

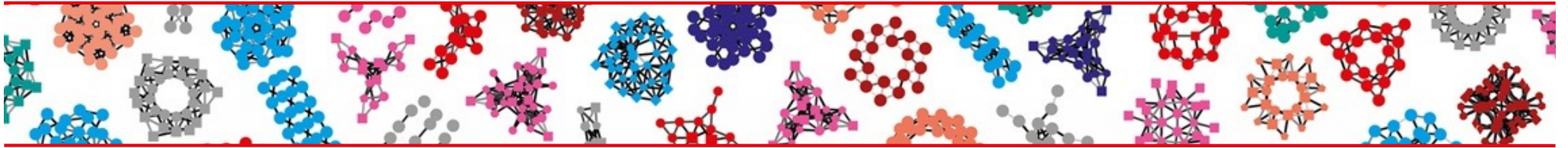
Swiss Institute of  
Bioinformatics

# Introduction to bioinformatics: Clinical Bioinformatics

Valérie Barbié, Director SIB Clinical Bioinformatics  
Zürich, 06 December 2022



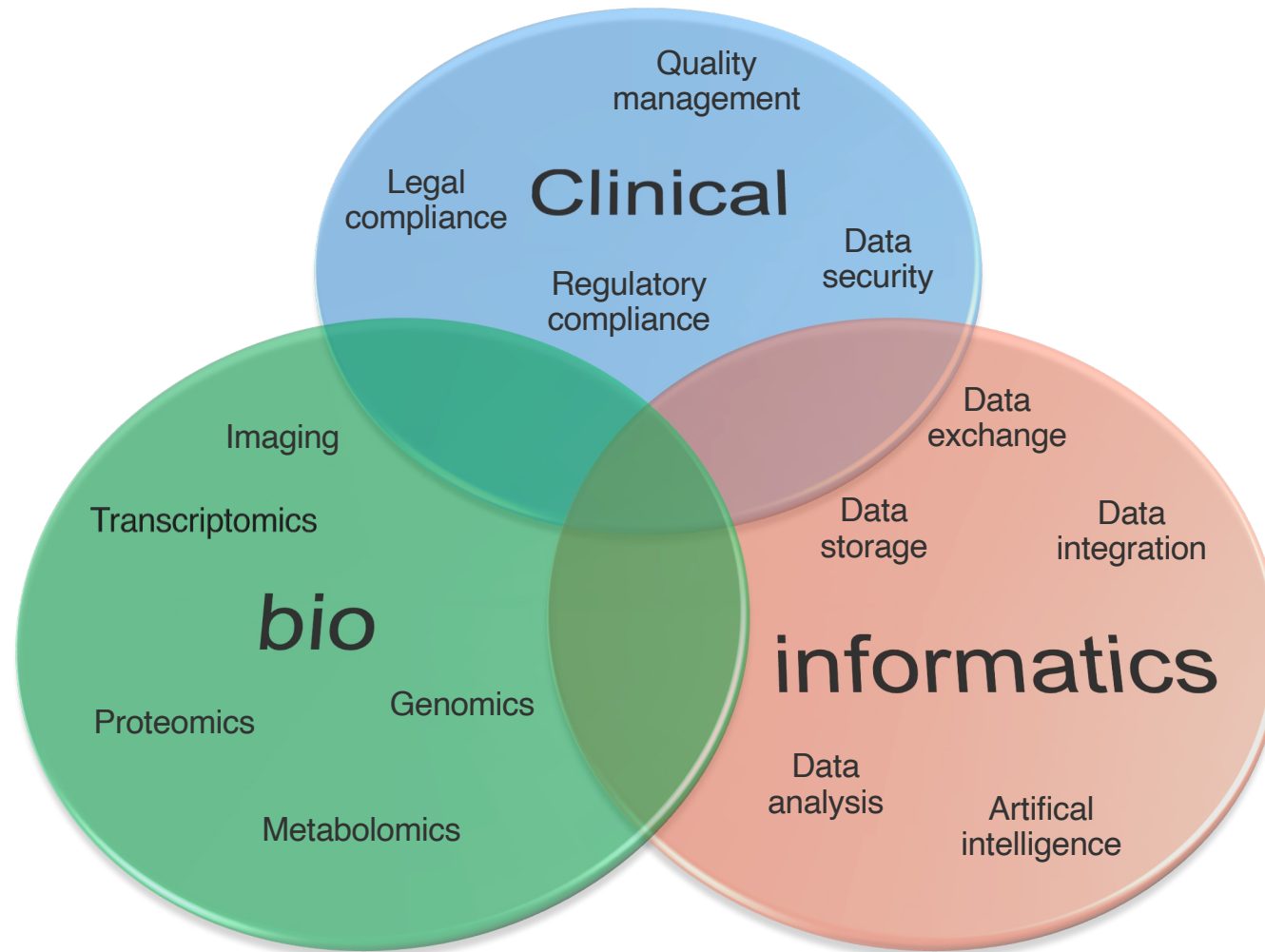
[www.sib.swiss](http://www.sib.swiss)

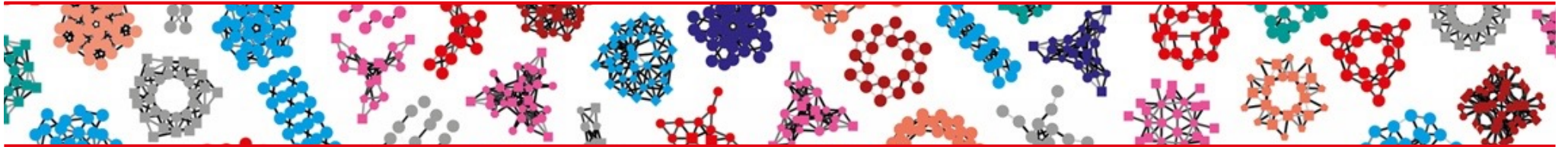


What is clinical bioinformatics?

