# Somatic Variant Calling

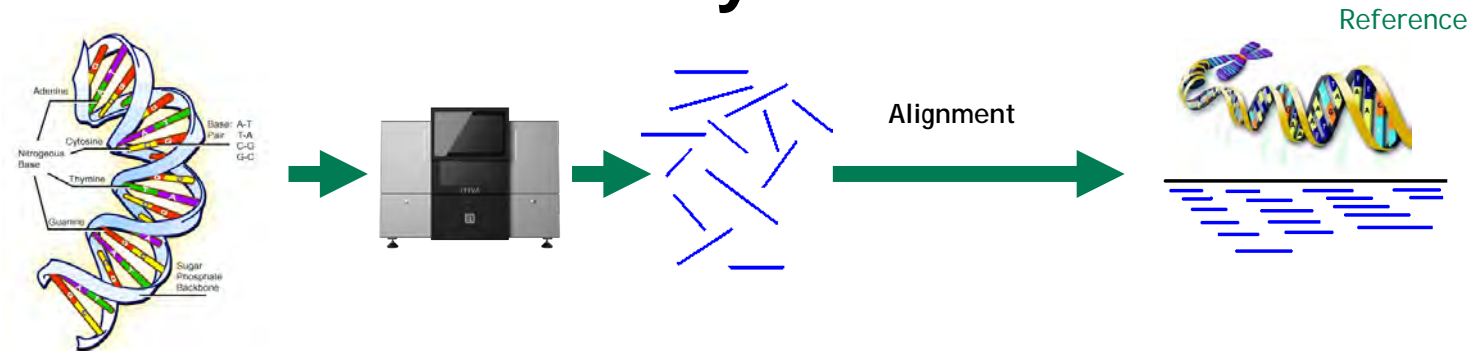**Tobias Rausch**

**European Molecular Biology Laboratory (EMBL)**
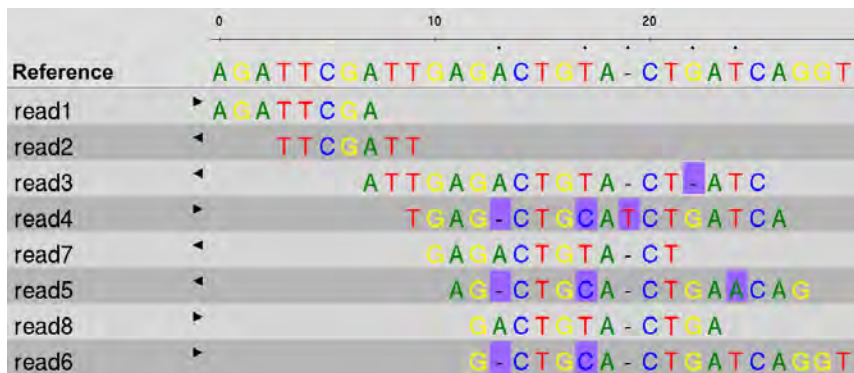
5 February 2026

EMBL

# Genome Variation Discovery
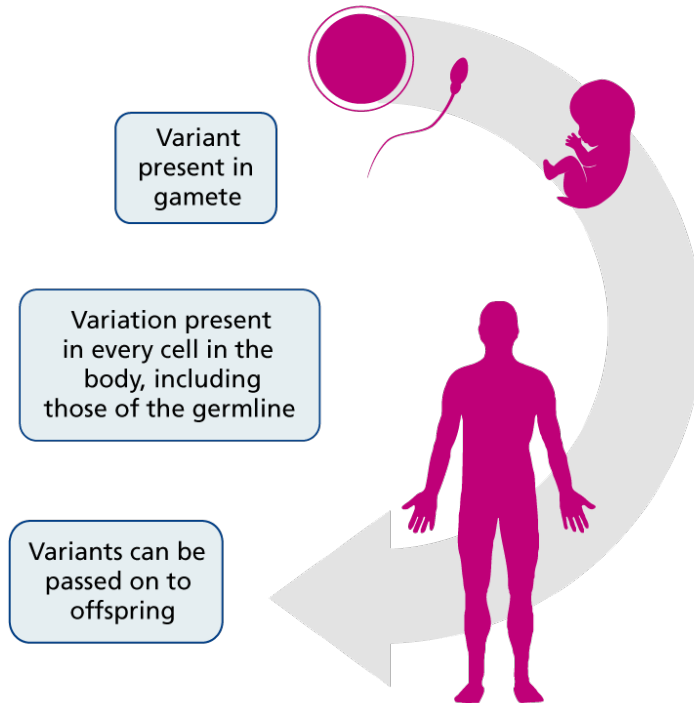


## Alignment



## Variants

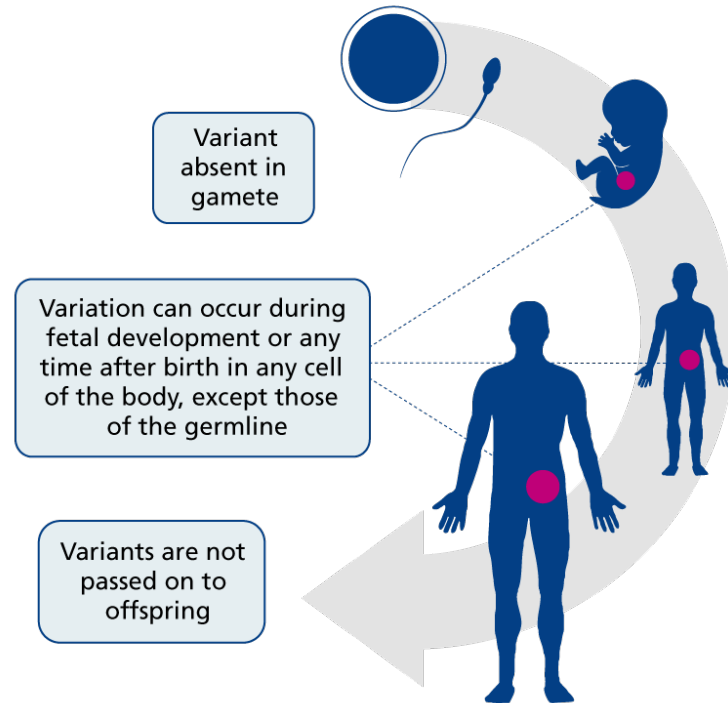| CHR | POS | ID | REF | ALT | GT |
|-----|-----|-----|-----|-----|-----|
| chr1 | 12 | . | GA | G | 0/1 |
| chr1 | 17 | rs123 | T | C | 0/1 |

Genotype (GT):
0/0: Homozygous reference
0/1: Heterozygous
1/1: Homozygous alternative

EMBL

# Germline and somatic variants
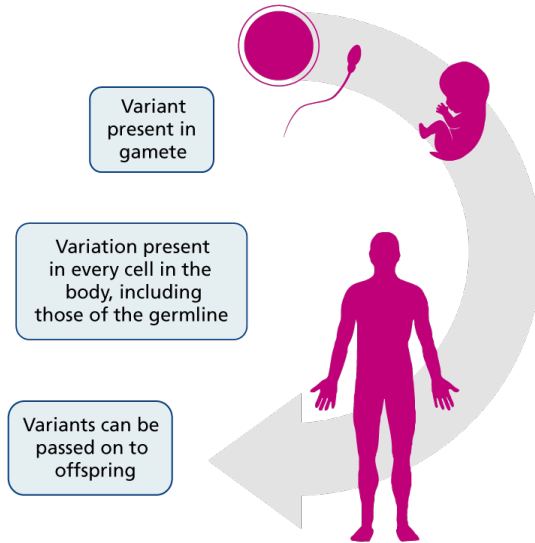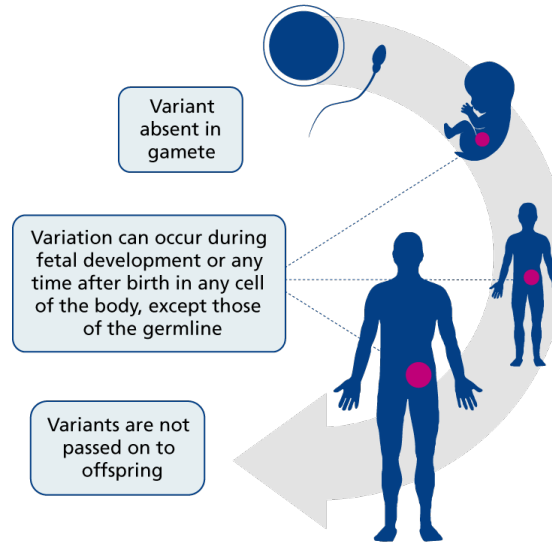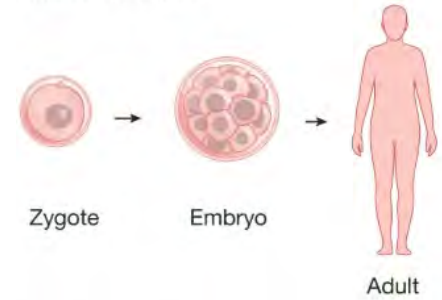


Germline Variants

Somatic Variants

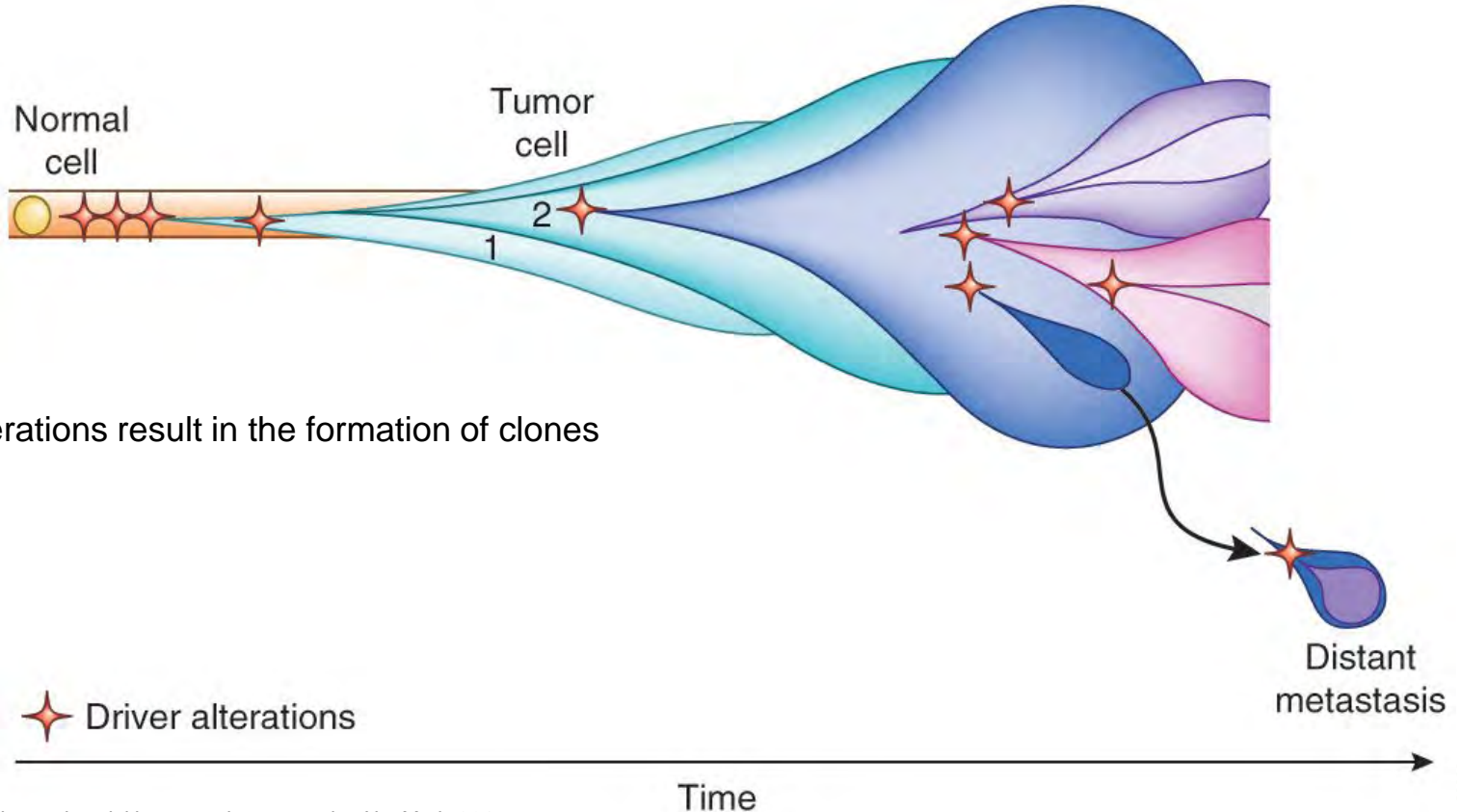Variant present in gamete

Variation present in every cell in the body, including those of the germline

Variants can be passed on to offspring

Variant absent in gamete

Variation can occur during fetal development or any time after birth in any cell of the body, except those of the germline

Variants are not passed on to offspring

EMBL

# Germline and somatic variants



Source: Coorens et al., The Somatic Mosaicism across Human Tissues Network. Nature. 2025

# Cancer as a disease of the genome



**Tumor evolution**
Somatic driver alterations result in the formation of clones

*Alizadeh et al. Toward understanding and exploiting tumor heterogeneity. Nat Med., 2015.*

# Agenda: Somatic variants as the driver of cancer

## Tumor Evolution



- Per genome
  - ~3-5 million germline variants
  - 100-100,000 somatic variants

## Germline and somatic variants



## Tumor heterogeneity

Schematic depiction of a bi-clonal tumor



normal    clone 1    clone 2

## Complex Rearrangements



Original chromosome sequence

Catastrophic chromosome breakage

Rearranged chromosome

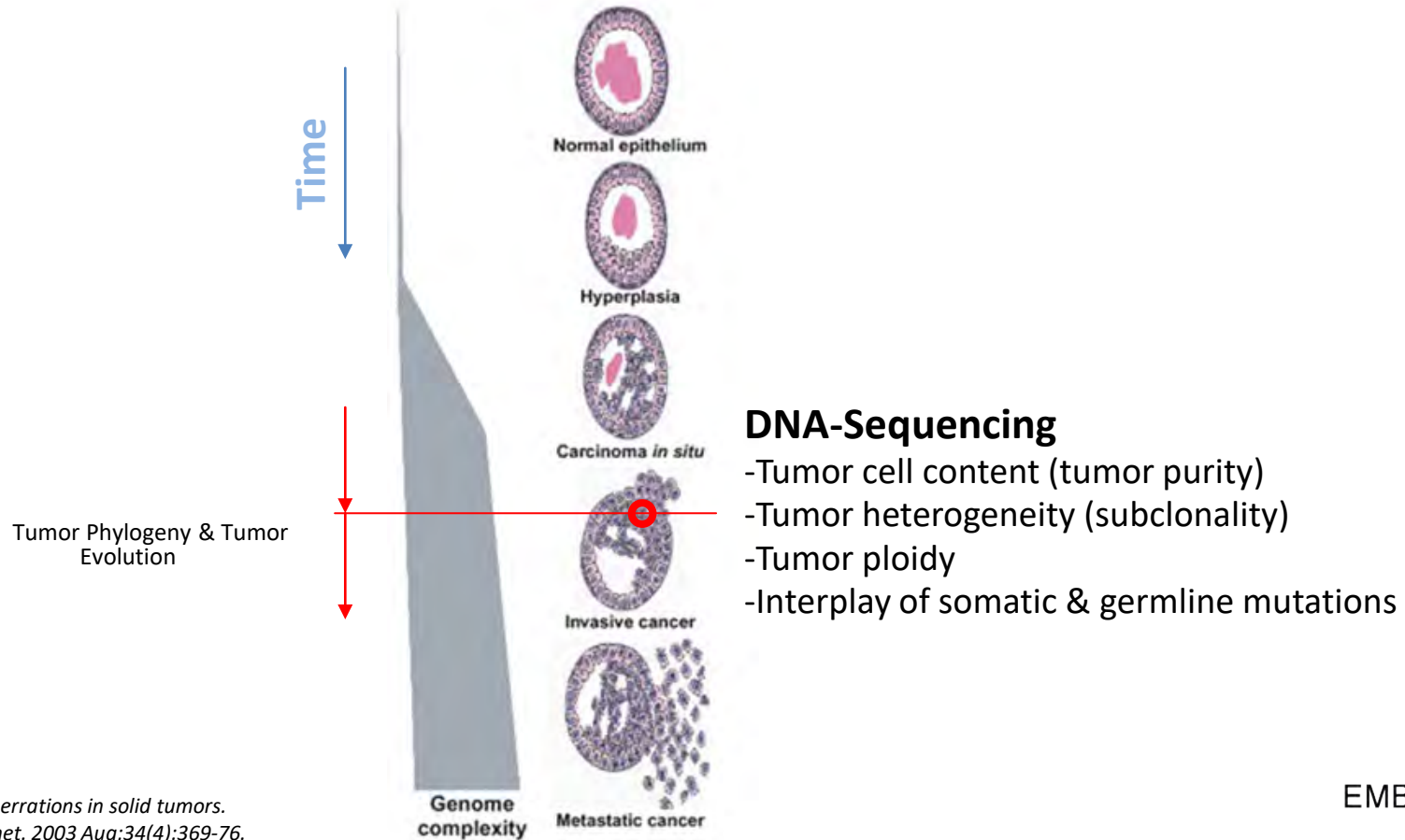Lost chromosomal material

Sources:
Genomic landscape of non-small cell lung cancer in smokers and never-smokers. Govindan et al., Cell. 2012 Sep 14;150(6):1121-34.
Evolution of the cancer genome. Yates and Campbell, Nat Rev Genet. 2012 Nov;13(11):795-806.
Genome sequencing of pediatric medulloblastoma links catastrophic DNA rearrangements with TP53 mutations. Rausch et al., Cell. 2012 Jan 20;148(1-2):59-71.
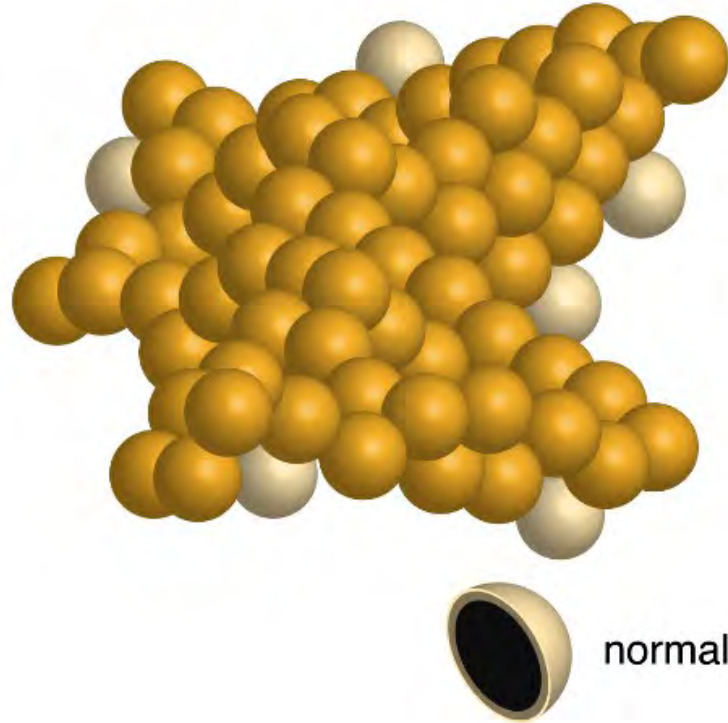
EMBL

