# Alignments

What is an alignment?

**Sequence alignments** are a way of arranging the sequences of DNA, RNA or proteins in order to **identify regions of similarity** that may be a consequence of functional, structural, or evolutionary relationships between the sequences.

ACGTCTTGACTGG -TTAAAATAC
AC - TCTTGACTGGATTAACATAC

# Alignments

Elements of an alignment

Alignment seeks to **reduce gaps and mismatches** and **maximize matches**.

ACGTTTTGCAGTAAATGCGGACTGA - T
ACGTTGTGCAGTAAATGCGGA -- GACT

mismatch
(SNVs)

match

gap
(indels)

# Alignments

Elements of an alignment

In the construction, each of these components has a penalty value associated. For gaps there is a penalty value for opening the gap and another for extending it.

ACGTTTGCAGTAAATGCGGACTGAT
ACGTTGTGCAGTAAATGCGGA-GACT     ➡ 1 gap 3 mismatches

ACGTTTGCAGTAAATGCGGACTGAT
ACGTTGTGCAGTAAATGCGGA -- GAC T     ➡ 1 extended gap 1 mismatch

ACGTTTGCAGTAAATGCGGACTGA -T
ACGTTGTGCAGTAAATGCGGA --GACT     ➡ 2 gaps

# Alignments

Objectives

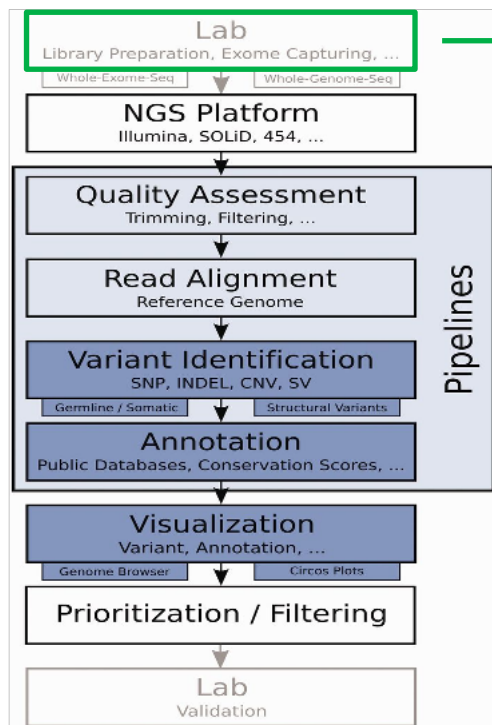The comparison between sequences in sequence alignment allows to:

1. Find homologous **positions**
2. Determine the **homology** degree
3. Identify **functional domains**
4. **Compare** the gene with its product
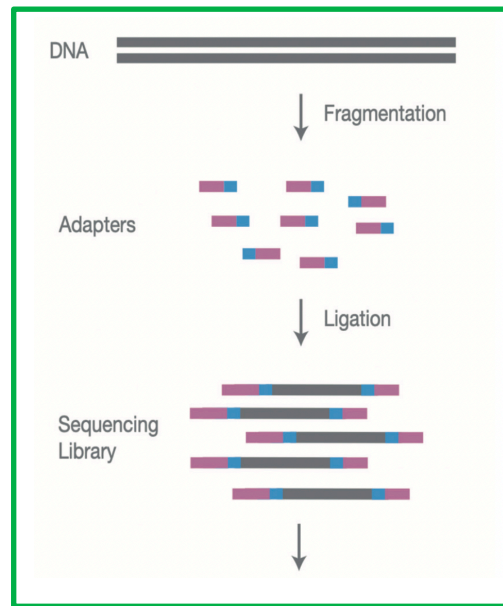5. Identify **differences**

Objective of Variant Calling in Next
Generation Sequencing (NGS)

