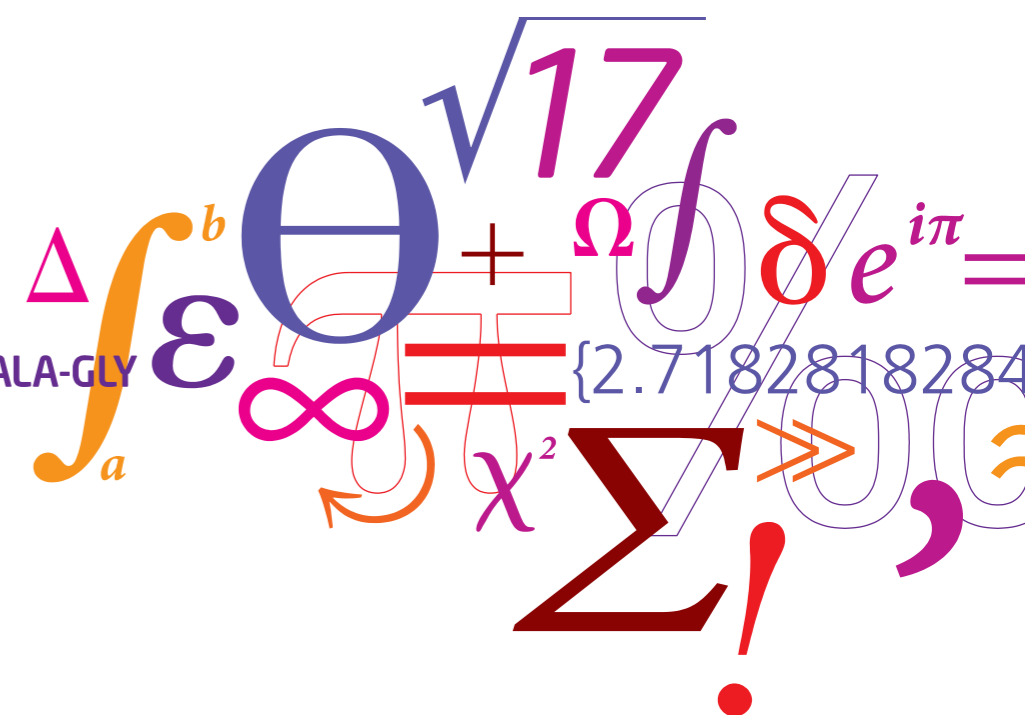# Understanding cancer genomics

## Adrian Otamendi Laspiur, Research Assistant (iCOPE)
## Original slides: Jose MG Izarzugaza

GCTGGT  > GCUGGU > ALA-GLY
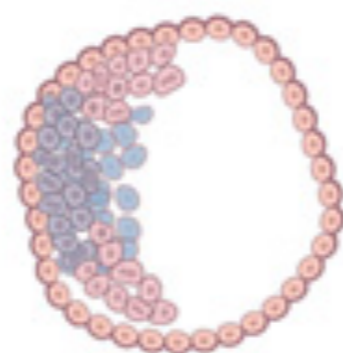CGACCA

# What is cancer?

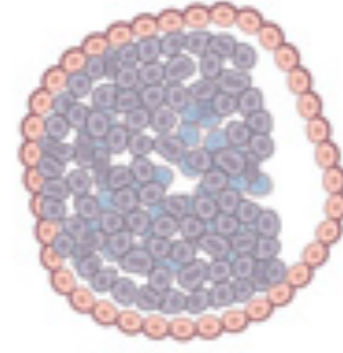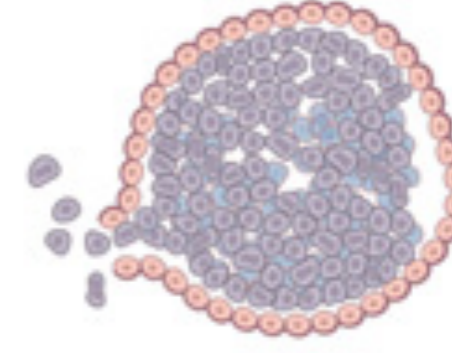**The disease caused by an uncontrolled division of *abnormal* cells in a part of the body**



