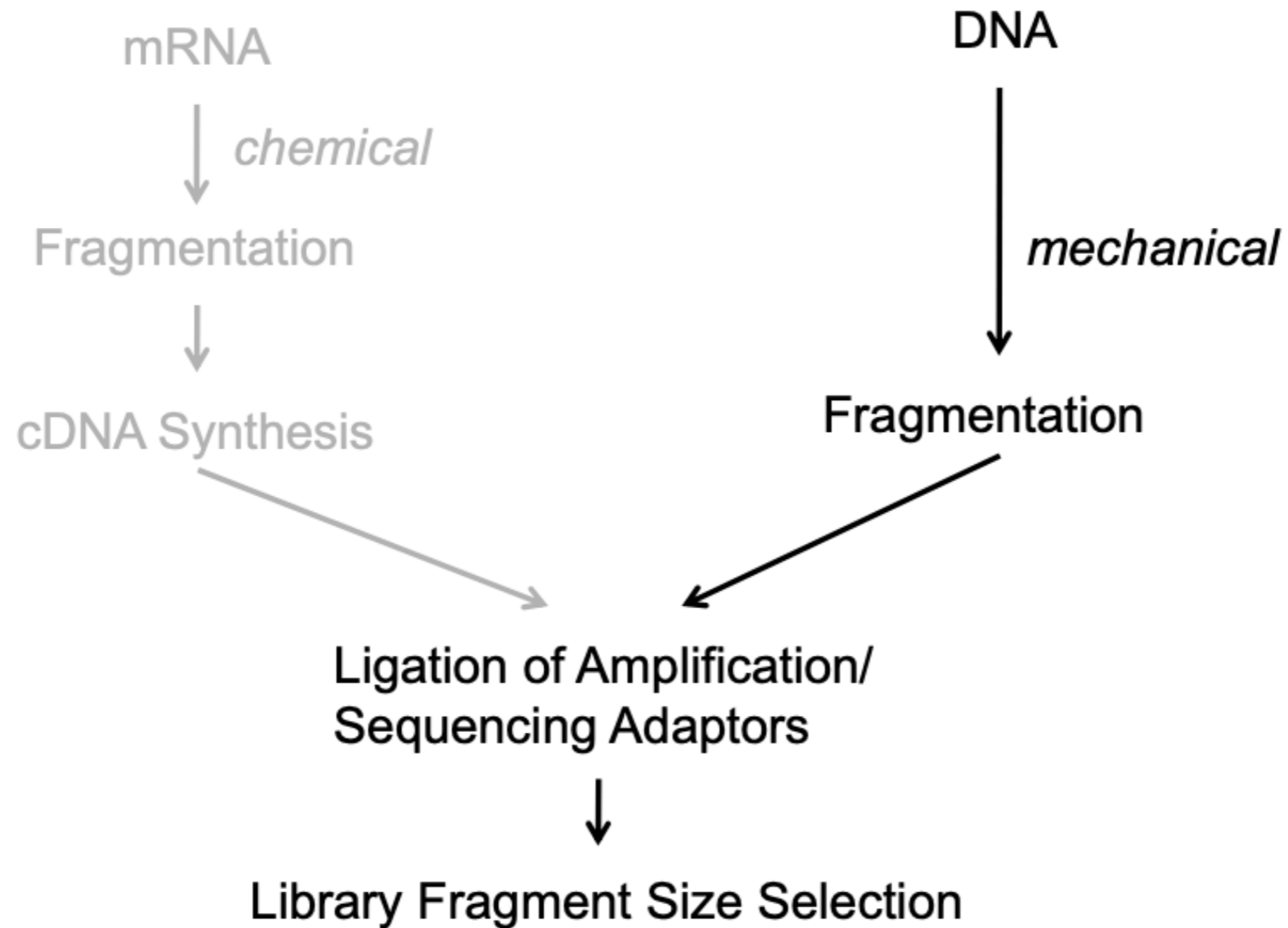


DNA and RNA sequencing

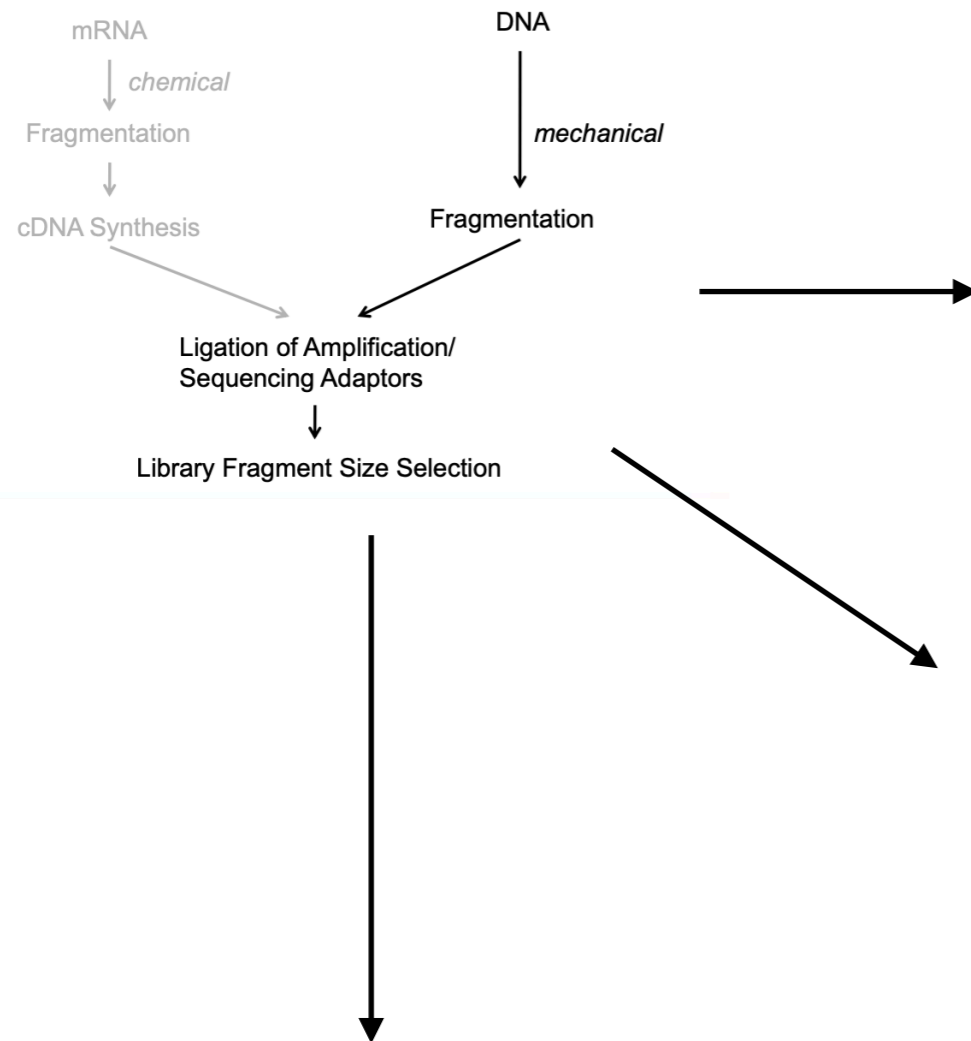
2024-04-04

DNA-sequencing

Key steps in sequencing



Main types of DNA sequencing

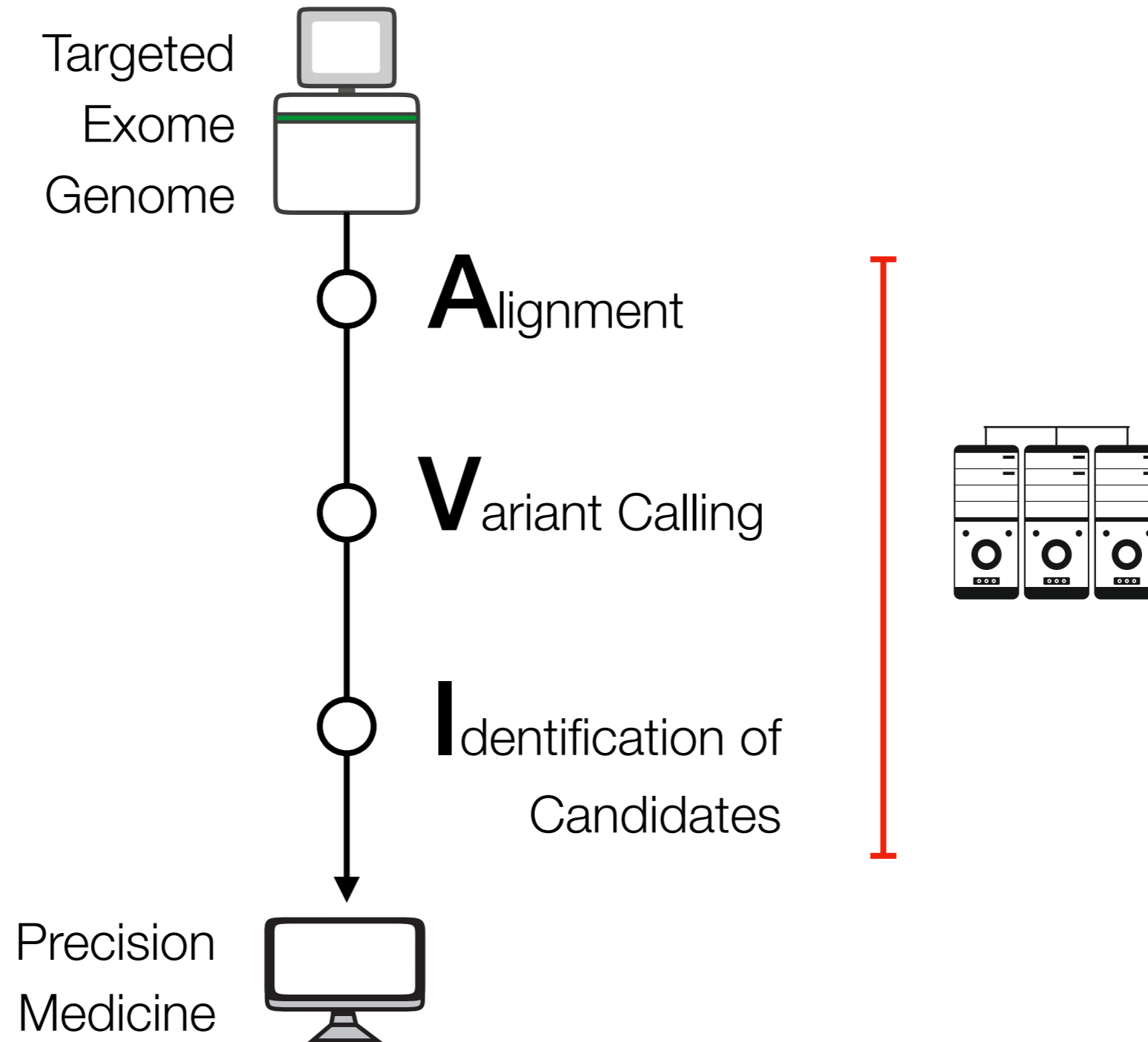


Targeted sequencing:
genes or regions of interest

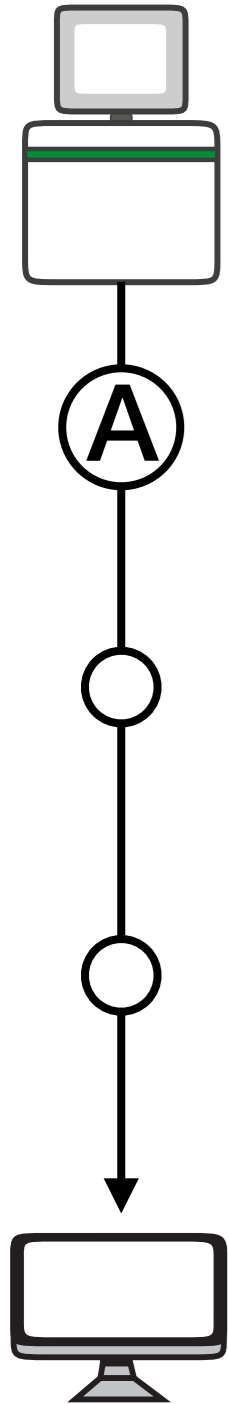
Whole EXome Sequencing
(WES/WXS): exome (coding sequences)

Whole Genome Sequencing
(WGS) : genome

Deciphering DNA-seq data



Alignment



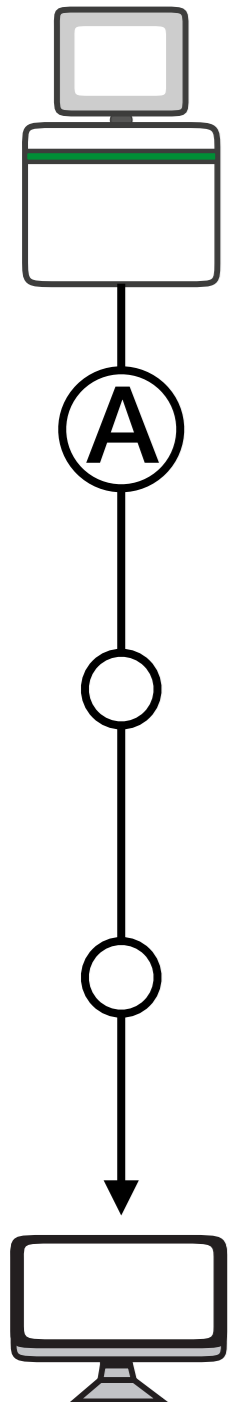
The first step in data analysis is alignment i.e. we need to understand where reads map on the genome.

Raw reads are usually found in **FASTQ format**, while the final output of the alignment is a **SAM/BAM** file.

The alignment requires a genome reference. The most recent release is GRCh38 (2013).

The FASTQ format

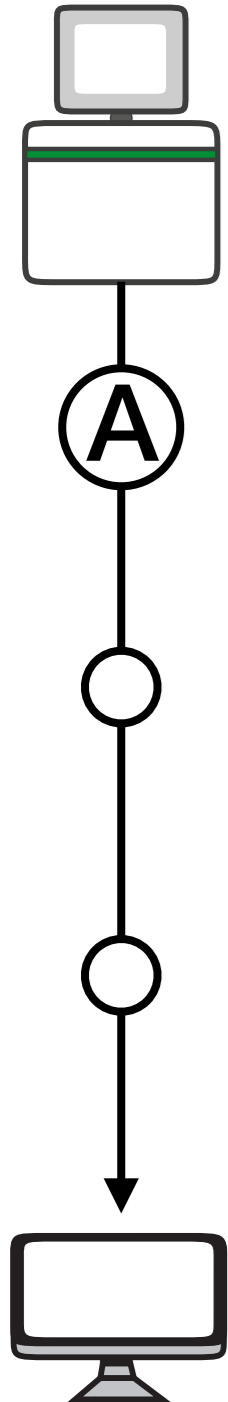
The FASTQ format keeps both read and read information in a unique file.



```
@A00959:190:HKH3CDRX2:1:2101:1118:1016 1:N:0:AACGCTTA
ANTGTGCTGGGCAAGGTTCTGGCATGAAGCCAGCACTCAATAAATGCAAGTTATTCTTTATGCAGATC
+
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFFFFFFFFFFFFFFFFFFFF
@A00959:190:HKH3CDRX2:1:2101:1353:1016 1:N:0:AACGCTTA
CNTGCGCTGGATGAGCAGGTAGAACACCACCTTCTGGTGGCCTGCTTCCTGGGCTGGCGCCGCTGGGTCC
+
F#FFFF:FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFFFFFFFFFFFF
@A00959:190:HKH3CDRX2:1:2101:1271:1031 1:N:0:AACGCTTA
GCCTAGCACATAGTAGGTGCTCAATAAATATGTGTTAAGGCCGGGCGCAGTGGCTCACGCCTATAGTCC
+
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFF
```

The FASTQ format

Each read is represented by four different fields (lines)



Sequence

Sequencing information

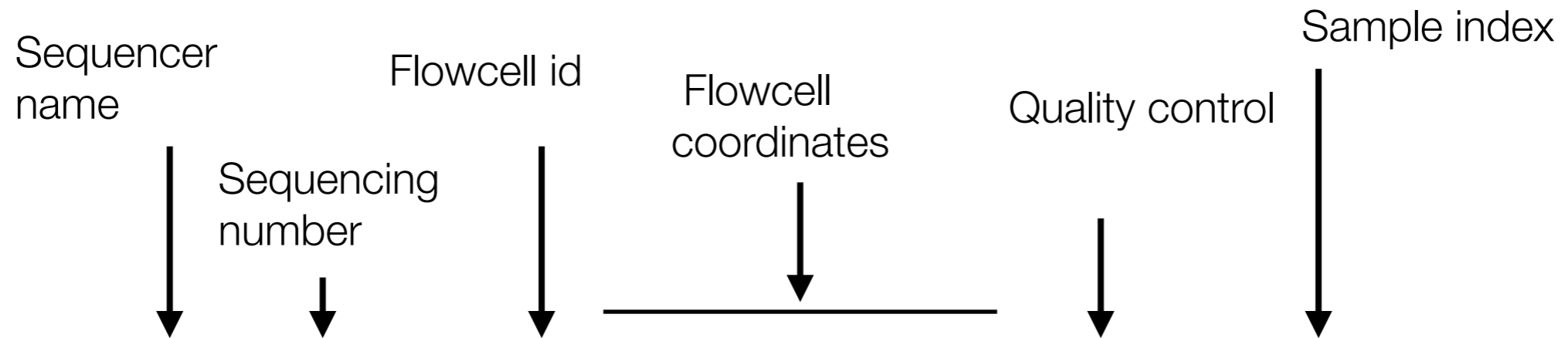
```
@A00959:190:HKH3CDRX2:1:2101:1353:1016 1:N:0:AACGCTTA
CNTGCGCTGGATGAGCAGGTAGAACACCACCTTCTGGTGGCCTGCTTCCTGGGCTGGCGCCGCTGGGTCC
+
F#FFFF:FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```

Null

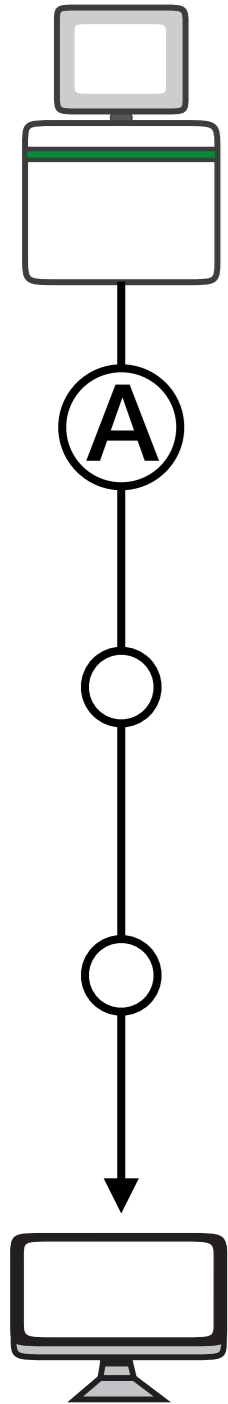
Base quality

The FASTQ format

Each read is represented by four different fields (lines)



```
@A00959:190:HKH3CDRX2:1:2101:1353:1016 1:N:0:AACGCTTA
CNTGCGCTGGATGAGCAGGTAGAACACCACCTTCTGGTGGCCTGCTTCCTGGGCTGGCGCCGCTGGGTCC
+
F#FFFF:FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```



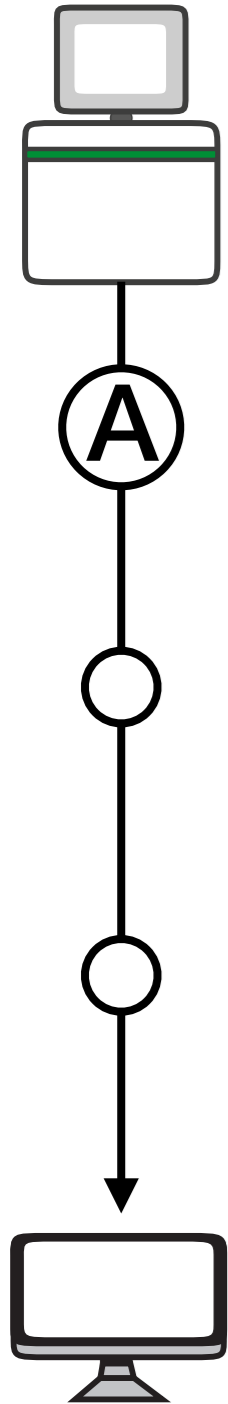
Base quality

The most used quality measure for sequencing data is the **Phred score**.

$$Q_{\text{PHRED}} = -10 \times \log_{10}(P_e)$$

| <u>Phred Quality Score</u> | <u>Probability of incorrect base call</u> | <u>Base call accuracy</u> |
|----------------------------|---|---------------------------|
| 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% |

In **fastq** format base quality is encoded in **ASCII**.