# What is clinical bioinformatics?

---



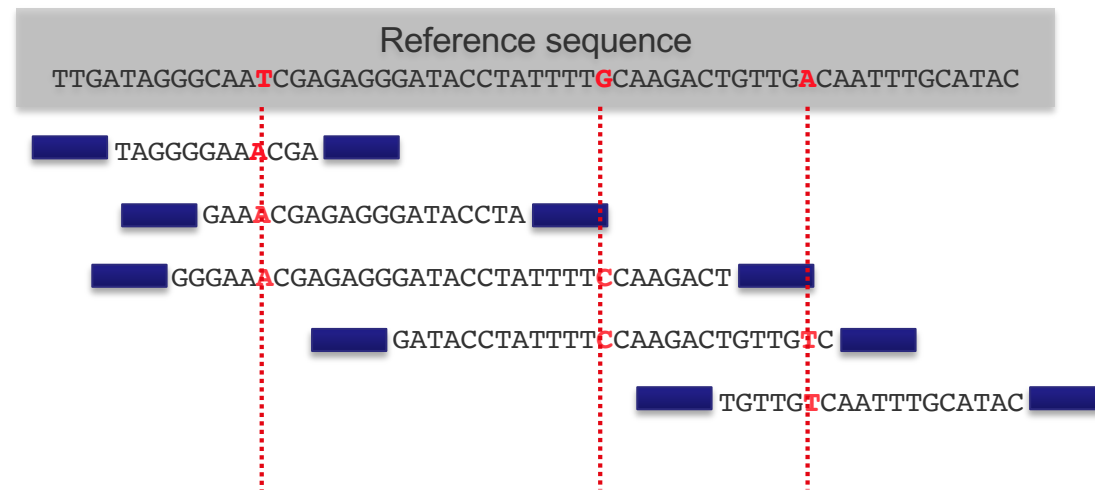


# Why clinical bioinformatics?

*The example of  
Next Generation Sequencing (NGS)  
in medical diagnosis*

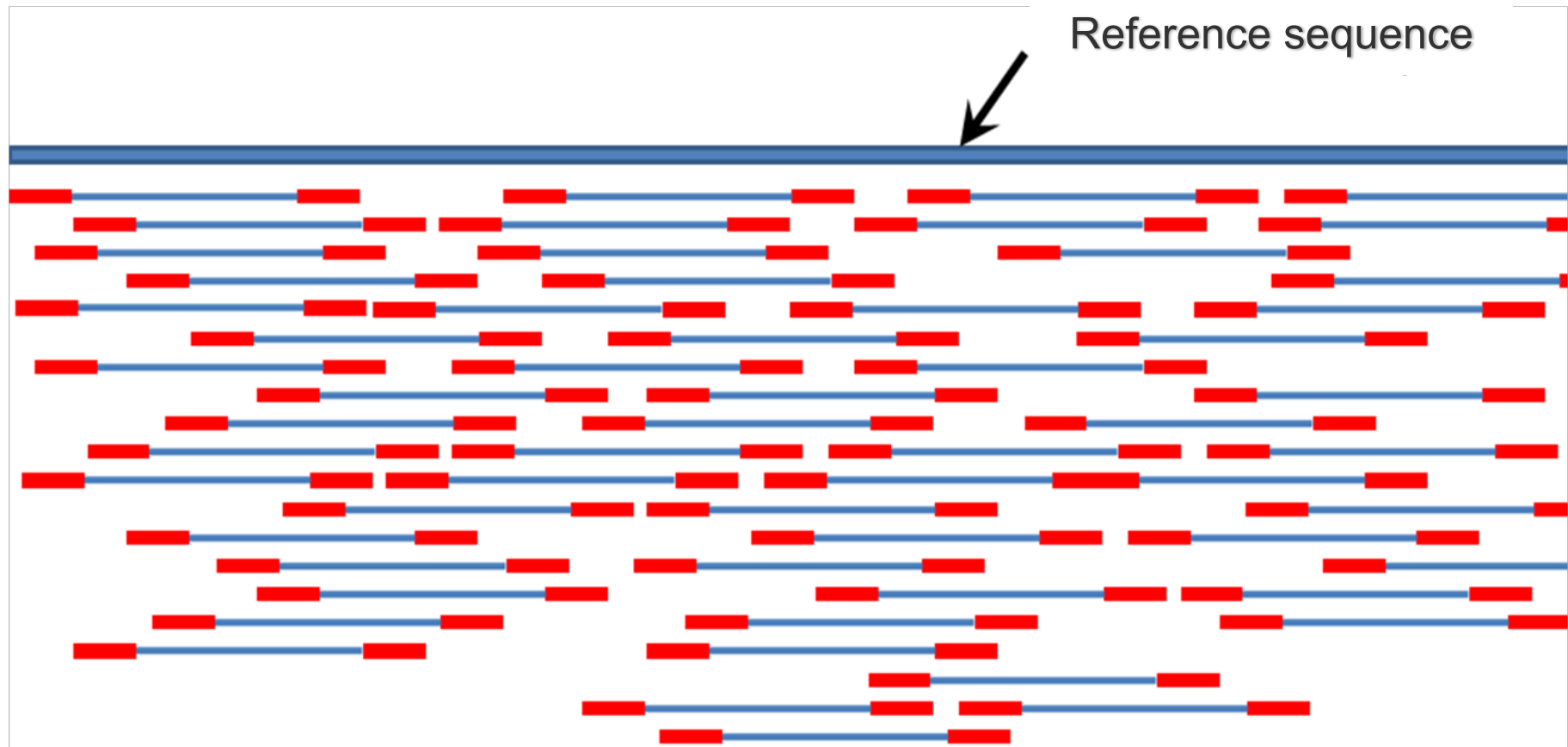
# Next Generation Sequencing principle

---



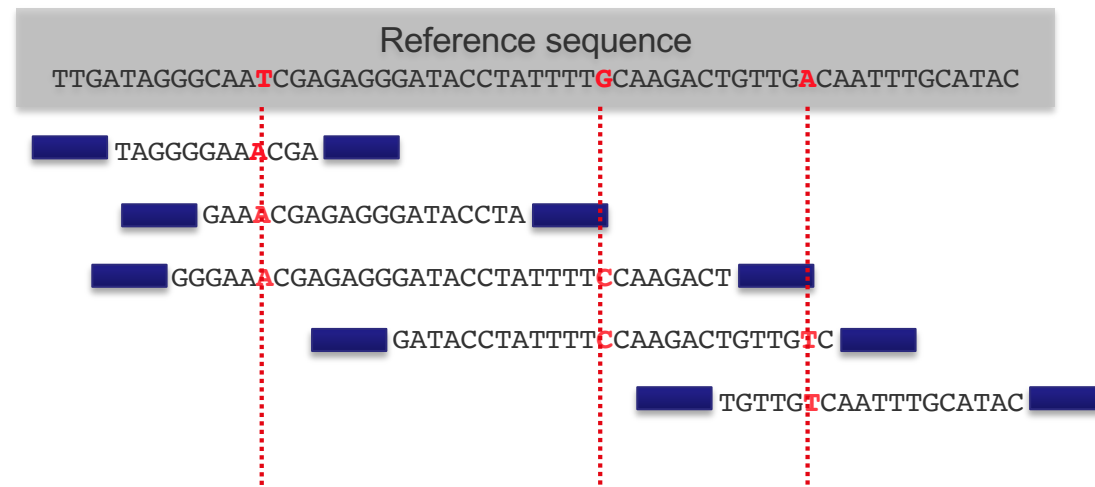
# Next Generation Sequencing principle

---



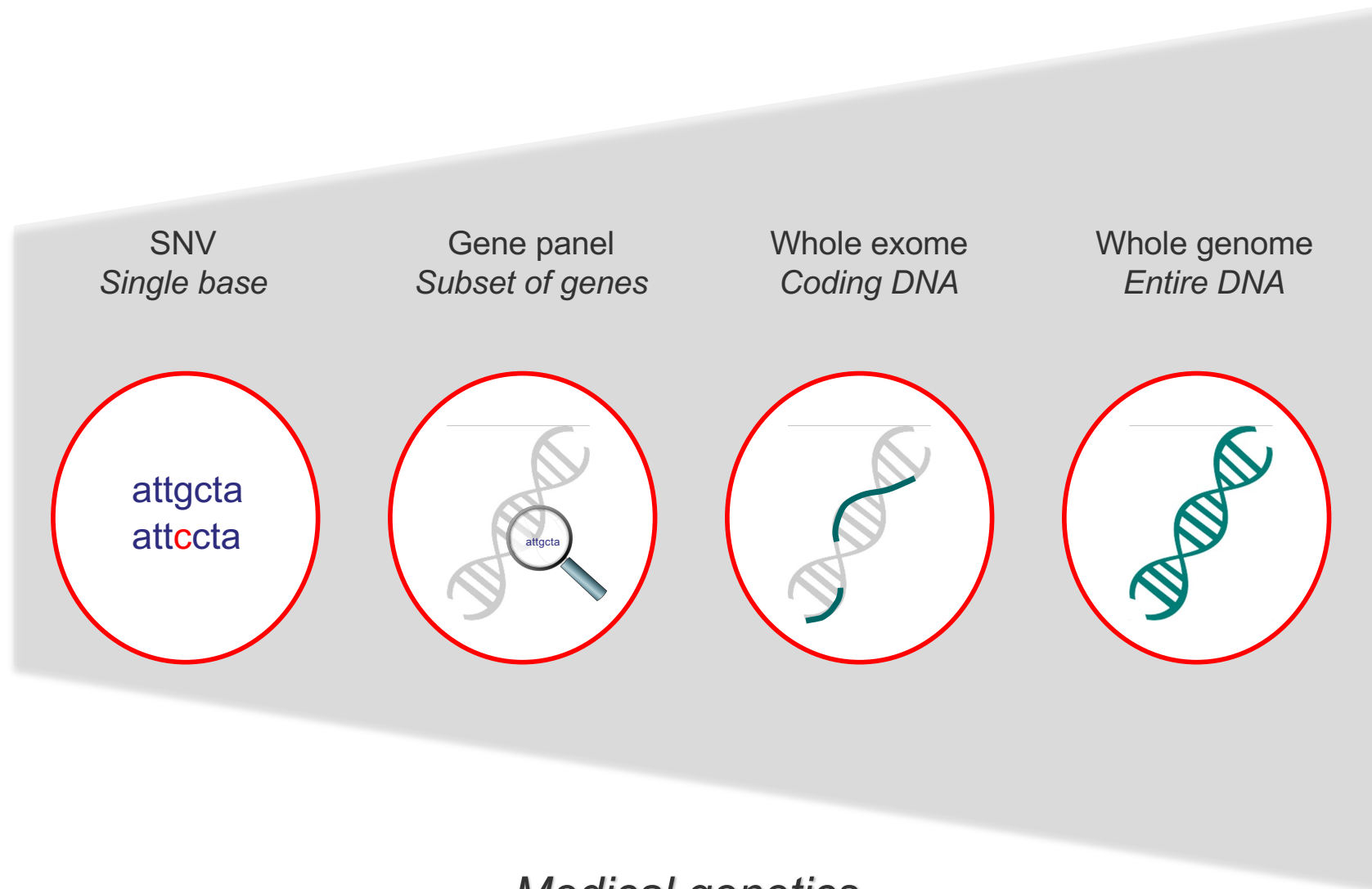
# Examples of NGS clinical applications

	Source DNA	Reference DNA
Oncology	Patient tumor or blood	Consensus human genome Germline
Microbiology	Patient	Pathogens genomes, resistance genes
Medical genetics	Patient	Family members, known defects
Pharmacogenetics	Patient	Drug-response or -sensitivity mutations



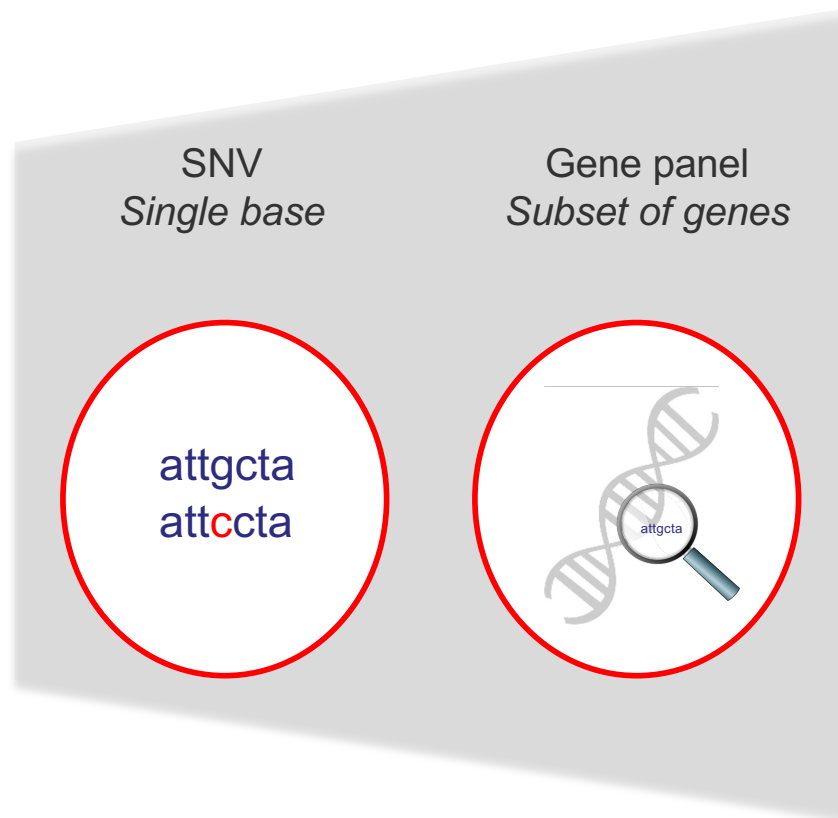
# Scale matters

---



# Scale matters

---



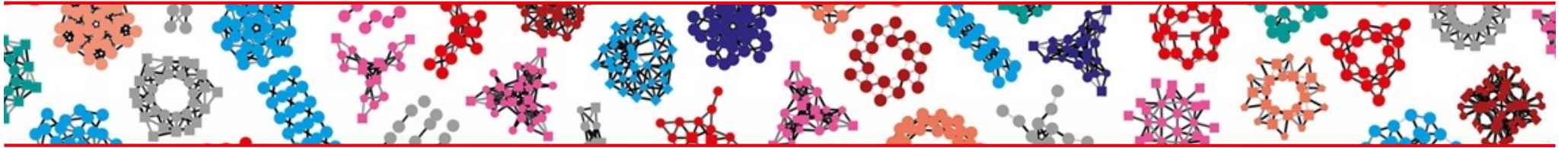
*Oncology*

Clinically-actionable variants

Reimbursement is limited

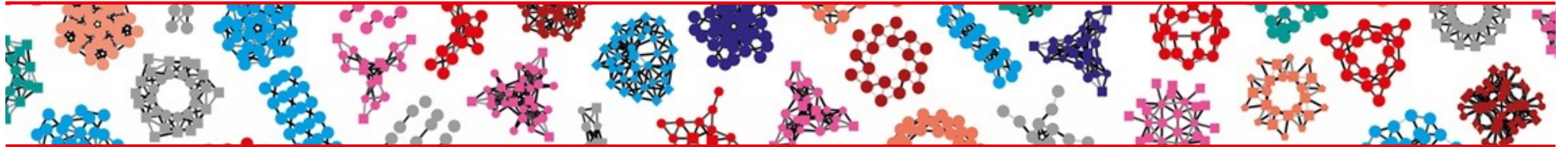
Incidental findings management

Results turn around time



# NGS in medical diagnosis

*Focus on oncology*



# PART I

## Overview of an NGS bioinformatics pipeline

# NGS in cancer diagnosis?

---

- Identify **single nucleotide variants (SNVs), insertions-deletions (indels)** to inform clinical management

at**t**cgggtcatgcccataagggg

Single Nucleotide  
Variant (SNV)