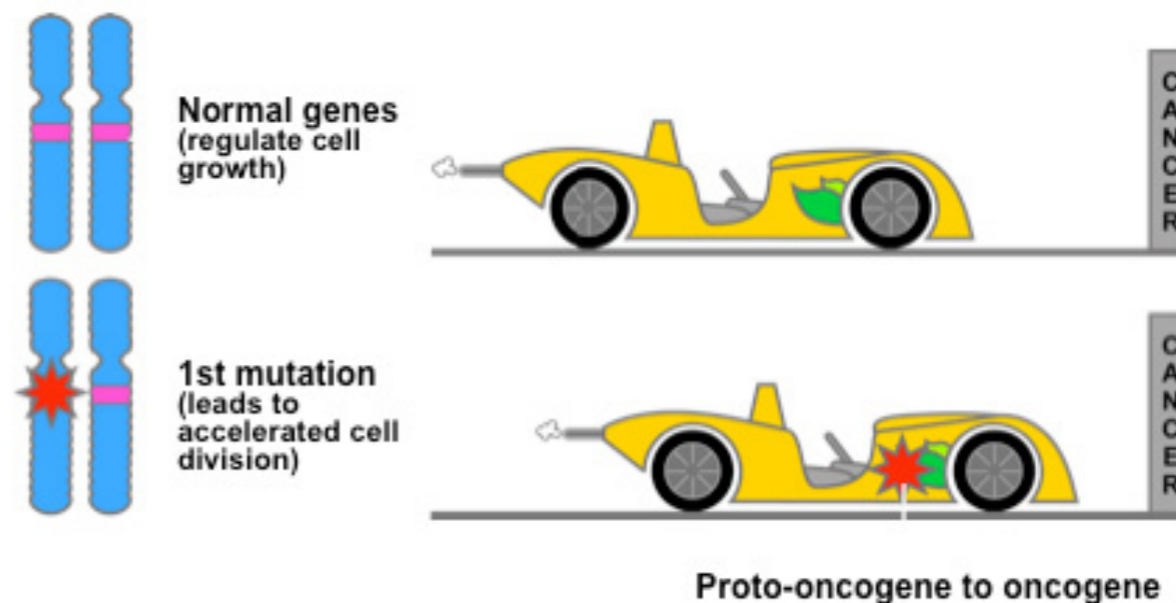Normal　　Hyperplasia　　Atypical hyperplasia　　Carcinoma in situ　　Microinvasive

**Oncogenes**:
- Mutated proto-oncogenes
- Turn abnormal cell growth on
- 70 protooncogenes
- gain of function genes
- primarily somatic activated
- [throttle pedal in a car]



Normal genes (regulate cell growth)

1st mutation (leads to accelerated cell division)

Proto-oncogene to oncogene

"Oncogenes are mutated genes whose PRESENCE can stimulate the development of cancer"