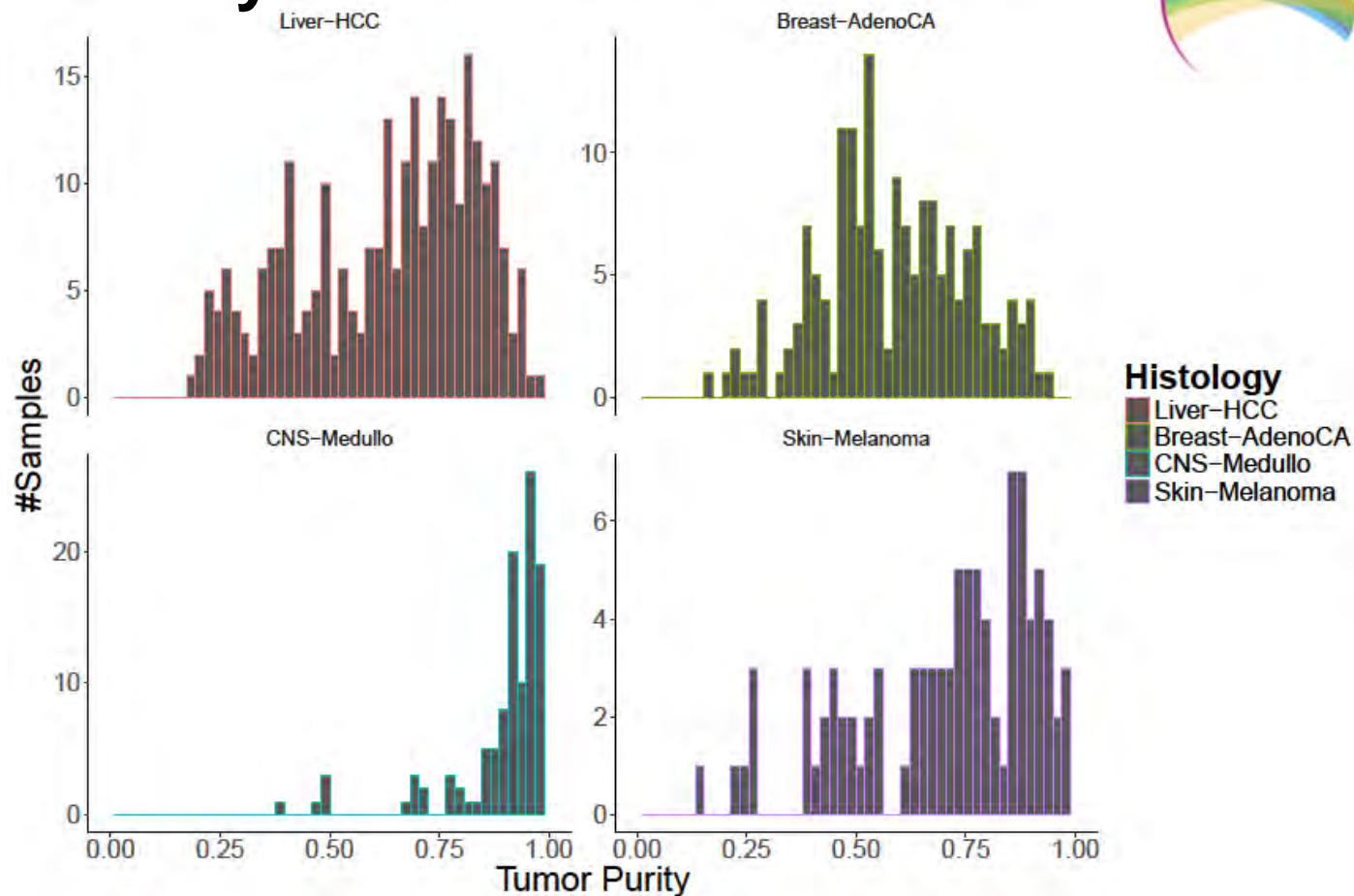# Cancer Genomics
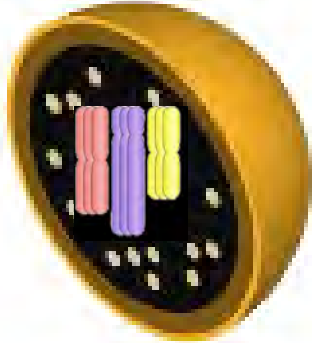## Sequencing provides a snapshot in time and space
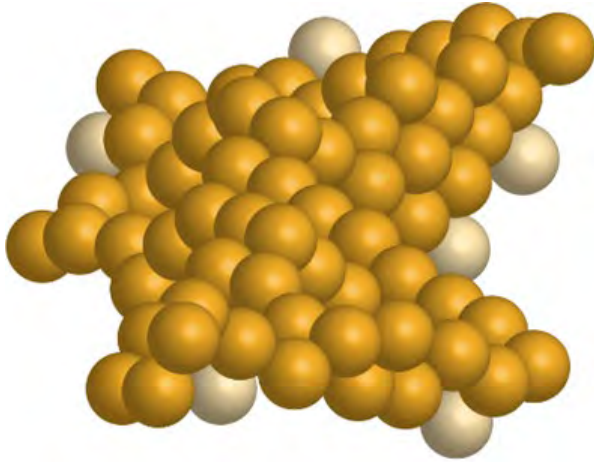


Time

Tumor Phylogeny & Tumor Evolution

Normal epithelium

Hyperplasia

Carcinoma *in situ*

Invasive cancer

Metastatic cancer

Genome complexity

**DNA-Sequencing**
-Tumor cell content (tumor purity)
-Tumor heterogeneity (subclonality)
-Tumor ploidy
-Interplay of somatic & germline mutations

EMBL

# Tumor Purity



Schematic depiction of a mono-clonal tumor

normal

EMBL

# Tumor Purity

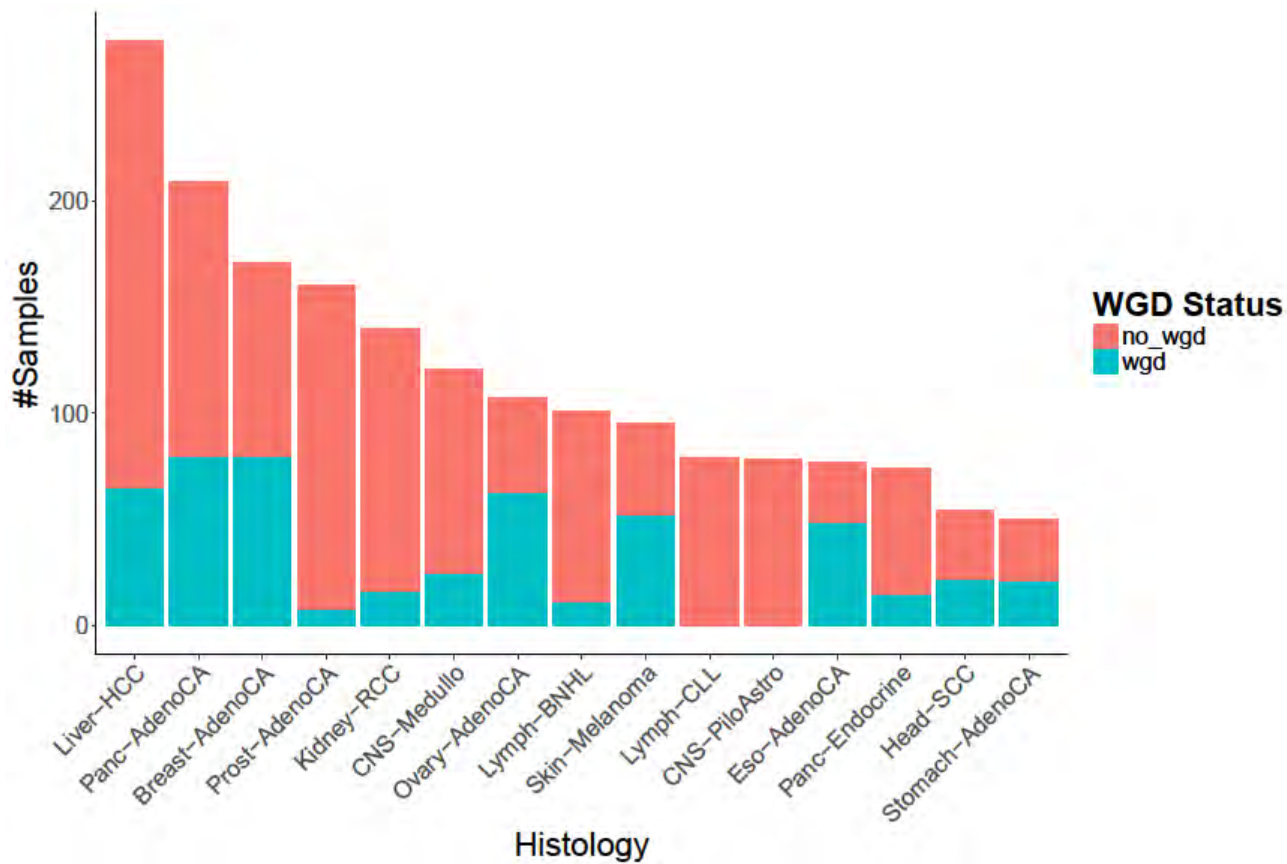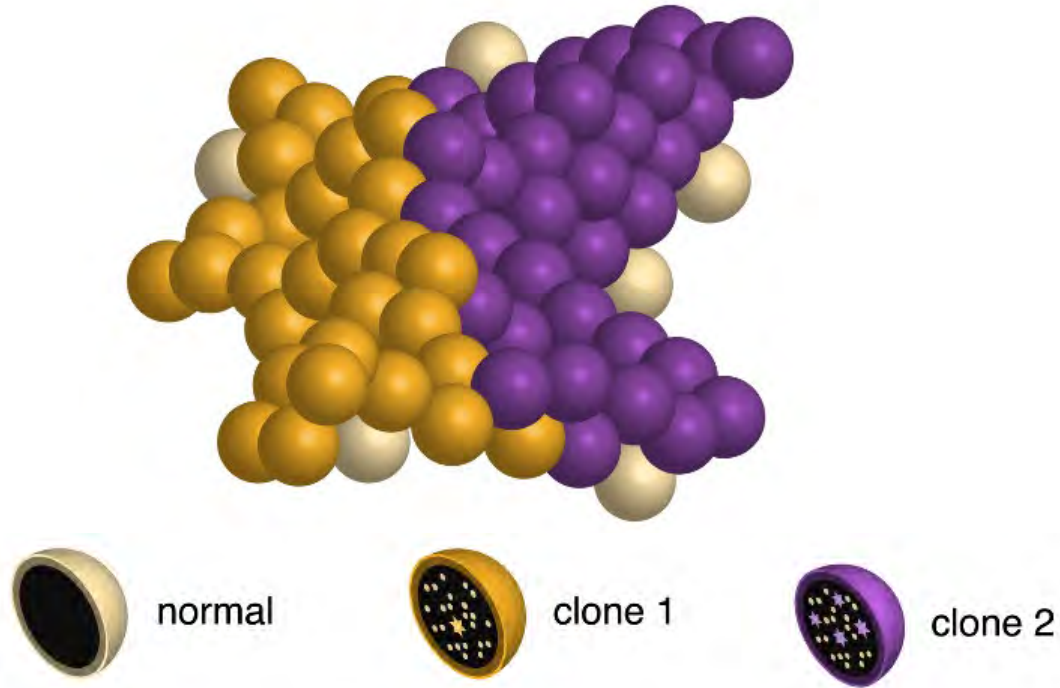# Tumor Ploidy

# Tumor Ploidy

# Inferred Whole-Genome Duplication (WGD) Status

# Tumor Clonality/Heterogeneity
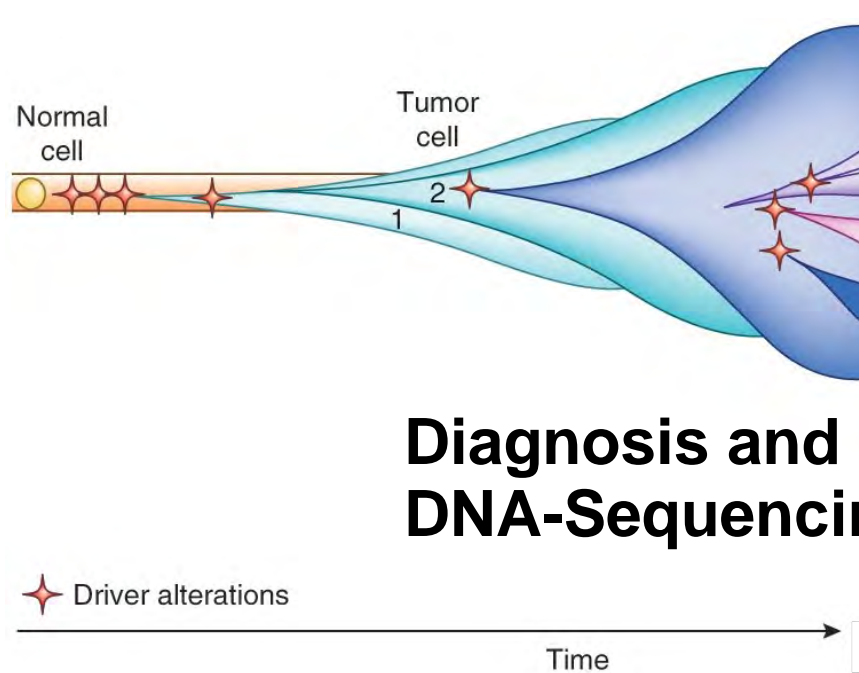
Schematic depiction of a bi-clonal tumor



normal          clone 1          clone 2

EMBL

# Tumor multi-clonality



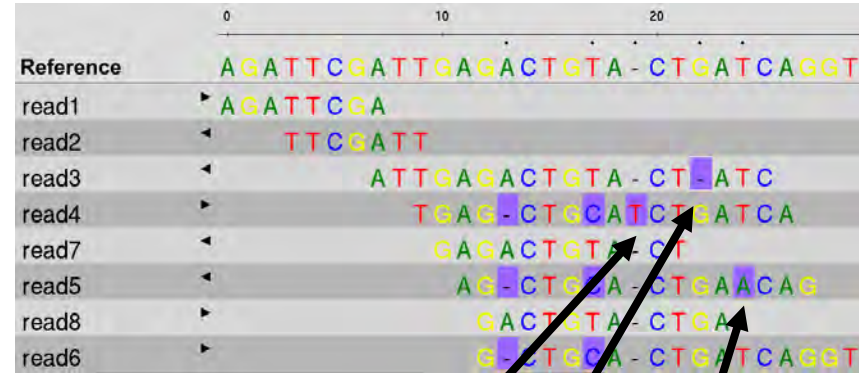Schematic depiction of a bi-clonal tumor

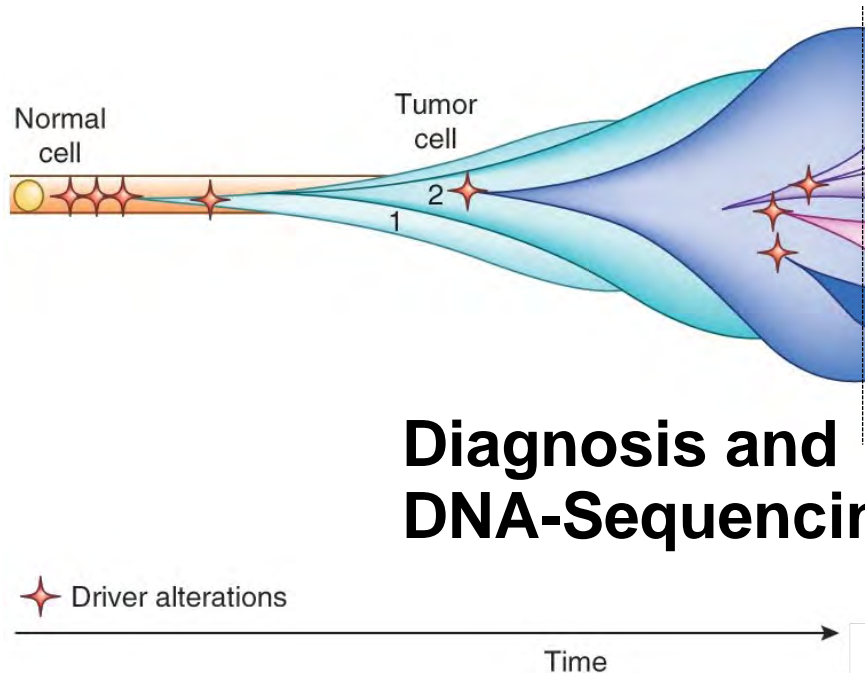## Diagnosis and DNA-Sequencing

# Tumor multi-clonality



**Diagnosis and DNA-Sequencing**

Sequencing errors?
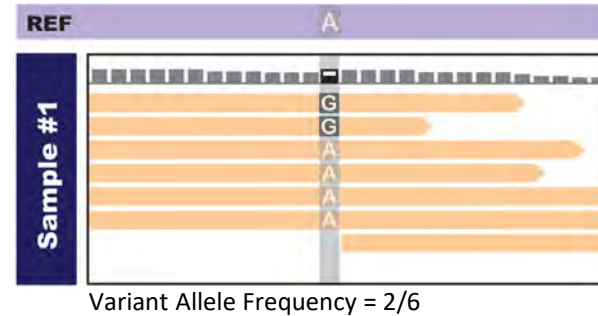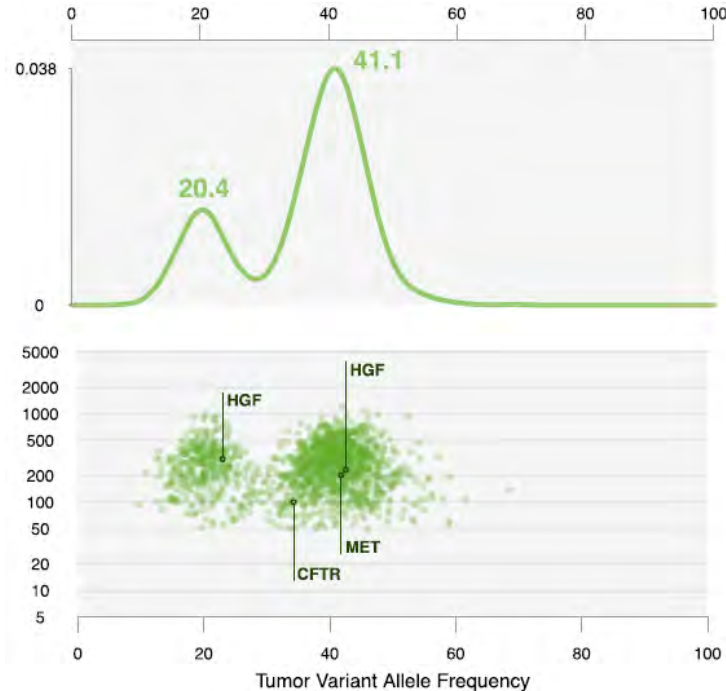
# Tumor multi-clonality



## Diagnosis and DNA-Sequencing

Sequencing errors?