at**g**cgggtcatgcccataagggg

Insertion

at**g**cgggtcat**cgtgtcc**gcccataagggg

Deletion

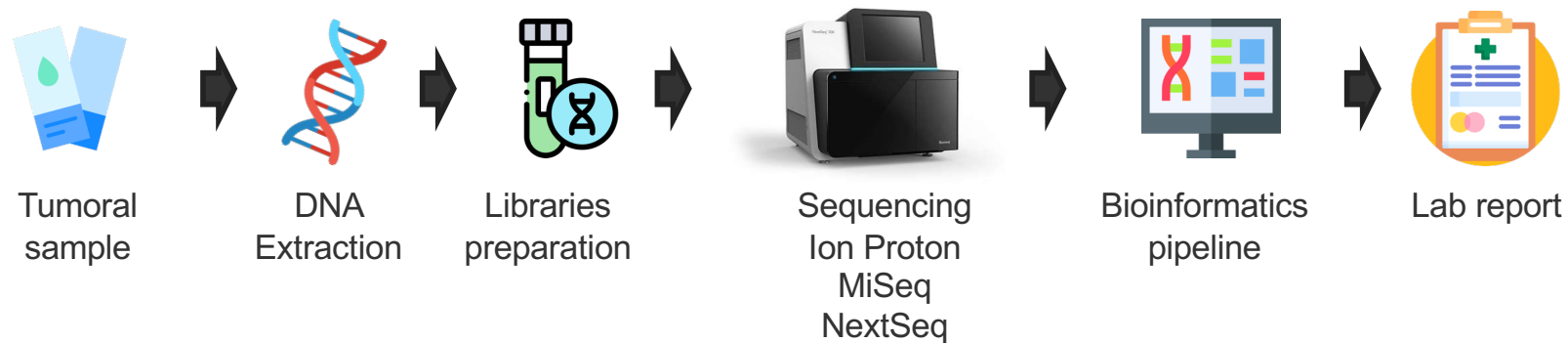
at**g**cgggtcatcgtgtccg....tagggg

**ccca**

---

# Overview of a NGS bioinformatics pipeline

---

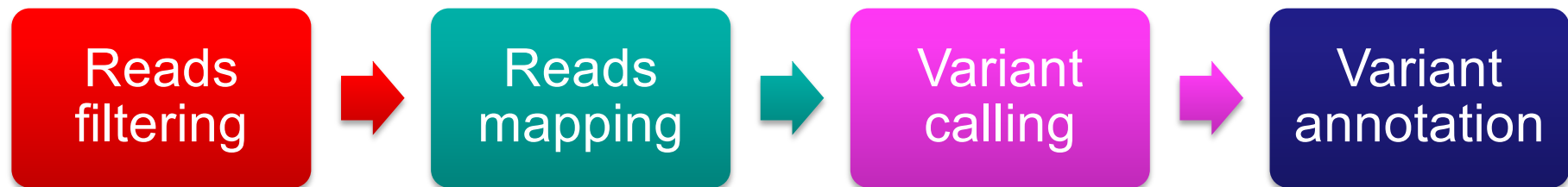
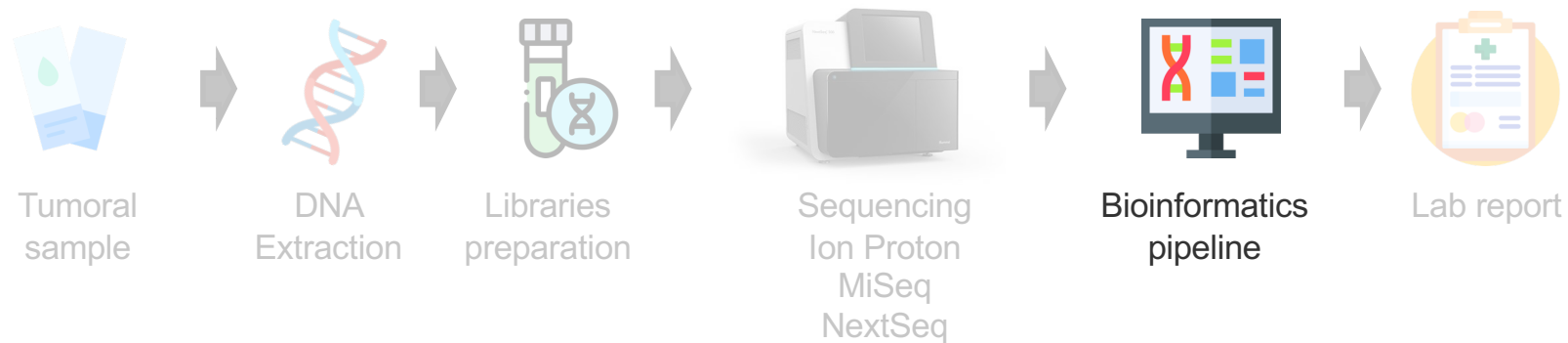


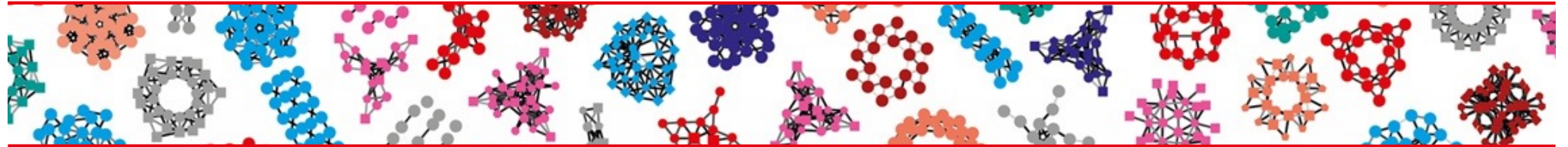
## ■ Gene panels analysis in clinical routine

- Identify **artifacts**: quality control
- Identify **somatic** vs. germline variants
- Variant **annotation**: does it provide clinically-useful information?

# Overview of a NGS bioinformatics pipeline

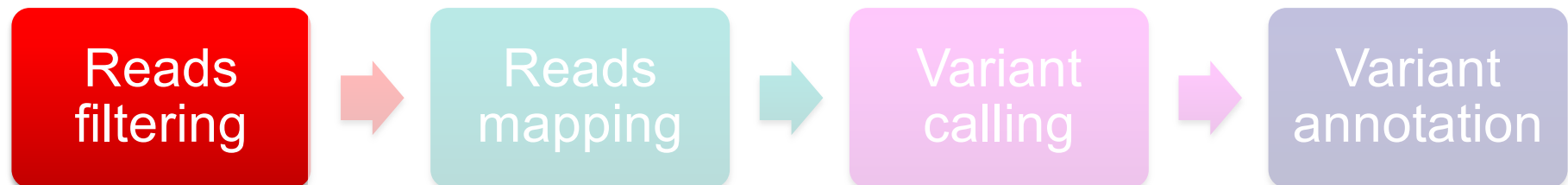
---





## PART II

### Quality control



# Out of the sequencer: FASTQ file

The diagram shows a FASTQ file format with four lines per record. Labels on the left point to specific lines: 'Identifier' points to the first line, 'Sequence' points to the second line, '+' sign points to the third line, and 'Quality scores' points to the fourth line. A red arrow points from the 'Quality scores' label to the 10th character 'h' in the first record's quality string. The first record's quality string is 'hhhhhhhhhhghhghhhhhfhhhhhfffffe'ee['X]b[d[ed'[Y[~Y'. The second record's quality string is 'hhhhghfhhcgghggfcffdhfehhhhcehdchhdhahehffffde'bVd'.

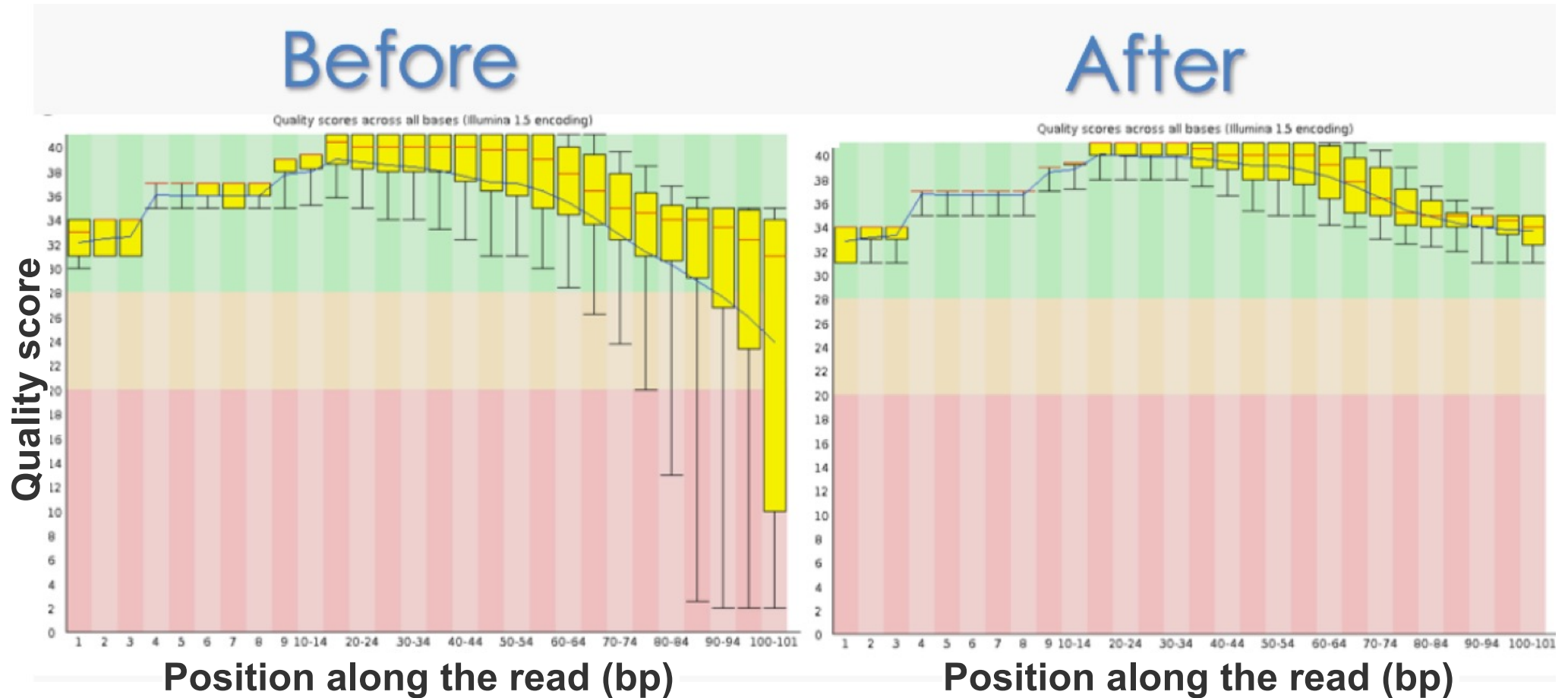
```
Identifier —● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence —● TTGCCTGCCTATCATTTTAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+' sign —● +
Quality scores —● hhhhhhhhhhhghhghhhhhfhhhhhfffffe'ee['X]b[d[ed'[Y[~Y
Identifier —● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence —● GATTTGATGAAAGTATACAACTAAACTGCAGGTGGATCAGAGTAAGTC
'+' sign —● +
Quality scores —● hhhhghfhhcgghggfcffdhfehhhhcehdchhdhahehffffde'bVd
```