Examples: HER-2/neu. RAS, MYC, SRC, hTERT

RAS, MYC, SRC are protein kinases ➔ Cell cycle regulation

*(Adapted from Nat. Cancer Inst.)*

**Tumour suppressor genes**:
- Stop the cell cycle, G1 phase
- Slow the cell cycle before S phase
- Can induce apoptosis

- primarily somatic de-activated
- loss-of-function mutations
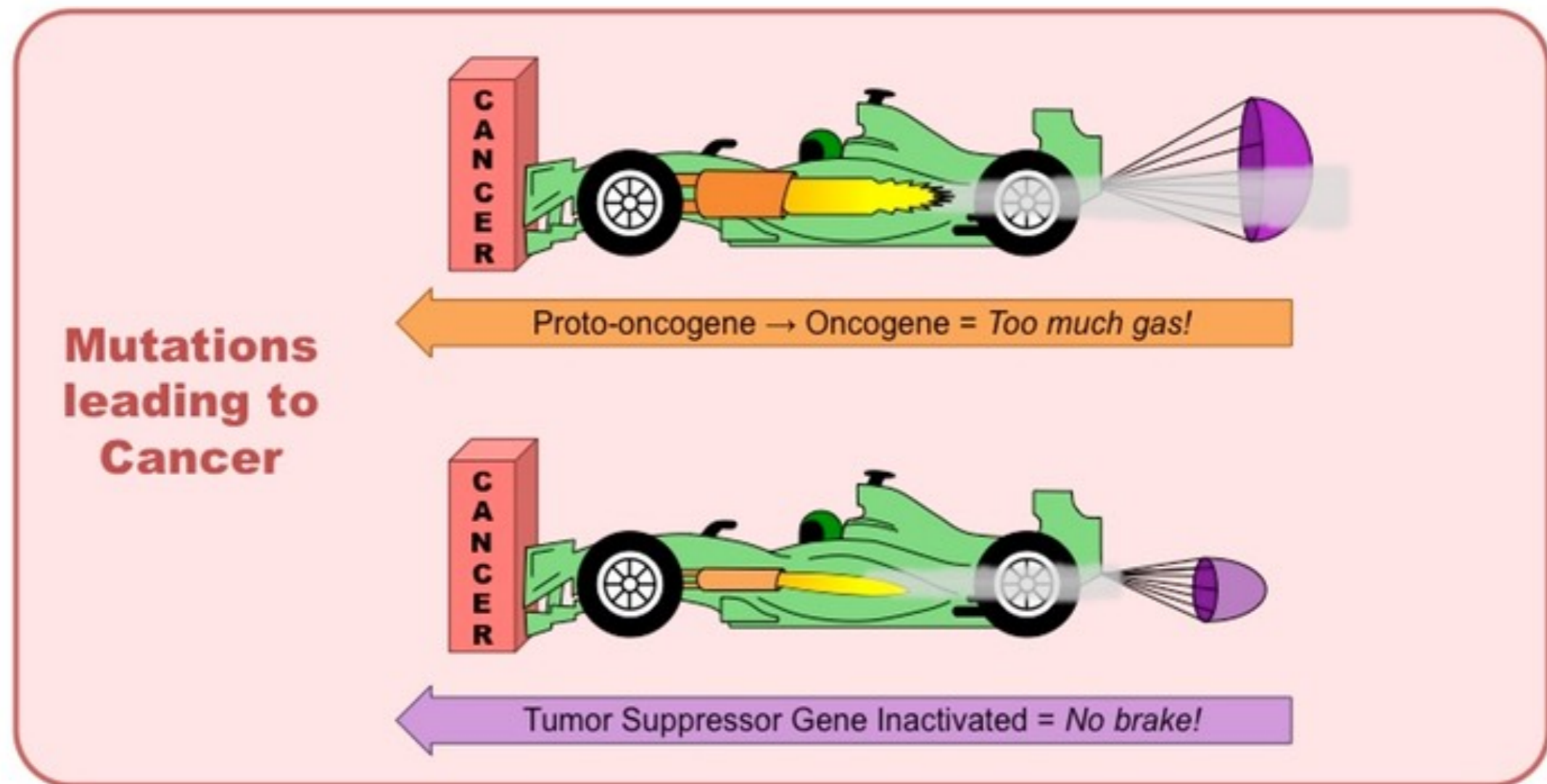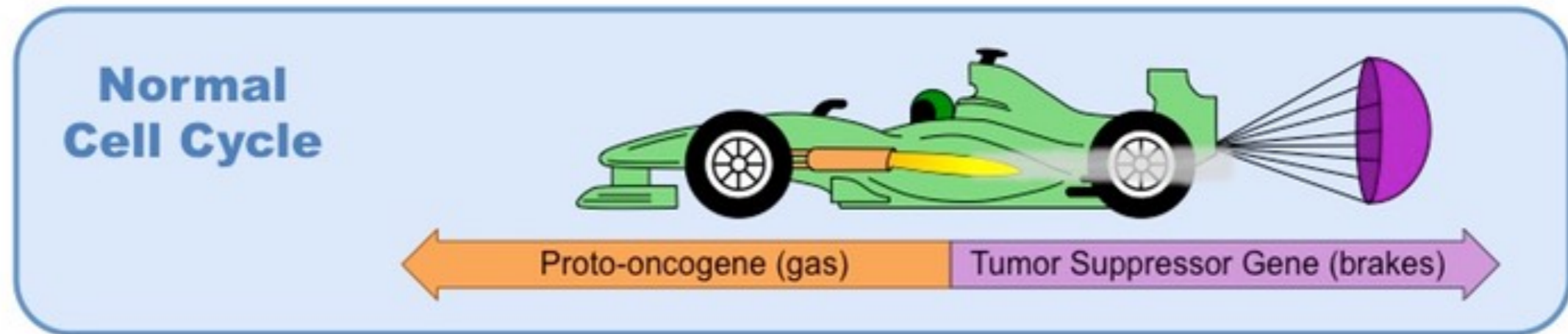- [brake pedal in a car]



"Tumour suppressors are normal genes whose ABSENCE can stimulate the development of cancer"