→ Cancer genomes are often sequenced to >60x because of tumor purity, tumor heterogeneity and many chromosomal aberrations
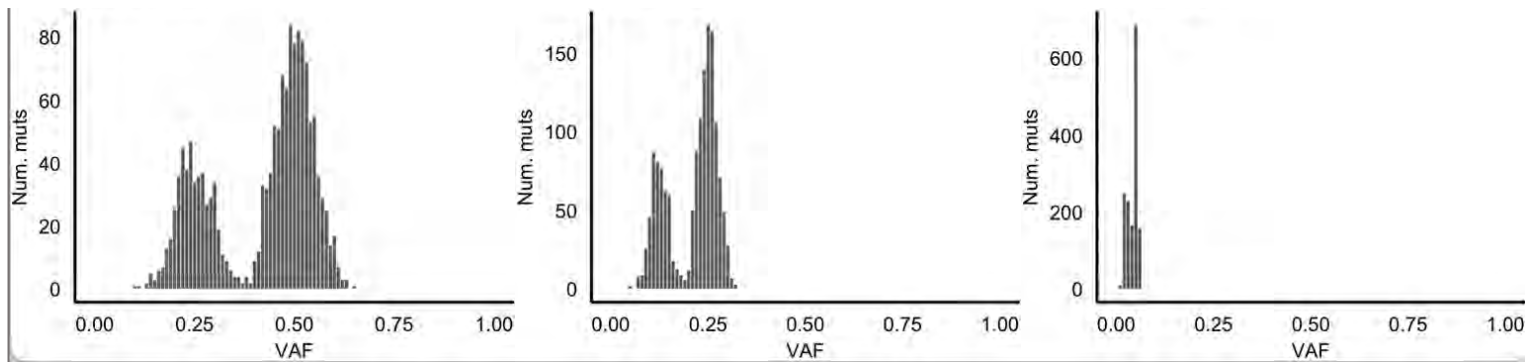
EMBL

# Tumor Clonality/Heterogeneity

- Requires high sequencing depth
  - Typically done with WES data (≥500x coverage)



Variant Allele Frequency = 2/6

Bi-clonal

82.2% tumor purity

*Source: Genomic landscape of non-small cell lung cancer in smokers and never-smokers. Govindan et al., Cell. 2012 Sep 14;150(6):1121-34.*
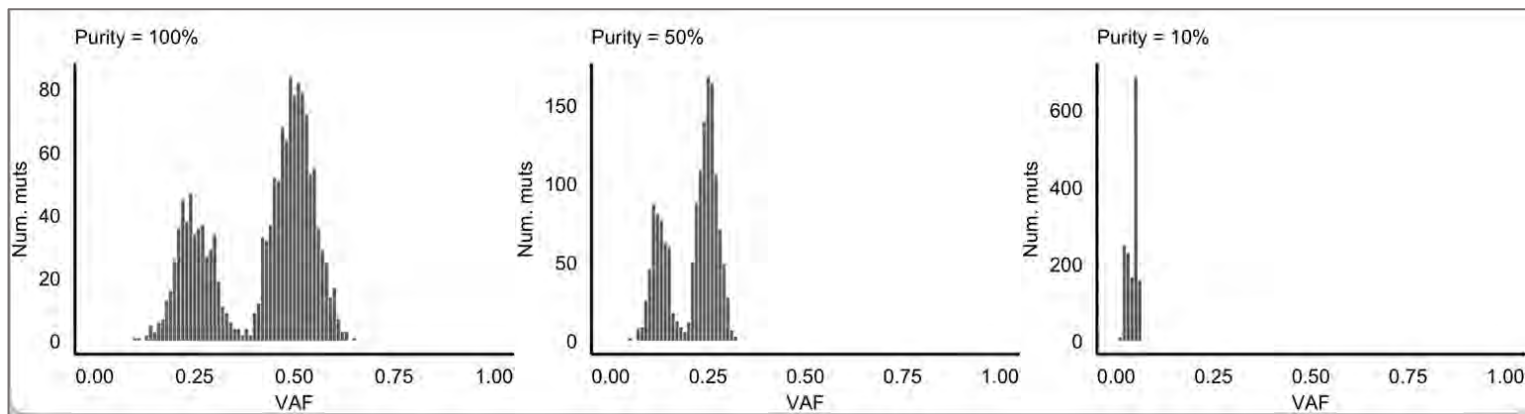
# Tumor purity

- Below are 3 tumor samples with varying levels of tumor cell content
- Can you guess the tumor purity based on the somatic VAF?
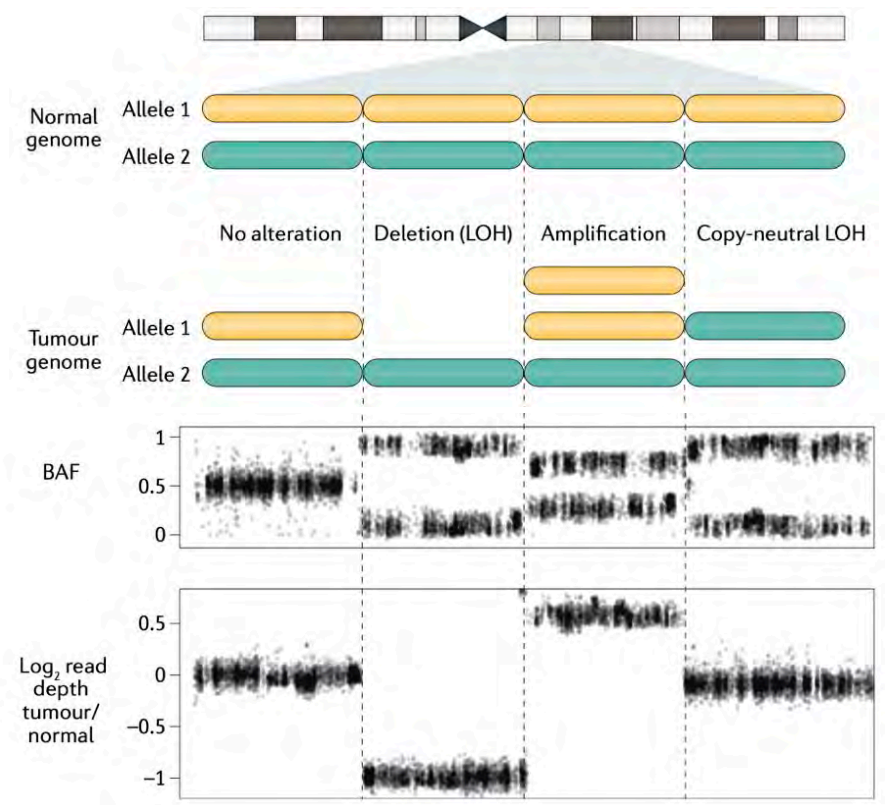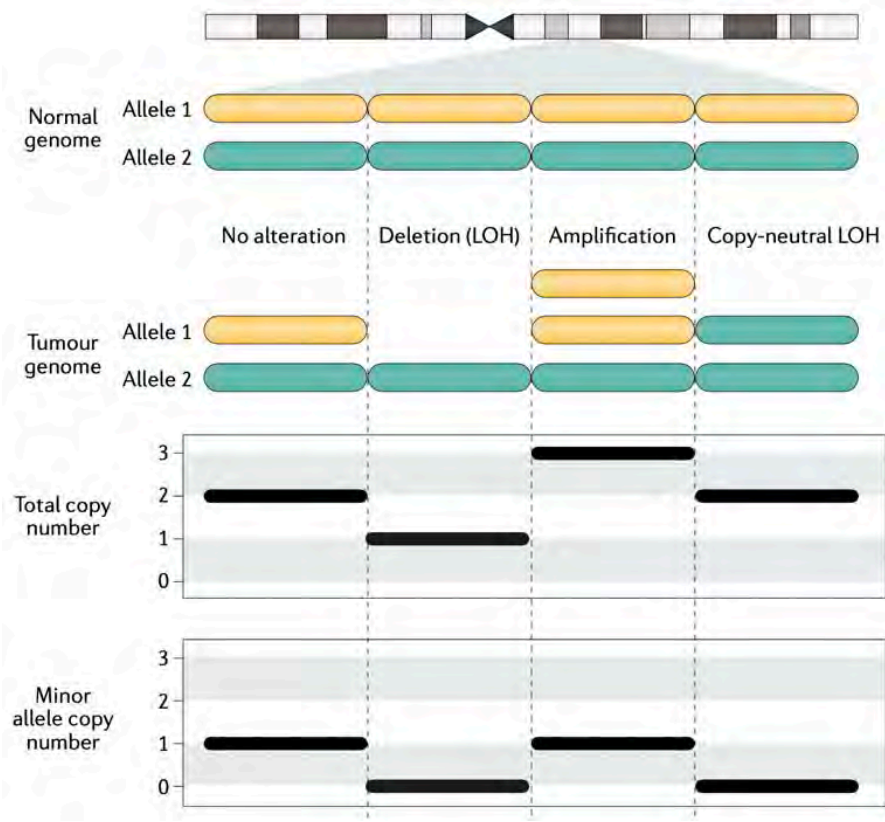


EMBL

# Tumor purity

- Below are 3 tumor samples with varying levels of tumor cell content
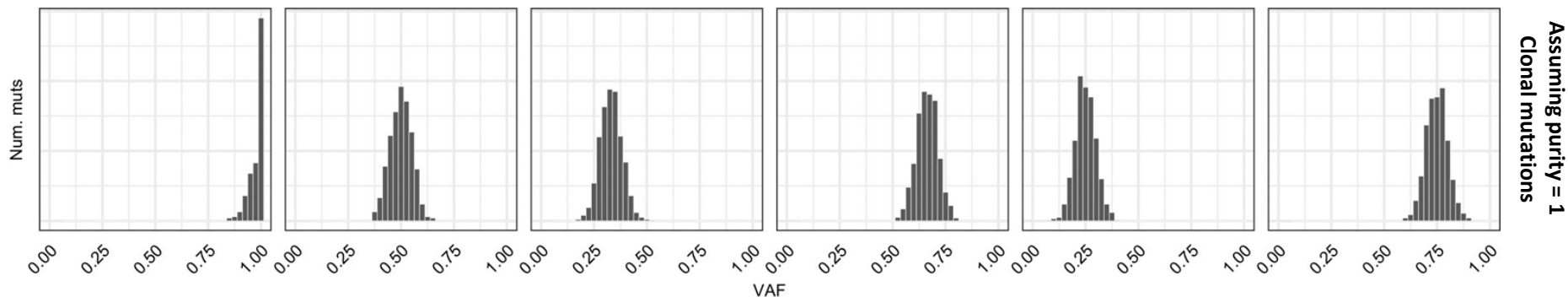- Can you guess the tumor purity based on the somatic VAF?
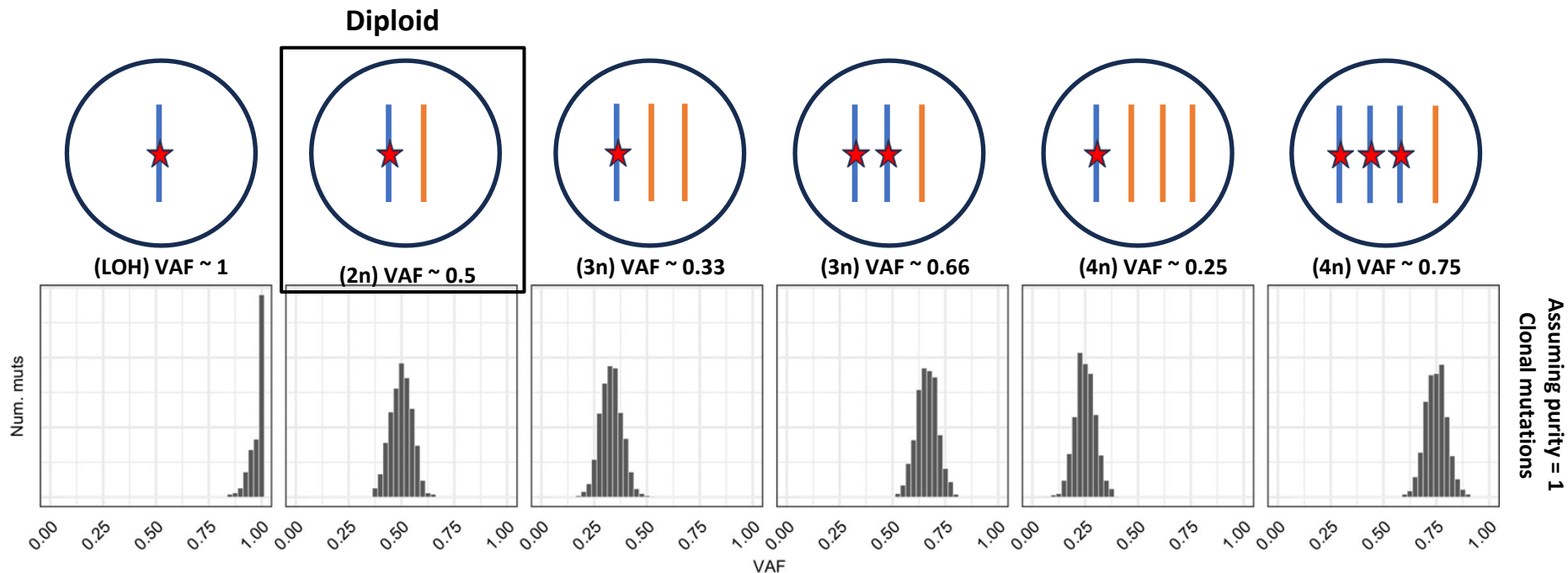
# Somatic copy-number variation (SCNA)

# Copy number changes affect the VAF distribution

- Assuming a 100% tumor purity and only clonal mutations
- Can you guess the total copy-number?

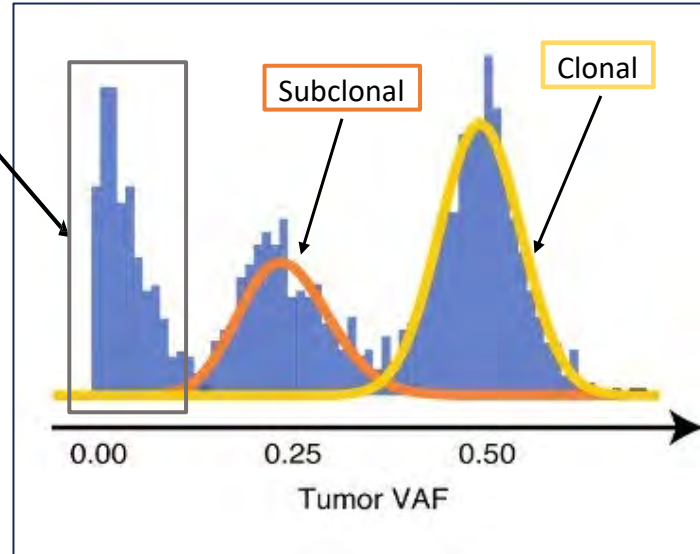# Copy number changes affect the VAF distribution

# Subclonal heterogeneity

- Bulk sequencing data provides an aggregated view of the genetic information
- Subclonal mutations that are present at low frequencies are difficult to detect
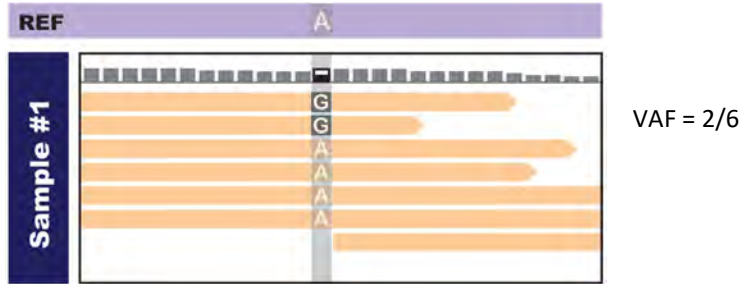  - → Underestimation of clonal diversity is common

Tail (not all from the same clone)
- Very hard to distinguish clones
- New mutations (majority passenger) occurring in every cell division
- The accuracy of the reconstruction **massively depends on the depth of coverage**
- Subclonal mutations at low VAF **heavily affected by sequencing errors**