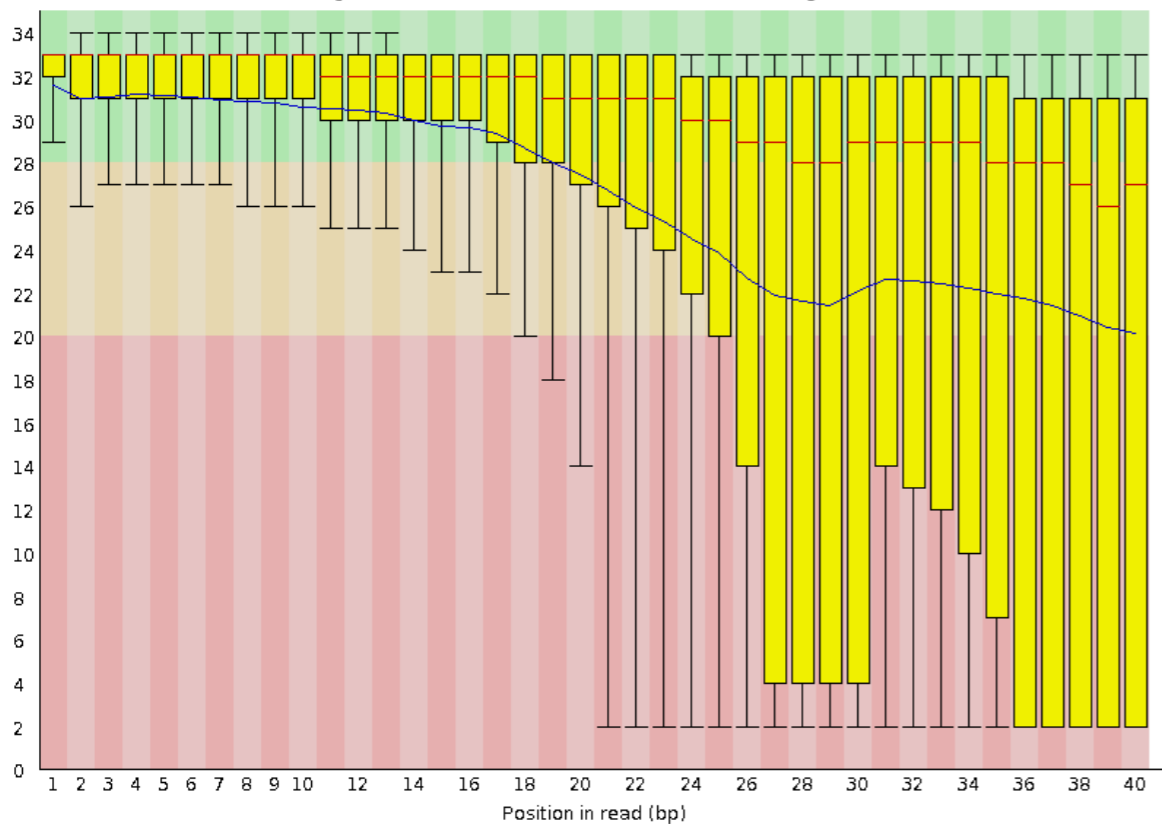
FastQC to do quality control on reads

FastQC Report

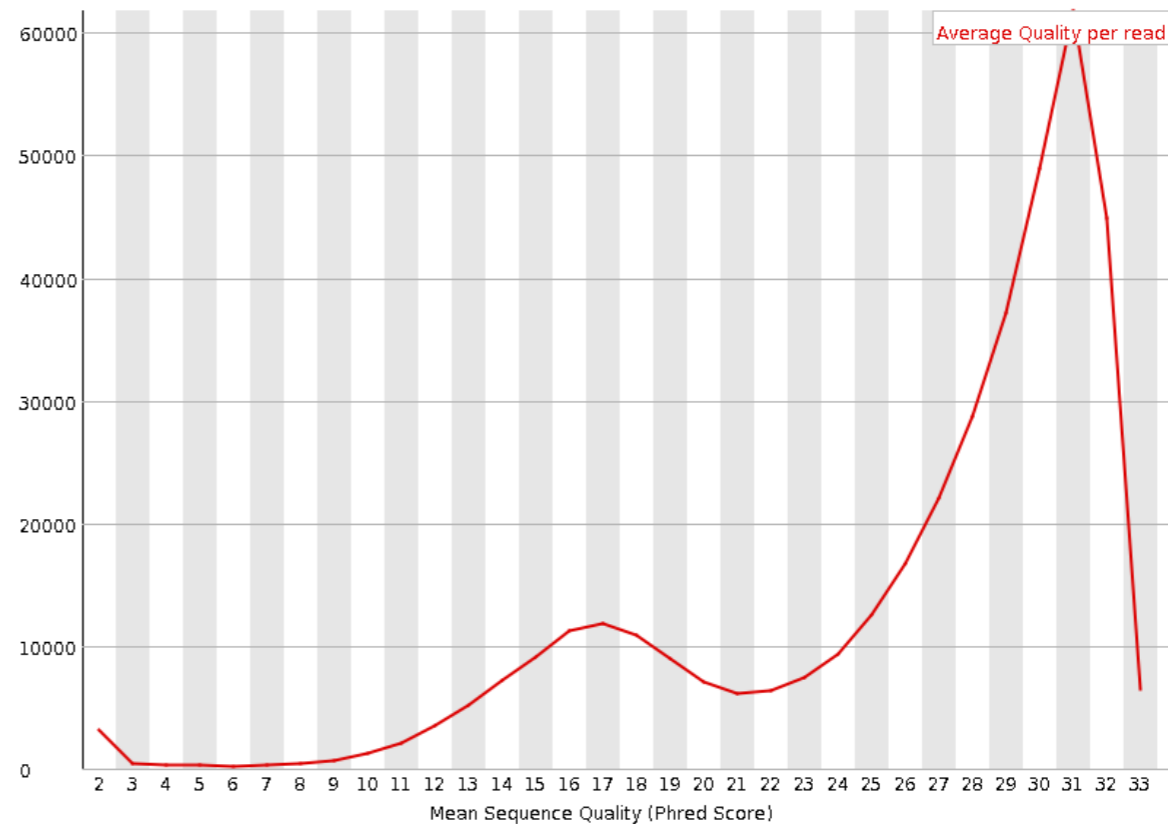
Summary

- ✓ [Basic Statistics](#)
 - ✓ [Per base sequence quality](#)
 - ✓ [Per tile sequence quality](#)
 - ✓ [Per sequence quality scores](#)
 - ✓ [Per base sequence content](#)
 - ✓ [Per sequence GC content](#)
 - ✓ [Per base N content](#)
 - ✓ [Sequence Length Distribution](#)
 - ✓ [Sequence Duplication Levels](#)
 - ✓ [Overrepresented sequences](#)
 - ✓ [Adapter Content](#)
- Base quality distribution along the length of the read
- Mean base quality distribution across reads
- Percentage of each base along the sequence
- Distribution of the percentage of GC along the sequence
- Distribution of the percentage of N bases (not properly called) along the sequence
- Fragment length distribution
- Duplication rate**
- If there are recurrent identical sequences

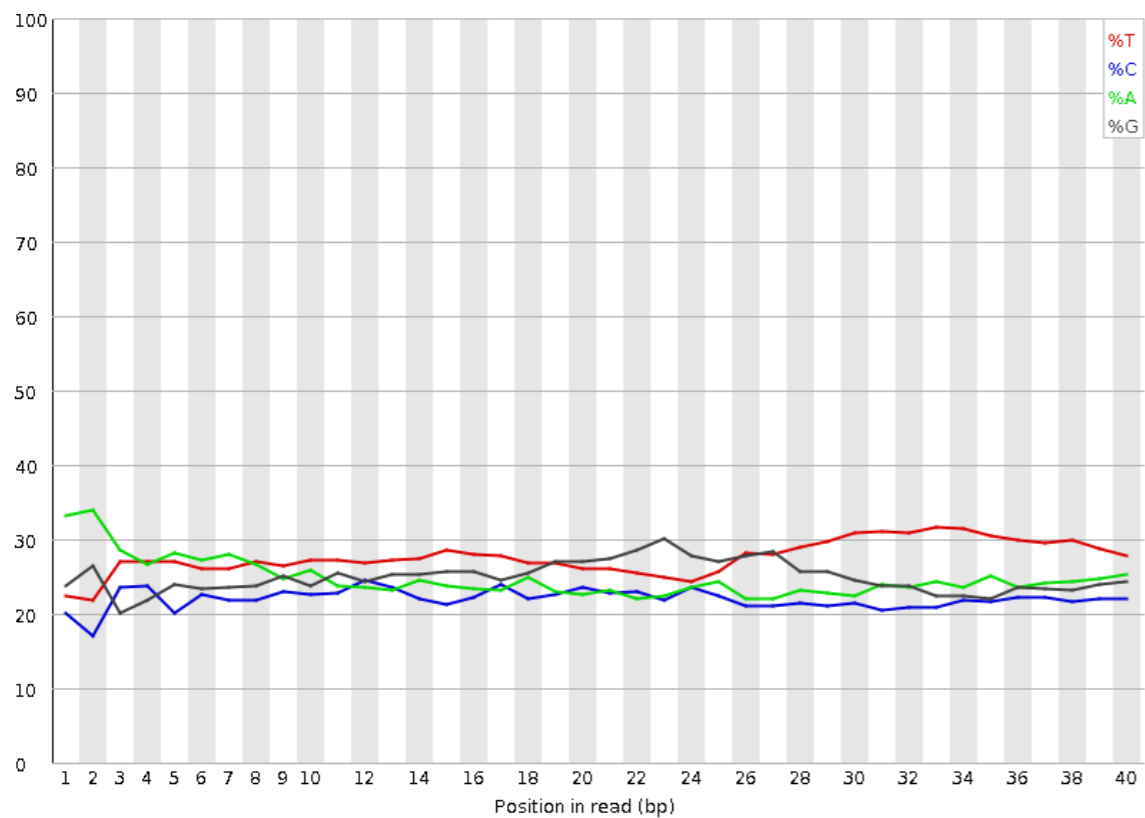
Base quality distribution along the read



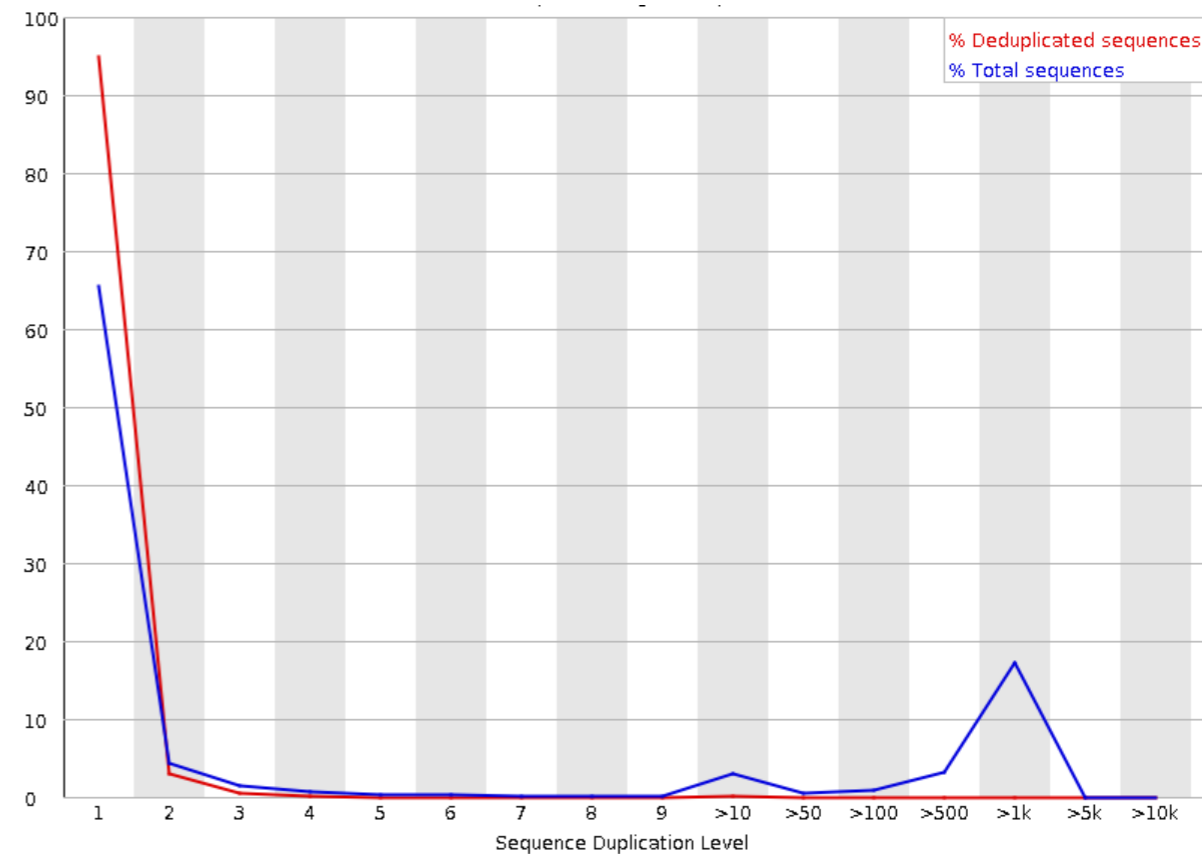
Mean base quality distribution



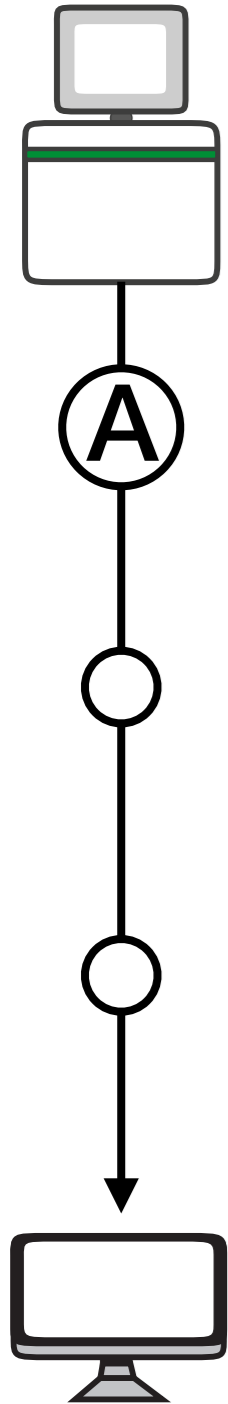
Percentage of each base along the sequence



Duplication rate



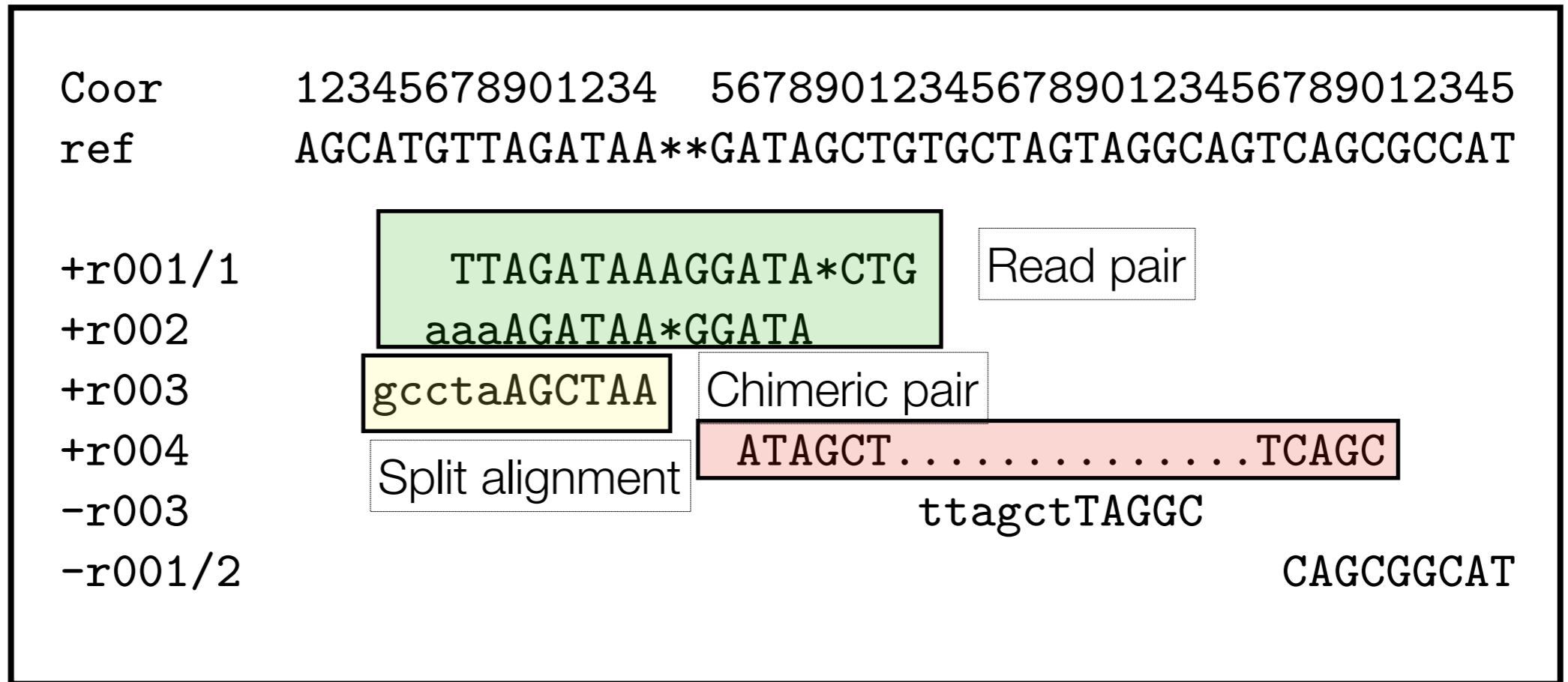
Come rappresentiamo l'allineamento?