Examples: p53, Rb, APC

Sometimes, a single functional copy (heterozygous) is enough to prevent cancer

*(Adapted from Nat. Cancer Inst.)*

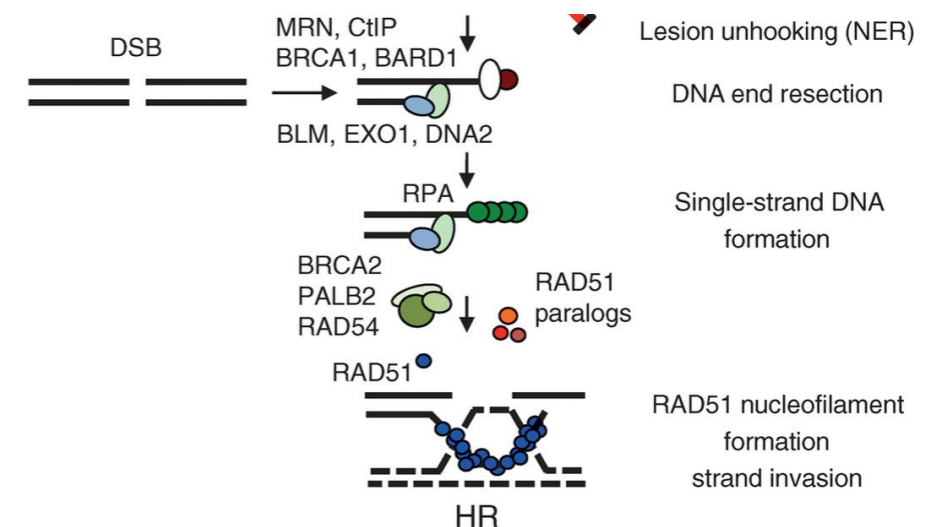# Need for speed: Oncogenes vs Tumour suppressors
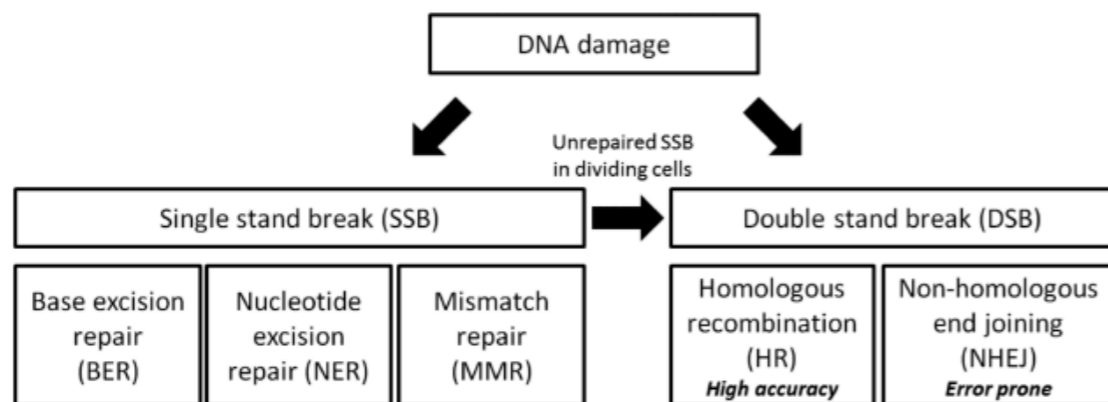


*(ib.bioninja.com.au)*

**DNA damage repair genes**
- Correct damage during DNA duplication
- Active in cell cycle, primarily G2
- After DNA replication, before Chr divides