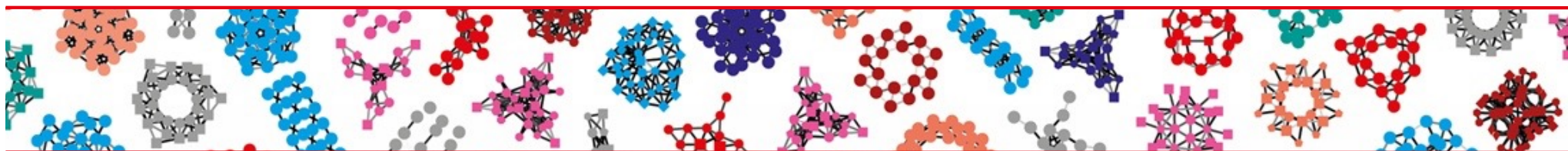
Each nucleotide has a **quality score (Phred score)** representing the probability that a base was miscalled by the sequencer

$$Q = -10 \log_{10} P$$

Phred Score	Prob. of incorrect base call	Base call accuracy	Code
10	1 in 10	90%	J
20	1 in 100	99%	T
30	1 in 1'000	99.9%	^
40	1 in 10'000	99.99%	h

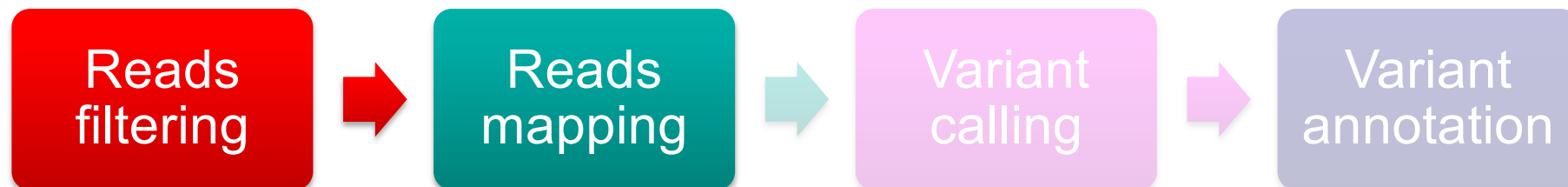
# Quality-based reads trimming





## PART III

# Variant identification



# Let's align the reads

---

Reference genome  
TCGCGCACAAG

! **Short reads** are likely to map at several positions along the reference genome

Reference genome  
CGTGGGACGAG

! **Mismatches** and **gaps** allowed  
→ algorithms have scoring functions

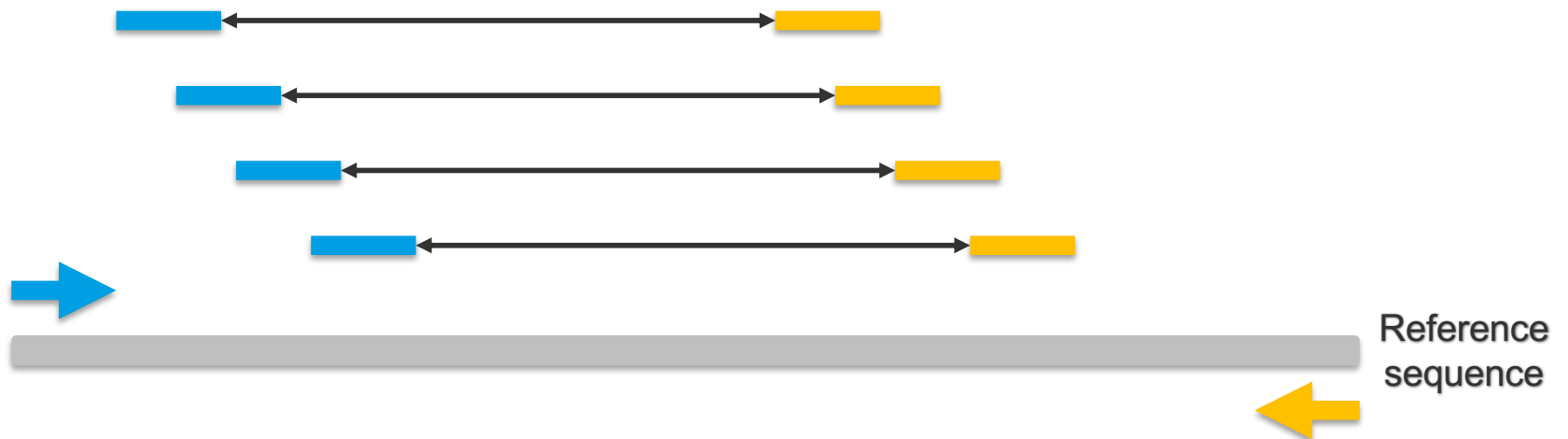
Reference genome  
TCGCGCACAAGACGTGGGACGAG

! **Longer reads** are less ambiguous  
→ but computationally more expensive

---

# Paired-end sequencing

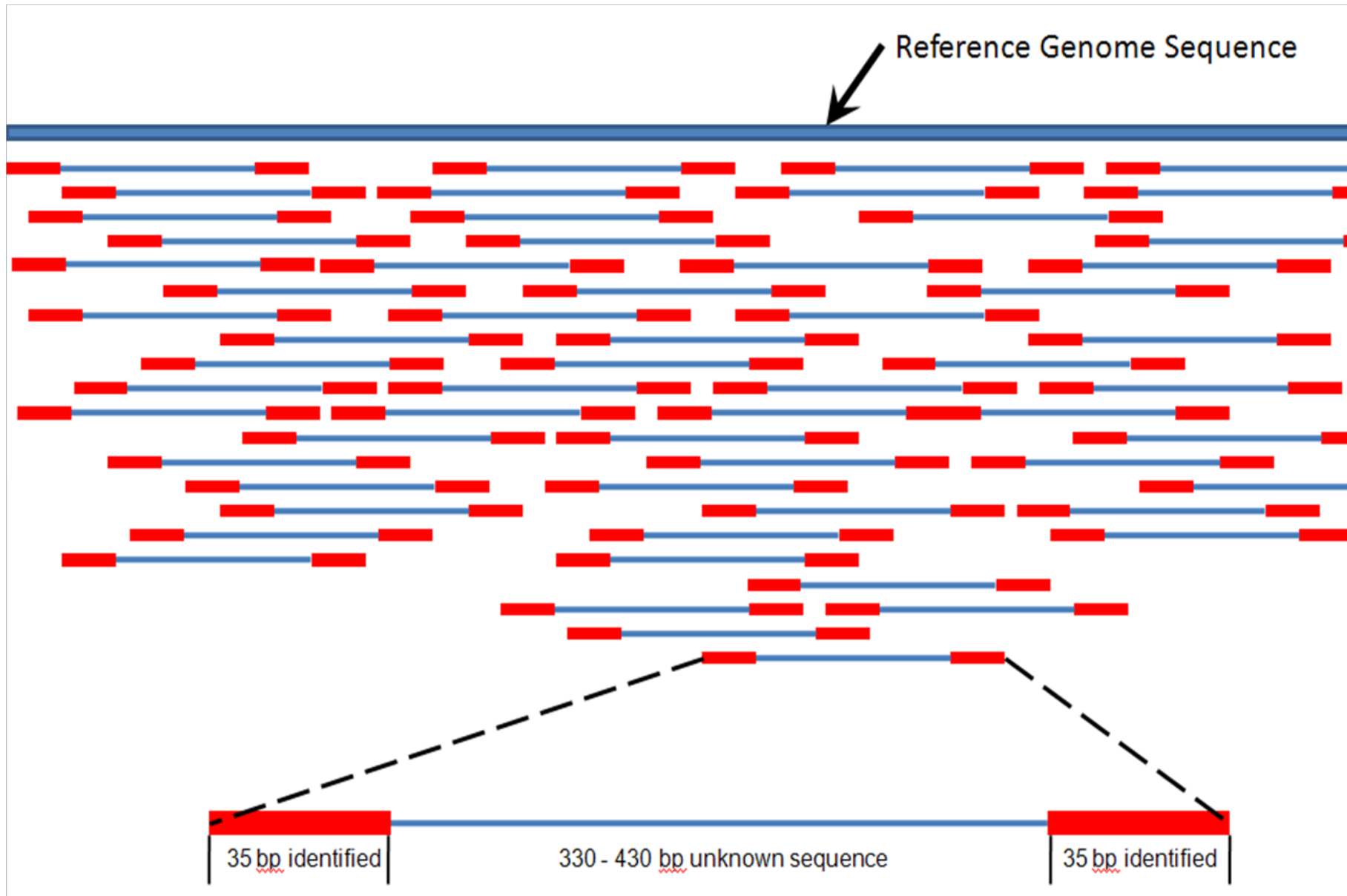
---

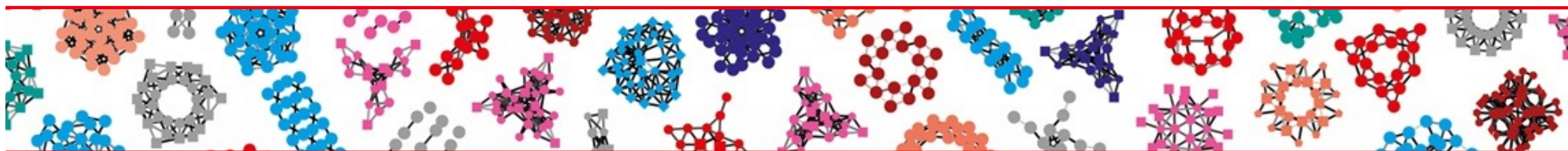


- Much better alignment on across regions difficult to sequence (e.g. repetitive regions)