EMBL

# Variant vs. Population Allele Frequency

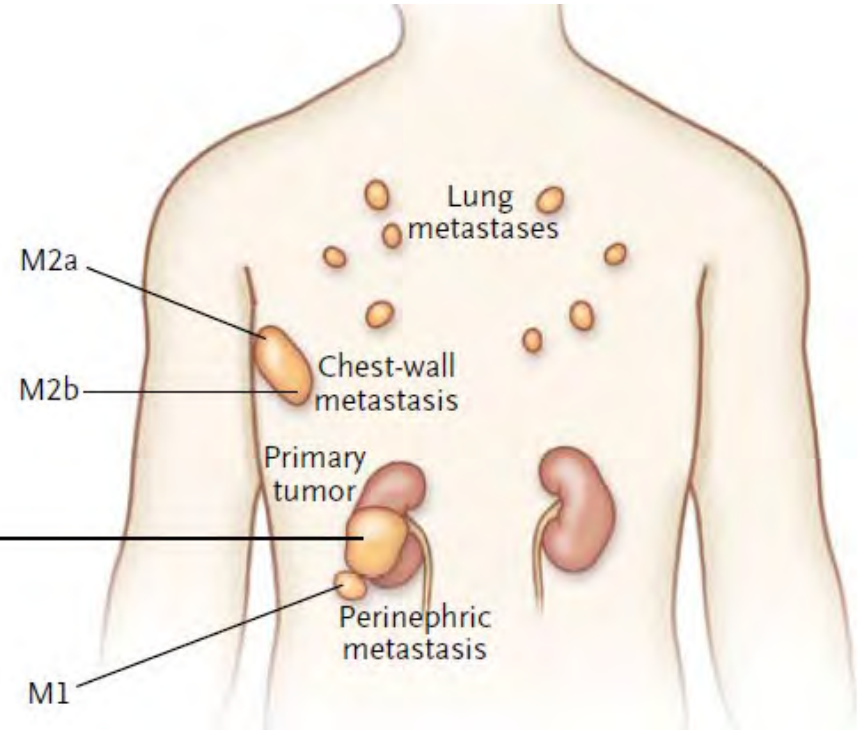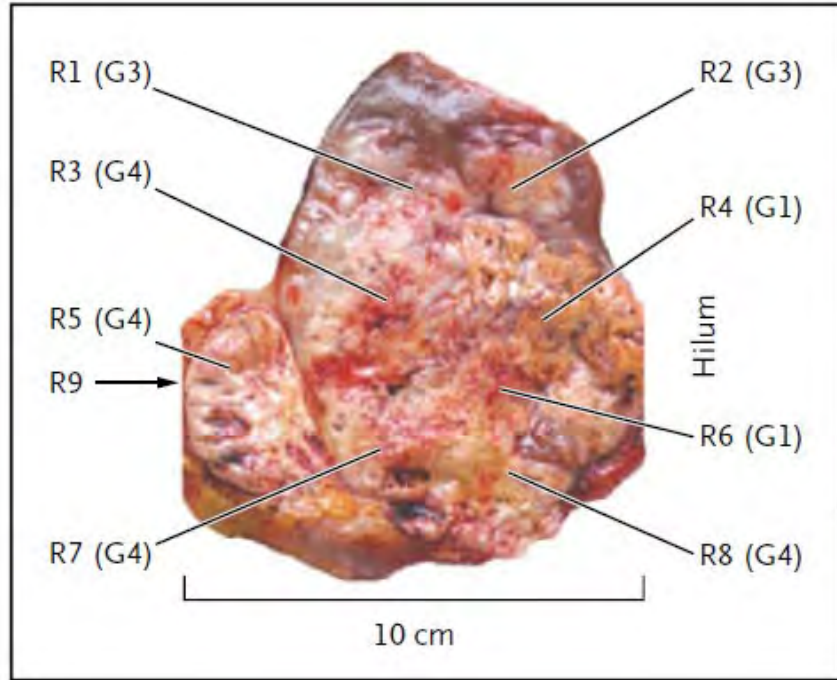- Somatic variant allele frequency



VAF = 2/6

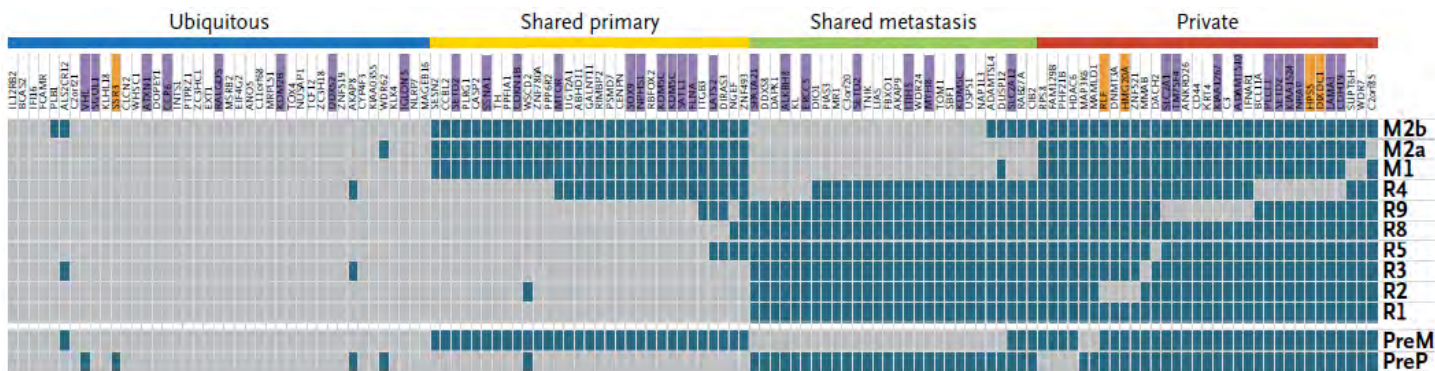- Population allele frequency of (potentially predisposing) germline variants
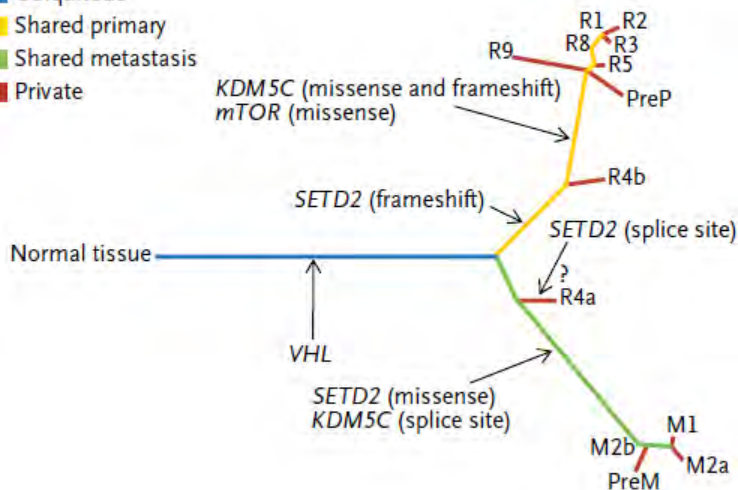  - rs6602666, A>G
    - African: ~26%
    - Finnish: 0%

# Tumor Heterogeneity & Tumor Evolution
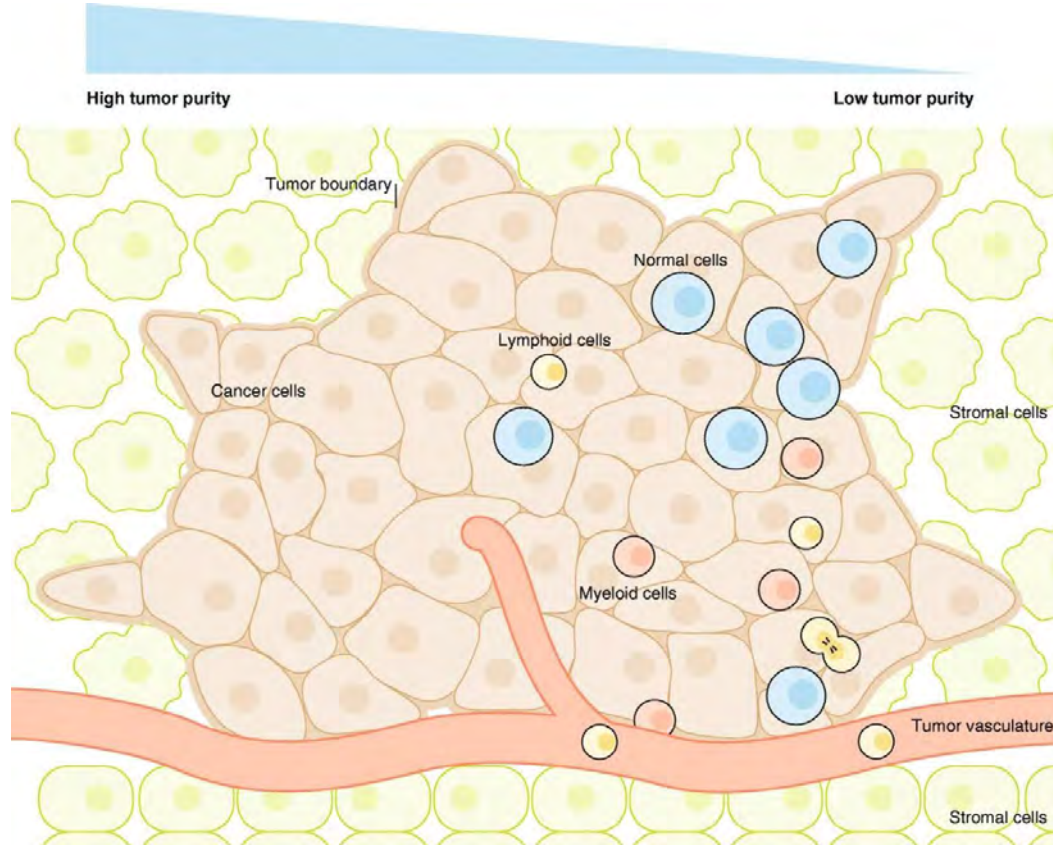
# Tumor Heterogeneity & Tumor Evolution



Tumor Phylogeny

# The hidden complexity in bulk sequencing

# Single-cell sequencing

# Single-cell (spatial) transcriptomics
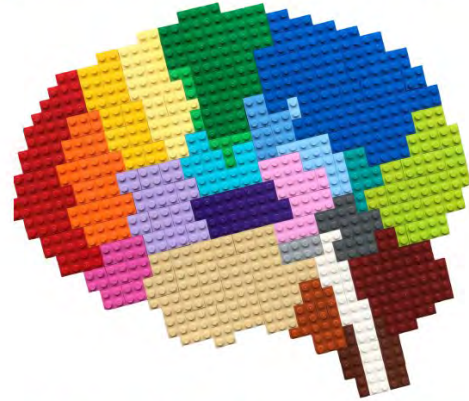


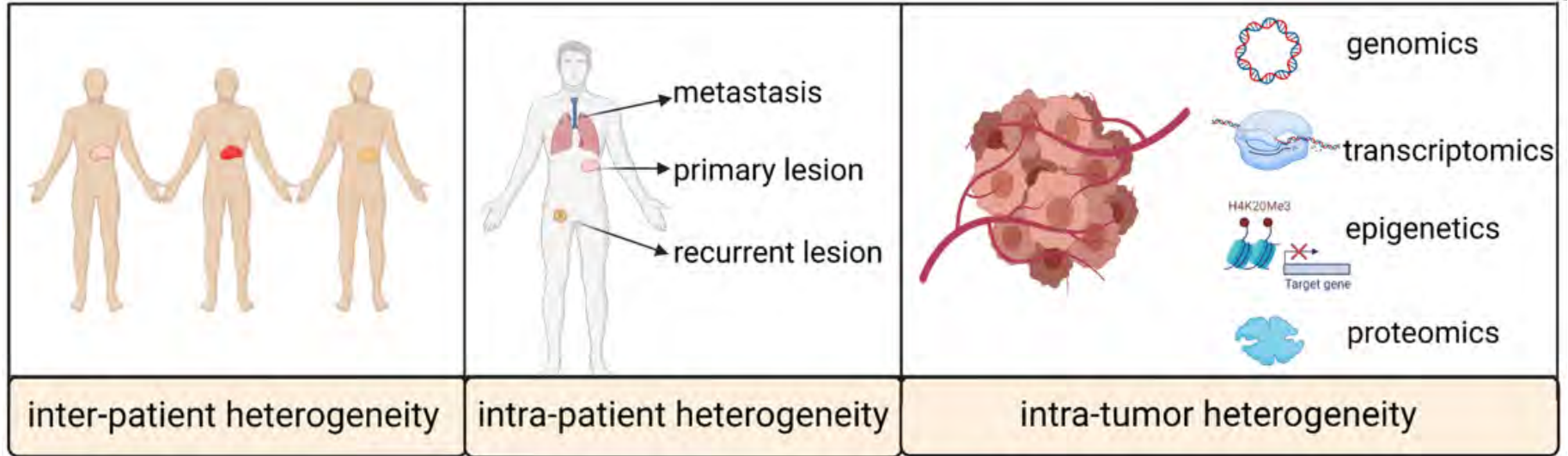**Bulk** RNA-sequencing

**Single-cell RNA-sequencing**

+ =

←— 6.5 mm —→

6.5 mm

**Spatial transcriptomics**

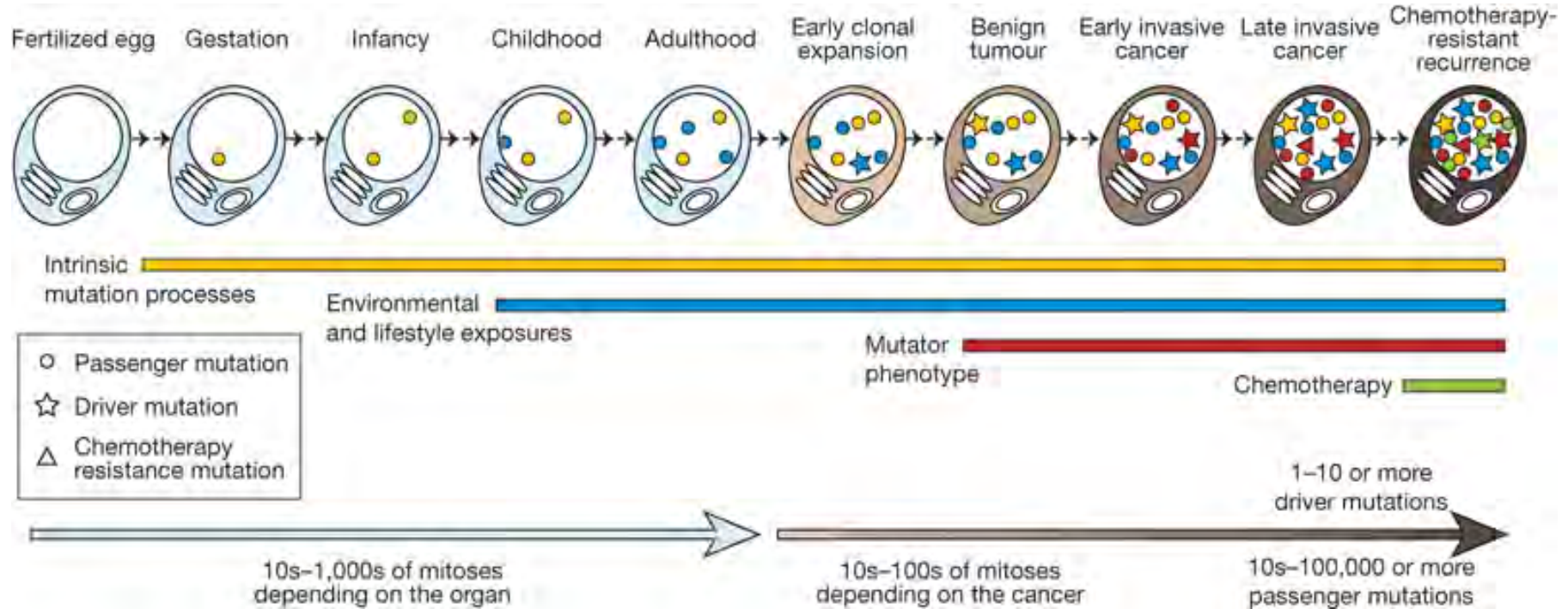**Spatial location of all cell types**

**Study cell-cell interactions**

EMBL

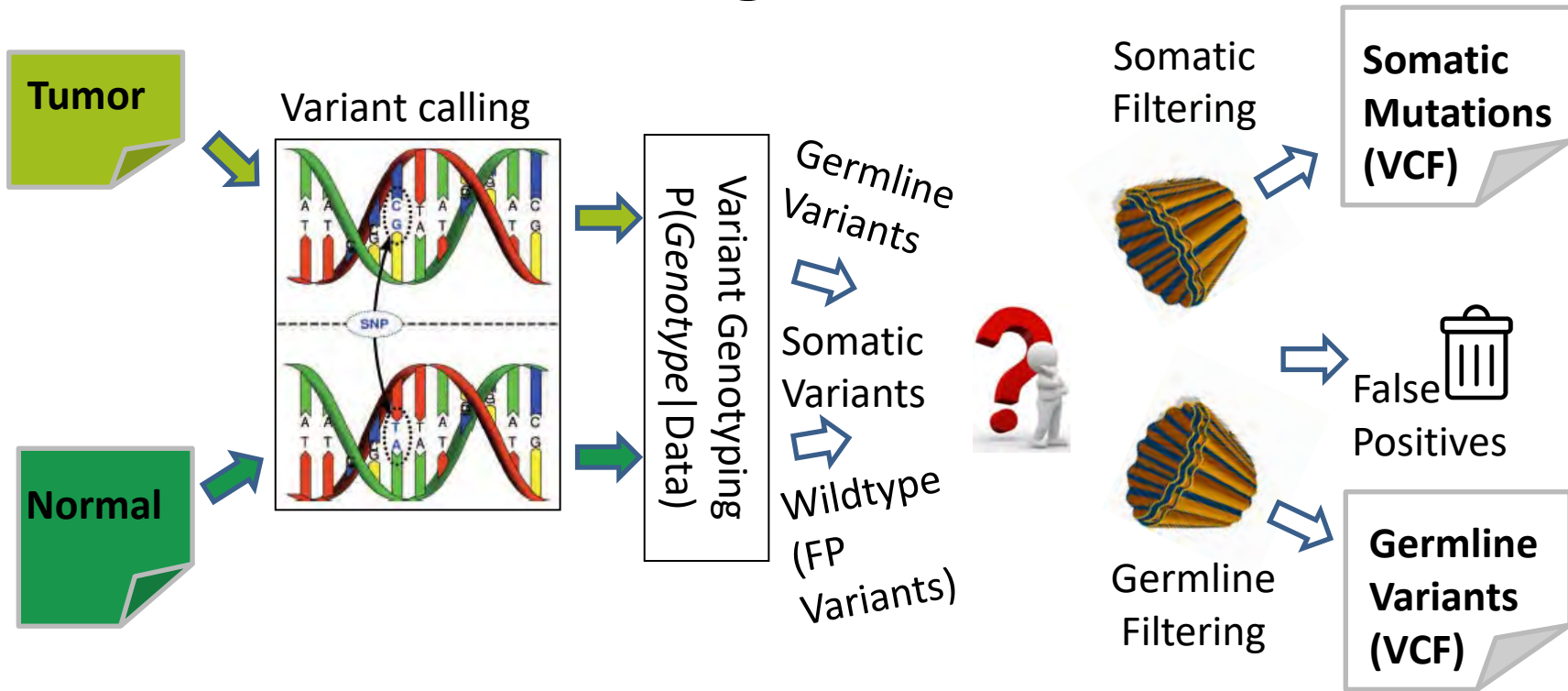# Tumor heterogeneity at the patient and tumor level

# Summary: Cancer as a genetic disease and evolutionary process



Fertilized egg | Gestation | Infancy | Childhood | Adulthood | Early clonal expansion | Benign tumour | Early invasive cancer | Late invasive cancer | Chemotherapy-resistant recurrence