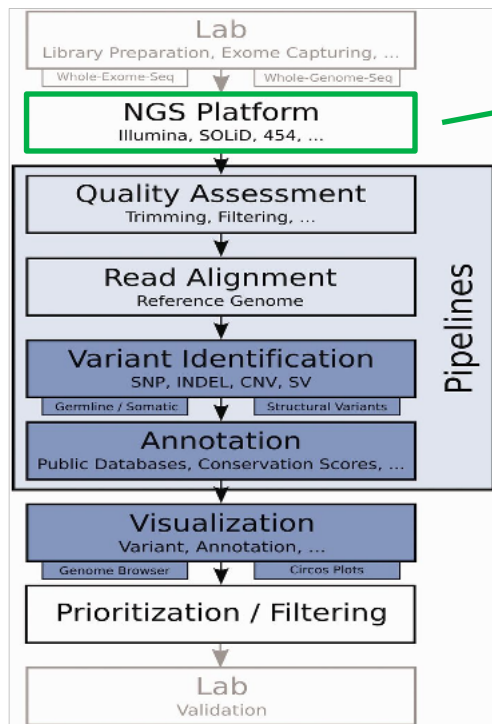# Next Generation Sequencing (NGS)

## Workflow



**Library Preparation:**
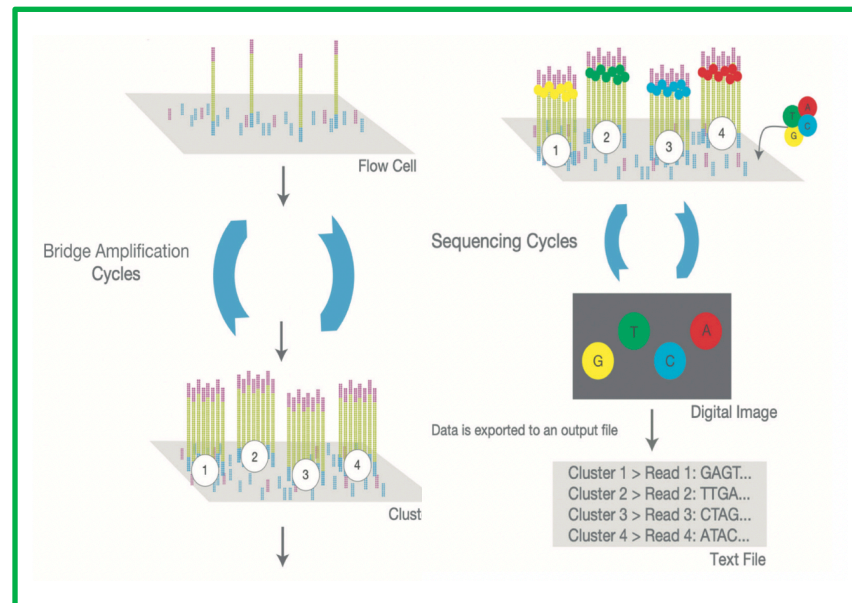DNA is fragmented into small pieces to form a **library.**



Stephan Pabinger et al. Brief Bioinform 2013;bib.bbs086

# Next Generation Sequencing (NGS)

## Workflow



Stephan Pabinger et al. Brief Bioinform 2013;bib.bbs086
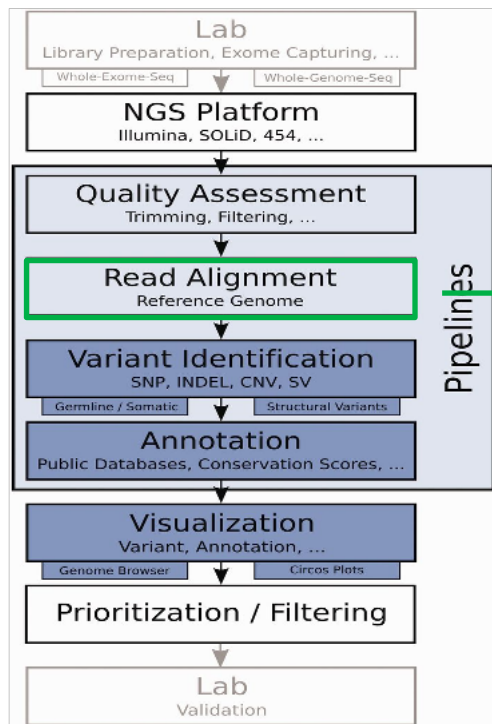
**Sequencing:**
Fragments are sequenced and stored as **reads** in text files.