# Mapping: finding the best position for each read

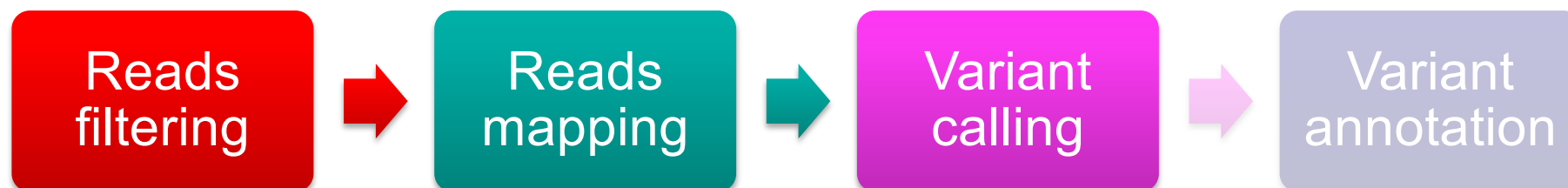
---





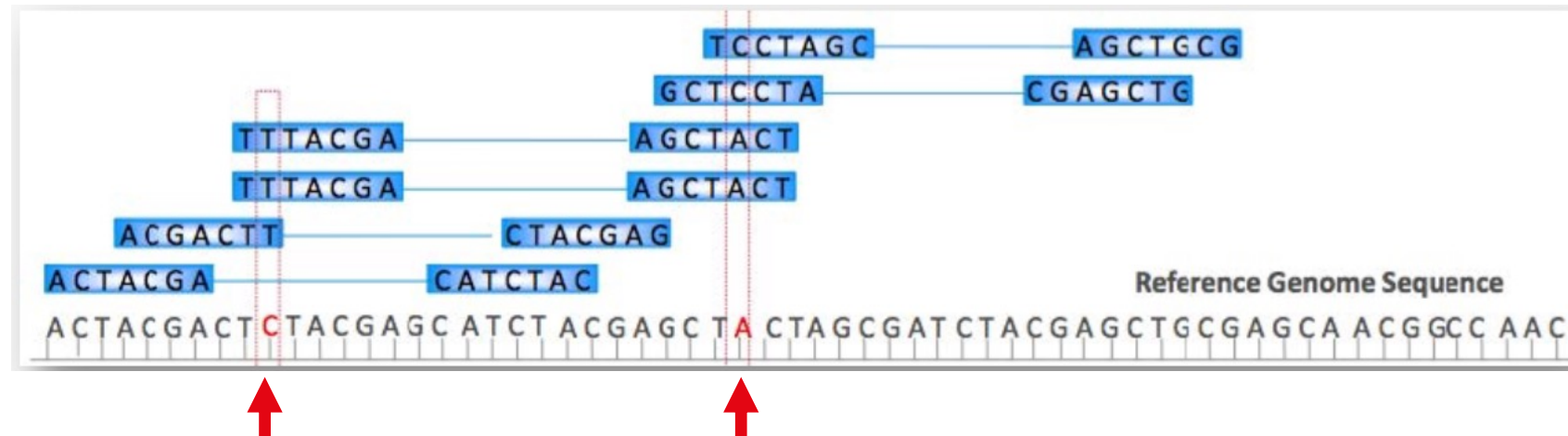
## PART III

# Variant identification



# Variant calling: putting it all together

---



*True variant or technical error?*

- Performed by the sequencer software or the bioinformatician
  - Germline vs somatic calling
    - Germline: constitutional genome analysis, where variants occur in **50%** (heterozygous) or **100%** (homozygous) of the reads.
    - Somatic: no ploidy assumption, low frequency alleles.
-

## VCF: Variant Call Format

## FORMAT meta-information

## Fixed fields

FORMAT	NORMAL	TUMOR
--------	--------	-------

GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51
GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3
GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2
GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51
GT:GQ:DP	0/1:35:4	0/2:17:2

# Output of the variant caller: VCF

---

## VCF: Variant Call Format

Fixed fields								
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	
BODY	20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2
	20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017
	20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;DB
	20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T
	20	1234567	microsat1	GTC	G,GTCTC	50	PASS	NS=3;DP=9;AA=G

---



---

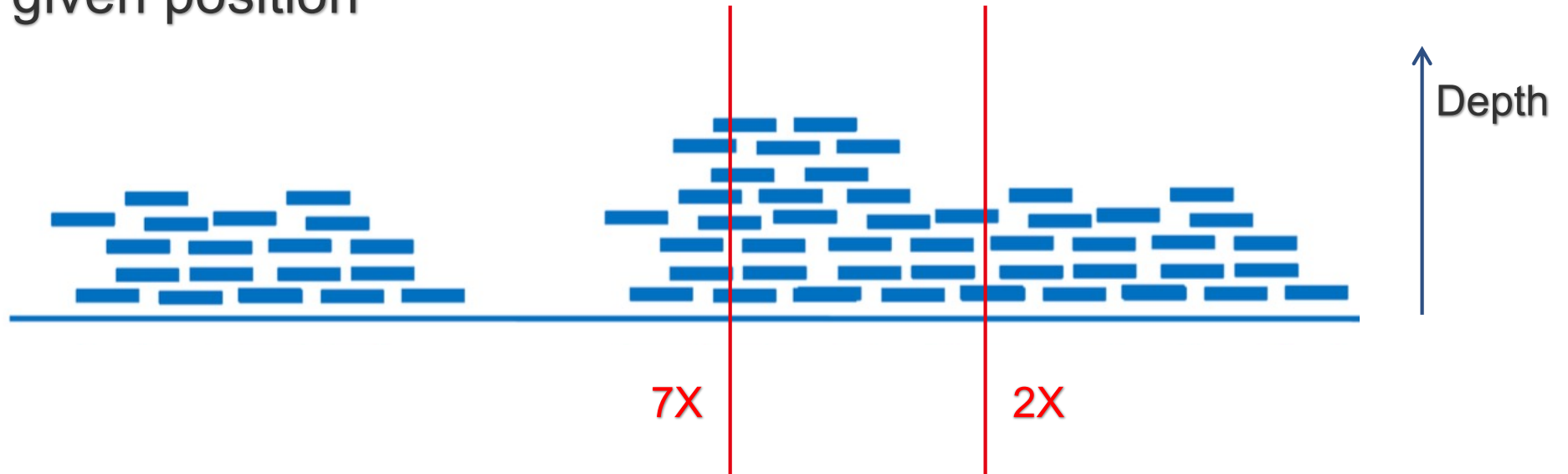
Things to watch out  
when assessing variant quality

---

# Depth

---

- **Depth:** nb of reads that include a given nucleotide, at a given position

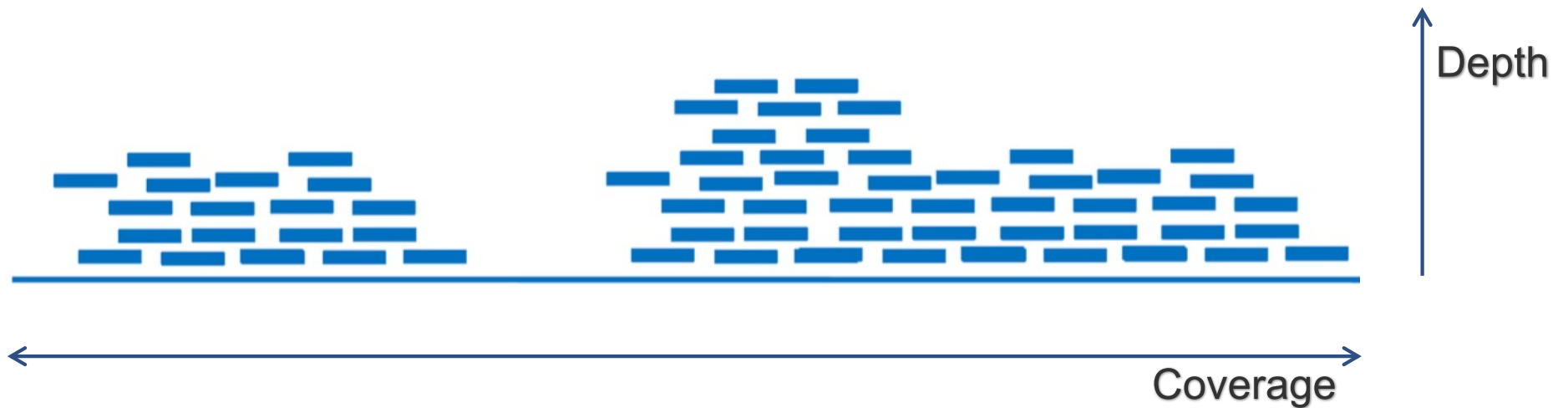


- Diagnosis: gene panel at 1500X, whole exome at 100X
  - In oncology, impossible to detect low frequency clones with exome analyses
-

# Coverage

---

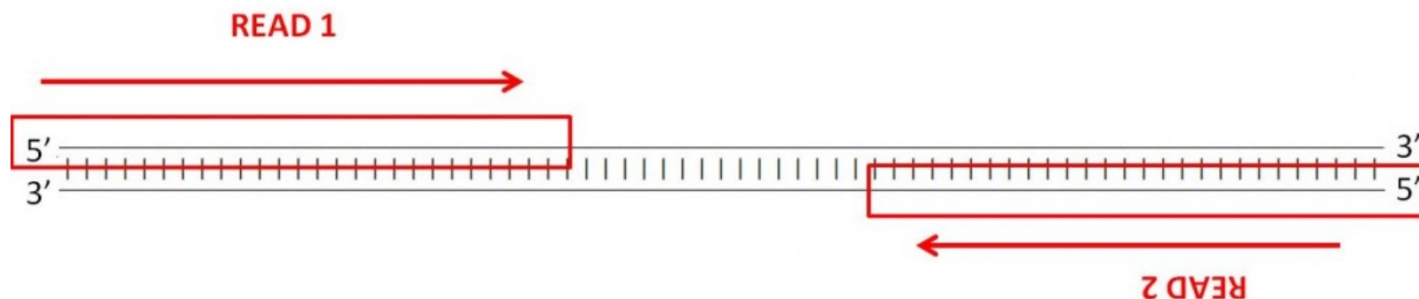
- **Coverage:** % or nb of bases of a reference genome that are covered with a certain depth, e.g. 90% at 5X

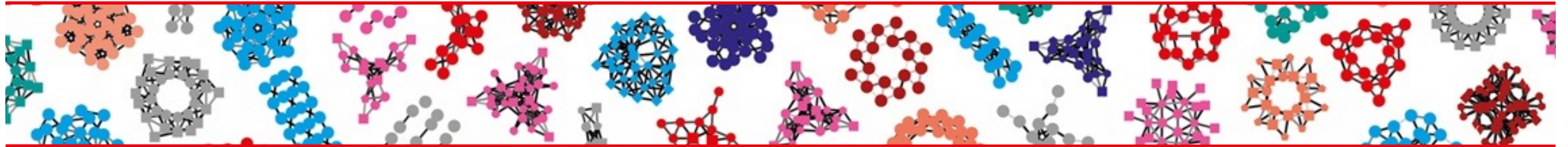


# Strand bias in paired-end sequencing

---

- **Both DNA strands are sequenced**
- Normal mutations should occur on both with equal frequencies





## PART IV

# Variant annotation and interpretation

# Medical genetics: focus on pathogenicity

© American College of Medical Genetics and Genomics | **ACMG STANDARDS AND GUIDELINES** | Genetics  
inMedicine

Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology

Sue Richards, PhD<sup>1</sup>, Nazneen Aziz, PhD<sup>2,16</sup>, Sherri Bale, PhD<sup>3</sup>, David Bick, MD<sup>4</sup>, Soma Das, PhD<sup>5</sup>, Julie Gastier-Foster, PhD<sup>6,7,8</sup>, Wayne W. Grody, MD, PhD<sup>9,10,11</sup>, Madhuri Hegde, PhD<sup>12</sup>, Elaine Lyon, PhD<sup>13</sup>, Elaine Spector, PhD<sup>14</sup>, Karl Voelkerding, MD<sup>15</sup> and Heidi L. Rehm, PhD<sup>15</sup>; on behalf of the ACMG Laboratory Quality Assurance Committee