FastQ

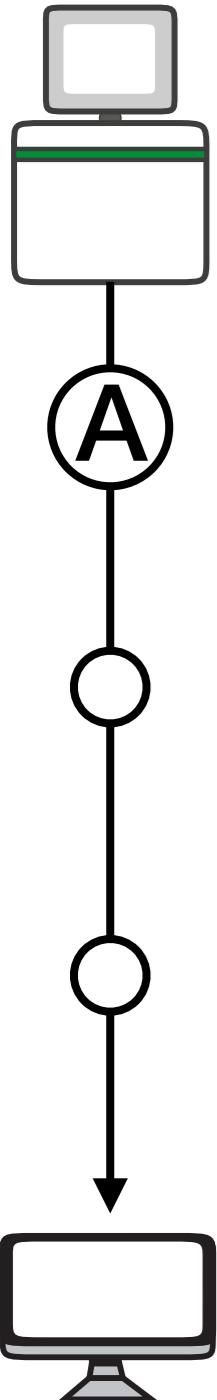
```
@A00959:190:HKH3CDRX2:1:2101:1271:1031 1:N:0:AACGCTTA
GCCTAGCACATAGTAGGTGCTCAATAAATATGTGTTAAGGCCGGGCGCAGTGGC
+
FFFFFFFFFFFFFFFF,FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```

↓ BWA-MEM to genome



↓ SAM/
BAM

Come rappresentiamo l'allineamento?



Intestazione

Allineamento

```

Coord      12345678901234 5678901234567890123456789012345
ref        AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1    TTAGATAAAGGATA*CTG
+r002      aaaAGATAA*GGATA
+r003      gcctaAGCTAA
+r004      ATAGCT.....TCAGC
-r003      ttagctTAGGC
-r001/2    CAGCGGCAT
    
```

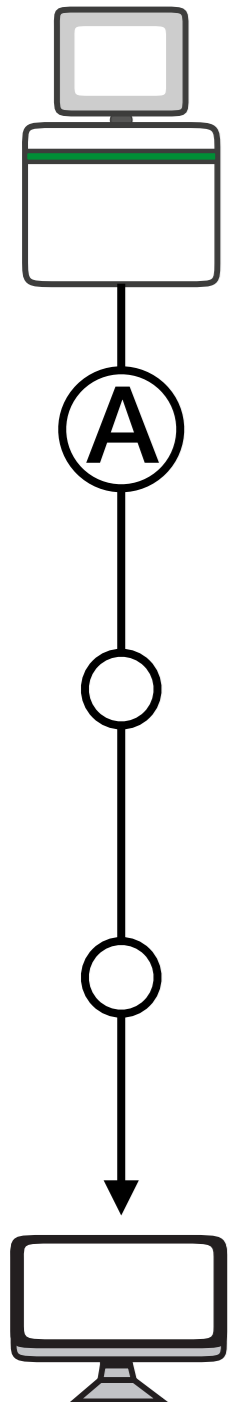
↓ SAM

```

@HD VN:1.6 S0:coordinate
@SQ SN:ref LN:45
r001  99 ref  7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002   0 ref  9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003   0 ref  9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004   0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001  147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
    
```

SAM è l'acronimo di Sequence Alignment/Map format. È un formato di testo delimitato da TAB costituito da una sezione di intestazione, che è facoltativa, e una sezione di allineamento.

Come rappresentiamo l'allineamento?



Dizionario delle sequenze di riferimento

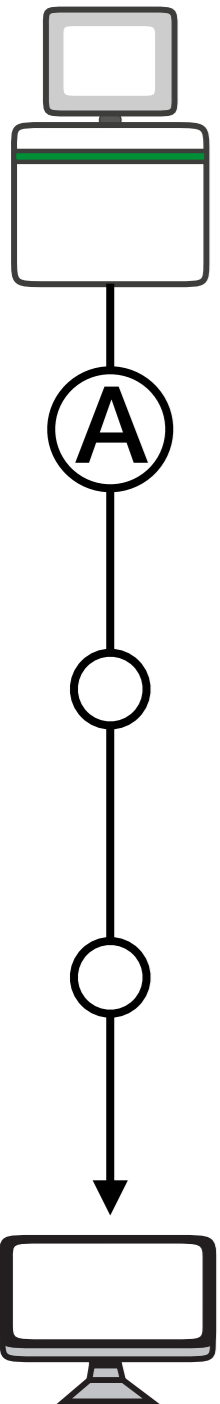
```
@HD VN:1.6 S0:coordinate
@SQ SN:ref LN:45
r001 99 ref 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

Versione

Tipo di ordinamento dell'allineamento

Nome della sequenza di riferimento

Lunghezza della sequenza di riferimento



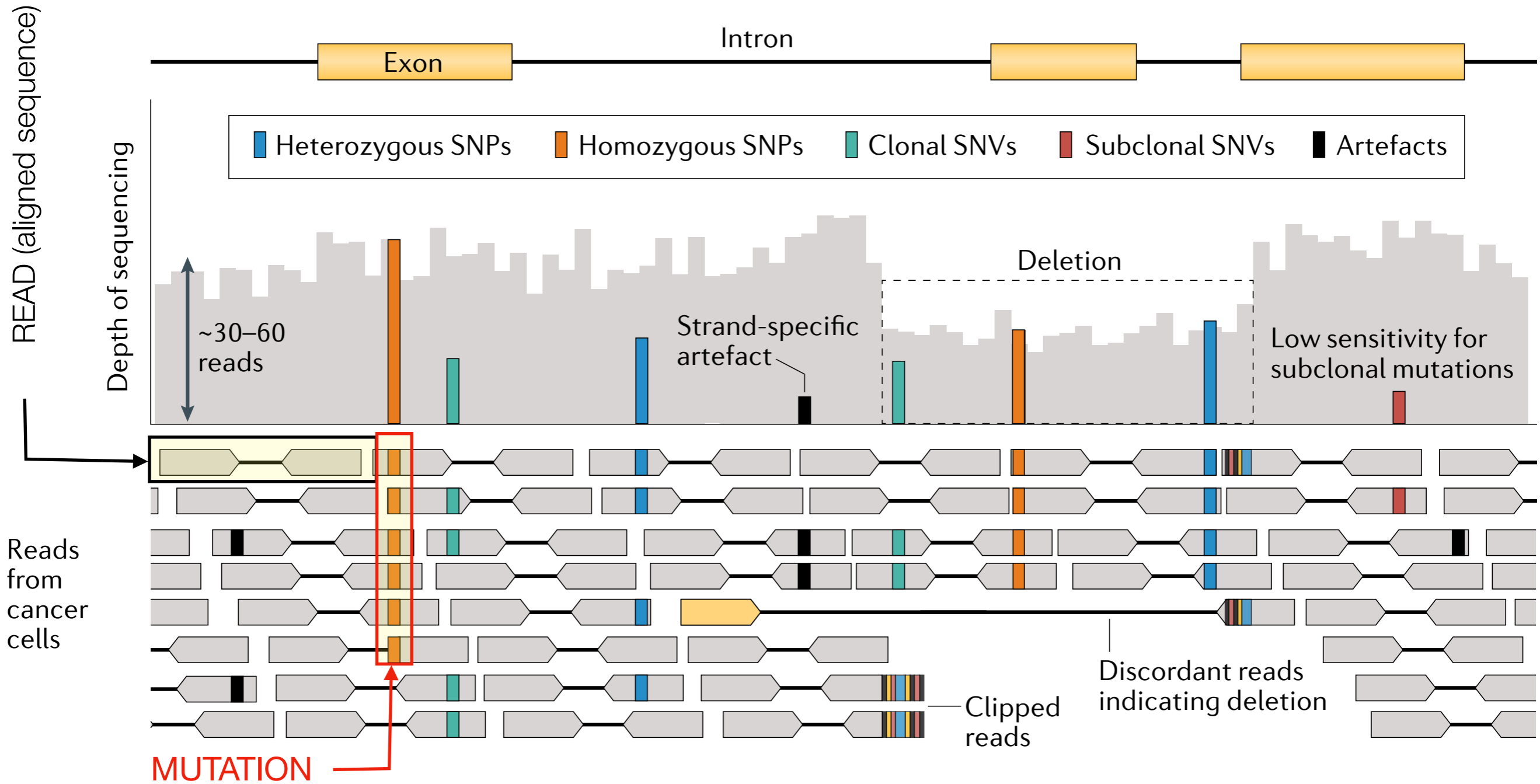
| Col | Field |
|-----|-------|
| 1 | QNAME |
| 2 | FLAG |
| 3 | RNAME |
| 4 | POS |
| 5 | MAPQ |
| 6 | CIGAR |
| 7 | RNEXT |
| 8 | PNEXT |
| 9 | TLEN |
| 10 | SEQ |
| 11 | QUAL |

| Op | Description |
|----|---|
| M | alignment match (can be a sequence match or mismatch) |
| I | insertion to the reference |
| D | deletion from the reference |
| N | skipped region from the reference |
| S | soft clipping (clipped sequences present in SEQ) |
| H | hard clipping (clipped sequences NOT present in SEQ) |
| P | padding (silent deletion from padded reference) |
| = | sequence match |
| X | sequence mismatch |

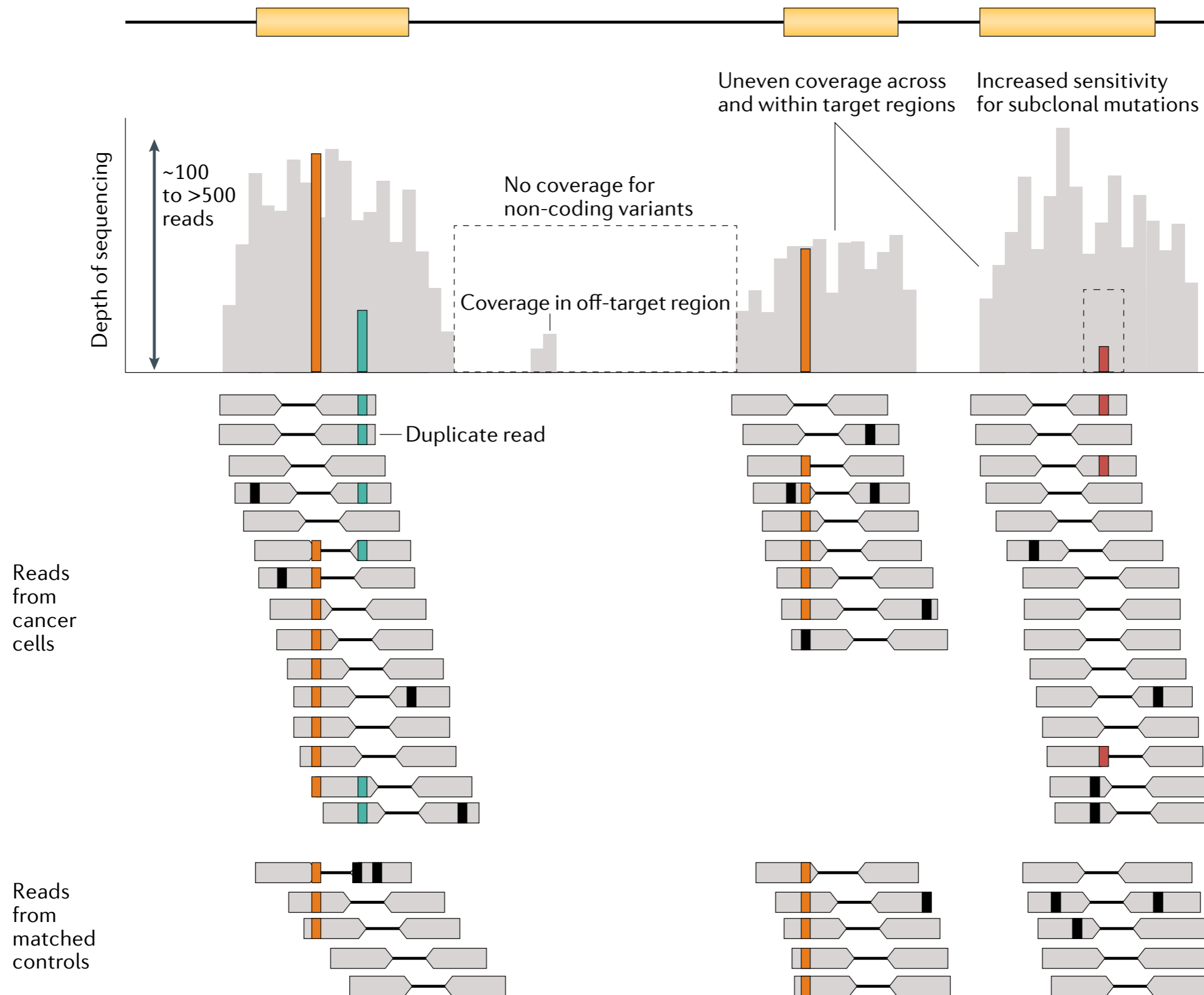
```

@HD VN:1.6 S0:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
  
```

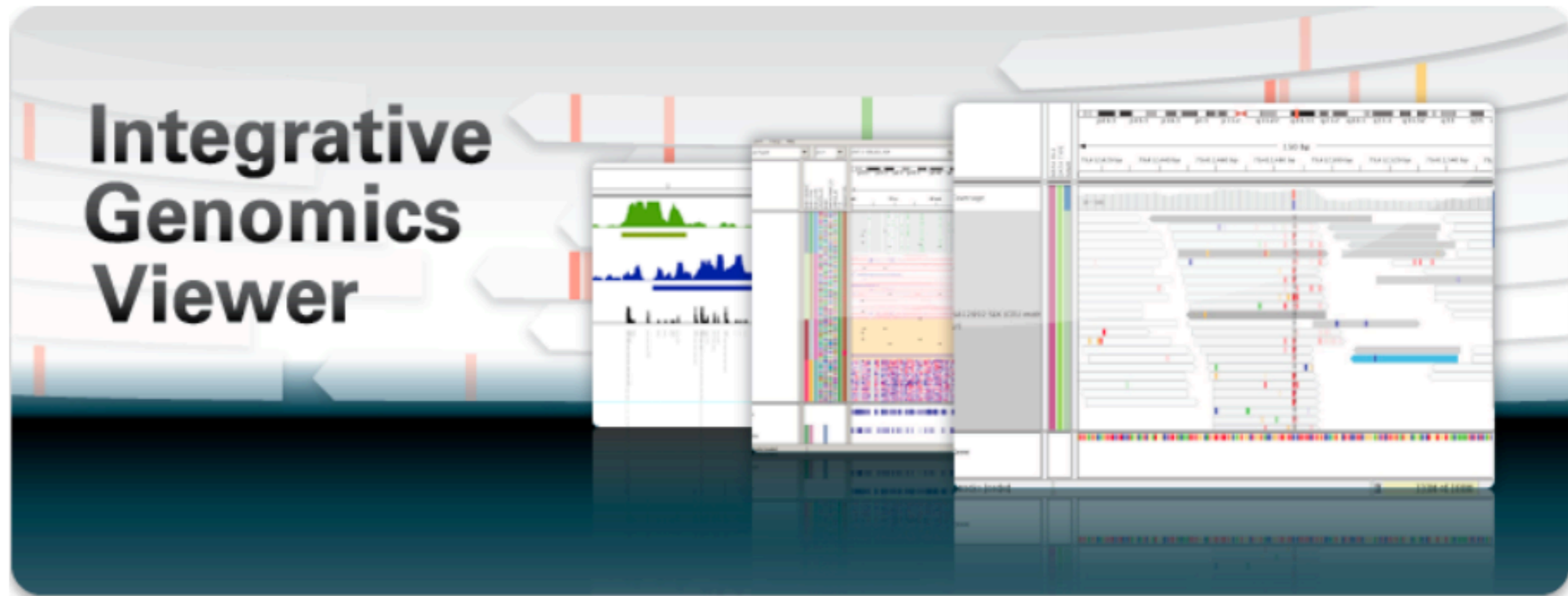
Whole-genome sequencing data



Whole-exome sequencing data



How can we see aligned data? Using IGV



According to the site ...“The **Integrative Genomics Viewer (IGV)** is a high-performance, easy-to-use, interactive tool for the visual exploration of genomic data. It supports flexible integration of all the common types of genomic data and metadata, investigator-generated or publicly available, loaded from local or cloud sources.

IGV is available in multiple forms, including:

- the original **IGV** - a Java desktop application,
- **IGV-Web** - a web application,
- **igv.js** - a JavaScript component that can be embedded in web pages (*for developers*)”

The Integrative Genomics Viewer (IGV)



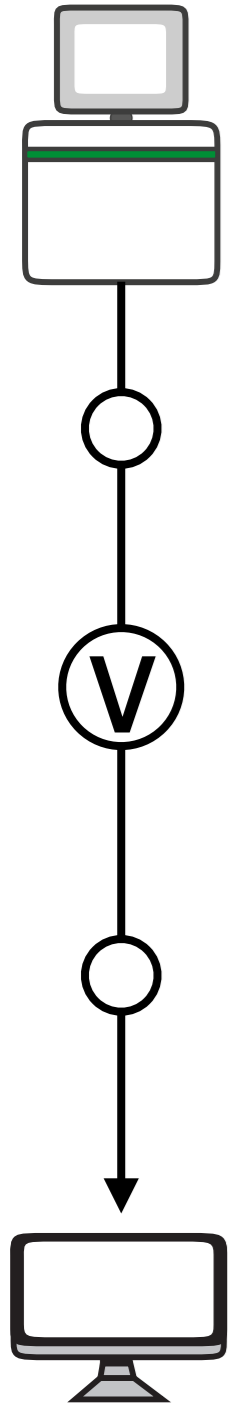
Mutations

Mutations identification

The main goal of variant identification is to evaluate if the alternative alleles supported by sequencing reads are true mutations or artefacts.

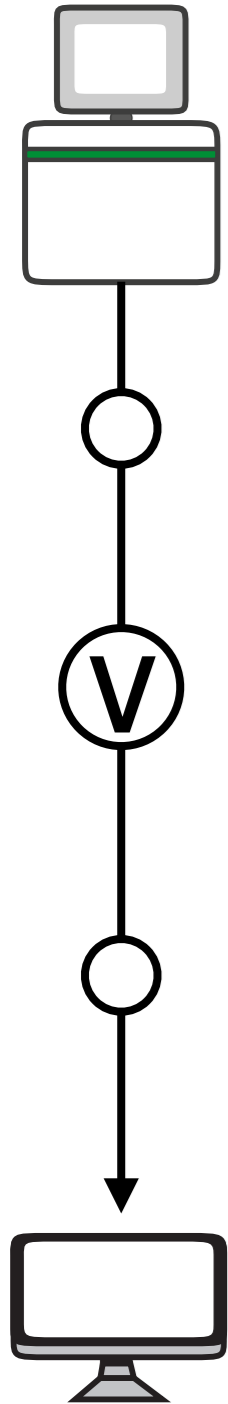
Vocabulary:

- **Single Nucleotide Polymorphisms (SNPs)**: is a germline substitution of a single nucleotide at a specific position in the genome
- **Single Nucleotide Variants (SNVs)**: a DNA sequence variation that occurs when a single nucleotide (adenine, thymine, cytosine, or guanine) in the genome sequence is altered. It is usually used for somatic mutations, nevertheless usually we can find the term “variant” for both somatic and germline mutations (Be aware of the context!!)
- **Small Insertions or Deletions (InDels)** (<50 bp)

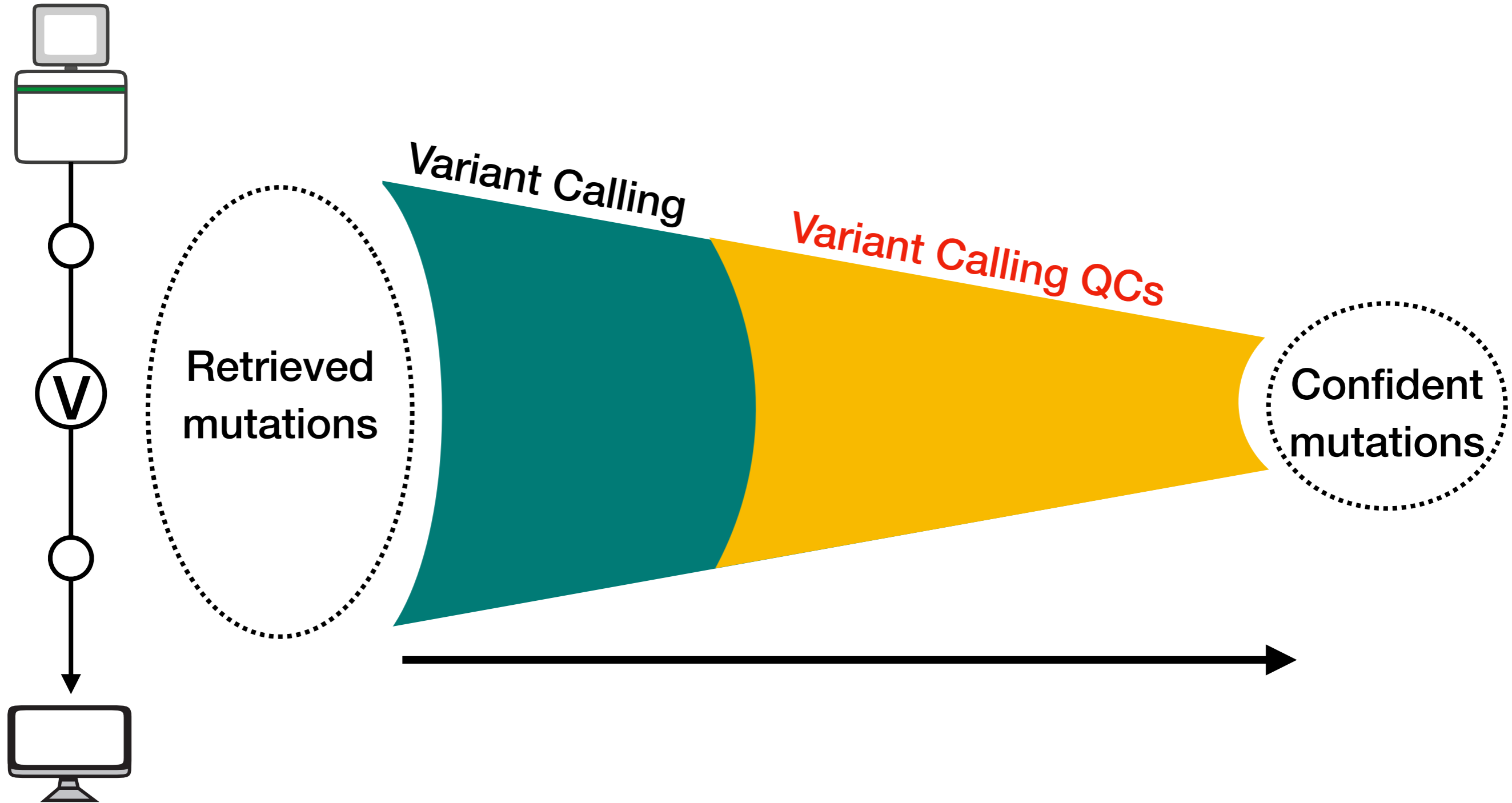


Mutations identification

A lot of different tools are available for variant identification:



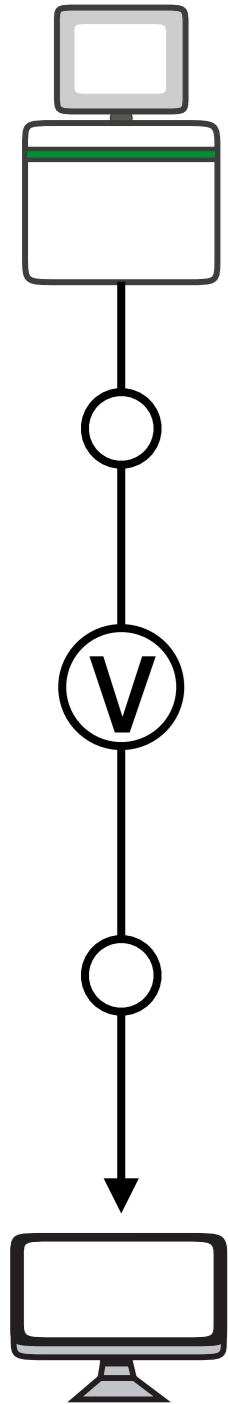
Error remotion



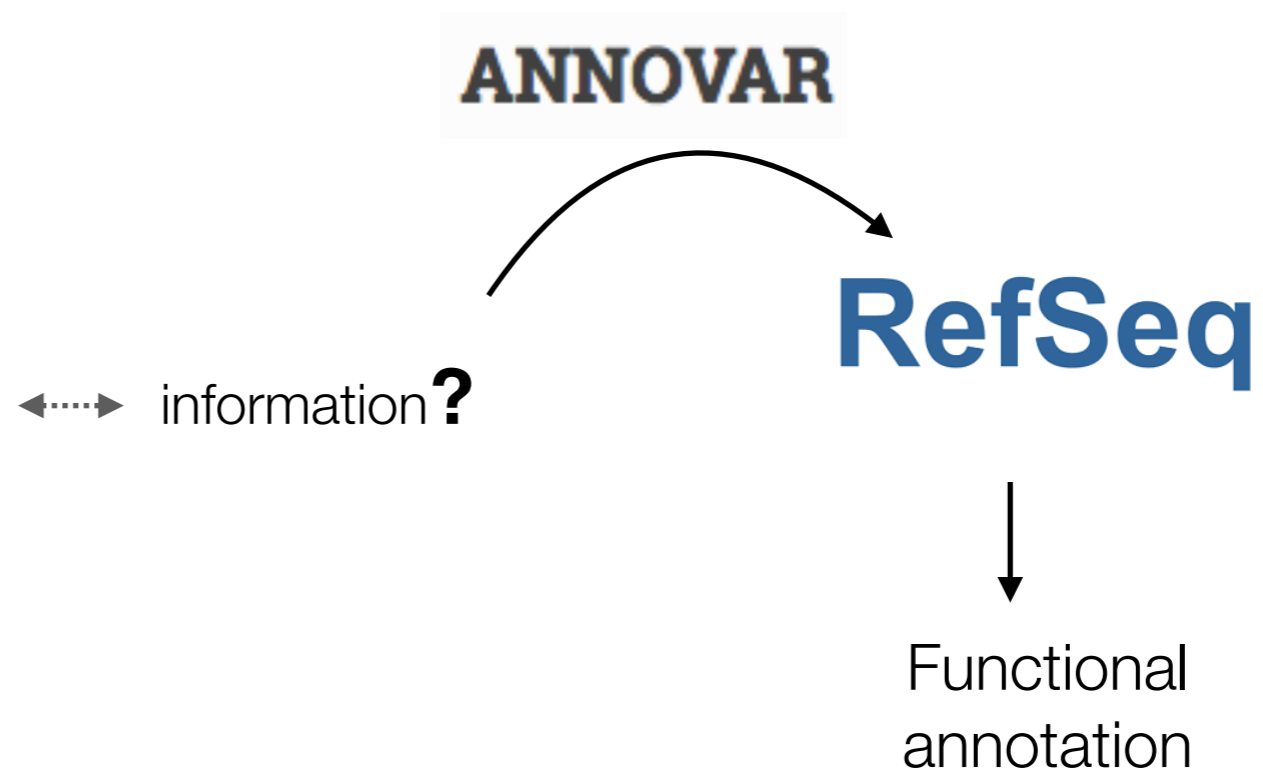
Functional annotation

Mutated genomic positions have to be **annotated** to understand their **biological meaning**.

Can the identified SNP/SNV/InDel cause changes in protein coding and interested amino acids?



| chrom | position | ref | var |
|-------|----------|-----|-----|
| chr17 | 17697102 | G | A |
| chr17 | 21319977 | A | G |
| chr17 | 26679861 | G | A |
| chr17 | 28890301 | G | A |
| chr17 | 36716758 | G | A |
| chr17 | 40369149 | G | A |



Functional annotation

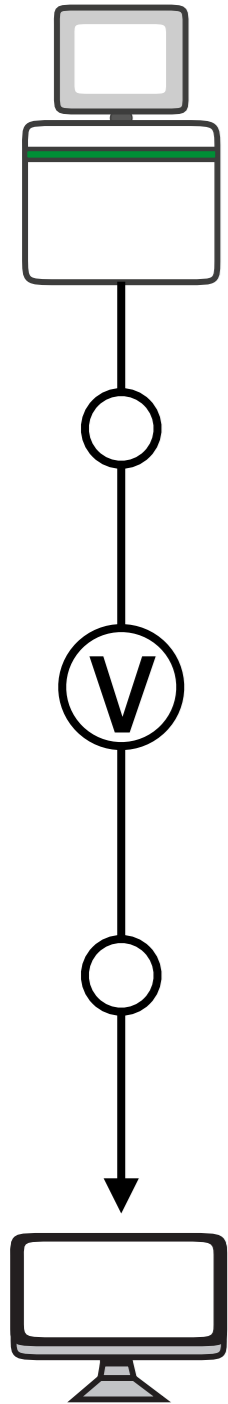
Functional
annotation



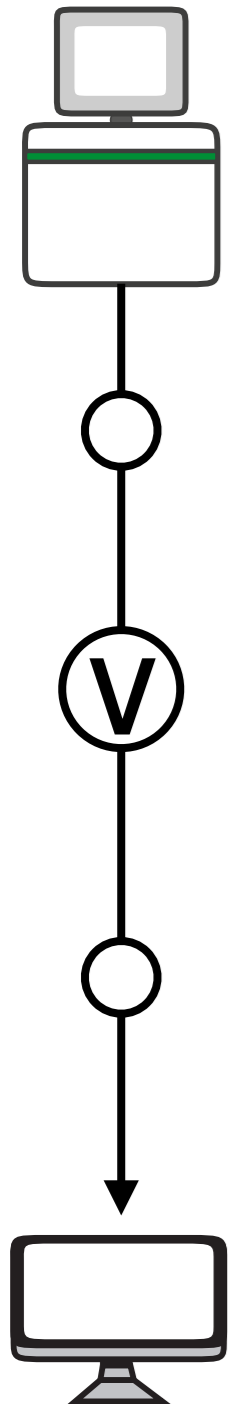
| chrom | position | ref | var | func | exonic.func | gene |
|-------|----------|-----|-----|----------|-------------------|---------|
| chr17 | 17697102 | G | A | exonic | synonymous SNV | RAI1 |
| chr17 | 21319977 | A | G | UTR3 | - | KCNJ18 |
| chr17 | 26679861 | G | A | intronic | - | POLDIP2 |
| chr17 | 28890301 | G | A | exonic | nonsynonymous SNV | TBC1D29 |
| chr17 | 36716758 | G | A | intronic | - | SRCIN1 |
| chr17 | 40369149 | G | A | intronic | - | STAT5B |



Exonic
functional
annotation

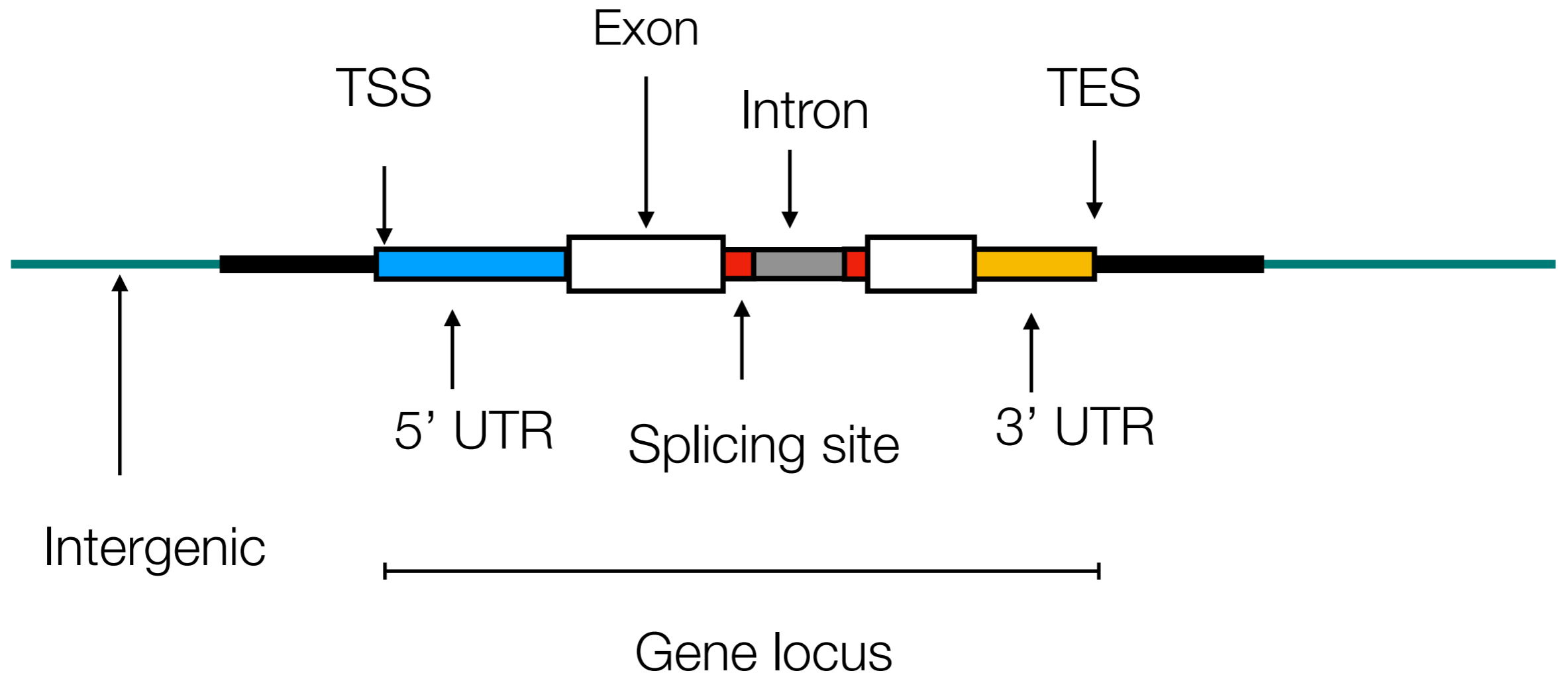
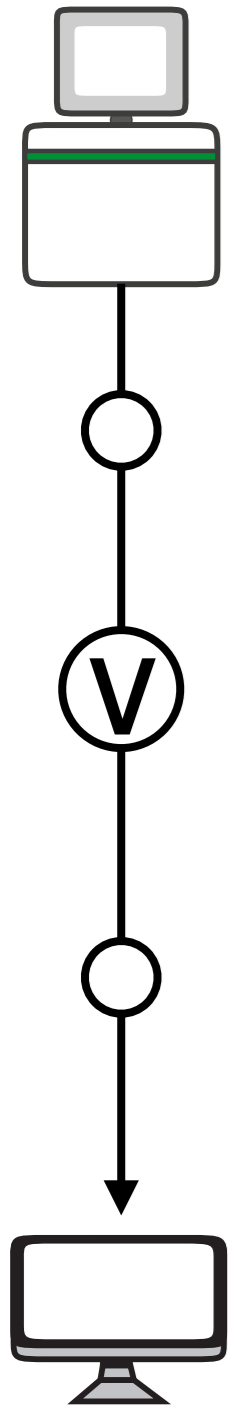


Functional annotation

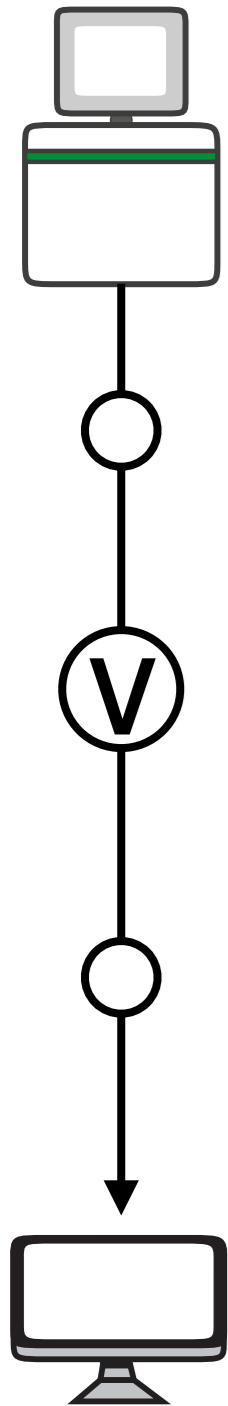


| Value | Rank | Explanation |
|------------|------|--|
| exonic | 1 | variant overlaps a coding |
| splicing | 1 | variant is within 2-bp of a splicing junction |
| ncRNA | 2 | variant overlaps a transcript without coding annotation in the gene definition |
| UTR5 | 3 | variant overlaps a 5' untranslated region |
| UTR3 | 3 | variant overlaps a 3' untranslated region |
| intronic | 4 | variant overlaps an intron |
| upstream | 5 | variant overlaps 1-kb region upstream of transcription start site |
| downstream | 5 | variant overlaps 1-kb region downstream of transcription end site |
| intergenic | 6 | variant is in intergenic region |

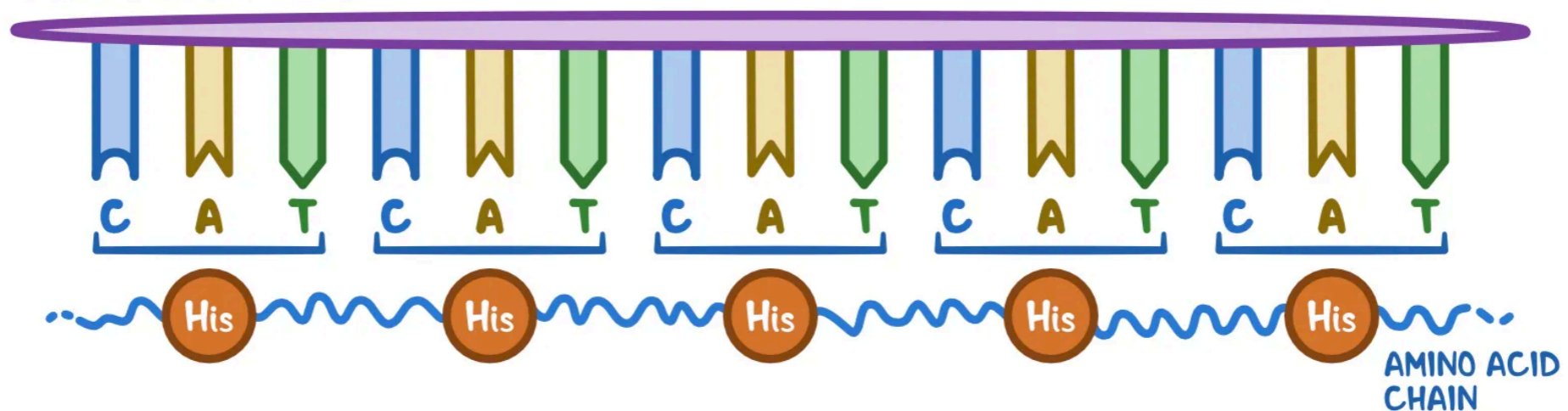
Functional annotation



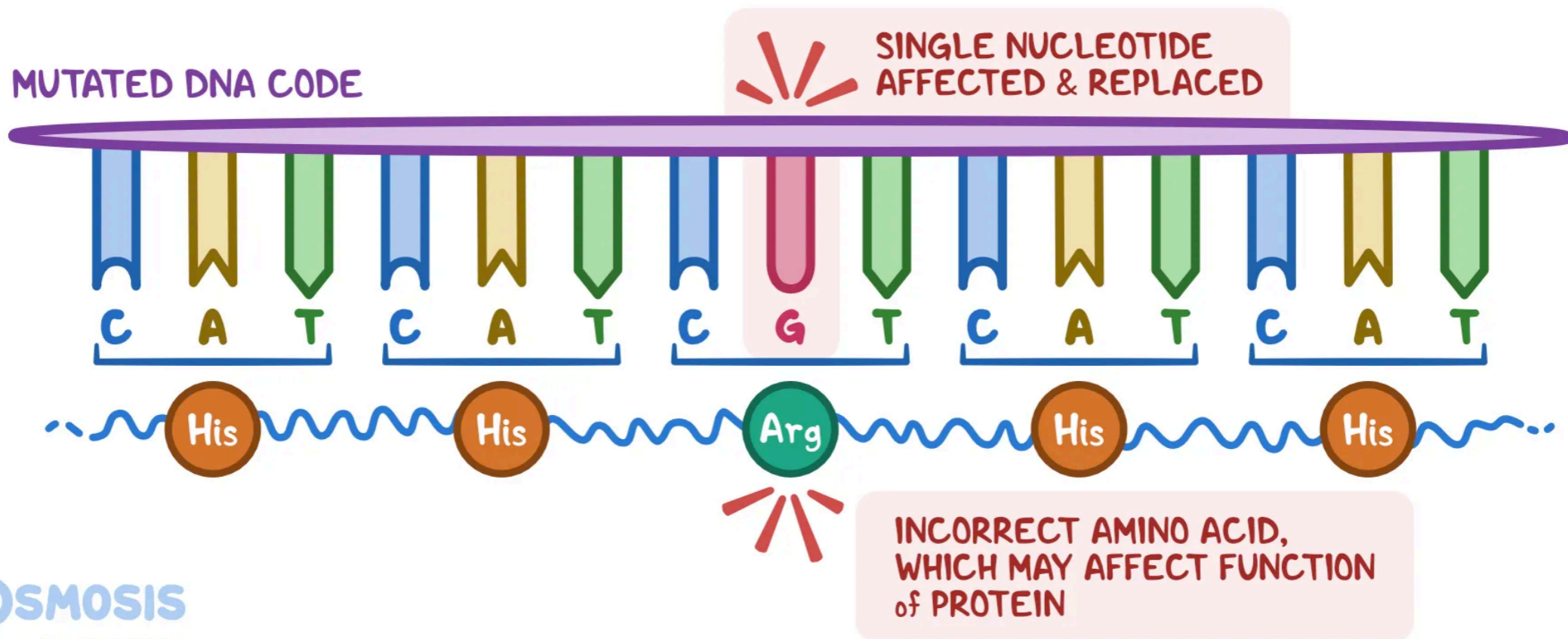
The effect of a mutation on protein



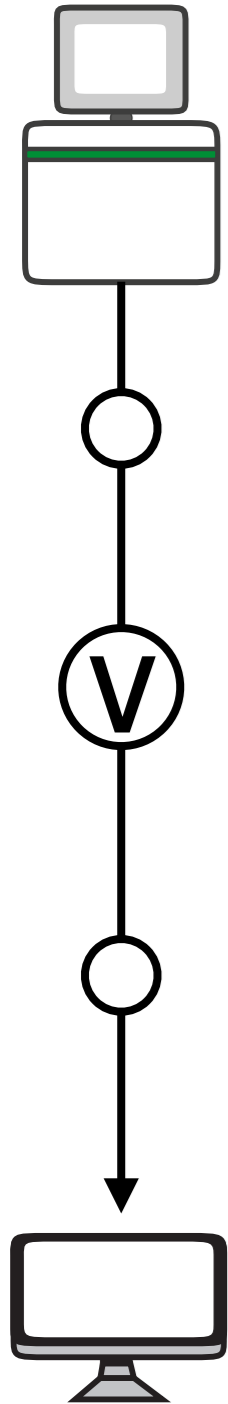
ORIGINAL DNA CODE



MUTATED DNA CODE

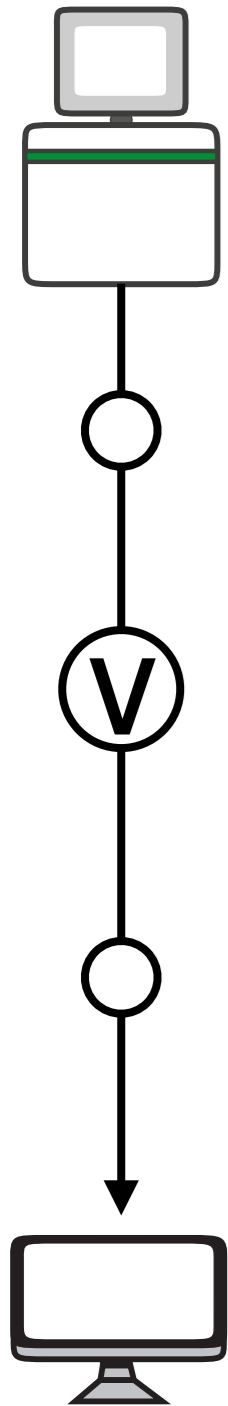


Exonic functional annotation



| Annotation | Rank | Explanation |
|---|------|---|
| frameshift insertion | 1 | an insertion of one or more nucleotides that cause frameshift changes in protein coding sequence |
| frameshift deletion | 2 | a deletion of one or more nucleotides that cause frameshift changes in protein coding sequence |
| frameshift block substitution | 3 | a block substitution of one or more nucleotides that cause frameshift changes in protein coding sequence |
| stopgain | 4 | a variant lead to the immediate creation of stop codon at the variant site. |
| stoploss | 5 | a variant that lead to the immediate elimination of stop codon at the variant site |
| nonframeshift insertion | 6 | an insertion of 3 or multiples of 3 nucleotides that do not cause frameshift changes in protein coding sequence |
| nonframeshift deletion | 7 | a deletion of 3 or multiples of 3 nucleotides that do not cause frameshift changes in protein coding sequence |
| nonframeshift block substitution | 8 | a block substitution of one or more nucleotides that do not cause frameshift changes in protein coding sequence |
| nonsynonymous SNV | 9 | a single nucleotide change that cause an amino acid change |
| synonymous SNV | 10 | a single nucleotide change that does not cause an amino acid change |
| unknown | 11 | unknown function (due to various errors in the gene structure definition in the database file) |