# Next Generation Sequencing (NGS)

Workflow



Stephan Pabinger et al. Brief Bioinform 2013;bib.bbs086

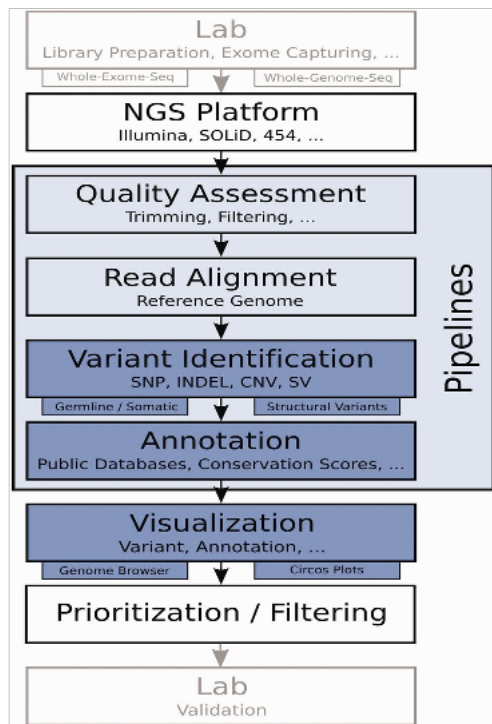**Alignment:**
Align the reads against the reference genome in order to find out where do they come from.

# Next Generation Sequencing (NGS)

## Workflow



Stephan Pabinger et al. Brief Bioinform 2013;bib.bbs086

Video: https://youtu.be/fCd6B5HRaZ8

# Short Read Aligners (SRA)

Challenges

Need to **align billions of reads** to a very **large reference genome:** SRA must be **extraordinarily efficient algorithms**
- Speed
- Memory use

As we need to align short reads, **a read may align in multiple positions**
- Either report multiple positions
- Or pick heuristically one of them

**Different NGS technologies** have different error profiles to take into account:
- **Roche 454:** Insertions or deletions in homopolymer runs
- **Illumina:** Increasing likelihood of sequence errors towards the end of the read

# Alignment
Reads

Short reads are stored in **text files** with the extension **FASTQ.**

Identifier ———— ● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence ———— ● TTGCCTGCCTATCATTTAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+' sign ———— ● +
Quality scores ———— ● hhhhhhhhhhghhghhhhhfhhhhhfffffe'ee['X]b[d[ed'[Y[^Y
Identifier ———— ● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence ———— ● GATTTGTATGAAAGTATACAACTAAAACTGCAGGTGGATCAGAGTAAGTC
'+' sign ———— ● +
Quality scores ———— ● hhhhgfhhcghghggfcffdhfehhhhcehdchhdhahehffffde'bVd

Hosseini, M.; Pratas, D.; Pinho, A.J. A Survey on Data Compression Methods for Biological Sequences. *Information*
**2016**, *7*, 56. doi: 10.3390/info7040056

# Alignment
Reference genome

The reference genome is stored in another **text file** with the extension **FASTA.**