**GENETICS in MEDICINE** | Volume 17 | Number 5 | May 2015

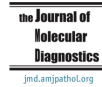
## Find **pathogenic** variants

*i.e. genetic alterations increasing an individual's susceptibility or predisposition to a certain disorder*



# Oncology: focus on clinical significance

The Journal of Molecular Diagnostics, Vol. 19, No. 1, January 2017



## SPECIAL ARTICLE

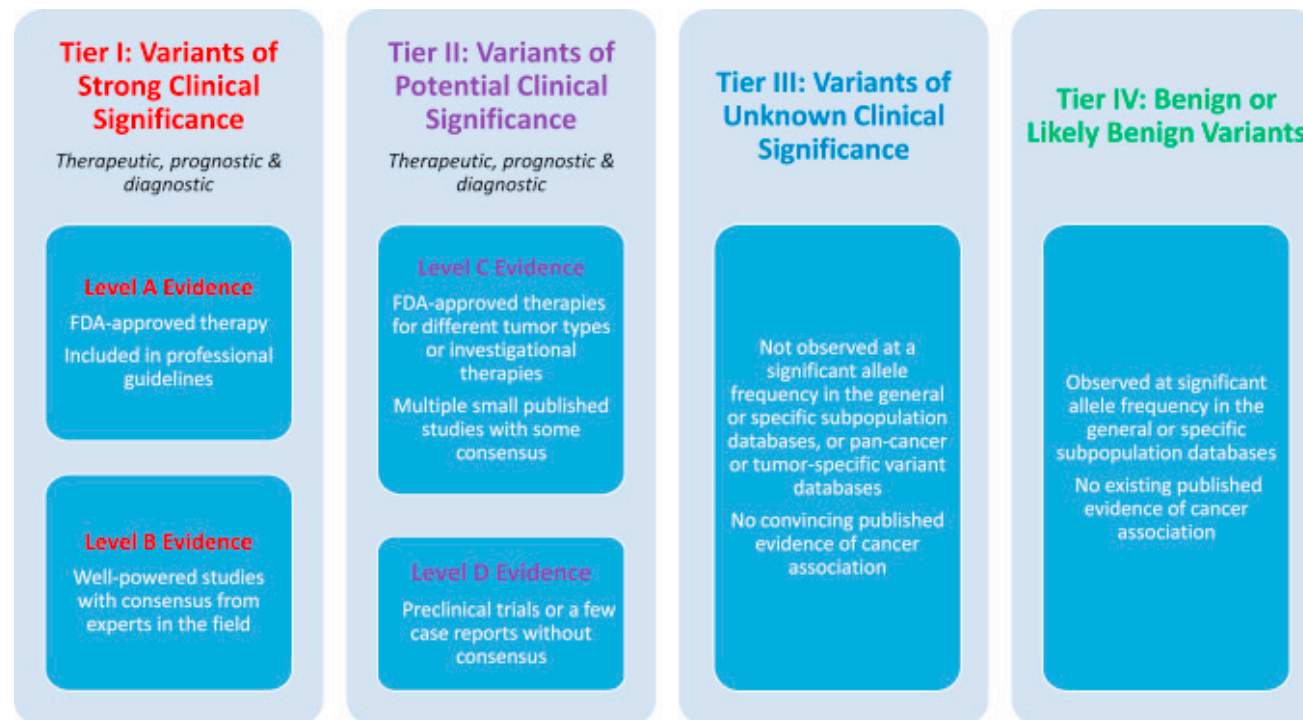
Standards and Guidelines for the Interpretation and Reporting of Sequence Variants in Cancer



*A Joint Consensus Recommendation of the Association for Molecular Pathology, American Society of Clinical Oncology, and College of American Pathologists*

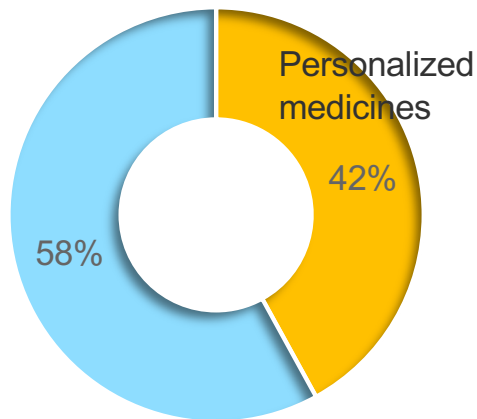
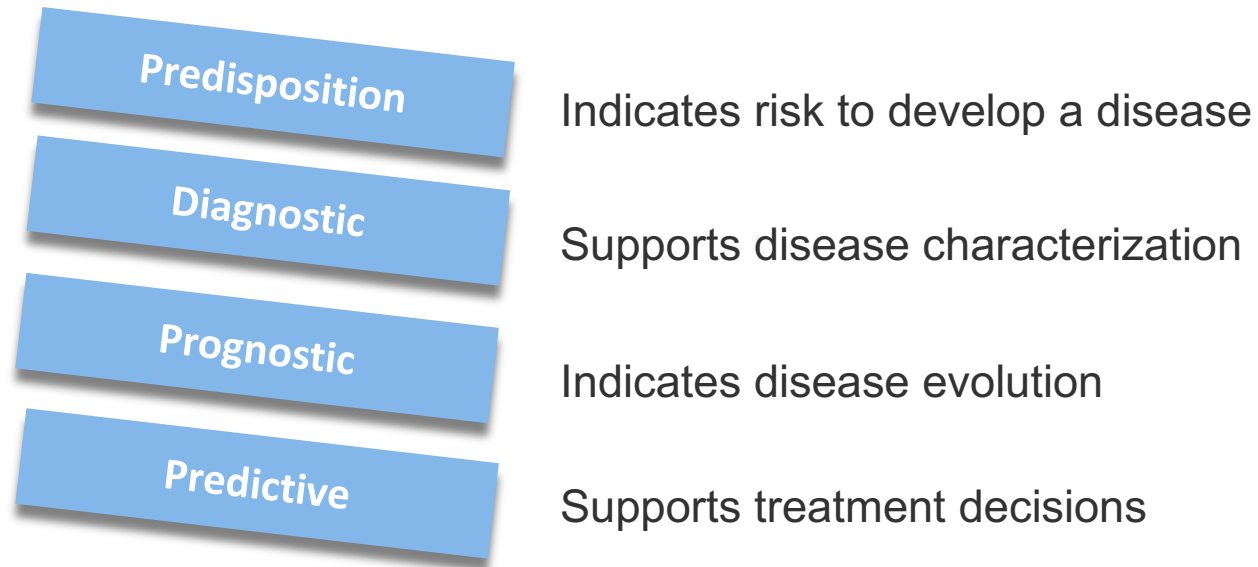
## Find **actionable** variants

*i.e. genetic alterations possibly having an impact on clinical care*

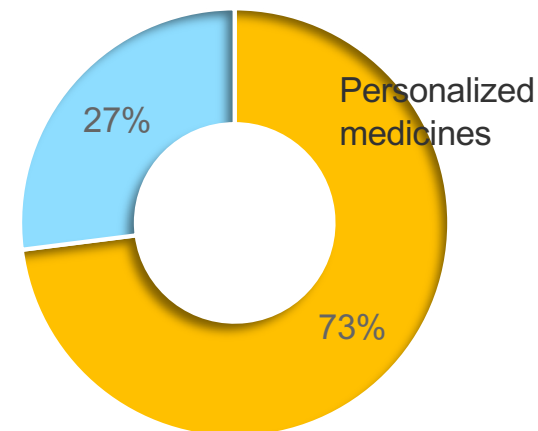


# Categories of markers

---

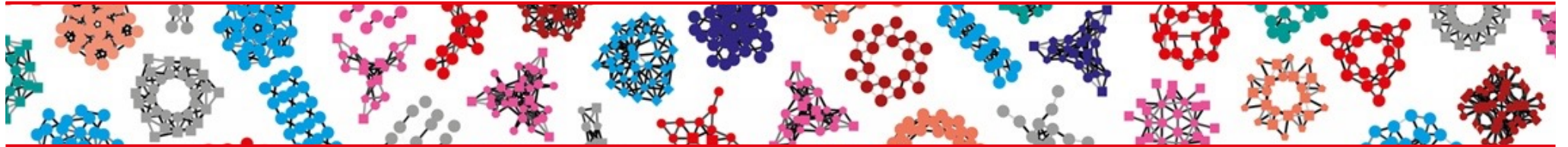


All drugs in development



Oncology drugs in development

---



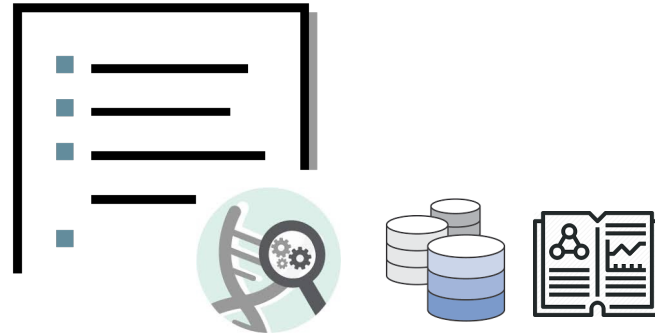
## PART IV

# Variant annotation and interpretation

... bioinformatics at the rescue

# Bioinformatics to the rescue...

---



- **Location** of the variant (e.g. intron, exon, regulatory region...)
- **Genes** and **transcripts** affected by the variant
- Predict **variant effect** (e.g. stop gained, missense...)