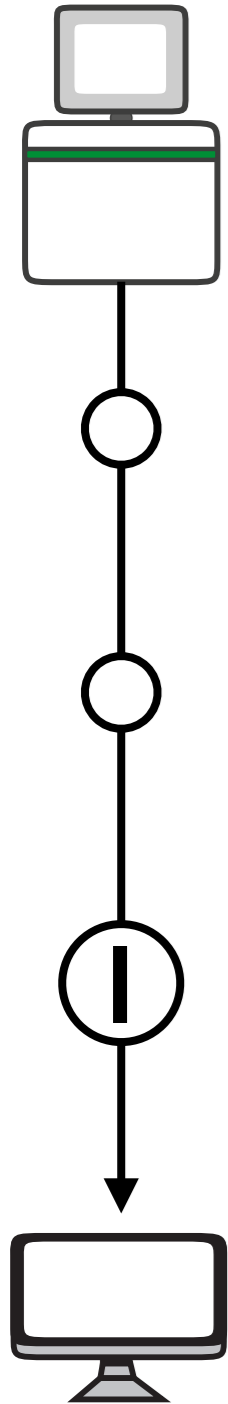
Exonic functional annotation (nonsilent mutations)



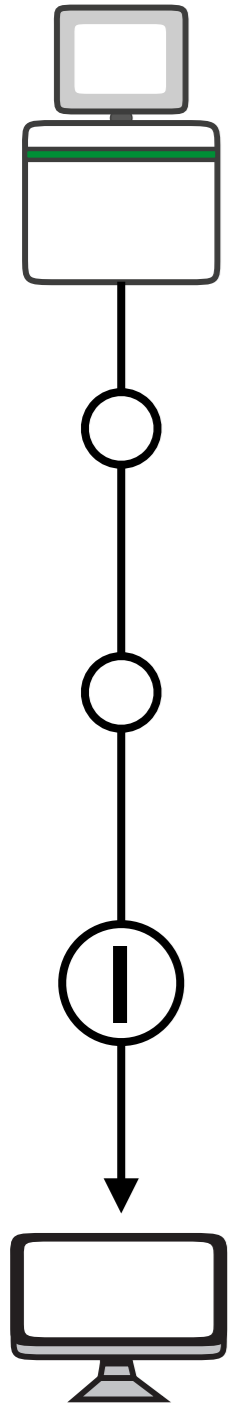
| Annotation | Rank | Explanation |
|---|------|---|
| frameshift insertion | 1 | an insertion of one or more nucleotides that cause frameshift changes in protein coding sequence |
| frameshift deletion | 2 | a deletion of one or more nucleotides that cause frameshift changes in protein coding sequence |
| frameshift block substitution | 3 | a block substitution of one or more nucleotides that cause frameshift changes in protein coding sequence |
| stopgain | 4 | a variant lead to the immediate creation of stop codon at the variant site. |
| stoploss | 5 | a variant that lead to the immediate elimination of stop codon at the variant site |
| nonframeshift insertion | 6 | an insertion of 3 or multiples of 3 nucleotides that do not cause frameshift changes in protein coding sequence |
| nonframeshift deletion | 7 | a deletion of 3 or multiples of 3 nucleotides that do not cause frameshift changes in protein coding sequence |
| nonframeshift block substitution | 8 | a block substitution of one or more nucleotides that do not cause frameshift changes in protein coding sequence |
| nonsynonymous SNV | 9 | a single nucleotide change that cause an amino acid change |
| synonymous SNV | 10 | a single nucleotide change that does not cause an amino acid change |
| unknown | 11 | unknown function (due to various errors in the gene structure definition in the database file) |

Clinical interpretation of variants

InterVar e **CancerVar** are curated databases for clinical interpretation. They contain catalogues of mutations previously pointed out to be pathogenetic or probably related to a disease.



Clinical interpretation of germline mutations



In 2015, the 'American College of Medical Genetics and Genomics (ACMG) e l'Association for Molecular Pathology (AMP) published standard criteria and updated guidelines for the clinical interpretation of sequence variants (relations between variants and human diseases).

InterVar generates an automatic interpretation based on 28 criteria, classifying mutations as '**Benign**', '**Likely benign**', '**Uncertain significance**', '**Likely pathogenic**' and '**Pathogenic**'

InterVar:Guideline

The American College of Medical Genetics and Genomics (ACMG) and the Association for Molecular Pathology (AMP) published in 2015 the updated standards and guidelines for the clinical interpretation of sequence variants, based on 28 criteria. However, variability between individual interpreters may be extensive due to lack of standard algorithms that implement these guidelines. This ACMG/AMP 2015 guideline is [at here](#)



Clinical Interpretation of genetic variants by ACMG/AMP 2015 guideline

InterVar is a bioinformatics software tool for clinical interpretation of genetic variants by the ACMG/AMP 2015 guideline. The input to InterVar is an annotated file generated from ANNOVAR, while the output of InterVar is the classification of variants into 'Benign', 'Likely benign', 'Uncertain significance', 'Likely pathogenic' and 'Pathogenic', together with detailed evidence code.

Search your **exonic** variants from pre-built wInterVar databases (updated *2022-June-13 17:57:28* with 100M sites):

If you already know the criteria of your variant, you can [click here](#) to interpret your variant directly.

This server is for exon variants interpretation only, if you have indels, you need to download the intervar tool from [github](#), then interpret your variant on local.

if you have cancer/somatic variant or CNV, you can click [CancerVar](#) to interpret your cancer variant directly.

if you have germline CNV, you can click [CNVinter](#) to interpret your copy number variation directly.

Please select the genomic version: hg19_updated.v.202107 ▾

Query by genomic coordinate

Chr: 1 ▾ POS: 115828756 Ref: G Alt: A

Query by dbSNP ID

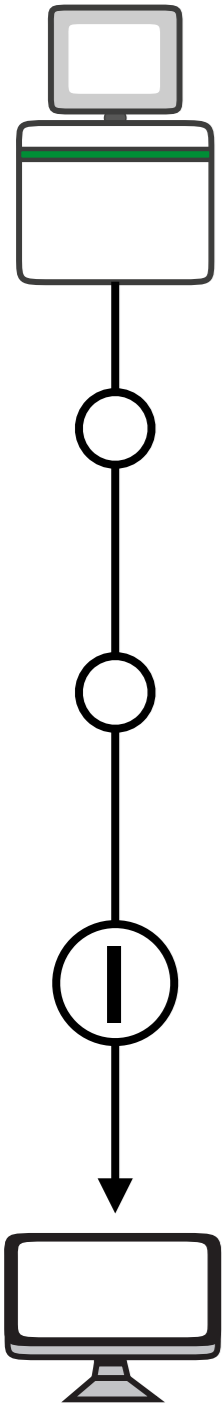
rs.: rs373849532

Query by HGNC gene symbol and cDNA

Gene: IL2RA cDNA change: c. C246A

Query by HGNC gene symbol and Protein Change

Gene: KRAS Protein change: p. G12C



Clinical Interpretation of genetic variants by ACMG/AMP 2015 guideline

InterVar is a bioinformatics software tool for clinical interpretation of genetic variants by the ACMG/AMP 2015 guideline. The input to InterVar is an annotated file generated from ANNOVAR, while the output of InterVar is the classification of variants into 'Benign', 'Likely benign', 'Uncertain significance', 'Likely pathogenic' and 'Pathogenic', together with detailed evidence code.

**Warning: All listed results were from the automated interpretation on default parameters!
Users are advised to examine detailed evidence and use prior knowledge on ethnicity/disease to perform manual adjustments.**

Database version:hg19_update

You searched by HGNC gene symbol with name as **KRAS**, and Protein change: **p.G12C**

Show/hide columns

Restore columns

Copy to clipboard

Download result as CSV

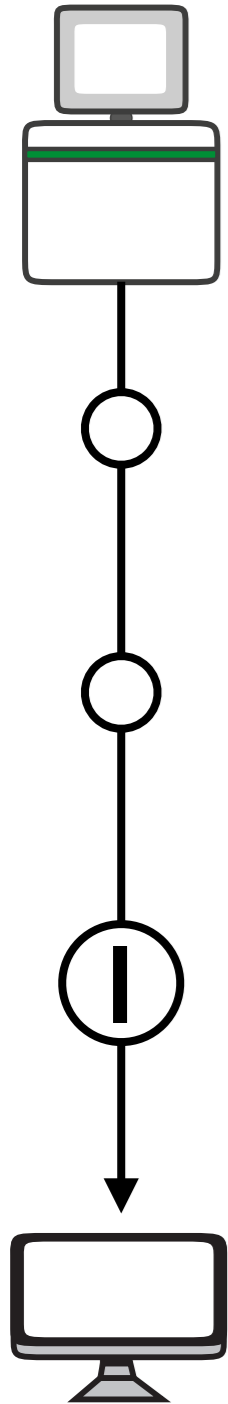
Search:

| Chr | Position | Ref | Alt | Gene (refGene) | Intervar | ExonicFunc (refGene) | SNP | Transcripts (Ref) | MAF in gnomAD_ALL(genome) | Disease in OrphaNet |
|-----|----------|-----|-----|----------------|---------------------------------------|----------------------|-----------------------------|--|---------------------------|--|
| 12 | 25398285 | C | A | KRAS | Pathogenic (Details&Adjust) | nonsynonymous SNV | rs121913530(details of MAF) | NM_004985 p.G12C NM_033360 p.G12C | (show in 7 POPs) | 1340 268114 648 519 1333 2612 26106 227535 |

Showing 1 to 1 of 1 entries

Previous 1 Next

Clinical interpretation of somatic mutations



CancerVar is a bioinformatic tool for the clinical interpretation of somatic variants based on guidelines of the AMP/ASCO/CAP/CGC 2017-2019

CancerVar classifies somatic variants as:

- Tier I/**Pathogenic**
- Tier II/**Likely pathogenic**
- Tier III/**Variants of Unknown Clinical Significance (VUS)**
- Tier IV/**Benign or Likely Benign Variants**

Clinical interpretation of somatic mutations

Tier I: Variants of Strong Clinical Significance

Therapeutic, prognostic & diagnostic

Level A Evidence

FDA-approved therapy
Included in professional guidelines

Level B Evidence

Well-powered studies with consensus from experts in the field

Tier II: Variants of Potential Clinical Significance

Therapeutic, prognostic & diagnostic

Level C Evidence

FDA-approved therapies for different tumor types or investigational therapies
Multiple small published studies with some consensus

Level D Evidence

Preclinical trials or a few case reports without consensus

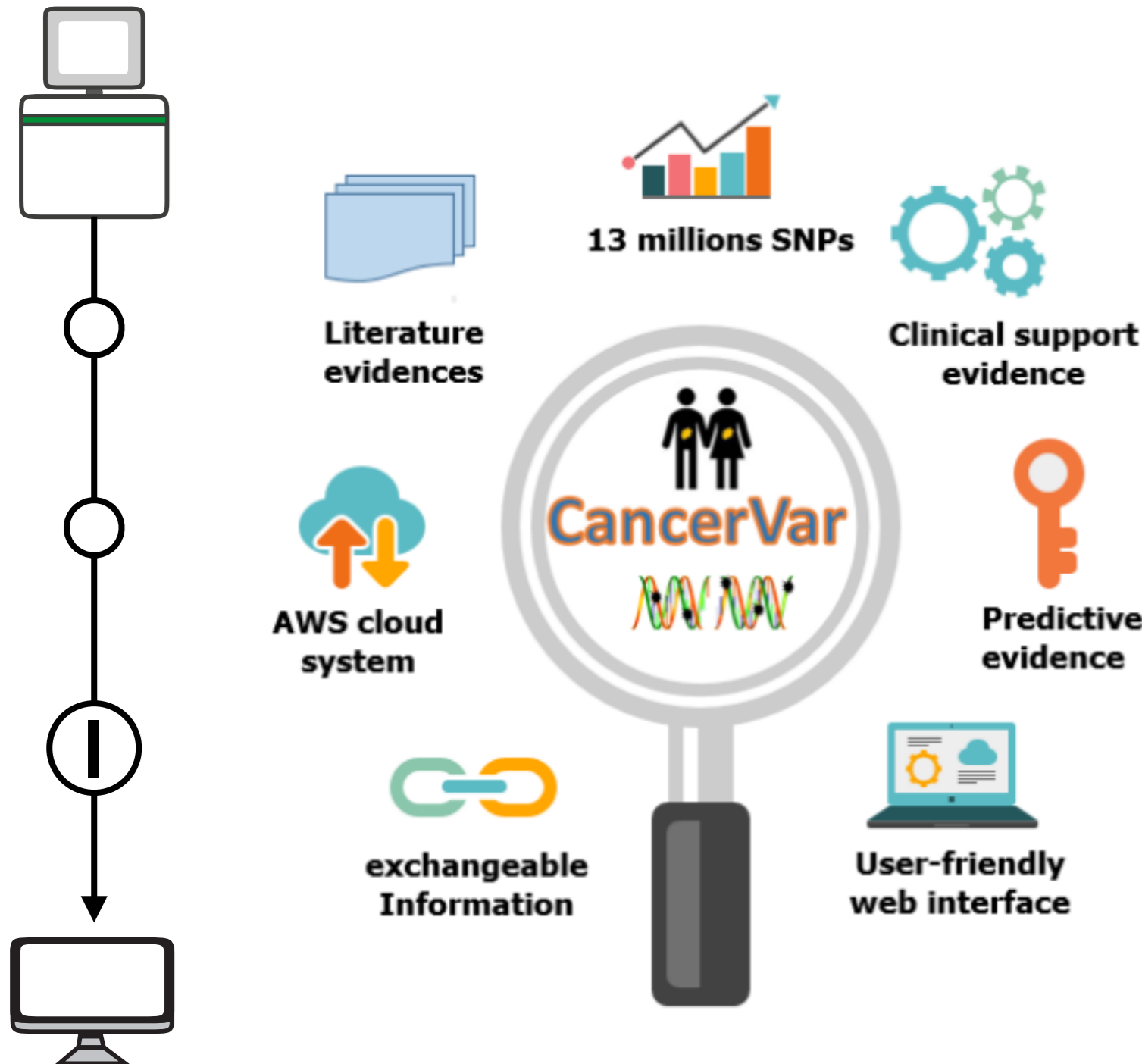
Tier III: Variants of Unknown Clinical Significance

Not observed at a significant allele frequency in the general or specific subpopulation databases, or pan-cancer or tumor-specific variant databases
No convincing published evidence of cancer association

Tier IV: Benign or Likely Benign Variants

Observed at significant allele frequency in the general or specific subpopulation databases
No existing published evidence of cancer association

Clinical interpretation of somatic mutations



clinical-based prediction

| | |
|-----------------------------|---|
| Therapeutic | FDA approved or investigational with strong evidence |
| Diagnostic | In Professional Guideline or reported evidence with consensus |
| Prognostic | In Professional Guideline or reported evidence with consensus |
| Mutation type | Activating, LOF (missense, nonsense, indel, splicing) |
| Population data | Absent or extremely low MAF in population databases |
| Somatic data | Most likely present in COSMIC and ICGC |
| Predictive Softwares | Most of Prediction showed pathogenic in |
| Pathway involvement | Involve disease-associated pathways or pathogenic pathways |
| Germline database | Present in <u>Clinvar</u> /HGMD as Pathogenic |
| Publication database | Publication reported pathogenic evidence with consensus |

CancerVar: a web server for improved AI and evidence-based clinical interpretation for cancer somatic mutations

CancerVar is a bioinformatics software tool for clinical interpretation of somatic variants.

Search your **exonic** variants from pre-built CancerVar databases(updated: with 13 Million mutations in 1911 cancer census gene):

This server is for exon variants, CNVs and some known indels interpretation only, if you have novel indels, you need to download the CancerVar tool from [github](#), then interpret your variant on local.

Please select the **genomic version**: and **cancer type**:

Query by **genomic coordinate**

Chr: POS: Ref: Alt:

Query by **dbSNP ID**

rs. :

Query by **HGNC gene symbol and cDNA**

Gene: cDNA change: c.

Query by **HGNC gene symbol and Protein Change**

Gene: Protein change: p.

Query by **HGNC gene symbol or Alternations**

Gene:

CancerVar Results

You searched by HGNC gene symbol with name as **KRAS**, cancer types: **All_types** and Protein change: **P.G12C**

GENOMIC VERSION

hg19



CANCER TYPE

All_types



ALTERNATIONS

Mutation



INTERPRETATION SUMMARY

Based on Evidence(CBPs), CancerVar assign the clinical significance of your somatic mutation in **KRAS** as:

Tier_I_strong and Score: **11(with sub-scores)**. Deep learning Score from OPAI: **0.99(Oncogenic)**.

This variant is **nonsynonymous SNV** in gene **KRAS**, located in chromosome **12:25398285**. There are some clinical and/or experimental evidence showed strong/potential clinical significance **in Therapeutic,Diagnosis,Prognosis**, From the population databases(gnomAD,1000 genome etc), this variant is **absent or extremely low minor allele frequency(MAF<0.1%)**. In Germline database of Clinvar/HGMD, this variant's clinical significane is **Pathogenic**. In most of the pathogenic or deleterious prediction softwares/algorithms, this variant was predicted as **Pathogenic**. When check the occurrence of somatic database in COSMIC or ICGC, this variant shows in **both** of them. From the KEGG pathway database, the gene of this variant **does** involve in disease-associated pathways or pathogenic pathways. Currently, searching the pubmed website, there are **some** publications from functional study, population study or other study as supporting evidence for clinical/biological significance.

Please review the cards below to get the detail of the interpretation.