Header — >VIT_201s0011g03530.1
Sequence — AATTAAGCATAAATACTCACTCTTACCCCCTTATTTTCTTATCTCTCATCACTTTTGGTGCGAAG
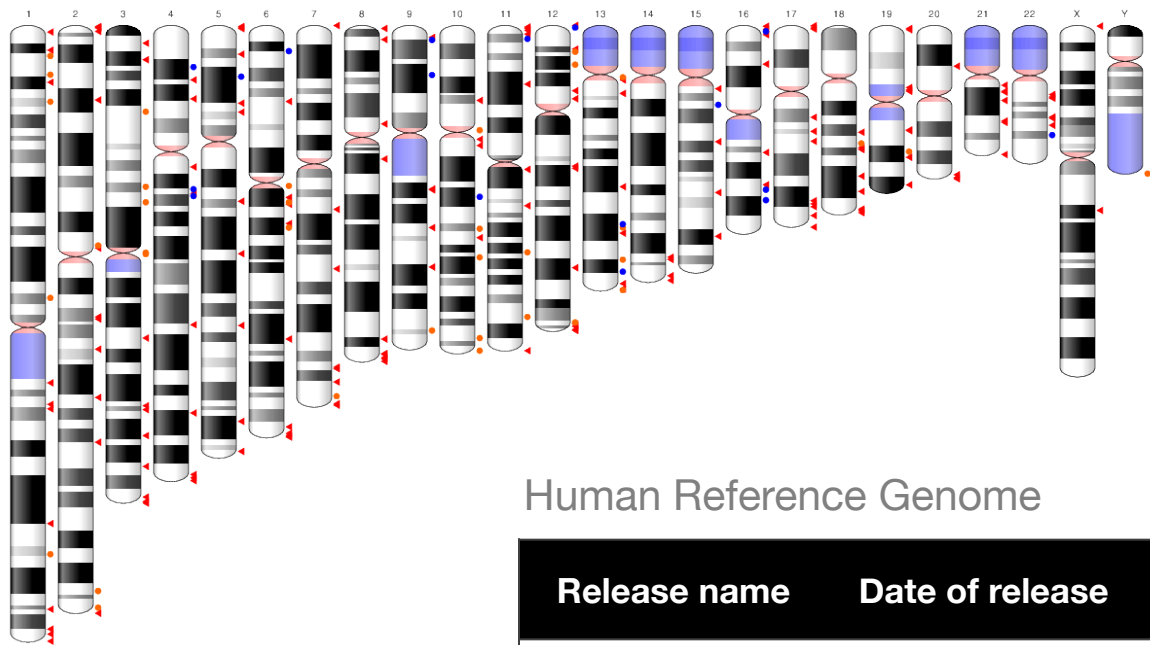GACCATGAGAACAAGCTGCAATGGGTGTAGGGTTCTTCGCAAGGCATGCAGCCAAGACTGCATCA
Header — >VIT_201s0011g03540.1
Sequence — CAGGTAGCGTGAAGTTAAACCCTAGCGCTTTAGACAAACAGCTGTAGTCACCGCCCACAAACACC
AGCCTCTGAGACACCACCTCAAACCTTTCCACTTAAATACACATCCCTCACACCCTTTTCAATTC
Header — >VIT_201s0011g03550.1
Sequence — CATGCAAAGCTGAACGCGATGCTGTGATTGGTGGTAAGTGGTAGTTGAGTAAATTTGACAGTGAA
GCCGAAATGGTAAAAGACTAAGGCTAGAAGTAGAATACCACTGTTCTTCTCATCACGTGGGCCCA