# Locating variants

---

- Convert **genomic coordinates** (chromosome, position) to the corresponding **cDNA/amino-acid coordinates**

- **HGVS nomenclature** (<http://varnomen.hgvs.org>)

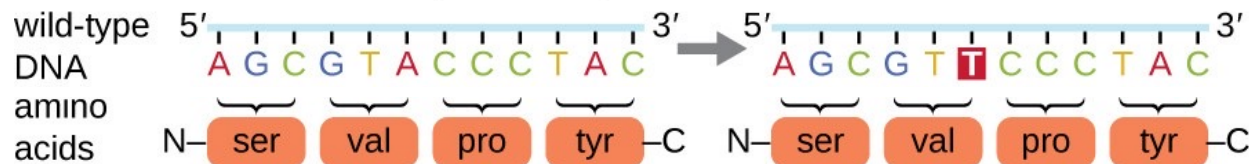
- Substitution c.76A>T
- Deletion c.76delA
- Insertion c.76\_77insG
- Genomic sequence g.476A>T
- Protein sequence p.Lys76Asn

- **Important to store for tracking**

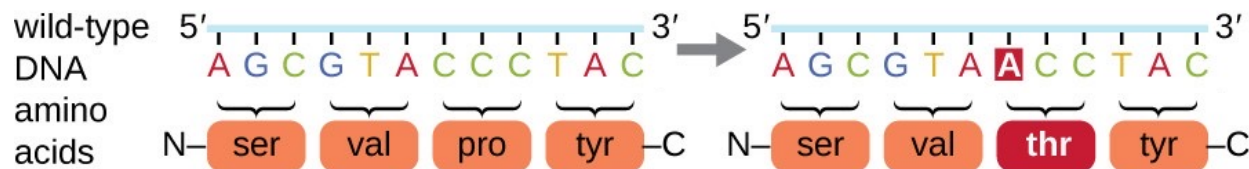
- Version of the human genome assembly
  - Accession and version of the mRNA transcripts
-

# Predicting variants effect on the protein

**silent:** has no effect on the protein sequence

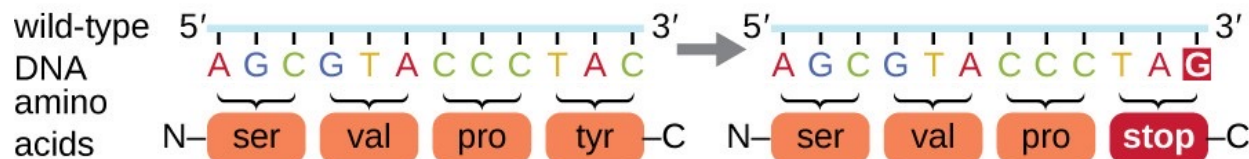


**missense:** results in an amino acid substitution

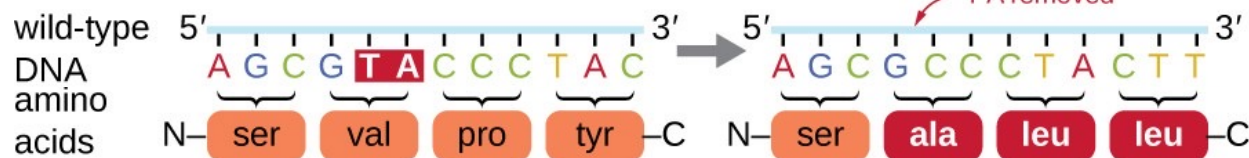


Point mutations  
(single base  
substitution)

**nonsense:** substitutes a stop codon for an amino acid



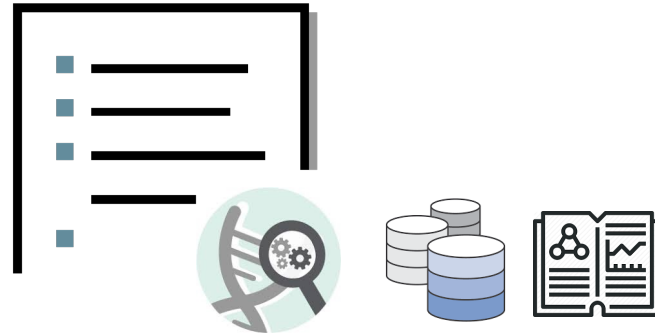
**Insertion or deletion** results in a shift in the reading frame.



Frameshift mutations  
(insertion or deletion of  
one or several bases)

# Bioinformatics to the rescue...

---



- ❑ **Location** of the variants (e.g. intron, exon, regulatory region...)
- ❑ **Genes** and **transcripts** affected by the variant
- ❑ Predict **variant effect** (e.g. stop gained, missense...)
- ❑ Predict **variant impact** on protein function, splicing

# Predicting variants impact: examples of tools

TOOLS	SnEff (ClinEff)	VEP	SIFT	PolyPhen-2	FATHMM
Variant effect and location (sequence ontology)	✓	✓			
Prediction of impact (score or category)	✓	←	✓	✓	✓
Features used for impact prediction	Rules based on variant effect (stop gained, lost...)		AA conservation in related seq.	AA conservation and structural features	AA conservation and protein tolerance to mutations

*Use a combination of tools and keep variants with consensus prediction.*

© American College of Medical Genetics and Genomics | **ACMG STANDARDS AND GUIDELINES** | Genetics in Medicine

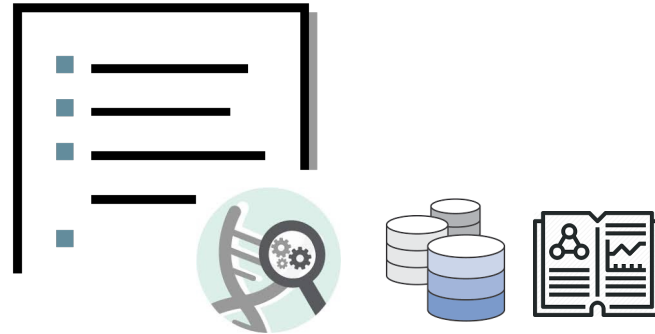
**Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology**

Sue Richards, PhD<sup>1</sup>, Nazneen Aziz, PhD<sup>2,16</sup>, Sherri Bale, PhD<sup>3</sup>, David Bick, MD<sup>4</sup>, Soma Das, PhD<sup>5</sup>, Julie Gastier-Foster, PhD<sup>6,7,8</sup>, Wayne W. Grody, MD, PhD<sup>9,10,11</sup>, Madhuri Hegde, PhD<sup>12</sup>, Elaine Lyon, PhD<sup>13</sup>, Elaine Spector, PhD<sup>14</sup>, Karl Voelkerding, MD<sup>15</sup> and Heidi L. Rehm, PhD<sup>15</sup>; on behalf of the ACMG Laboratory Quality Assurance Committee

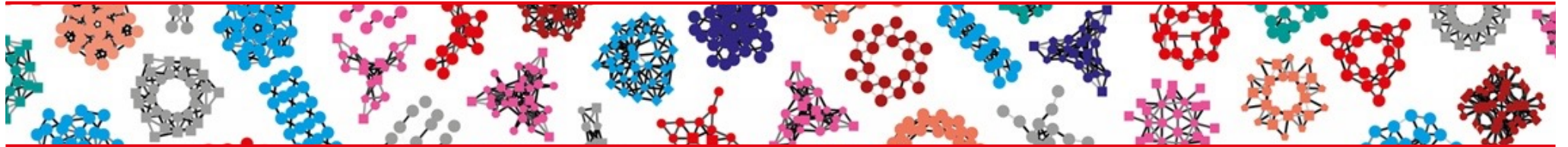
**GENETICS in MEDICINE** | Volume 17 | Number 5 | May 2015

# Bioinformatics to the rescue...

---



- ❑ **Location** of the variants (e.g. intron, exon, regulatory region...)
- ❑ **Genes** and **transcripts** affected by the variant
- ❑ Predict **variant effect** (e.g. stop gained, missense...)
- ❑ Predict **variant impact** on protein function, splicing
- ❑ Retrieve annotations from **public databases**



## PART IV

# Variant annotation and interpretation

... with knowledge-bases

# Important questions

---

- Is it prevalent in the cancer subtype of interest?
  - Is it known in other cancer subtypes or diseases?
  - Is it present in the general population?
  - Is it related to an ongoing clinical trial?
  - What is the evidence level? Observed vs. predicted
  - Are there other known variants in the same gene?
-

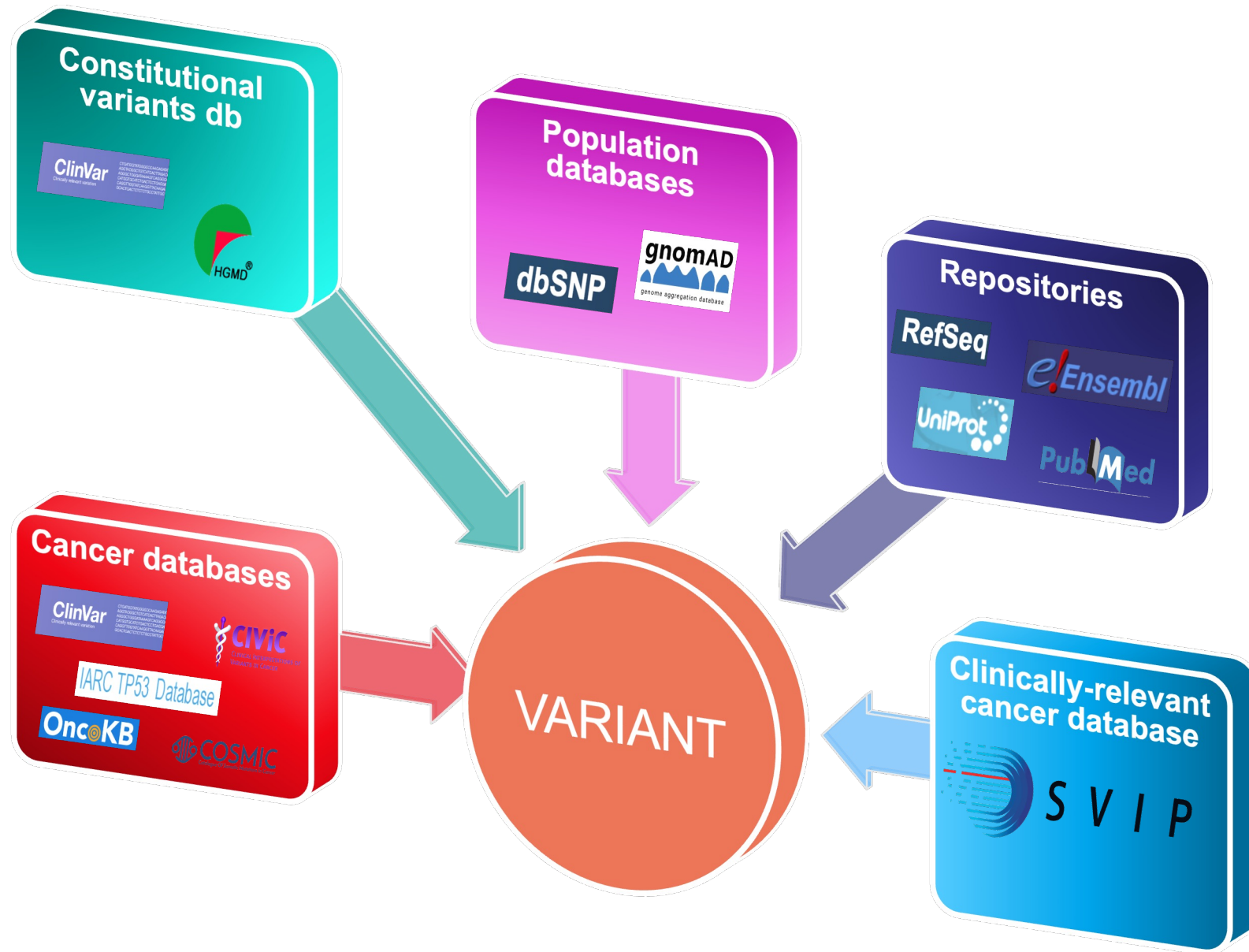
# Important questions

- ❑ Is the mutation in an **evolutionarily conserved** region accross species?