Evidence Overview

CBP1:Therapeutic: FDA approved or investigational with strong evidence.In total of **78** records (you specified cancer types as: **All_types**),most of the cancers are located in:**Lung 35%,Colorectal 33%,Cancer 8%**,and they have been mostly treated with:**Melphalan 5%,Cetuximab 4%,EGFR mAb inhibitor 4%**,the treatment of the drugs are:**Responsive 35%,Sensitivity/Response 28%,Resistance 21%**,if you need more information, please click [Detail...](#)

CBP2:Diagnostic: In Professional guideline or reported evidence with consensus.In total of **1** records (you specified cancer types as: **All_types**),most of the cancers are located in:**Lung 100%**, the diagnostic are:**Positive 100%**,if you need more information, please click [Detail...](#)

CBP3:Prognostic: In Professional guideline or reported evidence with consensus.In total of **4** records (you specified cancer types as: **All_types**),most of the cancers are located in:**Lung 75%,Other 25%**, the prognostic are:**Poor Outcome 100%**,if you need more information, please click [Detail...](#)

CBP4:Mutation type: Activating, LOF (missense, nonsense, indel, splicing), CNAs, fusions.

CBP5:Variant frequencies:Mostly mosaic. Need user's knowledge.

CBP6:Potential germline: Mostly nonmosaic. Need user's knowledge.

CBP7:Population databases: Absent or extremely low MAF. [MAF In GnomAD_genome . \(show in 7 POPs\)](#)

CBP8:Germline databases: may be present in HGMD/ClinVar **Pathogenic**.

CBP9:Somatic databases: Most present in COSMIC, ICGC, My Cancer Genome, TCGA.

CBP10:Predictive from: SIFT, PolyPhen2,MutationAssessor,MetaSVM,MetaLR,FATHM M,GERP++_RS, and mostly as **Pathogenic**.

CBP11:Pathway: involve in Disease-associated pathways or pathogenic pathways. [KEGG Pathway](#)

CBP12:Publications: Convincing evidence from Functional study, population study, other.

Evidence

Mutation Information

Chromosome: 12

Position: 25398285

Reference Allele: C

Alternative Allele: A

Function in refGene: nonsynonymous SNV

Minor allele frequency in(. means absent): ESP6500:.

1000 genome:.

gnomAD genome_ALL: [More](#)

ExAC:1.976E-5

Transcript in refGene: [NM_033360](#)

Exon location:exon2

Nucleotide change:c.34G>T

Residue change :p.G12C

Clinvar:Pathogenic/Likely_pathogenic

While the KRAS G12 region is a widely studied recurrent region in cancer, its impact on clinical action is still debated. Often associated with tumors that are wild-type for other drivers (EGFR and ALK specifically), the prognosis for patients with this mutation seems to be worse than the KRAS wild-type cohort in patients with colorectal and pancreatic cancer, however this hypothesis is in need of further validation. This mutation, along with the mutations affecting the neighboring G13 position, may result in a less responsive tumor when treated with first-generation TKI's like gefitinib. However, cetuximab treatment was shown to extend survival in a cohort of colorectal patients.

Deep learning Score: **0.99**(**OPA!**,**Oncogenic**)

CancerVar: **Tier_1_strong with score: 11**

Need to review scores or adjust and manually re-interpret? Please Click [Adjust!](#)

● Clinical significance ● Unknown ● Benign

Mutation

Gene Information

KRAS

Name: [Kirsten rat sarcoma viral oncogene homolog](#)

Location: [chr12:25357723-25403870\(Grch37\)](#)

Cytoband: [12p12.1](#)

GeneCards Summary:

KRAS (KRAS Proto-Oncogene, GTPase) is a Protein Coding gene. Diseases associated with KRAS include Oculoectodermal Syndrome and Noonan Syndrome 3. Among its related pathways are Oocyte meiosis and Oxytocin signaling pathway. Gene Ontology (GO) annotations related to this gene include GTP binding. An important paralog of this gene is NRAS. [More on GeneCards](#)

CIViC Summary:

Mutations in the RAS family of proteins are frequently observed across cancer types. The amino acid positions that account for the overwhelming majority of these mutations are G12, G13 and Q61. The different protein isoforms, despite their raw similarity, also behave very differently when expressed in non-native tissue types, likely due to differences in the C-terminal hyper-variable regions. Mis-regulation of isoform expression has been shown to be a driving event in cancer, as well as missense mutations at the three hotspots previously mentioned. While highly recurrent in cancer, attempts to target these RAS mutants with inhibitors have not been successful, and has not yet become common practice in the clinic. The prognostic implications for KRAS mutations vary between cancer types, but have been shown to be associated with poor outcome in colorectal cancer, non-small cell lung cancer, and others.

Gene

Clinical publications



Pathway



Domain



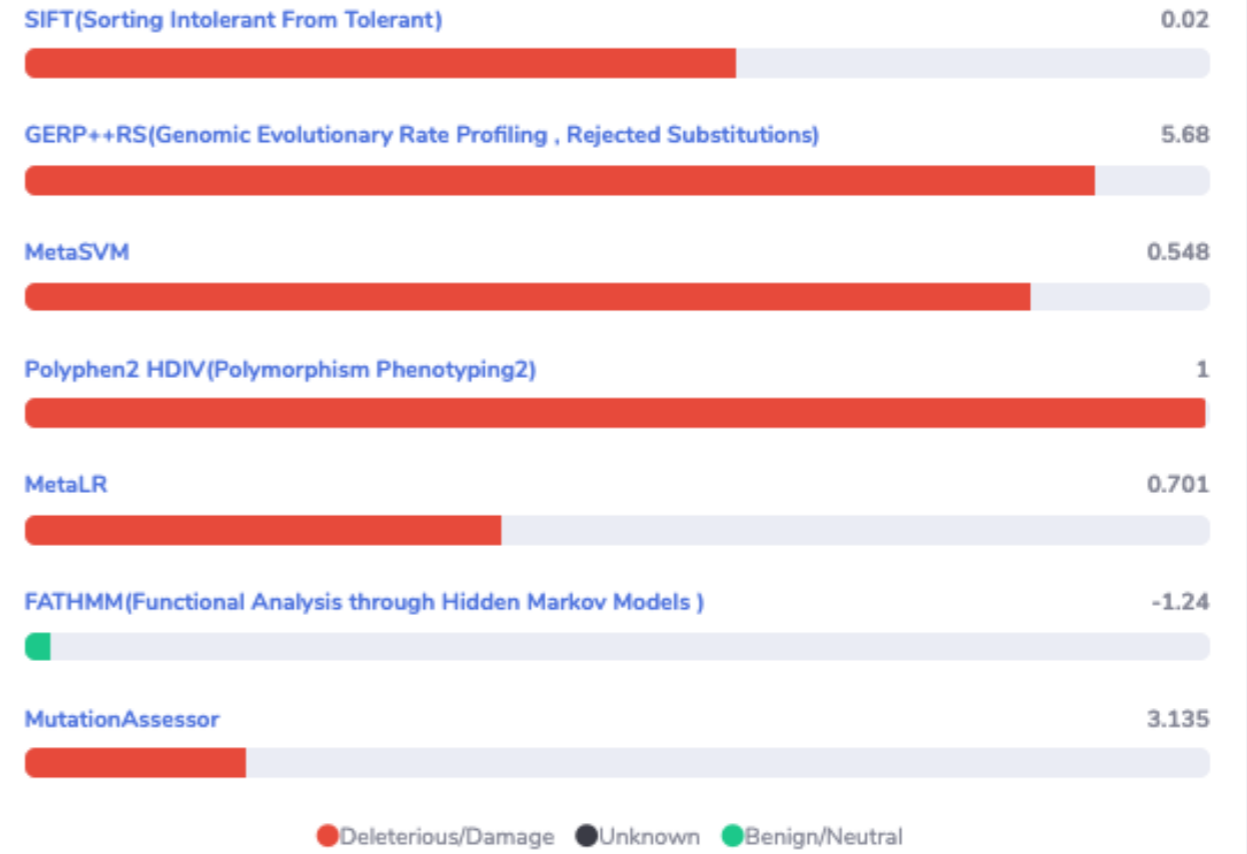
P-loop containing nucleoside triphosphate hydrolase;Small GTP-binding protein domain

More resources

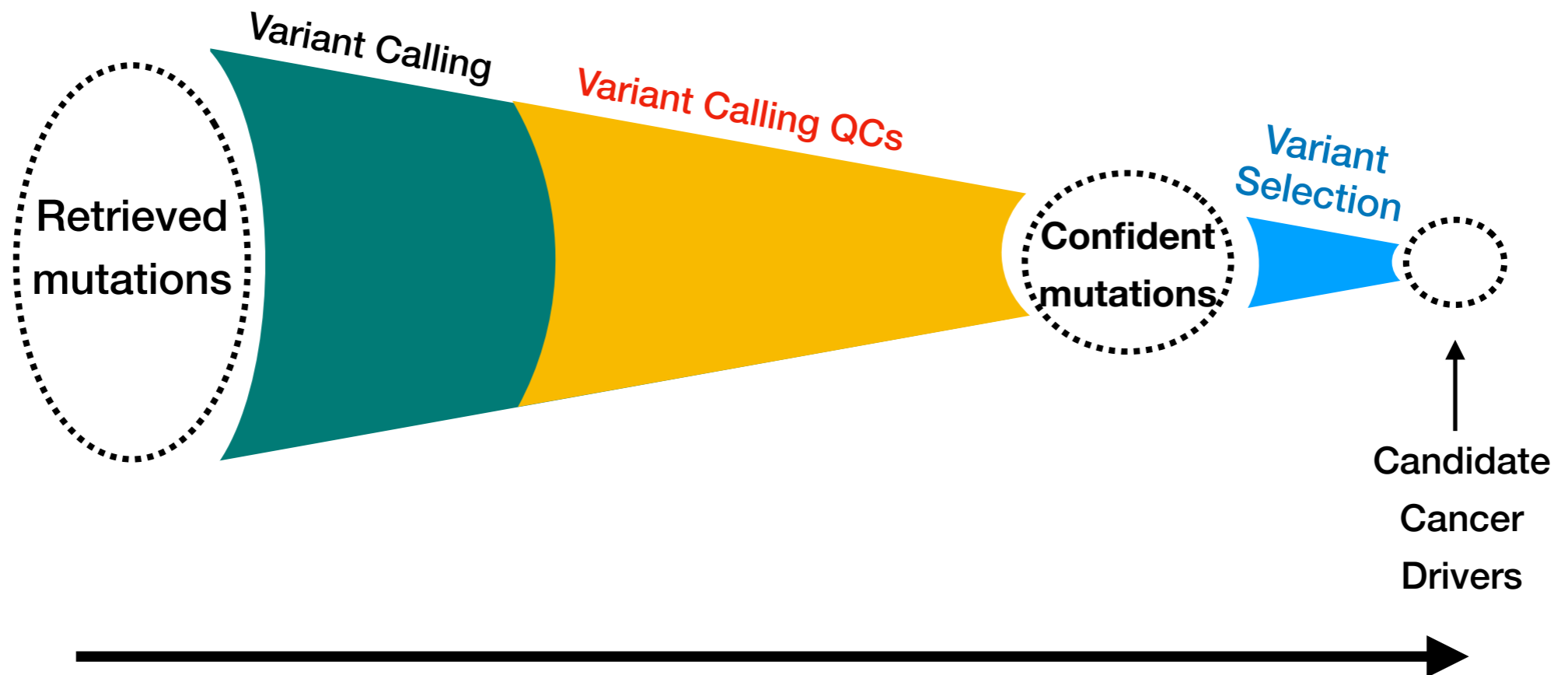
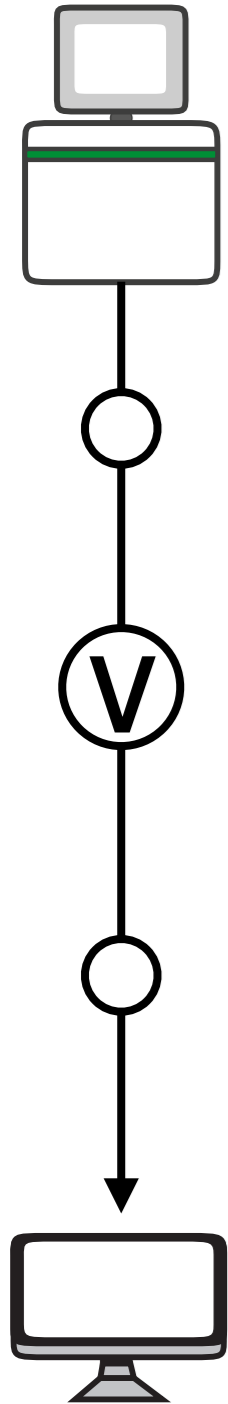


- [Gene on Oncokb](#)
- [Variant on Oncokb](#)
- [Clinvar](#)
- [Cosmic](#)

Deleterious Predictions from other softwares(. means no prediction)