Intrinsic mutation processes

Environmental and lifestyle exposures

Mutator phenotype

Chemotherapy

○ Passenger mutation
☆ Driver mutation
△ Chemotherapy resistance mutation

1–10 or more driver mutations

10s–1,000s of mitoses depending on the organ

10s–100s of mitoses depending on the cancer

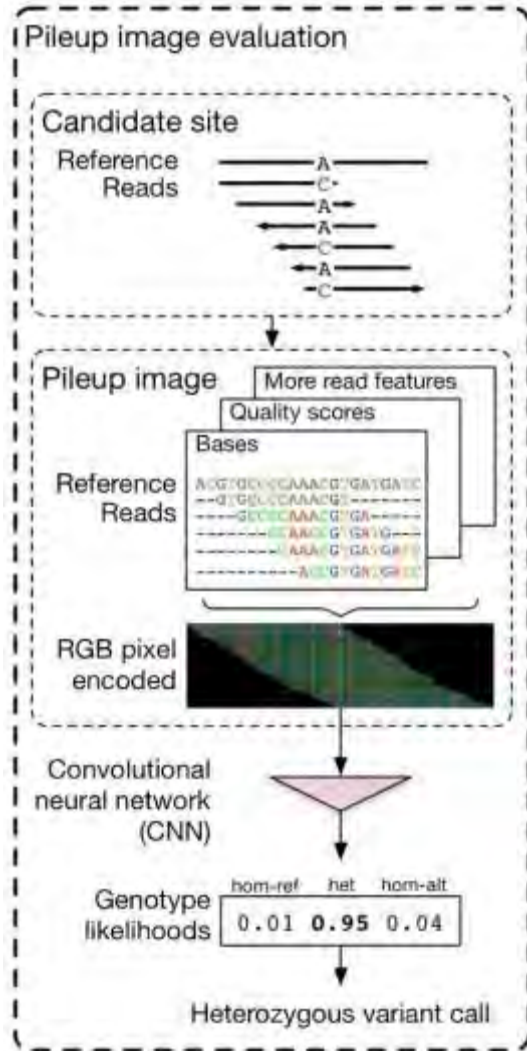10s–100,000 or more passenger mutations

EMBL

# Cancer Genome Data Analysis
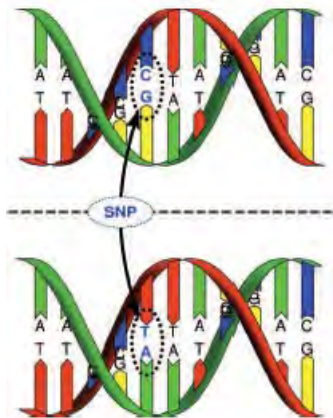
# Somatic variant calling

# Variant Calling Methods

- Four broad categories
  - Heuristic Methods
    - Hardly used anymore
  - Probabilistic methods
    - Bayesian methods
  - Machine Learning methods
    - Deep Learning methods, e.g., DeepVariant
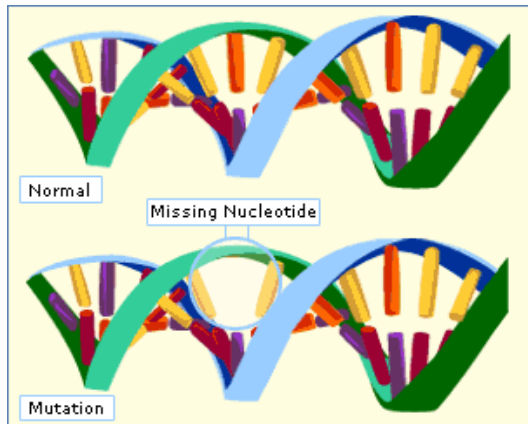  - Graph-based methods
    - Pangenome methods

# Types of Variants

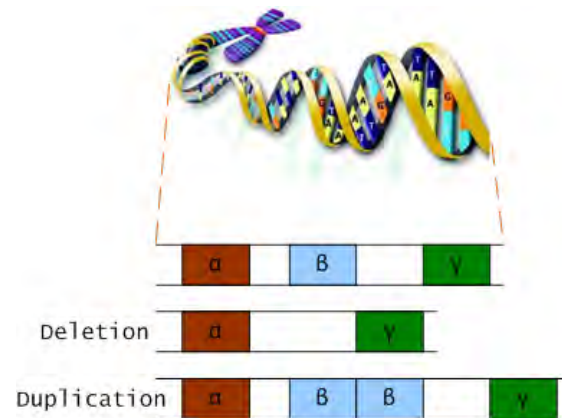Single-nucleotide variants (SNVs)

Short insertions & deletions (InDels)
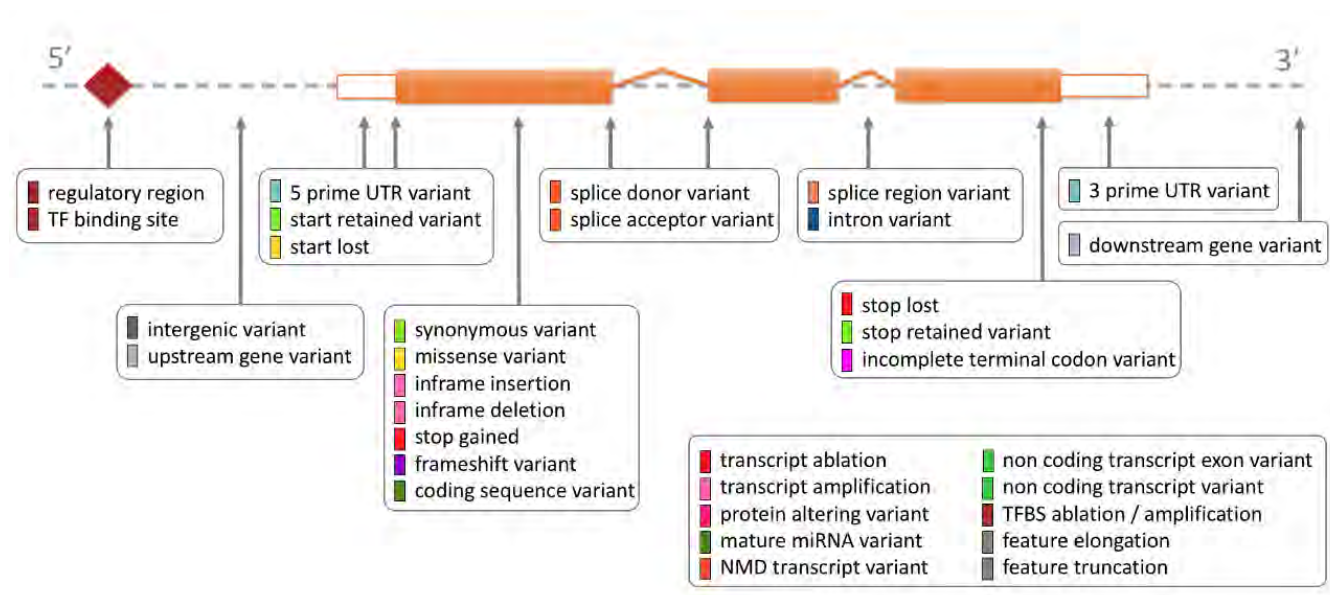
Copy-Number Variants (CNVs) & Structural Variants (SVs)



*Size: 1bp*

*1-50bp*

*>50bp*

Different methods are used to discover and genotype SNVs, InDels, SVs and CNVs.

Further information:
Methods in Genomic Variant Calling: https://www.youtube.com/watch?v=zO9WCOaq3aQ
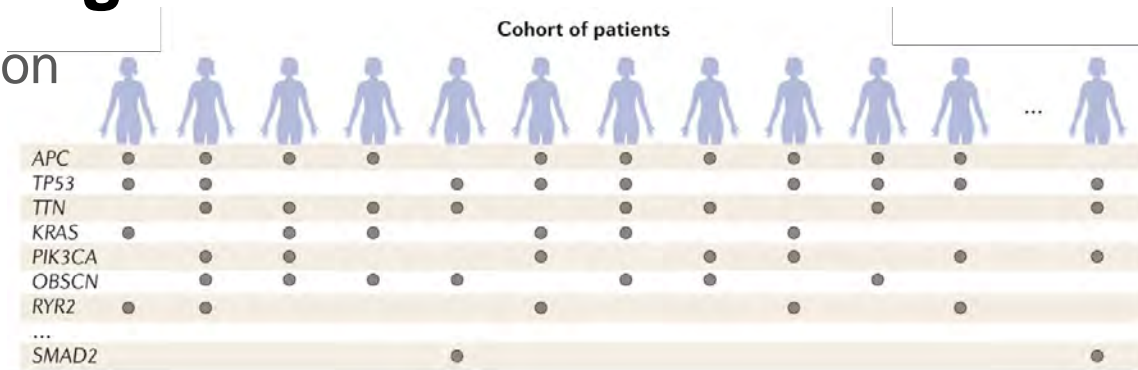
EMBL

# Interpreting Genomic Variants



Popular Tools: VEP, Annovar, snpEff

Mediocre support for annotation of copy-number and structural variants

*Source: https://m.ensembl.org/info/genome/variation/prediction/predicted_data.html*

# Driver versus Passenger

- Signals of positive selection

    - Mutation Recurrence

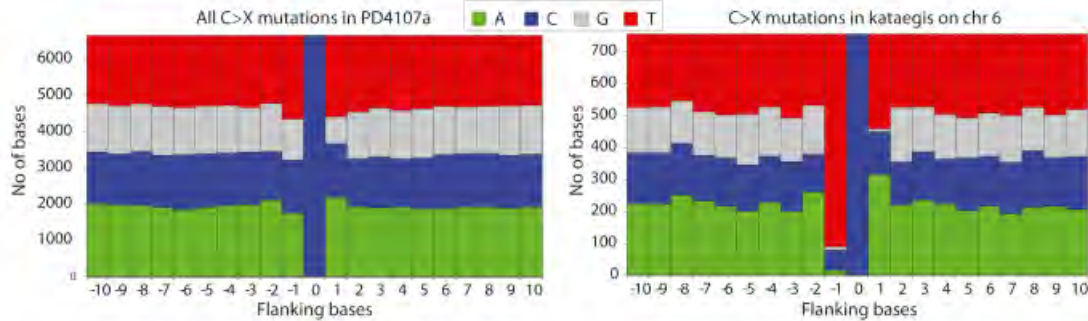    - Hotspot mutations

    - Protein domain clustering

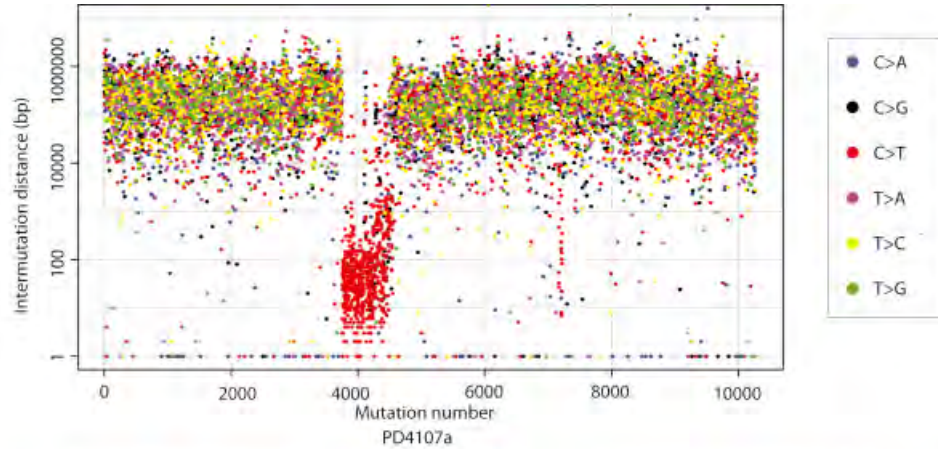    - Functional impact bias
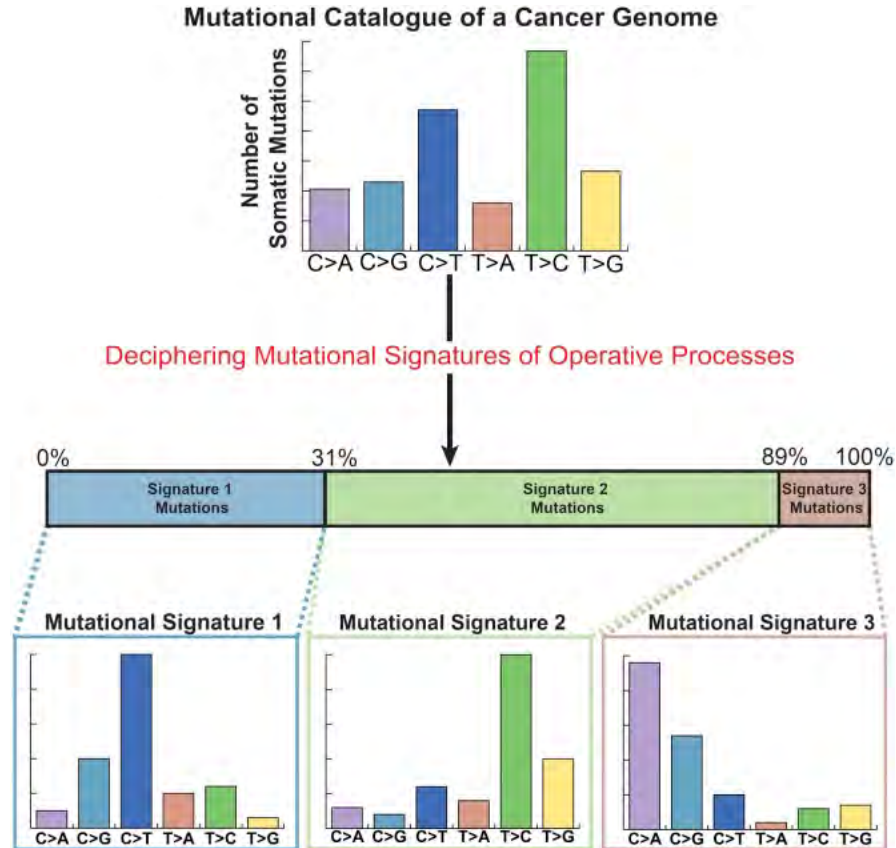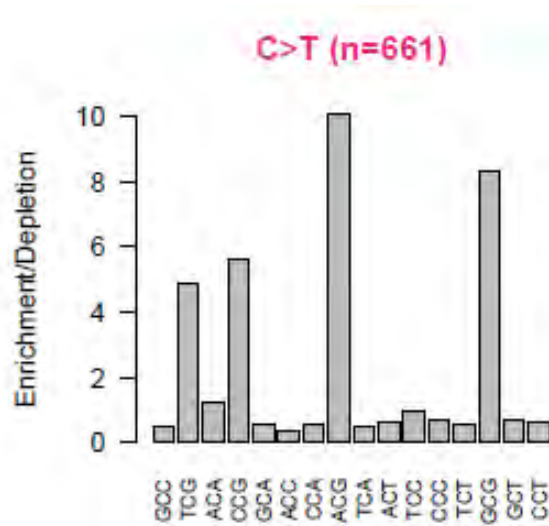
# Mutational Signatures
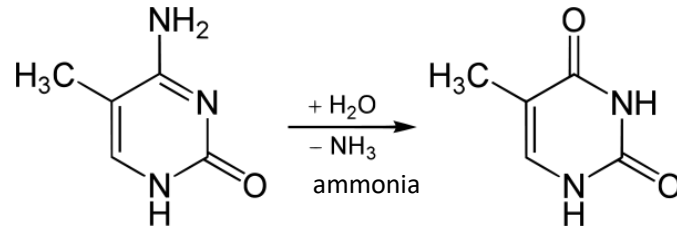
# Mutational Signatures



Tri-Nucleotide Context TpCpX

EMBL

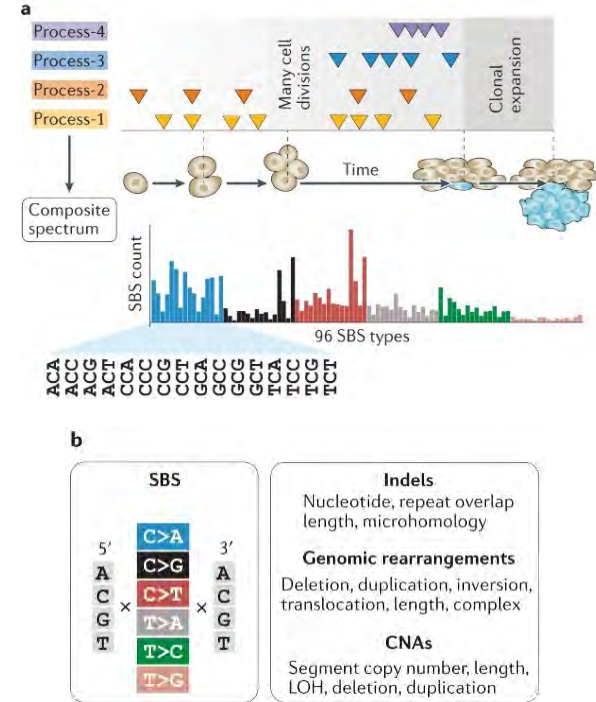# Mutational Signatures

# Tri-Nucleotide Sequence Context



C>T (n=661)

- Elevated **C>T** mutation rate at **XpCpG** trinucleotides
- Deamination of methylated cytosines to thymine (usually at XpCpGs)
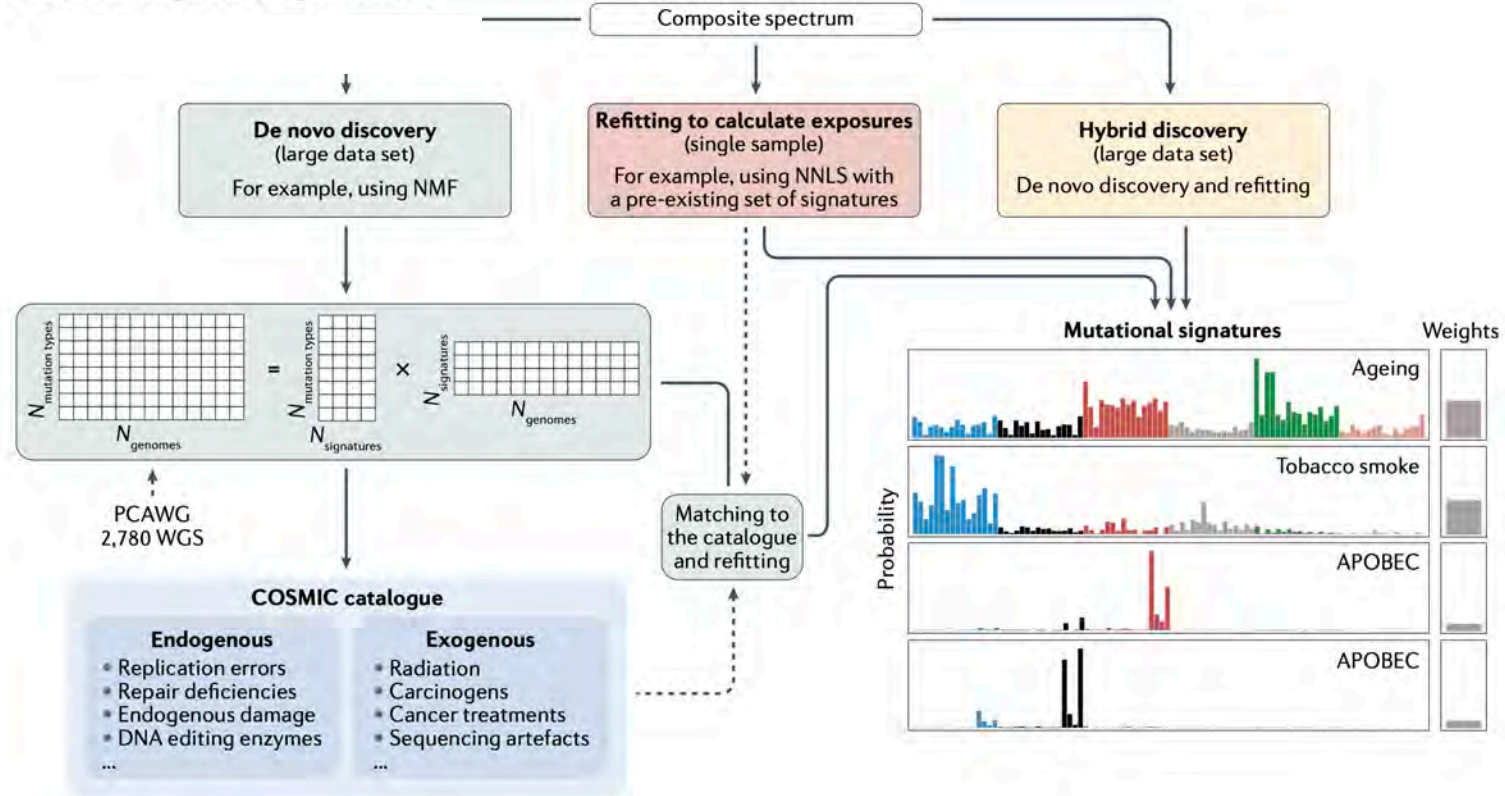


EMBL

# Mutational Signature

- Characterize patterns of DNA mutations

- These patterns can arise from various mutational processes that lead to changes in DNA sequences

- Mutations can occur due to a variety of factors
  - DNA replication errors
  - Exposure to mutagenic agents (e.g., ultraviolet radiation, chemical carcinogens)
  - Defects in DNA repair mechanisms.