- loss-of-function mutations ➔ increased mutation burden
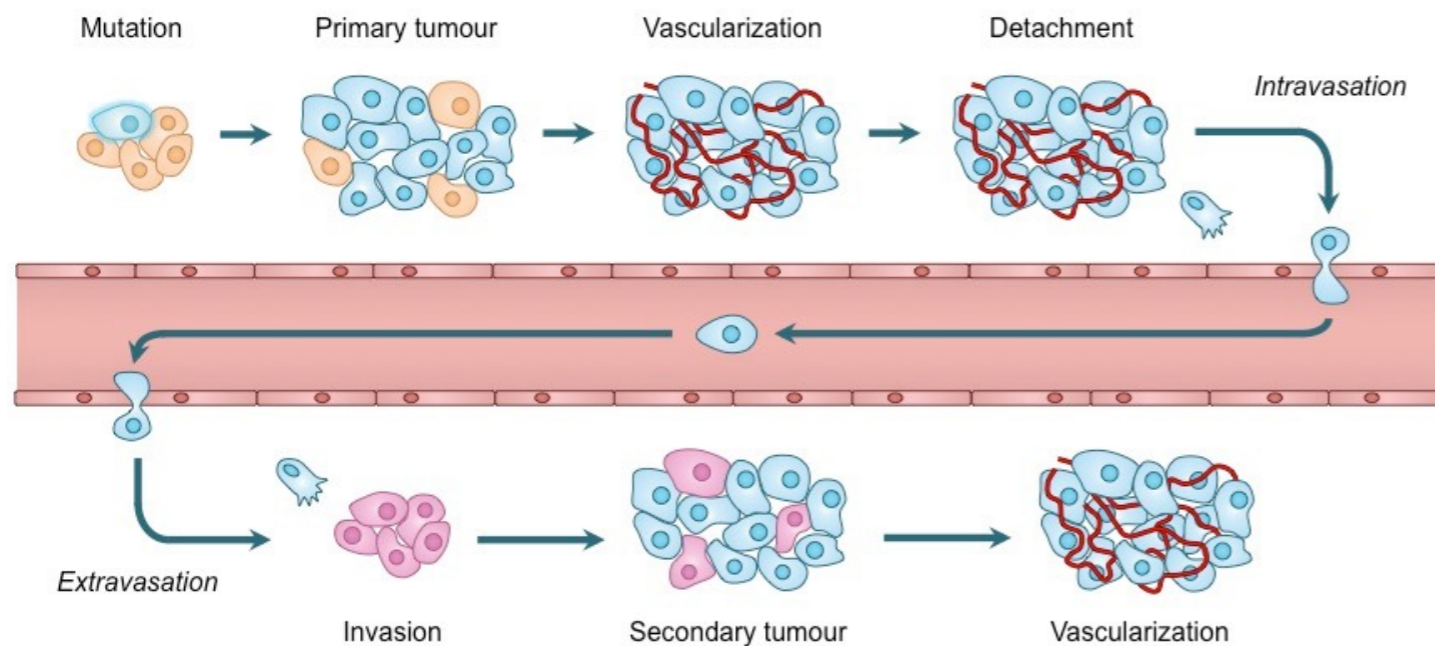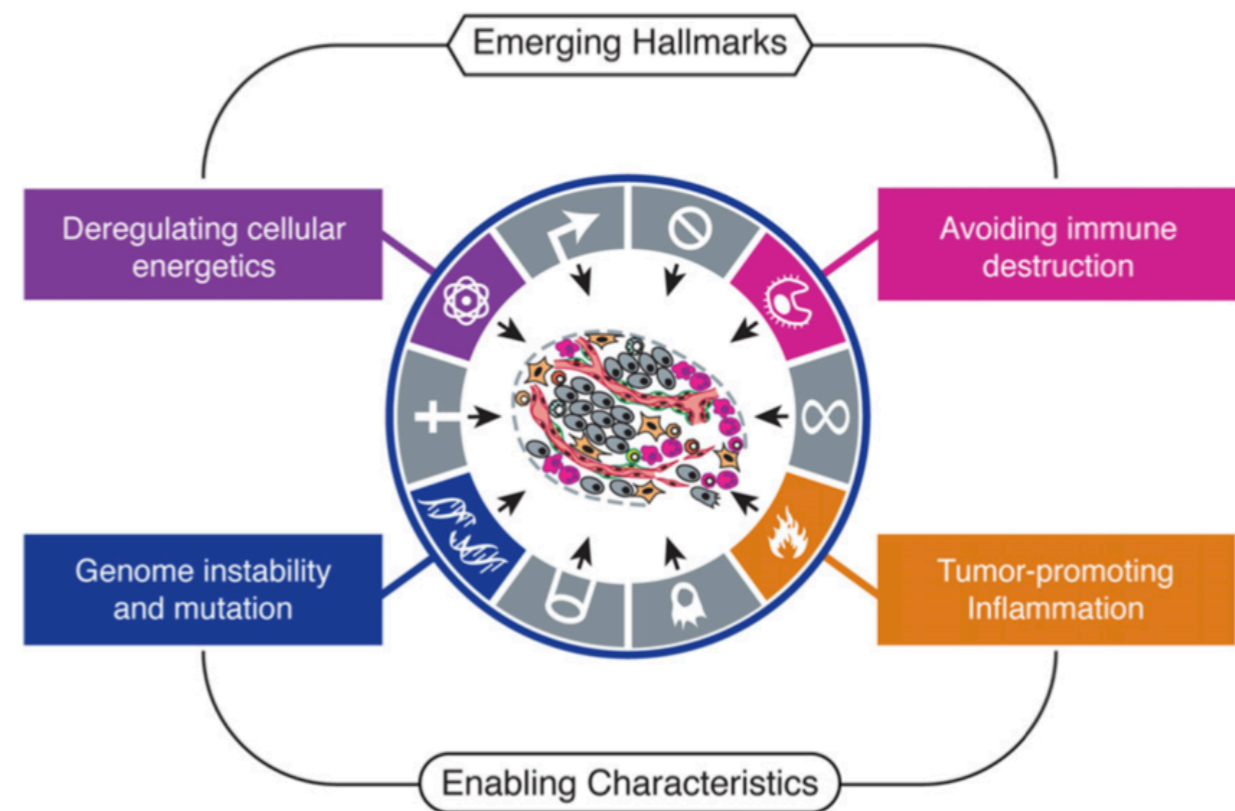Examples: BRCA1 and BRCA2 in breast cancer
          Also, mDDRG in hereditary colon cancer