# Knowledge bases

---



*Non exhaustive*

---

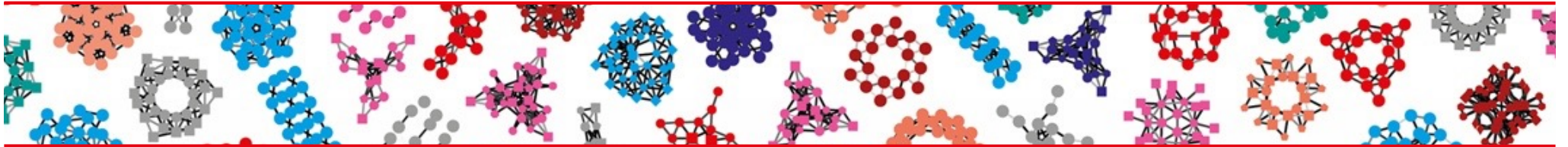
# I found a damaging mutation: is it always bad?

---

- Keep the mutation in context: what is the gene function?
  - **Tumor suppressor gene**  
Damaging mutations are pathogenic.
  - **Oncogene**  
Activating mutations are pathogenic.  
(beware: damaging mutation can be activating!)

**Keep the gene function in mind  
when interpreting its deleteriousness**

---



Other considerations...

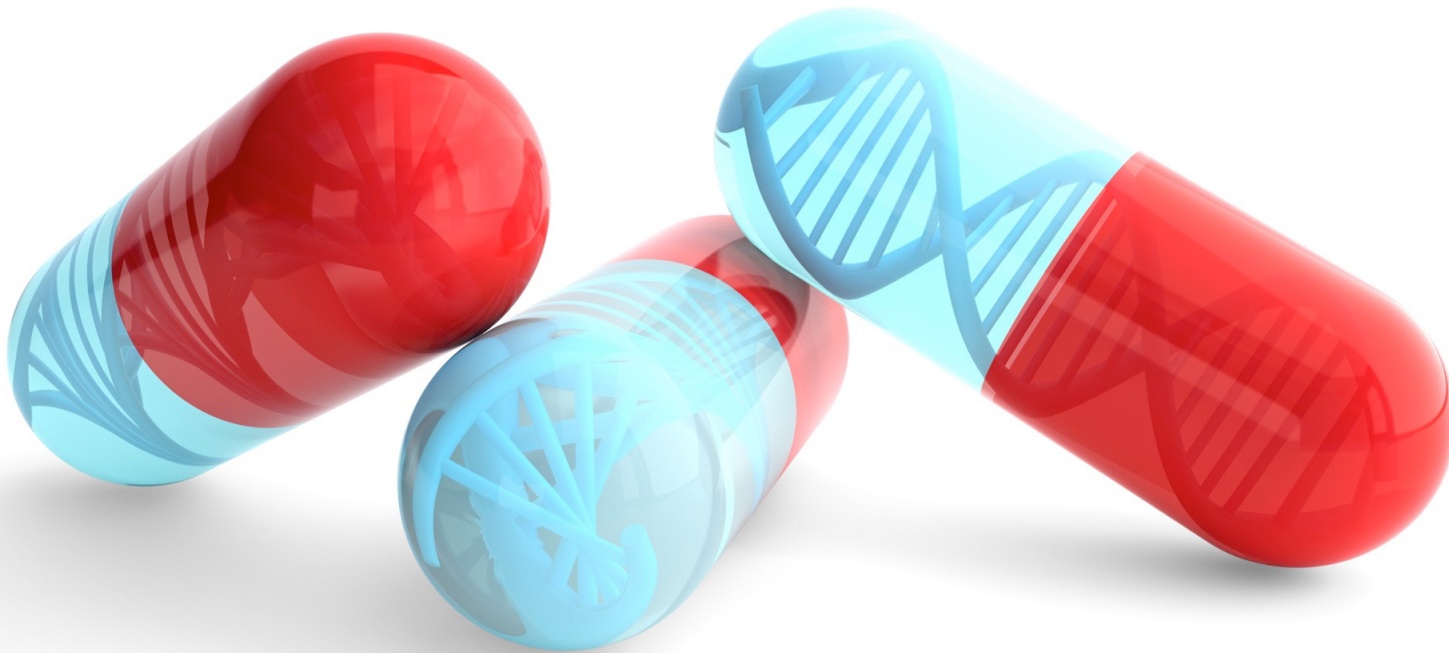
# Real-life constraints in the clinics

---



# Certificate of Advanced Studies (CAS) in Personalized molecular oncology

*[pmo.unibas.ch](http://pmo.unibas.ch)*

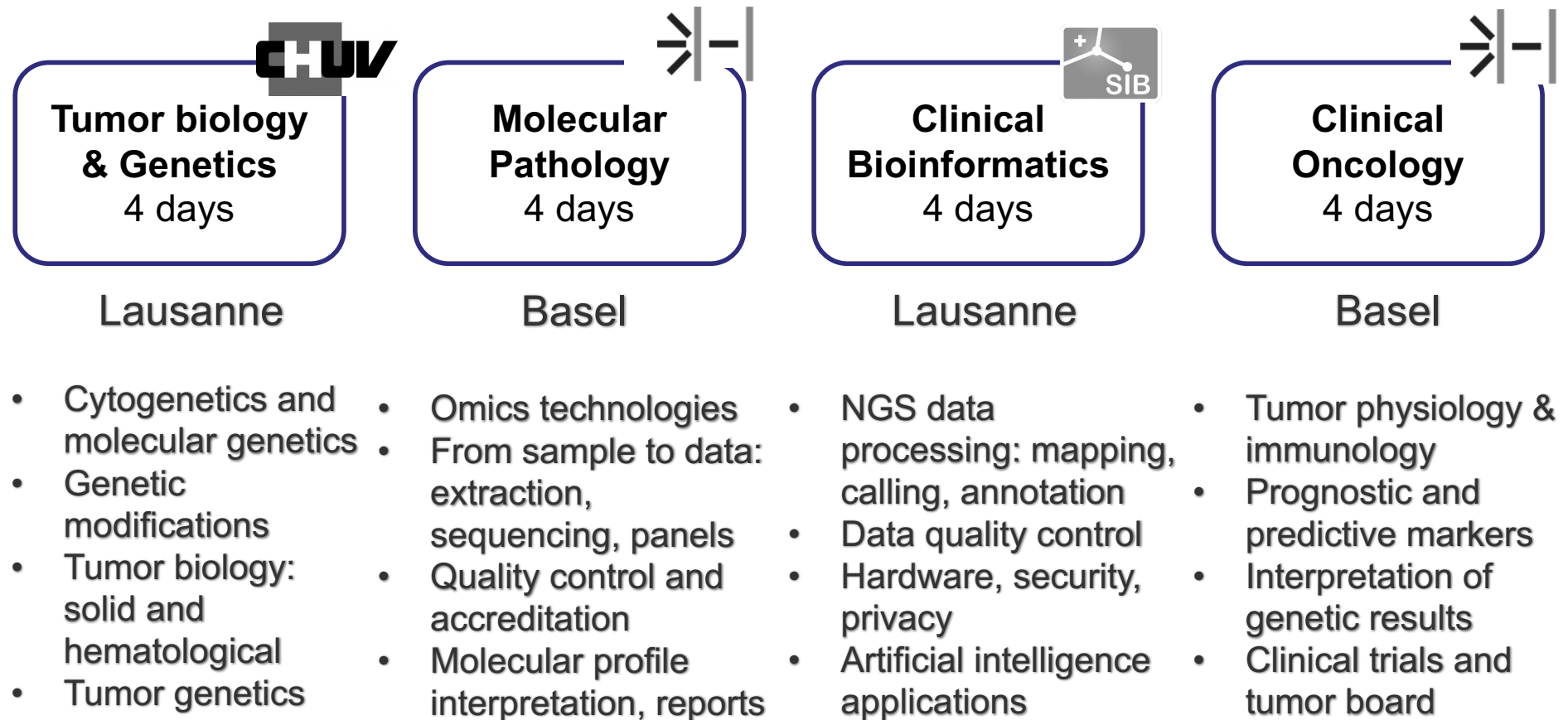


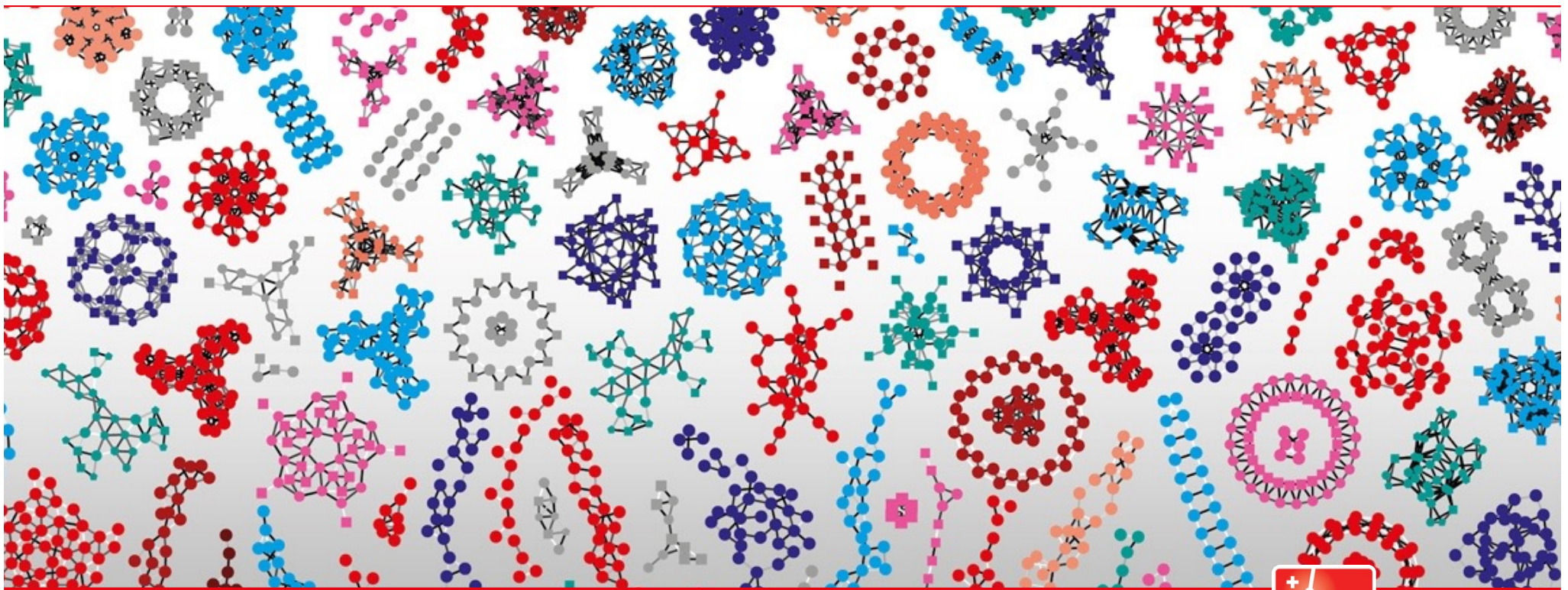
Swiss Institute of  
Bioinformatics



# CAS PMO: 4 modules and a mini-thesis

---





Swiss Institute of  
Bioinformatics

Thank You