*Cortes-Ciriano et al 2022*

# De novo discovery or re-fitting

# Mutational signature catalog

# Mutational Signatures across Human Cancer Types

# Classification of mutational signatures based on biological processes

- Cancer Aetiology Investigation
  - Mutational signatures provide insights into the underlying causes
  - https://cancer.sanger.ac.uk/signatures/

# Recurrently Mutated Cancer Genes

# Summary: Human Genome Variation



**A Cancer Genome**
**Somatic Variants**

**A Typical Genome**
**Germline Variants**

**Population of 2,504 peoples**
**Germline Variants**

### Origin of Variants

|  | Coding | Non-coding |
|---|---|---|
| Germ-line | 22K | 4.1 – 5M |
| Somatic | ~50 | 5K |

### Class of Variants

| SNP | 3.5 – 4.3M |
|---|---|
| Indel | 550 – 625K |
| SV | 6 – 8K |
| Total | 4.1 – 5M |

| SNP | 84.7M |
|---|---|
| Indel | 3.6M |
| SV | 60K |
| Total | 88.3M |

**Prevalence of Variants**

Passenger

Common
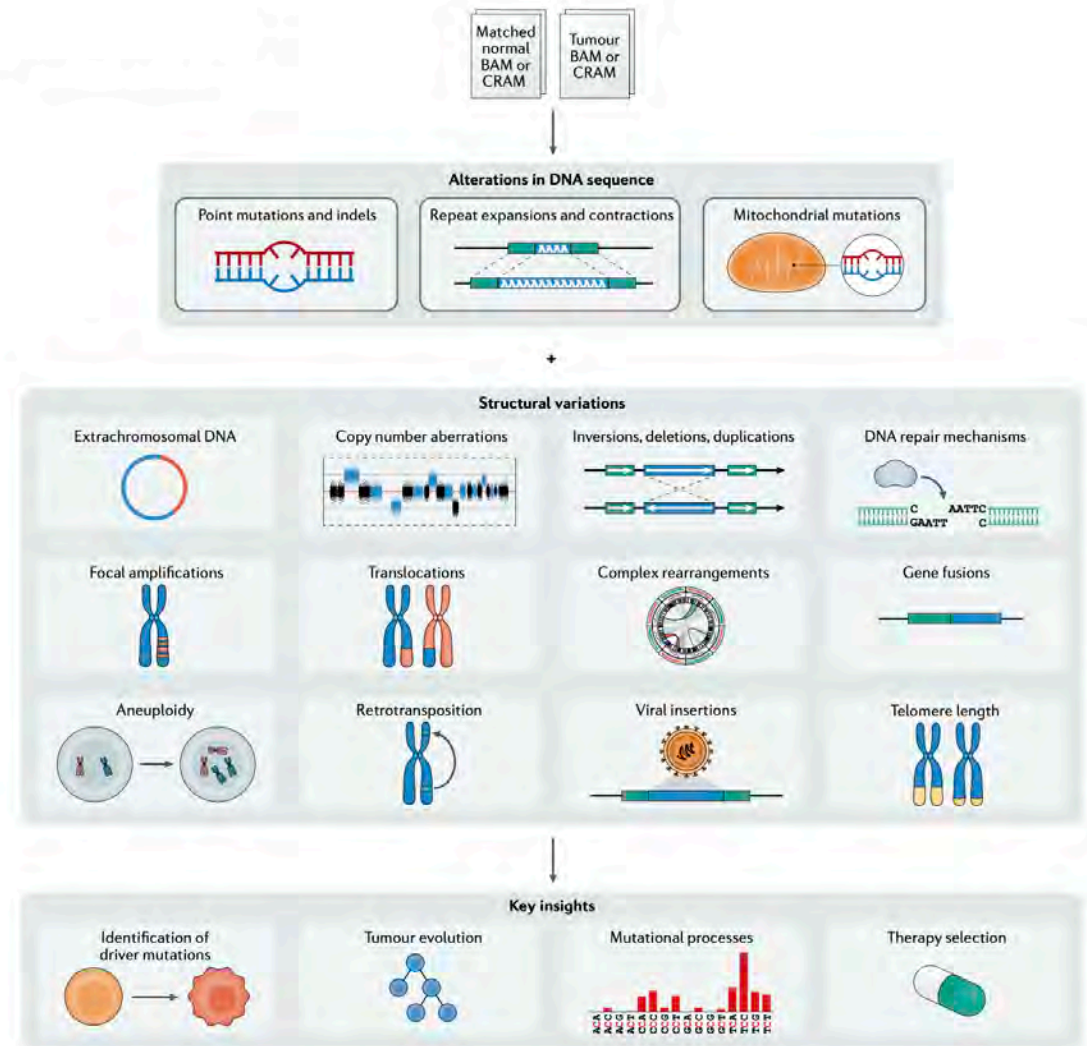
Common

Driver (~0.1%)

Rare* (1-4%)

Rare (~75%)

\* Variants with allele frequency < 0.5% are considered as rare variants in 1000 genomes project.

EMBL

# Tumor-only sequencing considerations

- Problem: A relatively small number of somatic variants is "hidden" in a large set of germline variants (4-5 Million)

- Study design options
  - Using an unmatched control
    - Rare variants appear as "somatic"
  - Using a panel of normal genomes
    - Fewer rare variants appear as "somatic"
  - Using germline variant catalogs (1000 Genomes, ExAC, gnomAD)
    - Catalogs are often highly curated, i.e., false positive variant calls from your analysis are likely not present and thus, still called as "somatic"

- Tumor-mutational burden (TMB) is often >2-fold over-estimated using tumor-only sequencing
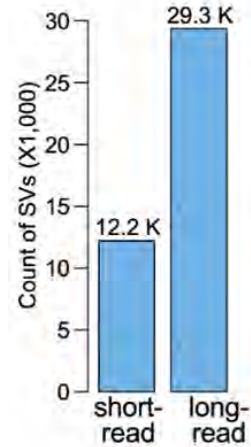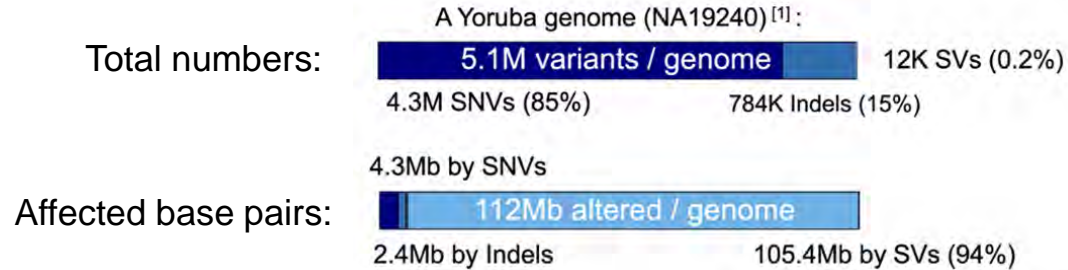
EMBL

# Cancer Genome Data Analysis

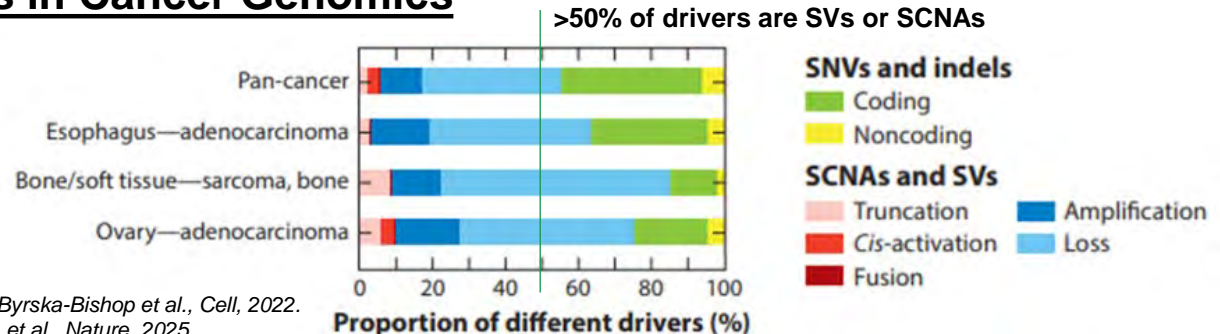# Structural and copy-number variants

# Structural Variants in Numbers

## Germline structural variants in Population Genomics



Total numbers:

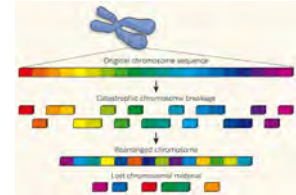A Yoruba genome (NA19240) [1]:

| 5.1M variants / genome | 12K SVs (0.2%) |

4.3M SNVs (85%)    784K Indels (15%)

4.3Mb by SNVs

Affected base pairs:

| 112Mb altered / genome |

2.4Mb by Indels    105.4Mb by SVs (94%)



## Somatic structural variants in Cancer Genomics



>50% of drivers are SVs or SCNAs



**SNVs and indels**
- Coding
- Noncoding

**SCNAs and SVs**
- Truncation
- Cis-activation
- Fusion
- Amplification
- Loss

*Sources: Pan-cancer analysis of whole genomes, Nature, 2020. Byrska-Bishop et al., Cell, 2022.*
*Cosenza et al., Annu Rev Genomics Hum Genet. 2022. Logsdon et al., Nature, 2025.*

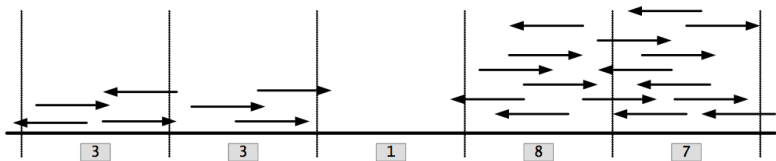# Wide Range of Chromosomal Aberrations in Cancer
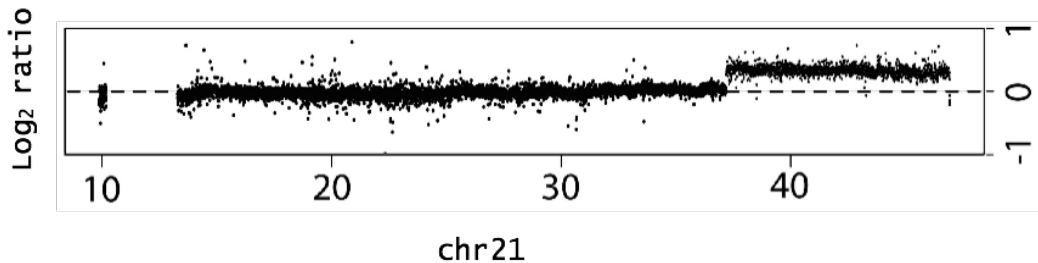
# Somatic structural variant calling

## Read-depth

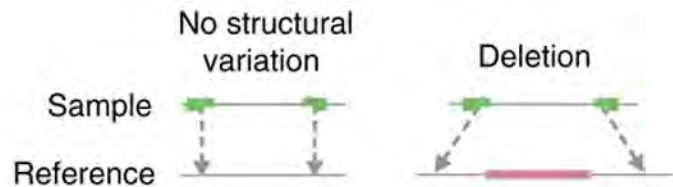- Read counting in windows for tumor and normal data



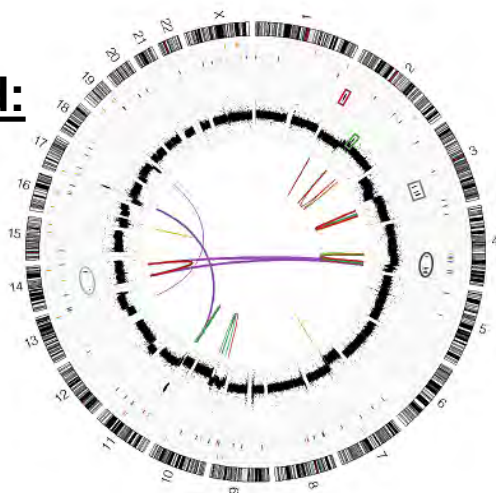- Log2 ratio for each window
- Chromosome-wide plot

$$\log_2 \frac{\#\text{Reads}_{Disease}}{\#\text{Reads}_{Normal}}$$



## Paired-end / Split-reads



No structural variation   Deletion
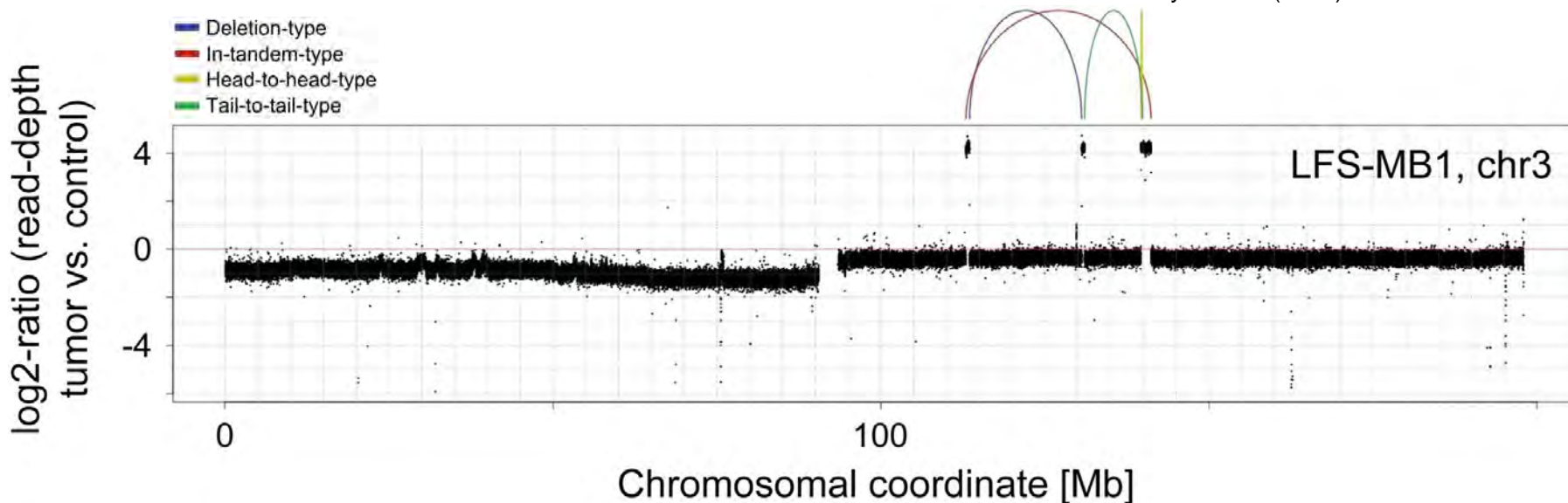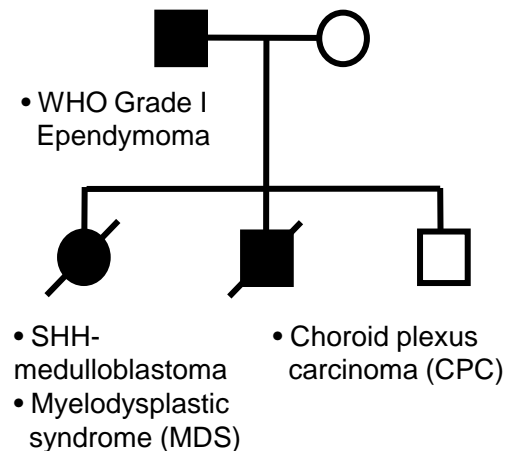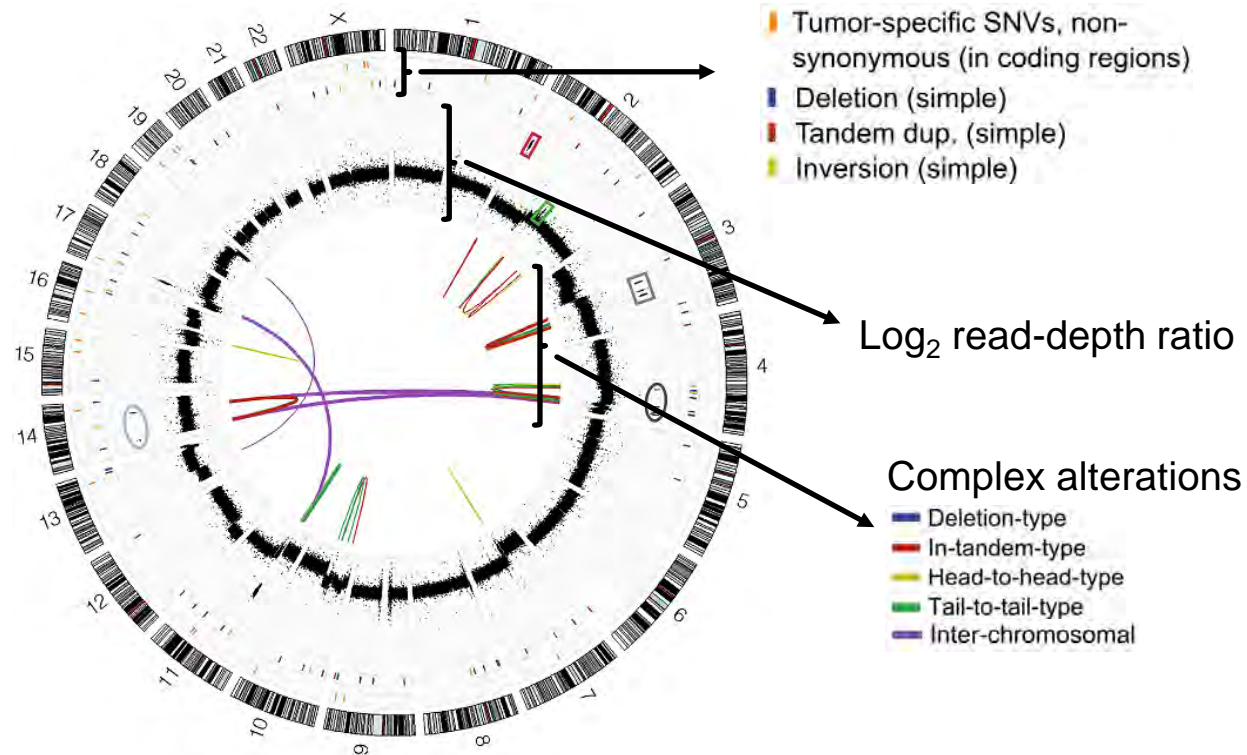
Sample