*(Adapted from Nat. Cancer Inst.)*

# The hallmarks of cancer
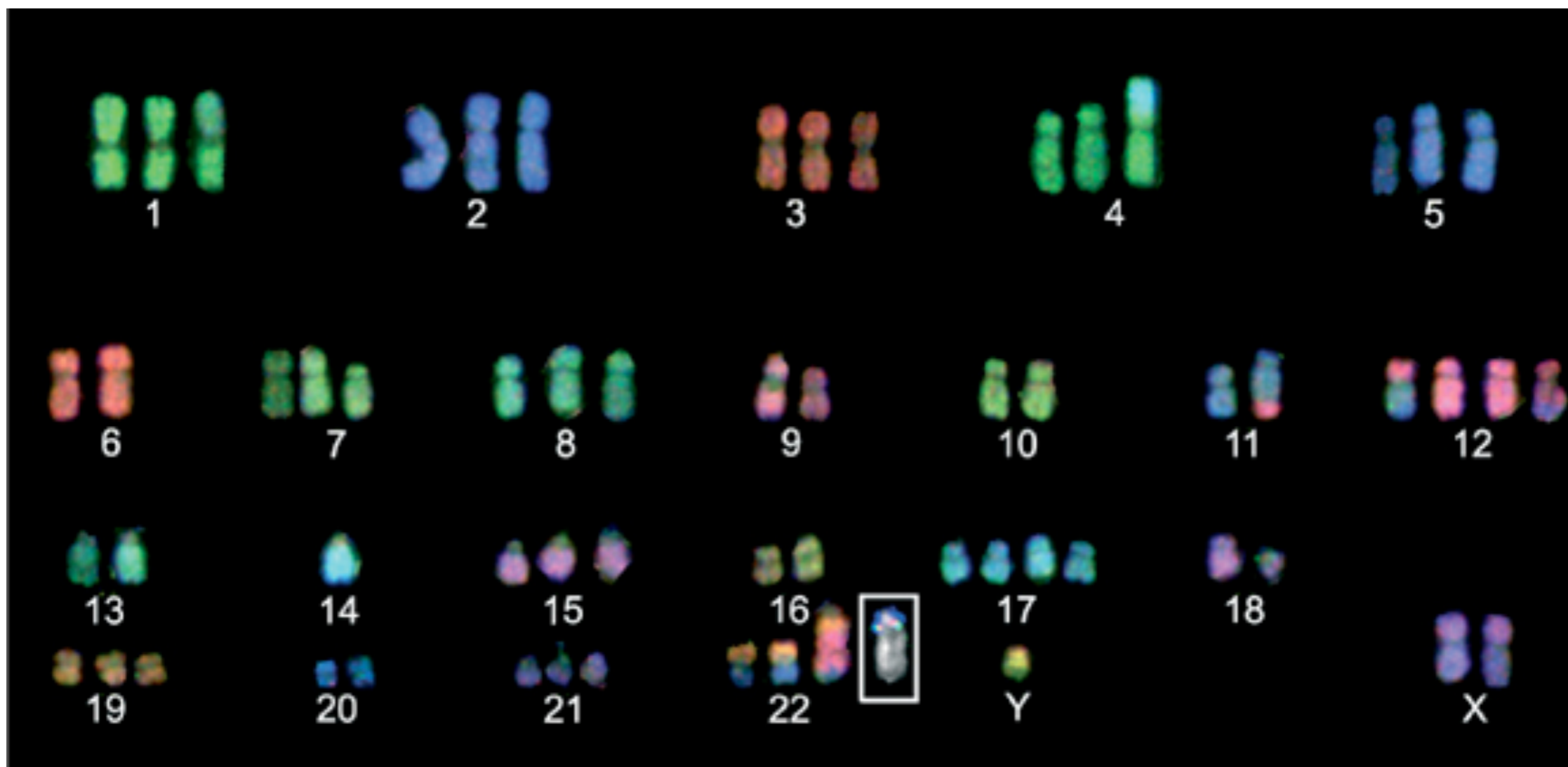
**Acquire functional capabilities**

- Sustaining proliferative signaling

- Evading growth suppressors

- Resisting cell death

- Enabling replicative immortality

- Inducing angiogenesis

- Activating invasion and metastasis

- Emerging Hallmarks

- Enabling characteristics



*The hallmarks o cancer*
*Hanahan and Weinberg, Cell 2011*



*(ib.bioninja.com.au)*

# What is cancer?

Cancer is a genetic disease: **chromosomal aberrations**