Candidate cancer driver

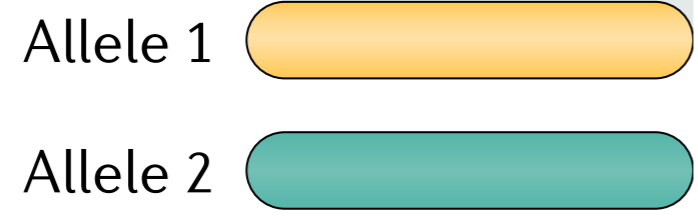


Copy Number Alterations

Allelic copies variation

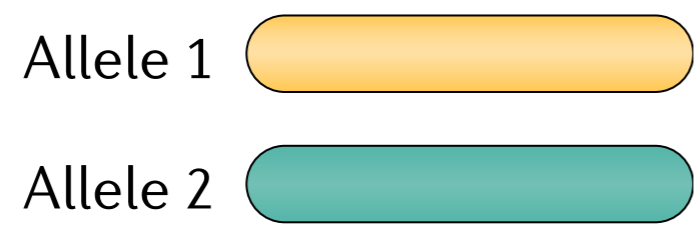


Normal genome

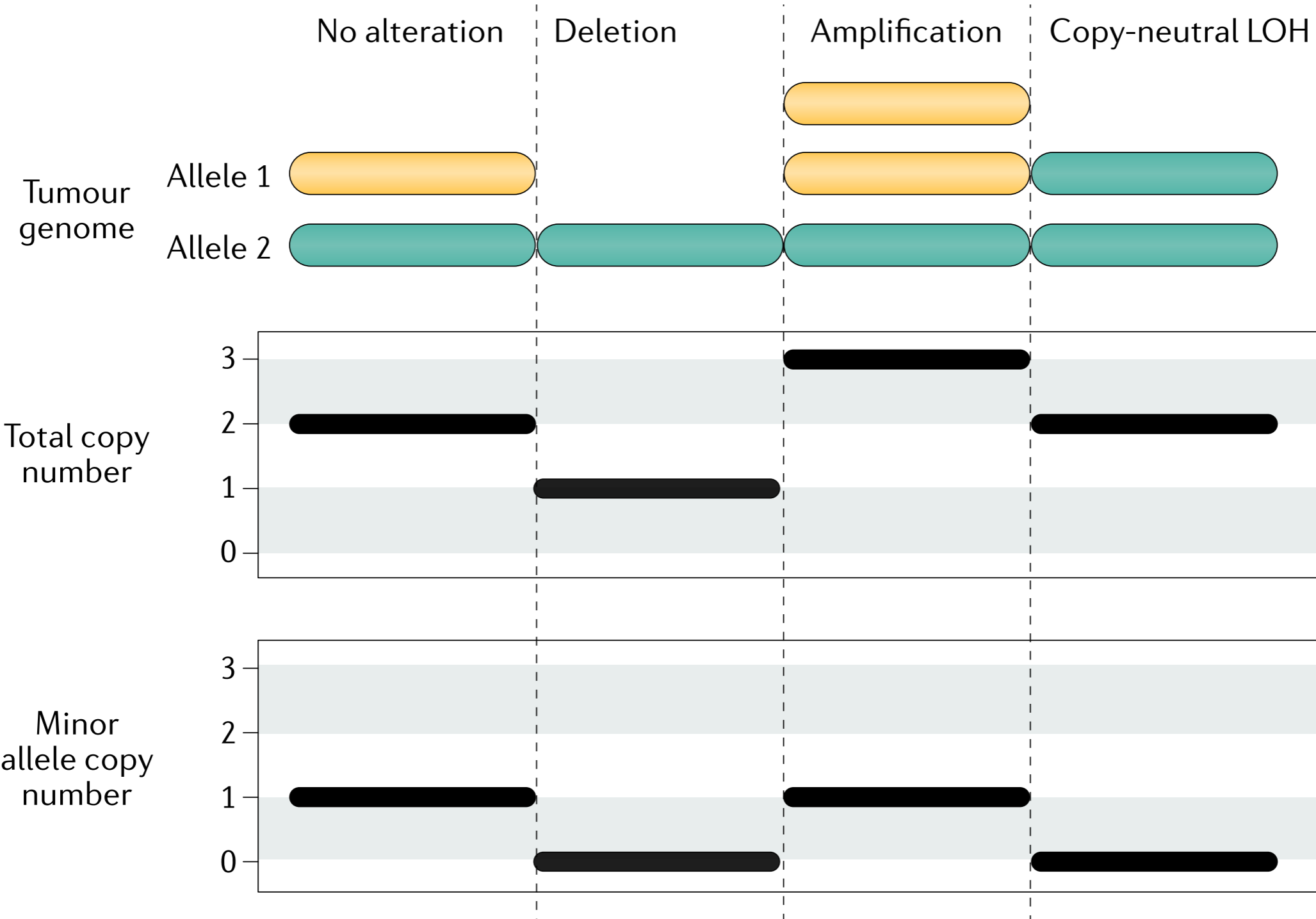


No alteration

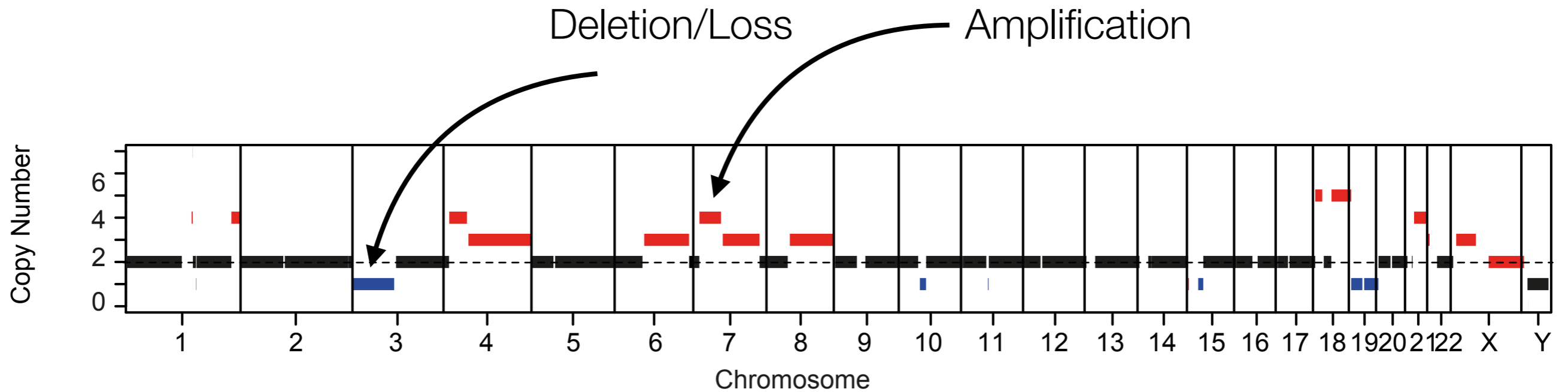
Tumour genome



Allelic copies variation

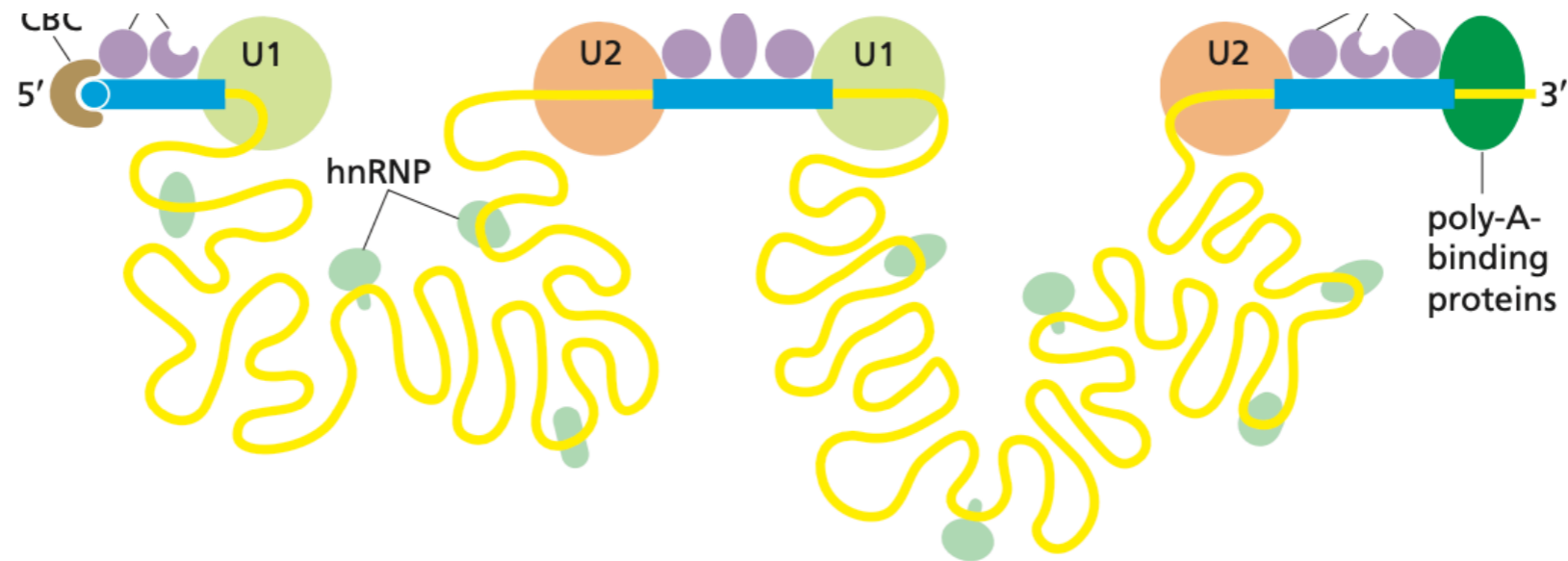


Genome CNV profile



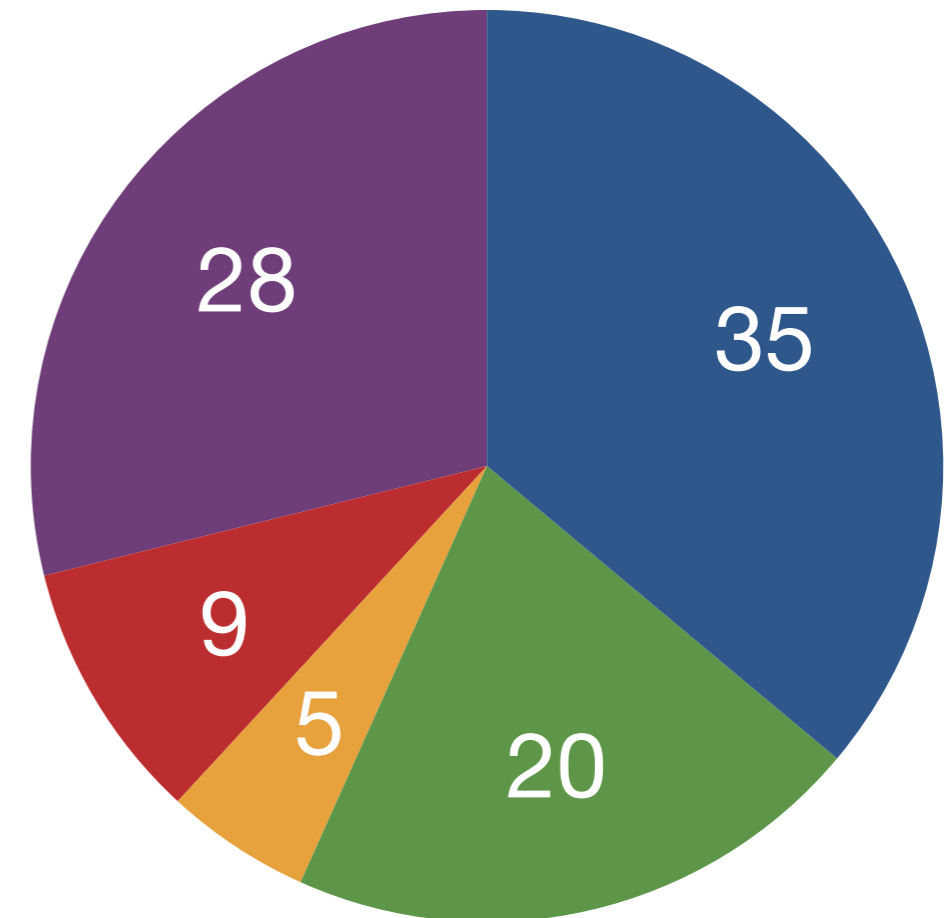
RNA-sequencing

Sequencing transcriptome



- Evaluate **expression** of genes/transcripts for:
 - All species of RNA
 - mRNA
 - small RNAs
- Evaluate expression levels of exons
 - Patterns of alternative splicing
- Evaluate transcriptional **alterations**
- Annotate **regions** and **functional elements**

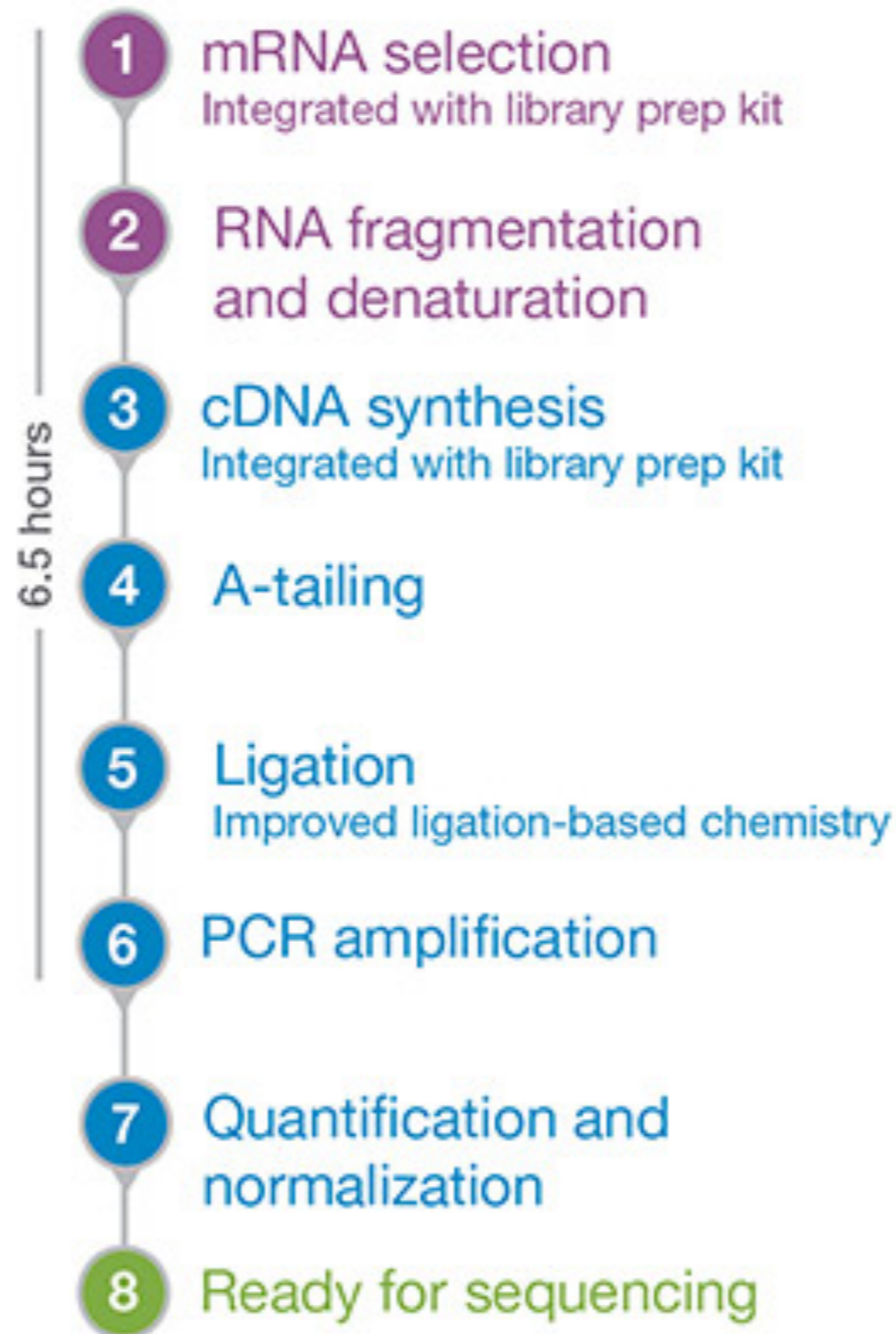
Protocolli di sequenziamento



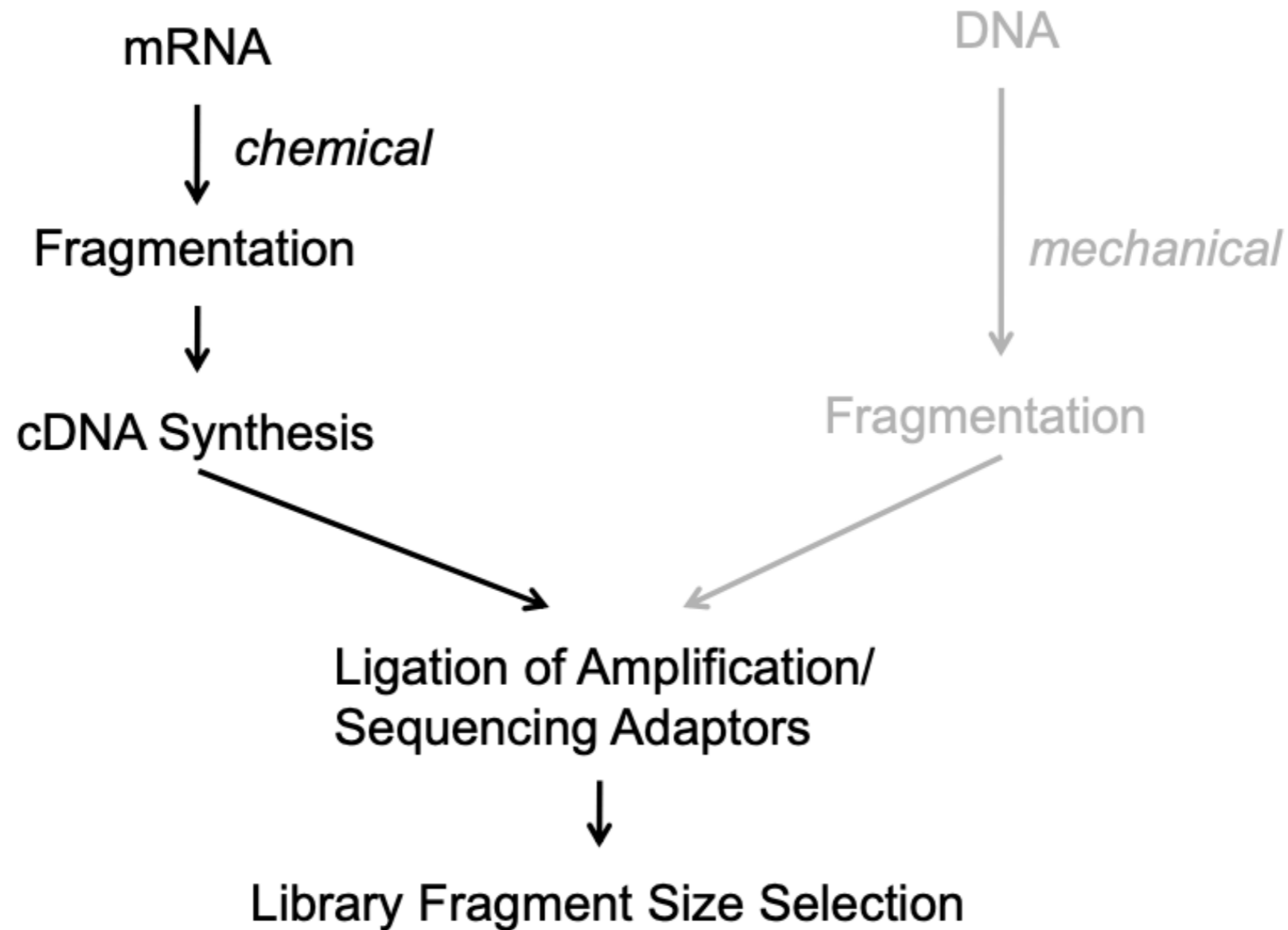
- RNA transcription
- RNA-Protein interactions
- RNA modifications
- RNA structure
- Low-level RNA detection

| RNA Transcription | RNA-Protein Interactions | RNA Modifications | RNA Structure | Low-Level RNA Detection |
|-------------------|--------------------------|-------------------|---------------|-------------------------|
| RNA-Seq | Ribo-Seq | MeRIP-Seq | SHAPE-Seq | scRNA-Seq |
| CaptureSeq | RIP-Seq | miCLIP-m6A | icSHAPE | SUPeR-Seq |
| RASL-Seq | CLIP-Seq | PSI-Seq | CIRS-Seq | UMI |
| ClickSeq | Pol II CLIP | Pseudo-Seq | SHAPE-MaP | Digital RNA Sequencing |
| 3Seq | miR-CLIP | ICE | DMS-Seq | MARS-Seq |
| cP-RNA-Seq | eCLIP | | SPARE | Quartz-Seq |
| 3P-Seq | irCLIP | | PARS-Seq | DP-Seq |
| 2P-Seq | PAR-CLIP | | Cap-Seq | Smart-Seq |
| 3'-Seq | iCLIP | | CIP-TAP | FRISCR |
| TIF-Seq | BrdU-CLIP | | | CEL-Seq |
| PEAT | AGO-CLIP | | | STRT-Seq |
| SMORE-Seq | PIP-Seq | | | TCR Chain Pairing |
| TL-Seq | hiCLIP | | | TCR-LA-MC PCR |
| TATL-Seq | RBNS | | | CirSeq |
| RARseq | TRIBE | | | TIVA |
| TAIL-Seq | HiTS-RAP | | | PAIR |
| PAL-Seq | TRAP-Seq | | | CLaP |
| FRT-S wcell | DLAF | | | CytoSeq |
| ChIRP | miTRAP | | | Drop-Seq: |
| CHART | CLASH | | | Hi-SCL |
| RAP | | | | InDrop |
| GRO-seq | | | | snRNA-Seq |
| Bru-Seq | | | | Nuc-Seq |
| BruChase-Seq | | | | Div-Seq |
| 5'-GRO-Seq | | | | SCRB-Seq |
| BruDRB-Seq | | | | G&T-Seq |
| 4sUDRB-Seq | | | | scM&T-Seq |
| PRO-Seq | | | | scTrio-seq |
| PRO-Cap | | | | |
| CAGE | | | | |
| 3'NT Method | | | | |
| NET-Seq | | | | |
| mNET-Seq | | | | |
| PARE-Seq | | | | |
| GMUCT | | | | |