Hosseini, M.; Pratas, D.; Pinho, A.J. A Survey on Data Compression Methods for Biological Sequences. *Information* **2016**, *7*, 56. doi: 10.3390/info7040056
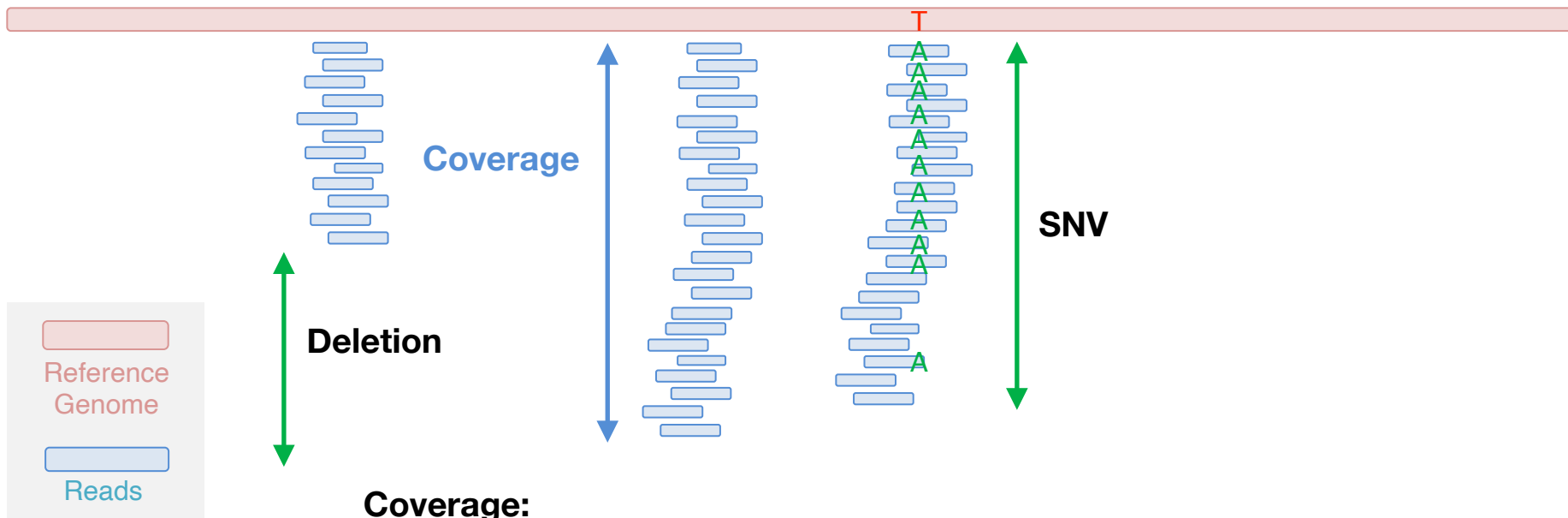
# Alignment

## Reference genome



Region containing alternate loci
Region containing fix patches
Region containing novel patches

## Human Reference Genome

| Release name | Date of release | Equivalent UCSC version | Base Pairs |
|---|---|---|---|
| GRCh38 | Dec 2013 | hg38 | 3,609,003,417 |
| GRCh37 | Feb 2009 | hg19 | 3,326,743,047 |

# Alignment

Reads are mapped against the reference genome in order to identify differences



**Coverage:**
Number of reads that align to the same region (60x, 100x,…)

# Alignment

Indexing

In order to optimize the alignment process, these algorithms use a strategy called **indexing. As in books,** indexing allows to organize information in a more easier and faster way to search.

# Alignment

Elements to consider

- Read type: DNA, RNA, etc.
- Read length
- Paired-end or single-end
- Computational requirements (number of processors, memory)
- Number of mismatches (limitation in allowed differences)
- Sequencing errors (can be platform dependent)

# Alignment

Elements to consider

- Read type: DNA, RNA, etc.



Mackenzie RJ. RNA-seq: Basics, Applications and Protocol. *Technology Networks Genomic Research*, **2018**. https://www.technologynetworks.com/genomics/articles/rna-seq-basics-applications-and-protocol-299461

# Alignment

Elements to consider

- Read type: DNA, RNA, etc.
- Read length
- Paired-end or single-end
- Computational requirements (number of processors, memory)
- Number of mismatches (limitation in allowed differences)
- Sequencing errors (can be platform dependent)

# Alignment

Elements to consider

- Read length

The longer the reads, the higher the likelihood of sequencing errors towards the end.