*Spectral karyotyping*

Chromosomal gain, loss
Translocation, inversion
Focal amplification

# What is cancer?

Cancer is a genetic disease: **point mutations**
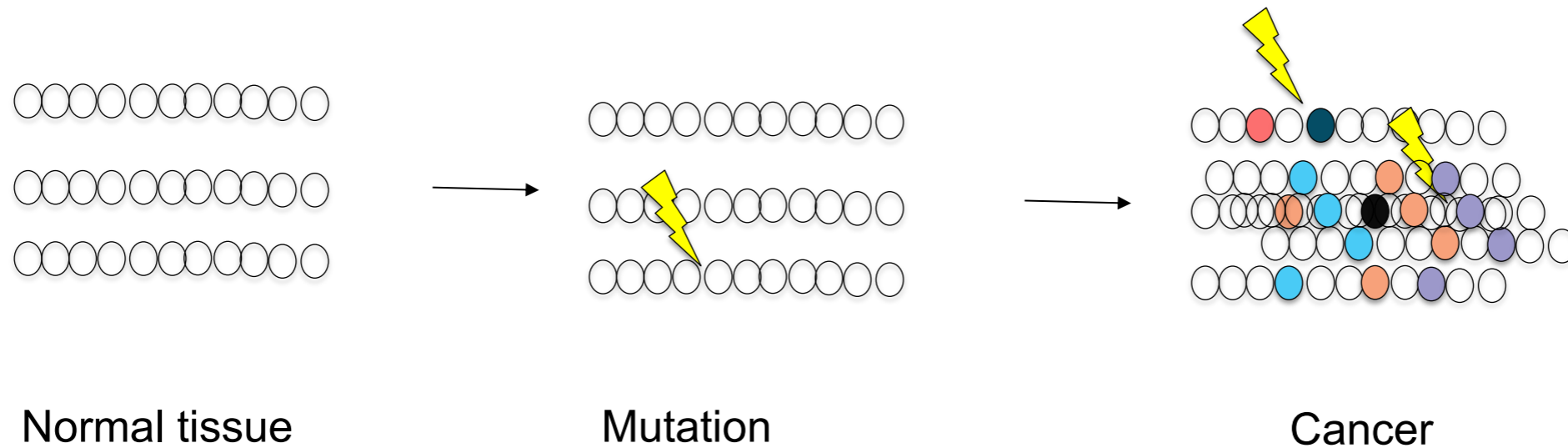
*Substitution*
*Insertion*
*Deletion*

**KRAS-wt**
```
ATGACTGAATATAAACTTGTGGTAGTTGGAGCTGGTGGCGTAGGCAAG...
-M--T--E--Y--K--L--V--V--V--G--A--G--G--V--G--K-...
```

**KRAS-G12D**
```
ATGACTGAATATAAACTTGTGGTAGTTGGAGCTGATGGCGTAGGCAAG...
-M--T--E--Y--K--L--V--V--V--G--A--D--G--V--G--K-...
```