Reference

## Combined:

# Childhood Brain Tumor Medulloblastoma

- Li-Fraumeni syndrome
  - Germline TP53 mutation

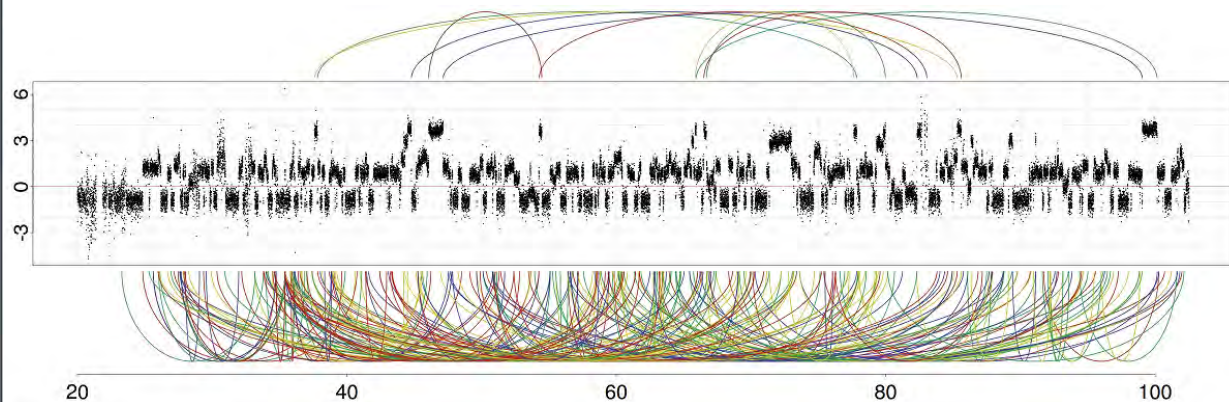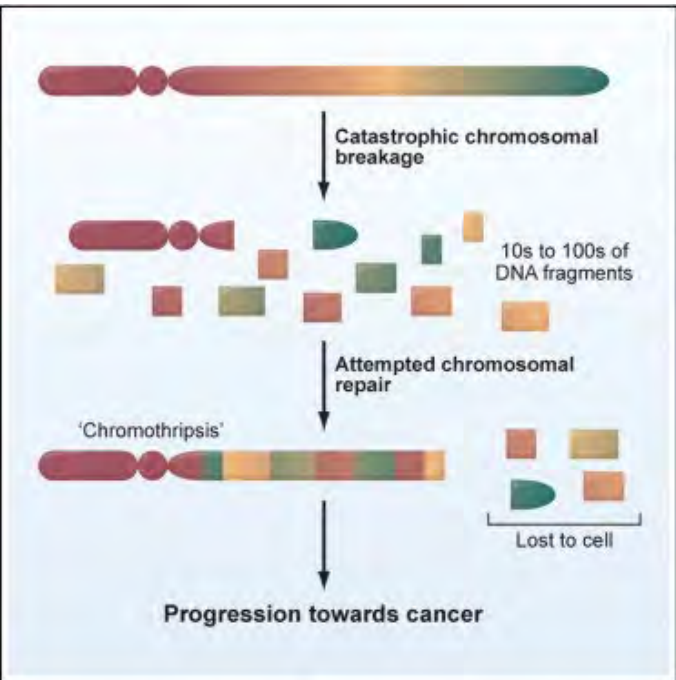# Somatic DNA alterations



Tumor-specific SNVs, non-synonymous (in coding regions)
Deletion (simple)
Tandem dup. (simple)
Inversion (simple)

Log$_2$ read-depth ratio

Complex alterations
Deletion-type
In-tandem-type
Head-to-head-type
Tail-to-tail-type
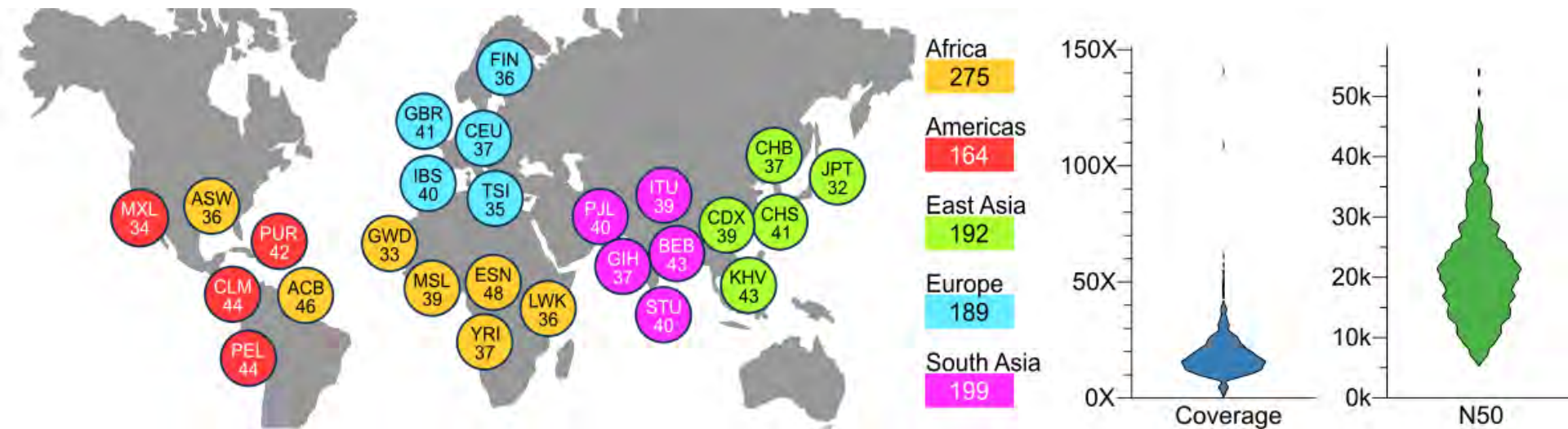Inter-chromosomal

EMBL

# Extra-chromosomal DNA (ecDNA)

# Chromothripsis

Structural variant calling using long-reads

# ONT Sequencing of 1,019 samples from the 1000 Genomes Project
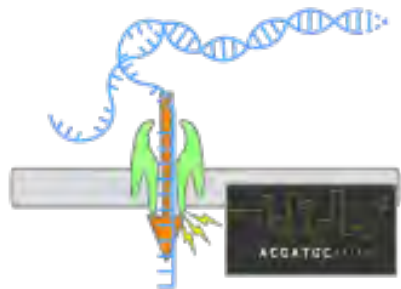


1,019 samples sequenced with ONT
- ~15x coverage
- **Structural variant calling** using long-reads

EMBL

# Pan-genomes and long-reads for SV discovery



Short-reads: 100bp-300bp

2 to 3-fold increase in detected SVs

Long-reads: N50 read length >15,000bp

Nanopore sequencing

? What have we missed with short-reads and a single reference genome?

Linear reference genome (GRCh38)

+ ~200Mbp

CHM13 (T2T)

+ ~100Mbp

Graph pan-genome

Human Pangenome Reference Consortium 44 samples

EMBL

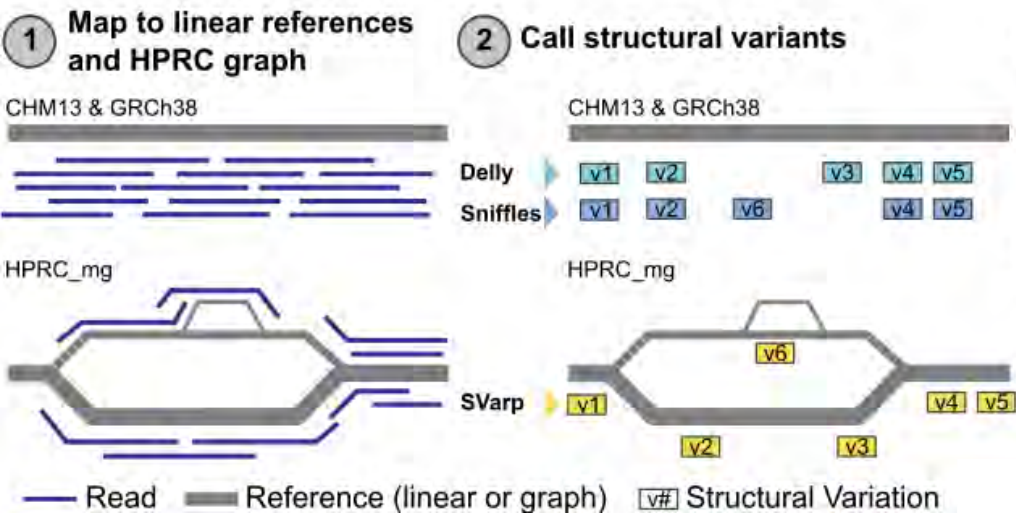# SV Analysis by Graph Augmentation (SAGA)
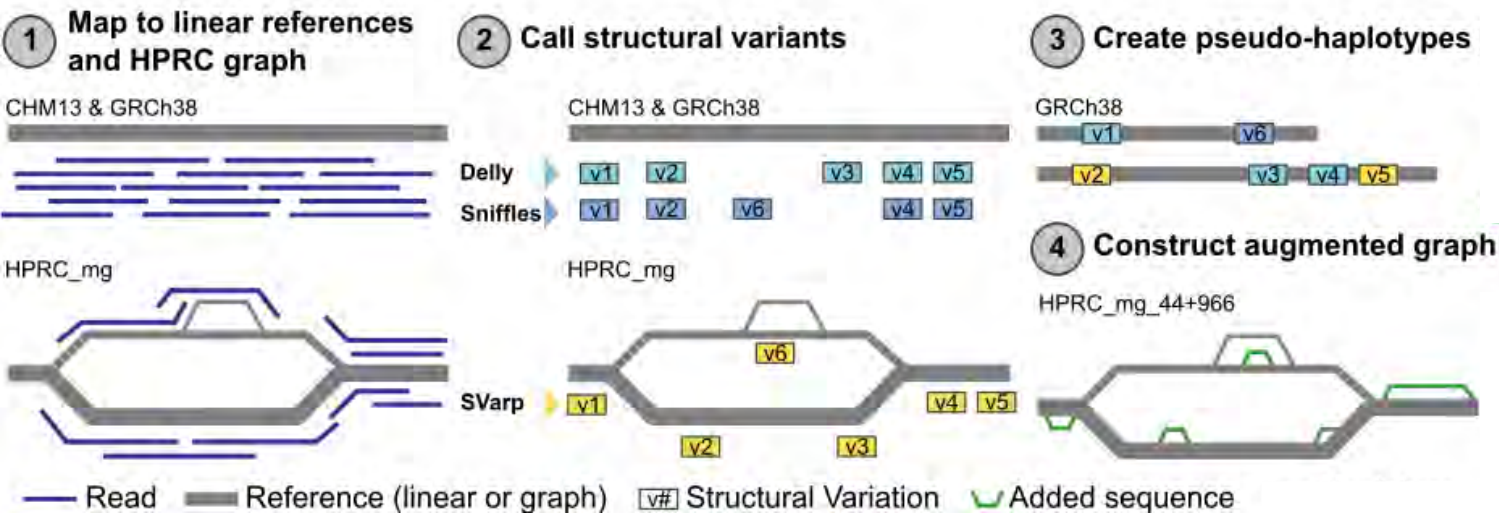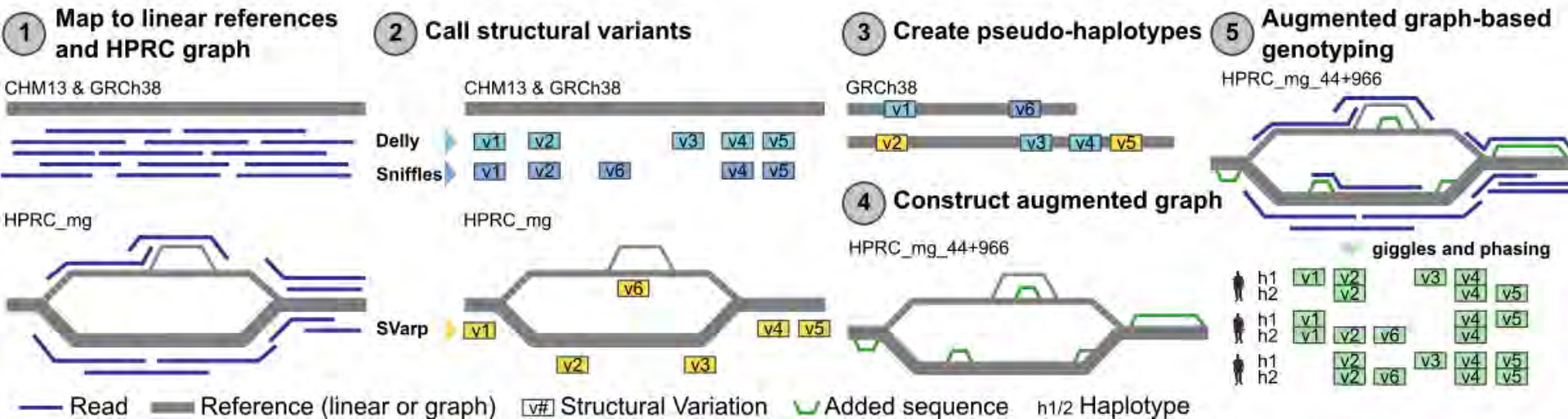
EMBL

# SAGA: SV Analysis by Graph Augmentation

# SAGA: SV Analysis by Graph Augmentation

# SAGA: SV Analysis by Graph Augmentation
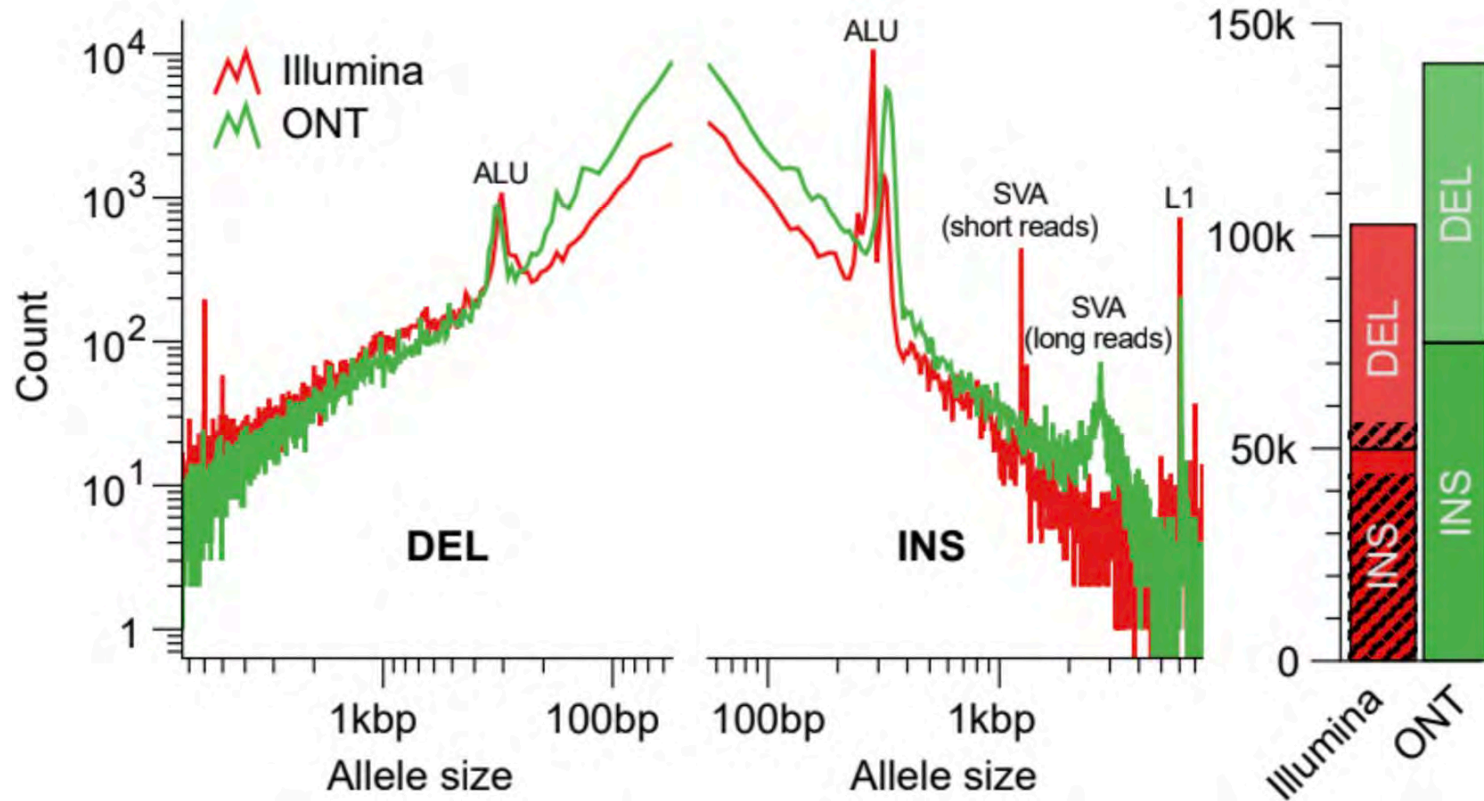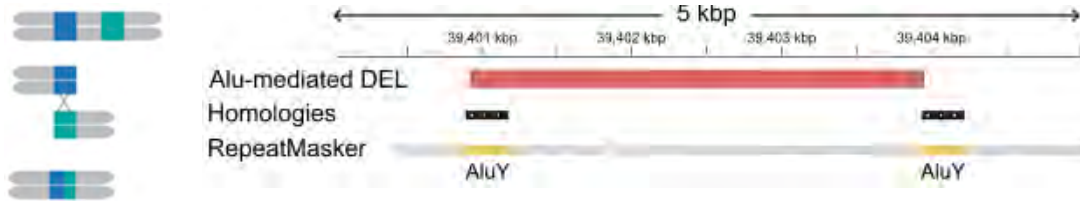
# SAGA: SV Analysis by Graph Augmentation