Illumina Stranded mRNA Prep

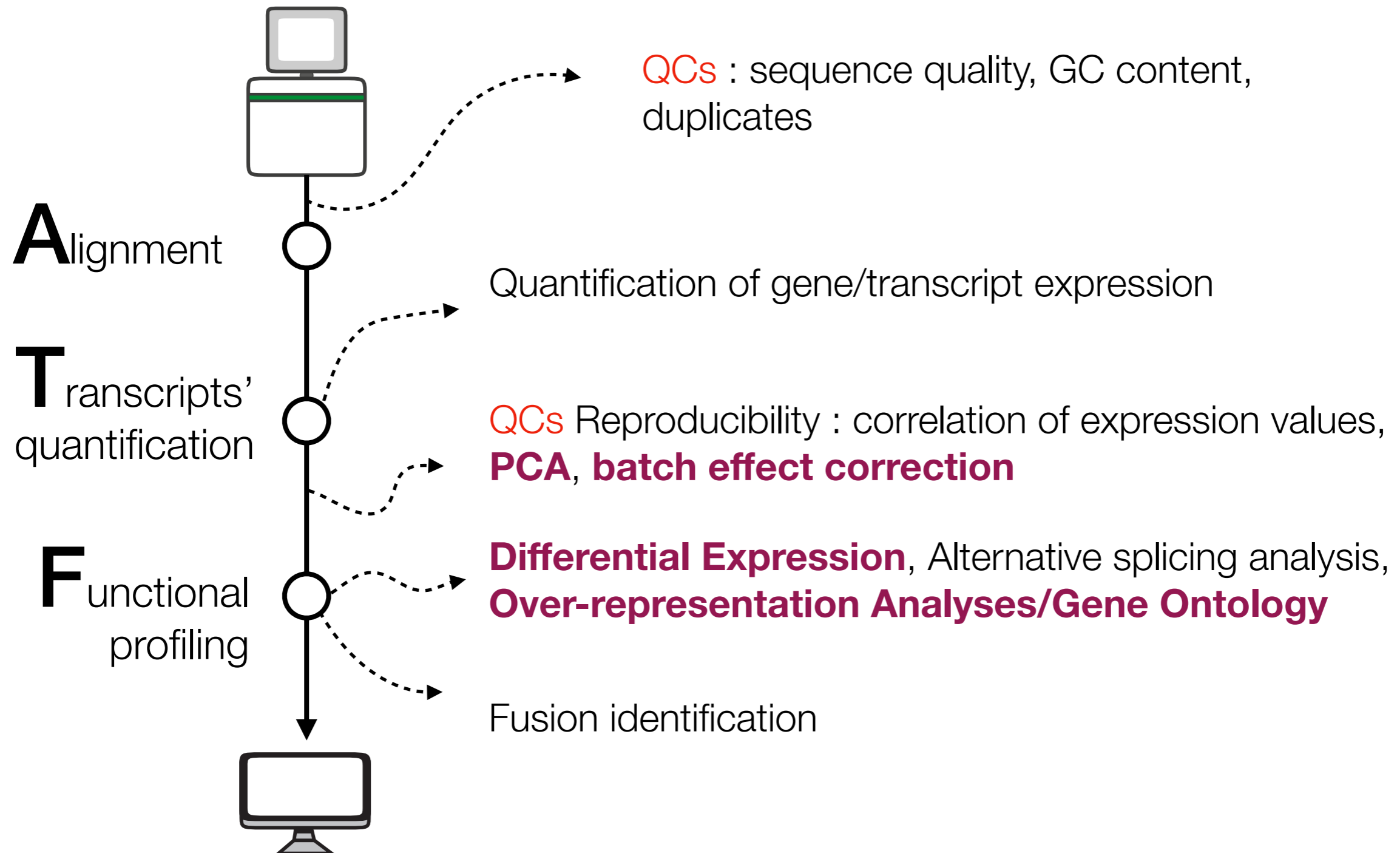


Key steps in sequencing

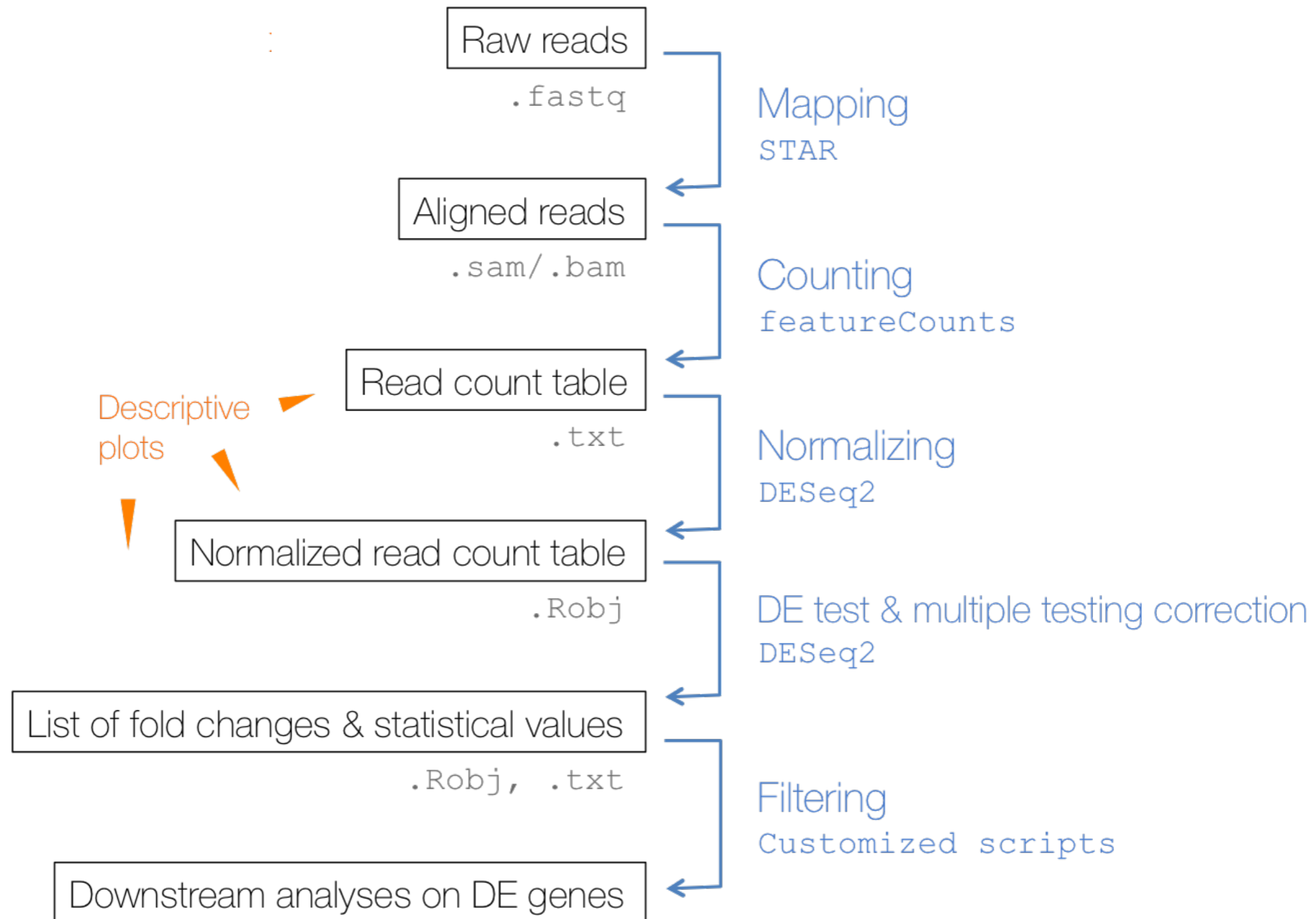


Deciphering gene expression

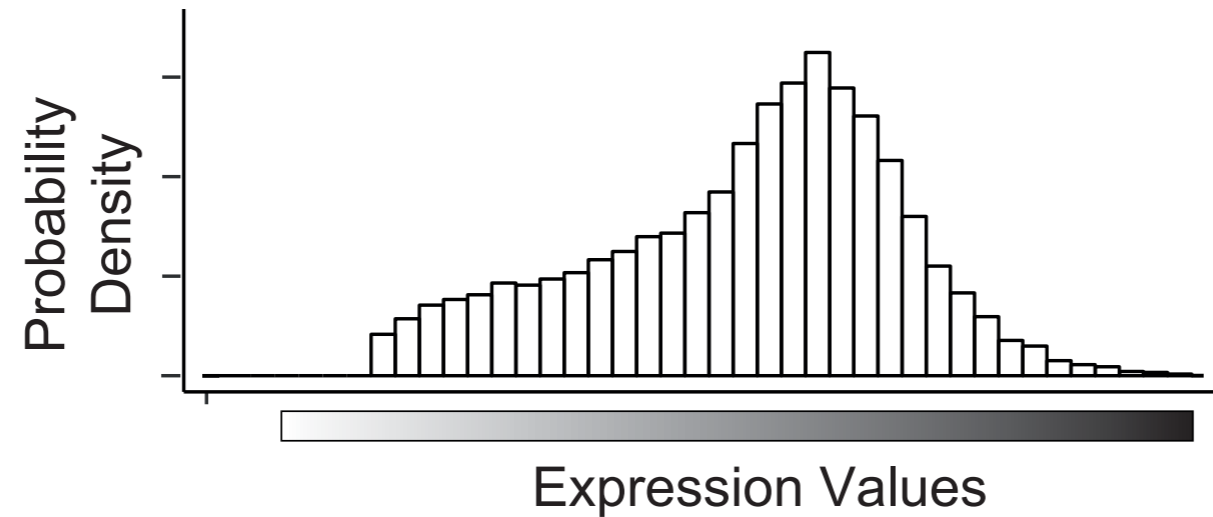
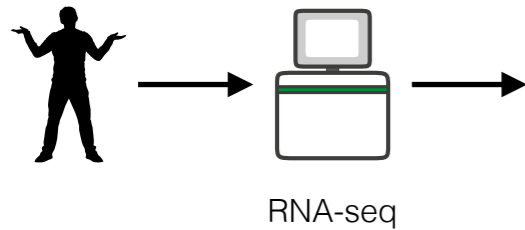
RNA-seq data analysis workflow:



Deciphering gene expression

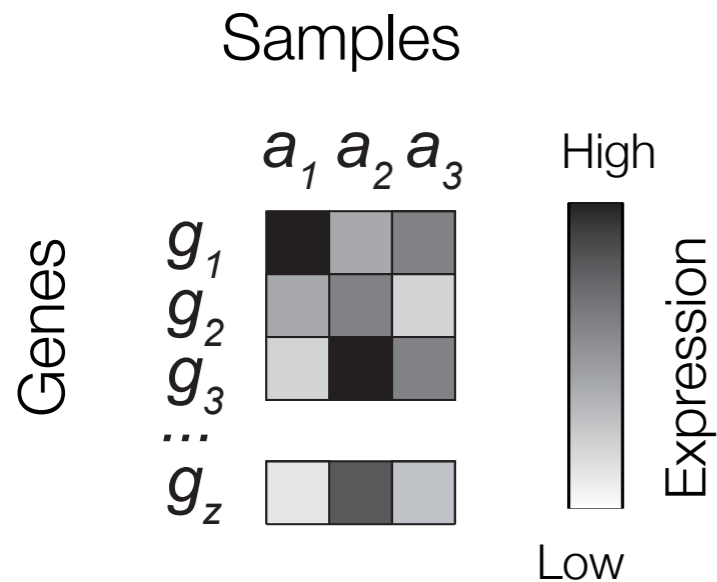


Gene expression level distribution



<https://academic.oup.com/nar/article/48/4/1730/5691219>

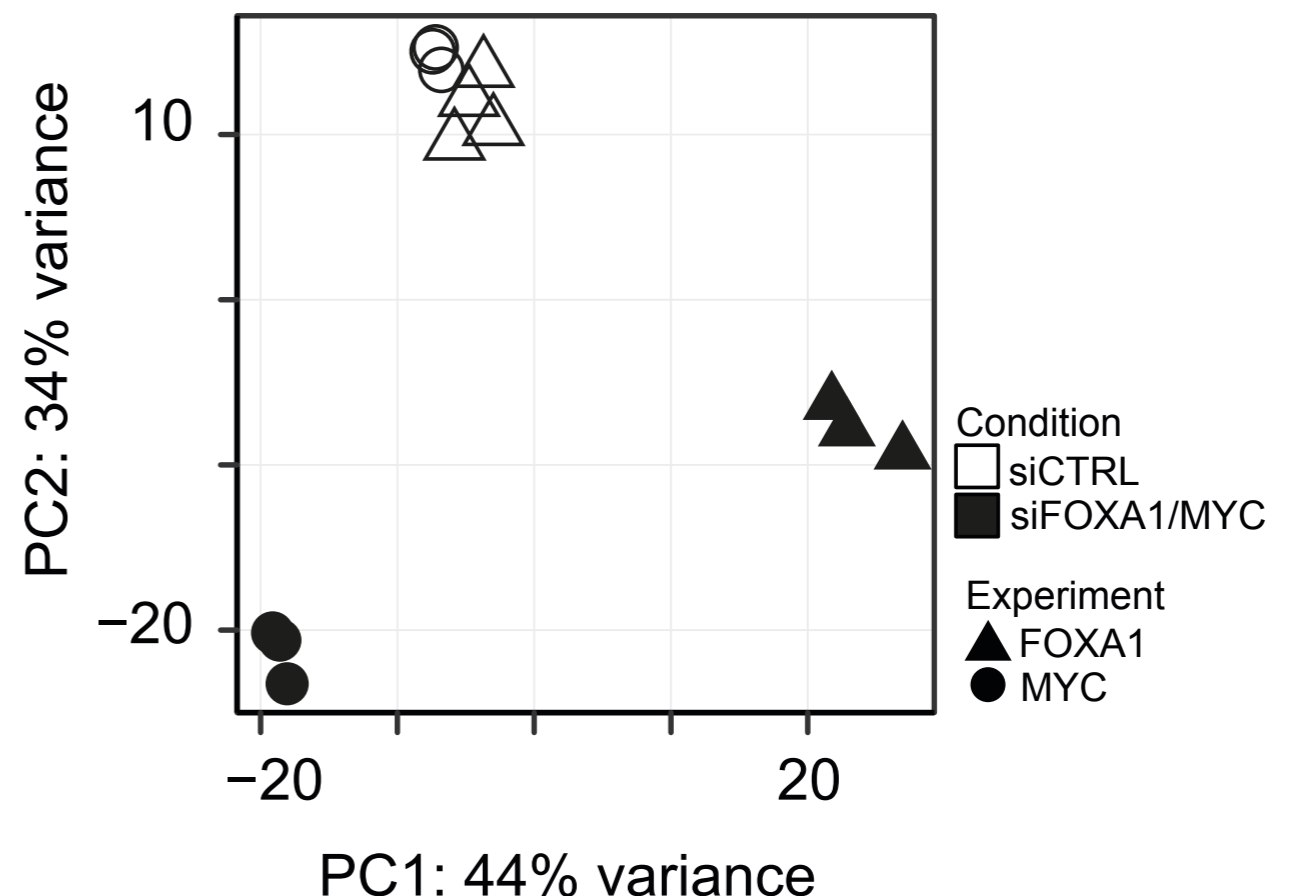
Gene expression level distribution



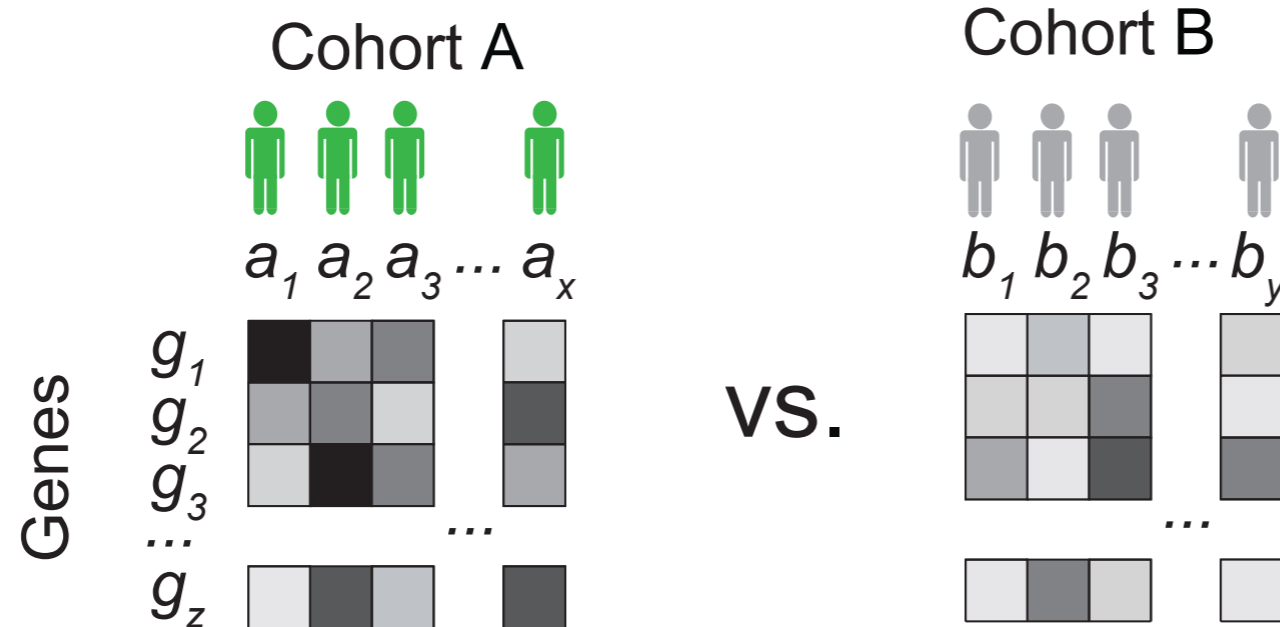
QC's Reproducibility

Principal Component Analysis (PCA)

PCA is a statistical technique for dimensionality reduction. We use PCA when a dataset presents a high number of features (genes in this case). It is like compressing information about ~20,000 in two dimensions or some more if we need it.



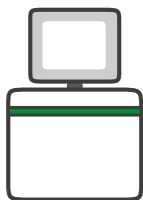
Differential expression



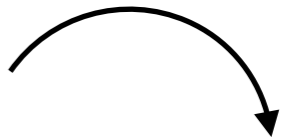
Two are the main goals of a differential expression (DE) analysis:

1. Estimate the **entity of variation** between the two conditions, i.e. calculate Fold Change (FC)
2. Estimate the **significance of the difference**, i.e. p-value, and correct it for multiple testing (p-adjusted).

DESeq2

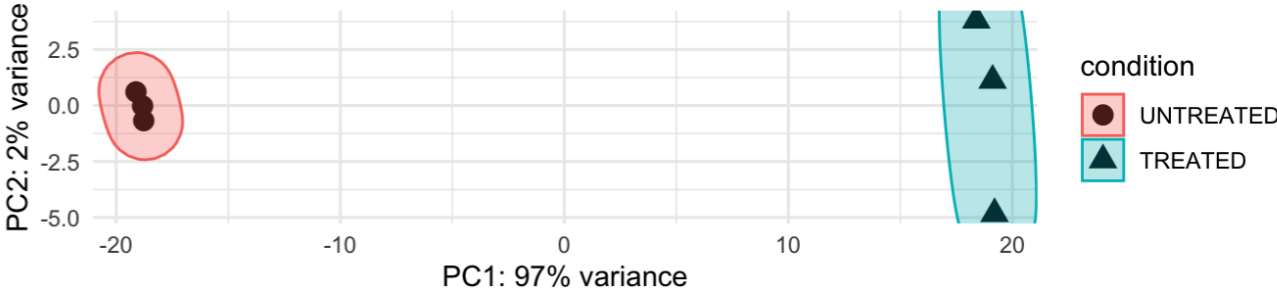
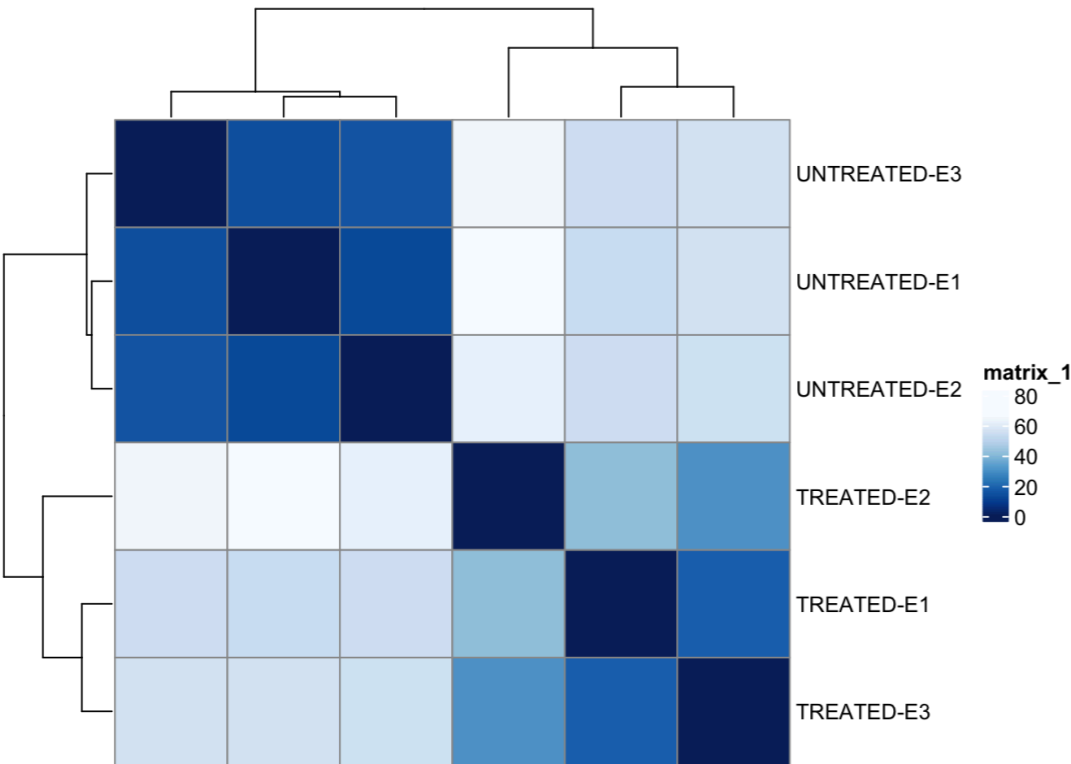


```
##  sampleName          fileName condition experiment
## 1  UNTREATED1  CTRL-1.star.nsrt.count  UNTREATED      E1
## 2  UNTREATED2  CTRL-2.star.nsrt.count  UNTREATED      E2
## 3  UNTREATED3  CTRL-3.star.nsrt.count  UNTREATED      E3
## 4   TREATED1  TREATED-1.star.nsrt.count   TREATED      E1
## 5   TREATED2  TREATED-2.star.nsrt.count   TREATED      E2
## 6   TREATED3  TREATED-3.star.nsrt.count   TREATED      E3
```



Raw Count
Normalisation

QC
reproducibility



Differential
expression analysis

Normalization

Normalising data is fundamental. If we skip this step we introduce biases in our analysis.

https://docs.gdc.cancer.gov/Data/Bioinformatics_Pipelines/Expression_mRNA_Pipeline/#mrna-expression-transformation

https://hbctraining.github.io/DGE_workshop/lessons/02_DGE_count_normalization.html

Functional annotation

Once identified differentially expressed genes, we can ask if they belong to some particular groups of genes, i.e. if they have common functionalities.

We can perform a gene ontology/over-representation analysis/gene set enrichment analysis



Molecular Function

Molecular-level activities performed by gene products. Molecular function terms describe activities that occur at the molecular level, such as “catalysis” or “transport”. GO molecular function terms represent activities rather than the entities (molecules or complexes) that perform the actions, and do not specify where, when, or in what context the action takes place. Molecular functions generally correspond to activities that can be performed by individual gene products (*i.e.* a protein or RNA), but some activities are performed by molecular complexes composed of multiple gene products. Examples of broad functional terms are *catalytic activity* and *transporter activity*; examples of narrower functional terms are *adenylate cyclase activity* or *Toll-like receptor binding*. To avoid confusion between gene product names and their molecular functions, GO molecular functions are often appended with the word “activity” (a *protein kinase* would have the GO molecular function *protein kinase activity*).

Cellular Component

The locations relative to cellular structures in which a gene product performs a function, either cellular compartments (*e.g.*, *mitochondrion*), or stable macromolecular complexes of which they are parts (*e.g.*, the *ribosome*). Unlike the other aspects of GO, cellular component classes refer not to processes but rather a cellular anatomy.

Biological Process

The larger processes, or ‘biological programs’ accomplished by multiple molecular activities. Examples of broad biological process terms are *DNA repair* or *signal transduction*. Examples of more specific terms are *pyrimidine nucleobase biosynthetic process* or *glucose transmembrane transport*. Note that a biological process is not equivalent to a pathway. At present, the GO does not try to represent the dynamics or dependencies that would be required to fully describe a pathway.

<http://geneontology.org/docs/ontology-documentation/>

Functional annotation

Once identified differentially expressed genes, we can ask if they belong to some particular groups of genes, i.e. if they have common functionalities.

We can perform a gene ontology/over-representation analysis/gene set enrichment analysis

Molecular Signatures Database

Human Collections

H **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.

C5 **ontology gene sets** consist of genes annotated by the same ontology term.

C1 **positional gene sets** corresponding to human chromosome cytogenetic bands.

C6 **oncogenic signature gene sets** defined directly from microarray gene expression data from cancer gene perturbations.

C2 **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.

C7 **immunologic signature gene sets** represent cell states and perturbations within the immune system.

C3 **regulatory target gene sets** based on gene target predictions for microRNA seed sequences and predicted transcription factor binding sites.

C8 **cell type signature gene sets** curated from cluster markers identified in single-cell sequencing studies of human tissue.

C4 **computational gene sets** defined by mining large collections of cancer-oriented microarray data.



<https://www.gsea-msigdb.org/gsea/msigdb/>