[https://github.com/spabinger/medizinische\\_genomanalysen\\_2021](https://github.com/spabinger/medizinische_genomanalysen_2021)

# Assignments during Lectures

- 30 mins per assignment
  - Description of assignment
  - Planning of tasks
- Accept Github classroom assignment
  - Link will be sent
  - Details on the next slides
- Work on assignment
- Send your name and the name of Github repository to  
**stephan.pabinger@gmail.com**
- Deadline for submission of all 3 assignments:  
**16.05.2021**

# Github Classroom – Accept assignment

GitHub Classroom

GitHub Education



FHMedGen2021-Class

## Accept the assignment — assignment1

Once you accept this assignment, you will be granted access to the  
assignment1- name repository in the FHMedGen2021 organization  
on GitHub.

---

Accept this assignment

# Github Classroom – After accepting

GitHub Classroom

GitHub Education



You accepted the assignment, **assignment1**. We're configuring your repository now. This may take a few minutes to complete. Refresh this page to see updates.



Join the GitHub Student Developer Pack

Verified students receive free GitHub Pro plus thousands of dollars worth of the best real-world tools and training from GitHub Education partners — for free. [Learn more](#)

Apply

# Github Classroom – Refresh page

You're ready to go!

You accepted the assignment, **assignment1**.

Your assignment repository has been created:



<https://github.com/FHMedGen2021/assignment1> name

We've configured the repository associated with this assignment ([update](#)).

# Github Classroom – Assignment

main ▾ 1 branch 0 tags Go to file Add file ▾ Code ▾

github-classroom GitHub Classroom Autograding Workflow 714cfb 28 seconds ago 3 commits

File	Description	Time
.github	GitHub Classroom Autograding Workflow	28 seconds ago
README.md	Initial commit	30 seconds ago
assignment2.py	Initial commit	30 seconds ago
assignment2_test.py	Initial commit	30 seconds ago
chr21_new.vcf	Initial commit	30 seconds ago

Test your code using pytest

assignment2\_test.py → README.md

Read first!

Use this file  
(no need to copy it)

# Github Classroom – Autograder

The screenshot shows the GitHub Classroom interface. At the top, there are navigation links: Code, Issues, Pull requests, Actions (which is highlighted with a red underline and has a black arrow pointing to it from the left), Projects, Security, and Insights. Below this, the 'Workflows' section is visible, with a 'New workflow' button and a 'GitHub Classroom Workflow' entry. The 'Actions' tab is active, displaying 'All workflows' with a filter bar. The main content area shows '2 workflow runs': 1) 'Update assignment1.py' (status: success, pushed by 'name' via 'GitHub Classroom Workflow #2', 1 minute ago, 1m 19s) and 2) 'GitHub Classroom Autograding Workflow' (status: error, pushed by 'github-classroom bot' via 'GitHub Classroom Workflow #1', 17 minutes ago, 56s). A large black arrow points upwards from the bottom-left towards the 'Actions' tab.

All workflows

New workflow

GitHub Classroom Workflow

All workflows

Filter workflows

2 workflow runs

	Event	Status	Branch	Actor
✓ Update assignment1.py GitHub Classroom Workflow #2: Commit 823c2d5 pushed by name	main	Success	1 minute ago	1m 19s
✗ GitHub Classroom Autograding Workflow GitHub Classroom Workflow #1: Commit deaf5eb pushed by github-classroom bot	main	Error	17 minutes ago	56s

Look at the results