*Schloissnig et al., Nature, 2025.*

# Long-reads facilitate the discovery of sequence-resolved insertions



Byrska-Bishop et al. Cell 2022
for comparison

*Source: Schloissnig et al., Nature, 2025.*
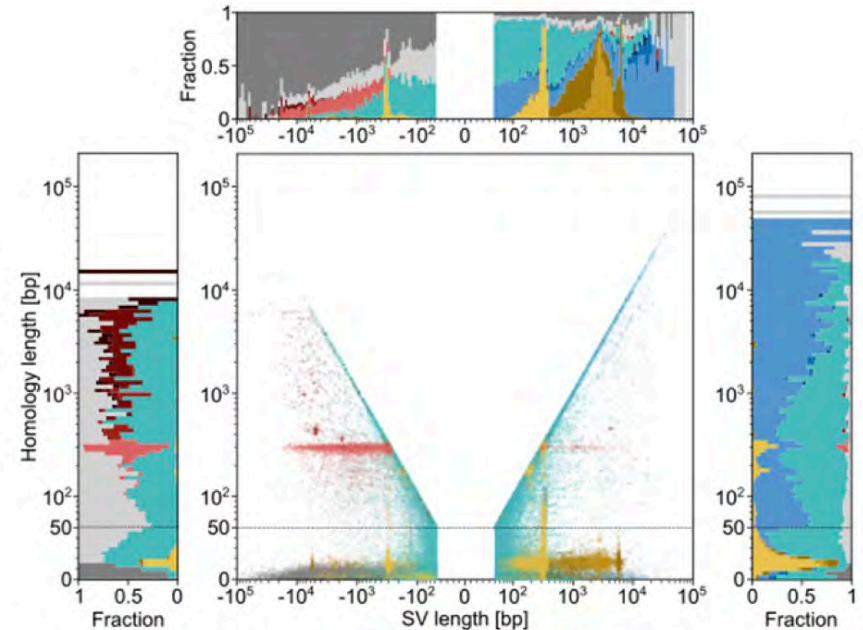
EMBL

# Repeat mediated SVs

- NAHR – Non-allelic homologous recombination

# Cancer Genomics using long-reads
# Deciphering haplotype-resolved complex rearrangements



Source: Rausch et al., Cell Genomics, 2023.

# Haplotype-resolved genome analysis



Heterozygous alleles



- Applications
  - Analyze compound heterozygotes in rare diseases
  - Measure allele-specific expression, methylation, TF binding, etc.
  - Determine how combinations of variants uniquely situated on each haplotype may affect gene function

EMBL

# Allele-specific methylation (ASM)



PTCH1

- Het. deletion in the promoter of *PTCH1*
- Methylated in the relapsed tumor

EMBL

# Long-read WGS in bladder cancer patients

- Multi-omics and spatial analyses

# Frequent somatic LINE-1 (L1) insertions

- Up to >500 somatic L1 insertions per tumor
- Evidence in some samples for L1 multi-jumps



Unmethylated L1 promoter

EMBL

# Somatic L1s are linked with downstream genomic rearrangements and chromosomal instability
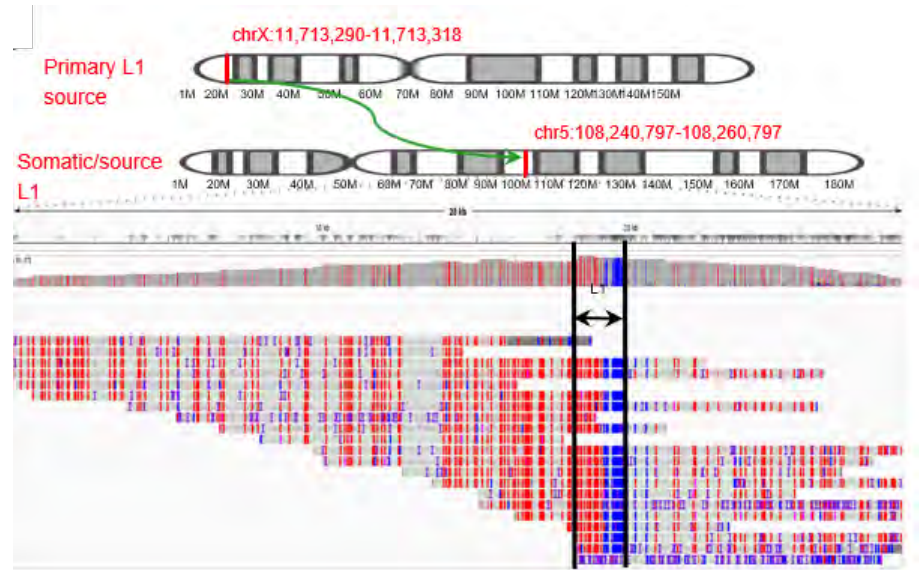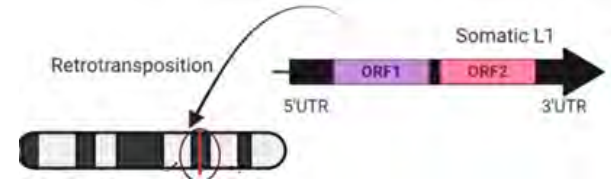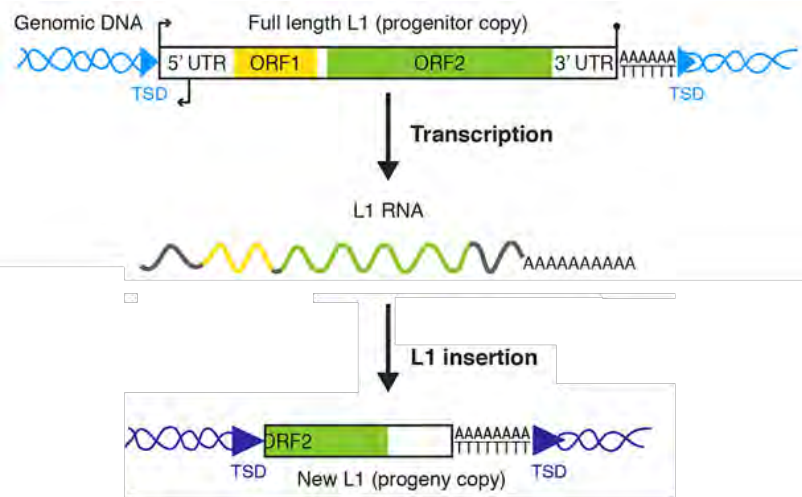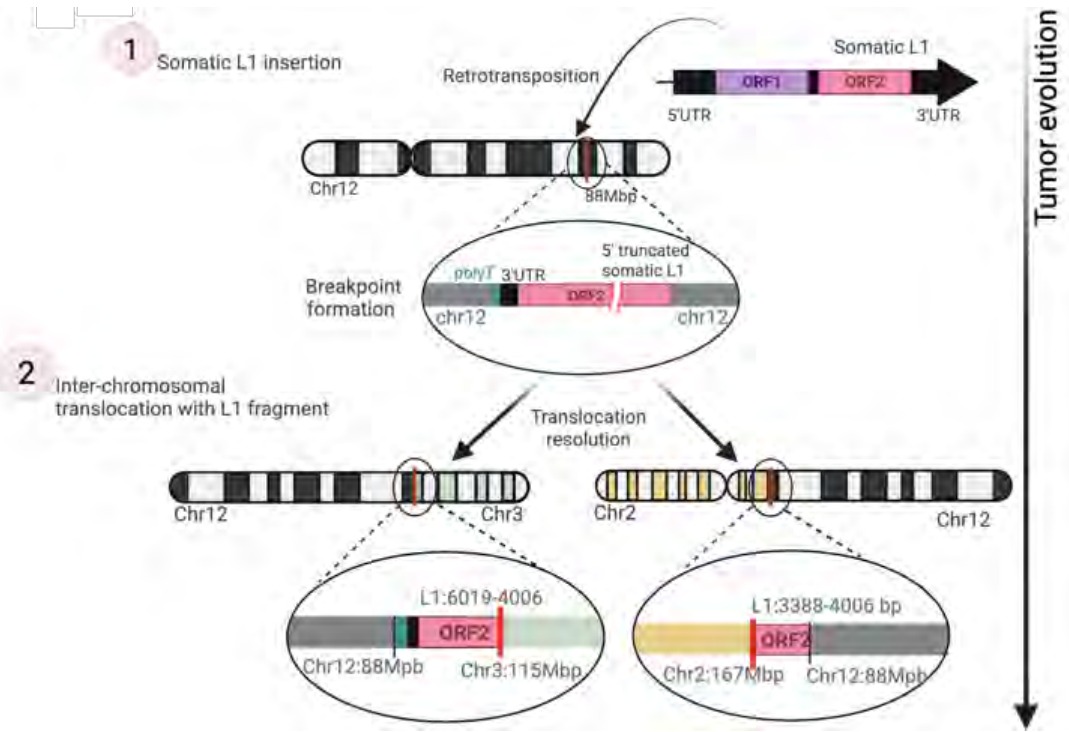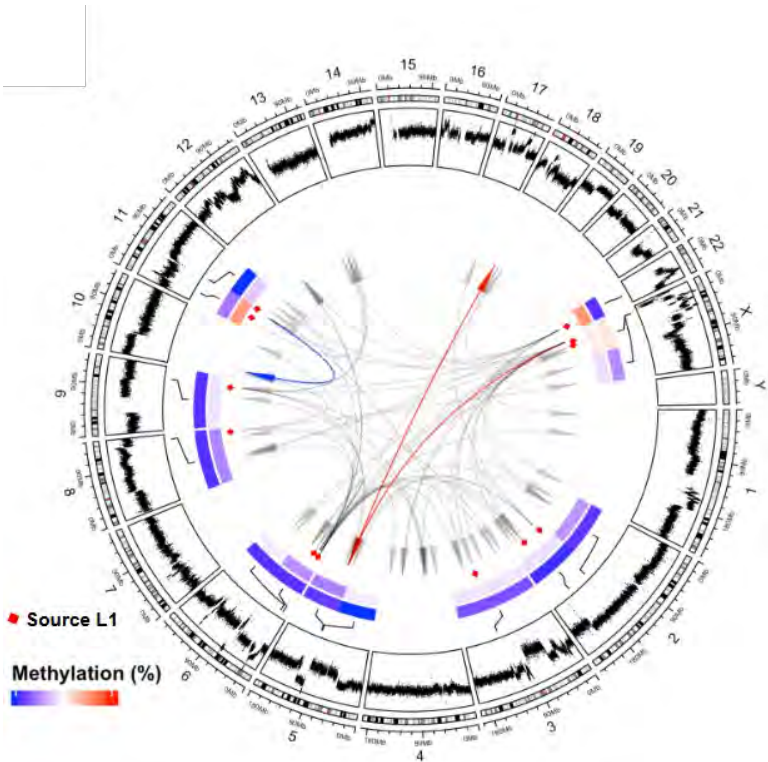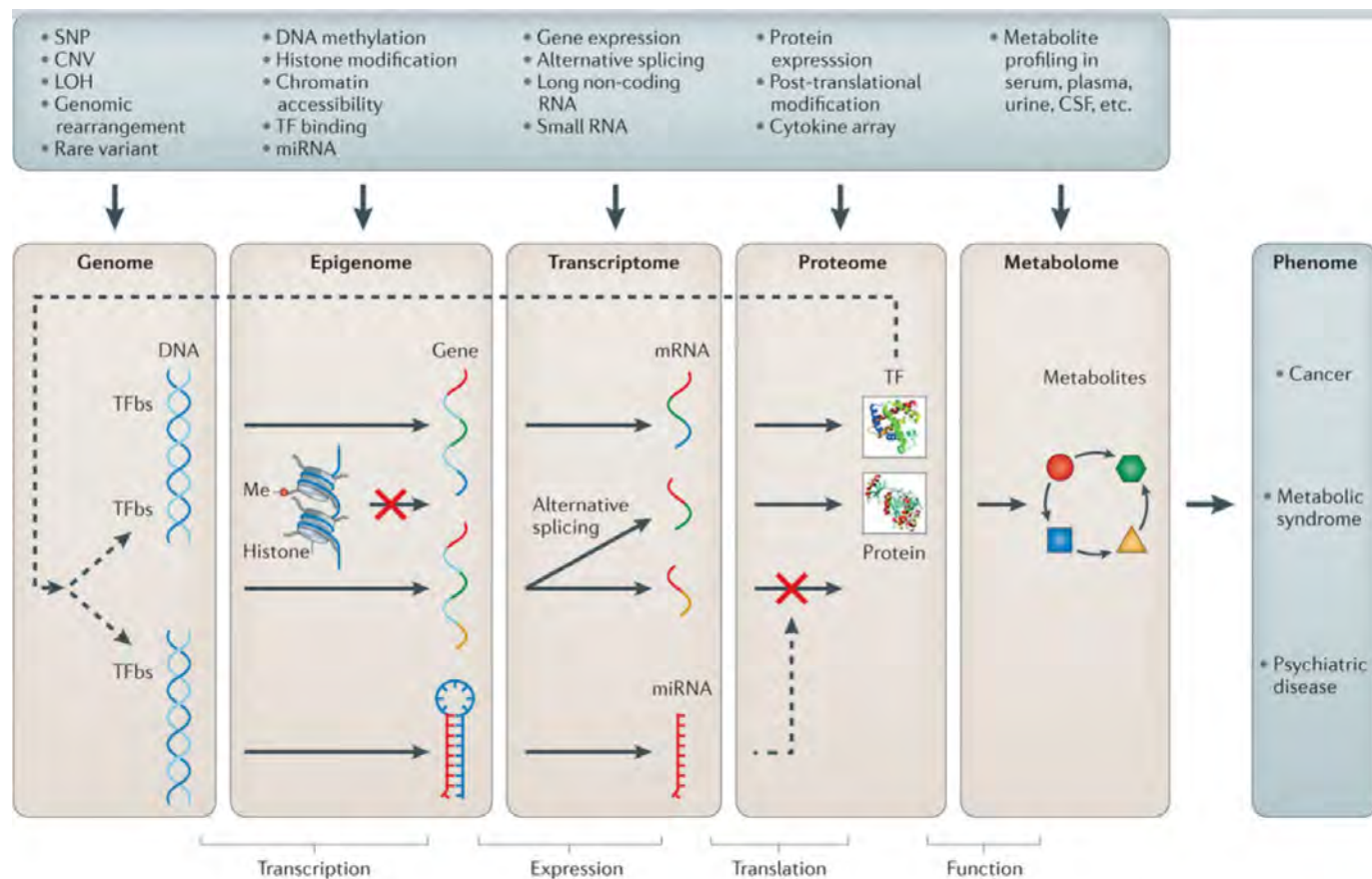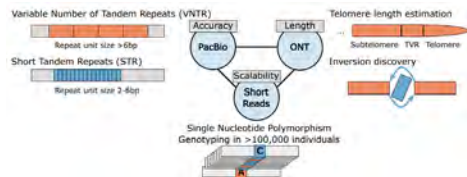
# Beyond somatic driver variants



Nature Reviews | Genetics

# Thank you!

The impact of long-read sequencing on human population-scale genomics
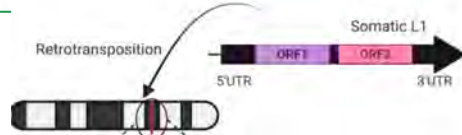
Tobias Rausch [1], Tobias Marschall [2,3], Jan O Korbel [1]

Variable Number of Tandem Repeats (VNTR)
Repeat unit size >6bp

Short Tandem Repeats (STR)
Repeat unit size 2-6bp

Accuracy
PacBio
Length
ONT
Scalability
Short Reads

Telomere length estimation
Subtelomere TVR Telomere

Inversion discovery

Single Nucleotide Polymorphism Genotyping in >100,000 individuals

Structural variation in 1,019 diverse humans based on long-read sequencing

Multi-allelic bubble
Bi-allelic bubble
Segment
Link
Path
Repeat Variation
...ACGAACTGACA...
Alu
Graph Construction
...ACGAACT GACA...
Augmentation
Alignment
SV Allele Disovery
= Alu

Image: Joana Carvalho/EMBL

Integrative spatial and multi-omic profiling in bladder cancer links L1 retrotransposition to extrachromosomal DNA, genomic instability, and viral mimicry response

Sophia J. Pribus, Ivana Osredek, Jan Otonicar, Milena Simovic-Lorenz, Michael Scherer, Sergio Manzano-Sanchez, Andreas Kienzle, Urja Parekh, Vladimir Benes, Pooja Sant, Philipp Mallm, Karsten Brand, Angelika B. Riemer, Holger Sültmann, Christoph Plass, Mladen Stankovic, Jan O. Korbel, Tobias Rausch, Aurélie Ernst

Retrotransposition
Somatic L1
ORF1  ORF2
5'UTR  3'UTR

Chromosome 4
Chromosome 5
Chromosome 7
Templated Insertion Thread

Long-read sequencing of diagnosis and post-therapy medulloblastoma reveals complex rearrangement patterns and epigenetic signatures

GeneCore  EMBL