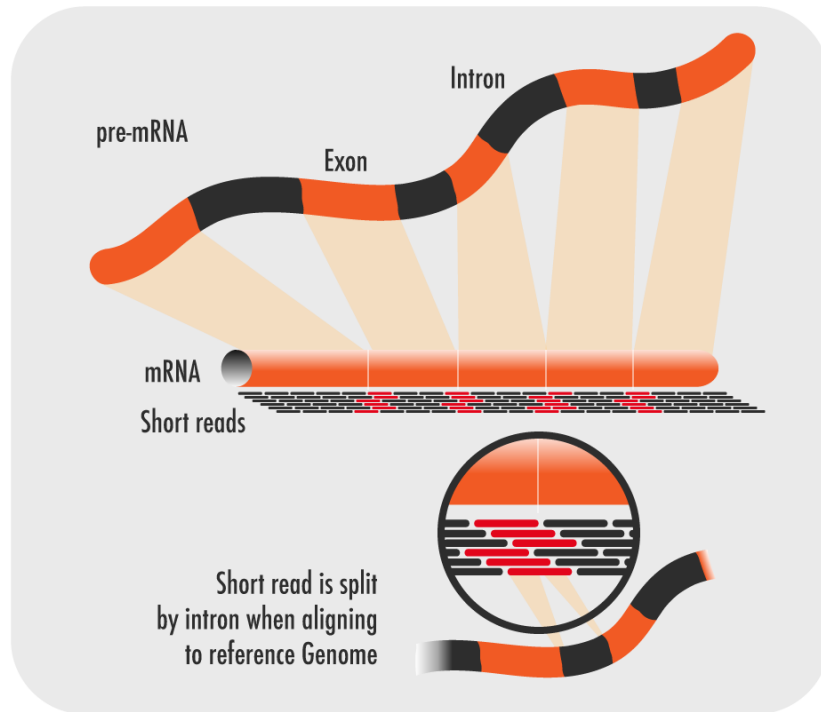# Alignment
Elements to consider

- Read type: DNA, RNA, etc.
- Read length
- Paired-end or single-end
- Computational requirements (number of processors, memory)
- Number of mismatches (limitation in allowed differences)
- Sequencing errors (can be platform dependent)

# Alignment
Elements to consider

- Paired-end or single-end

In a **paired-end** experiment, **both ends** of the fragment are sequenced.

The **distance between** the each paired read **is known (insert size)**, thus the aligners can **map** the reads **more precisely.** Important for correctly mapping **repetitive regions** of the genome.

Single-end reads

reference sequence

Paired-end reads

reference sequence

sequenced fragment    unknown sequence    sequenced fragment

200 - 1000bp

NGS overlapping reads. Biostars. Post 241139.

# Alignment

Elements to consider

- Read type: DNA, RNA, etc.
- Read length
- Paired-end or single-end
- Computational requirements (number of processors, memory)
- Number of mismatches (limitation in allowed differences)
- Sequencing errors (can be platform dependent)
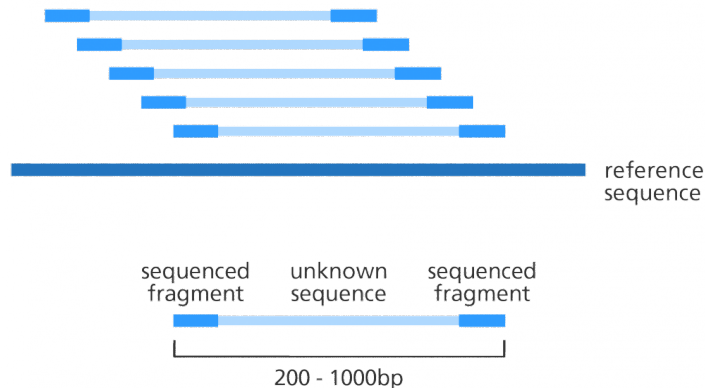
# Alignment

Elements to consider

- Read type: DNA, RNA, etc.

- Read length

- Paired-end or single-end

- Computational requirements (number of processors, memory)

- Number of mismatches (limitation in allowed differences)

- Sequencing errors (can be platform dependent)

# Alignment

Elements to consider

- Read type: DNA, RNA, etc.

- Read length

- Paired-end or single-end

- Computational requirements (number of processors, memory)

- Number of mismatches (limitation in allowed differences)

- Sequencing errors (can be platform dependent)

# Alignment
Errors and biases

- Sequencing errors:
  - Increases mismatches
  - Higher at the end of the reads
  - Technology-dependent
- Different regions in the DNA sequence cause aligning biases:
  - Repetitive regions:
    - Found in different locations
    - Place of sequencing errors
    - Place of real mutations and structural variants
  - Difficulties in the alignment of indels (gaps)

# Alignment

Solutions

Post-alignment Quality Control and local realignment of indels.

**More about indel realignment:** https://qcb.ucla.edu/wp-content/uploads/sites/14/2016/03/GATKwr12-3-IndelRealignment.pdf

# Thanks!