- Frequent d**river mutation** for tumors of the lung, colon, etc.
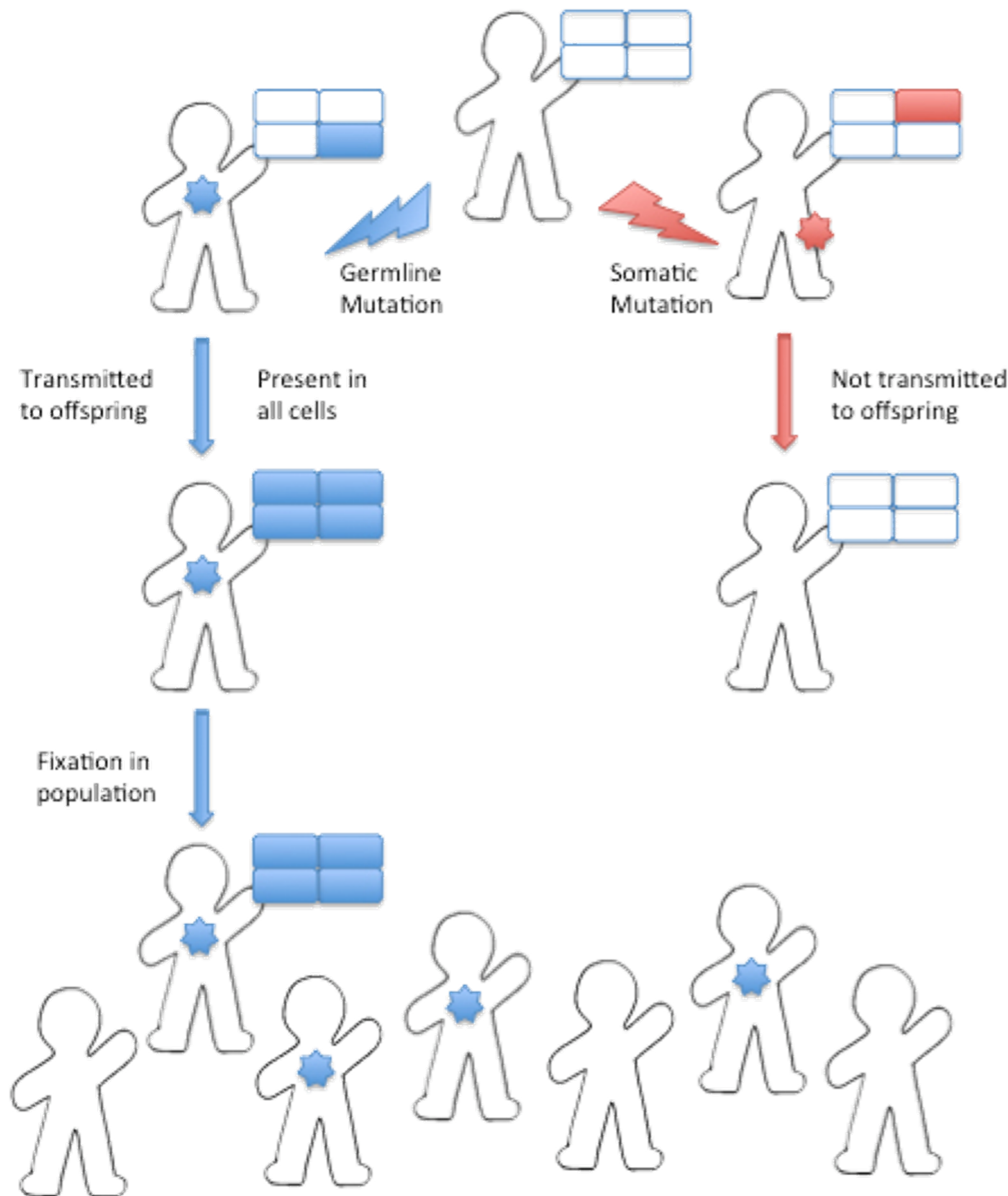- **Predicts lack of benefit** from EGFR inhibitors

# The drivers and passengers



Normal tissue          Mutation                    Cancer

⚡  Driver mutation

🔵🟠  Passenger mutations

Driver       ➔ Confers selective advantage
                Disease associated, pathogenic
Passenger ➔ Present in the clonal progenitor

# Germline vs Somatic mutations



Germline Mutation

Somatic Mutation

Transmitted to offspring

Present in all cells

Not transmitted to offspring

Fixation in population

**Germline Mutations**
Present in **all** cells
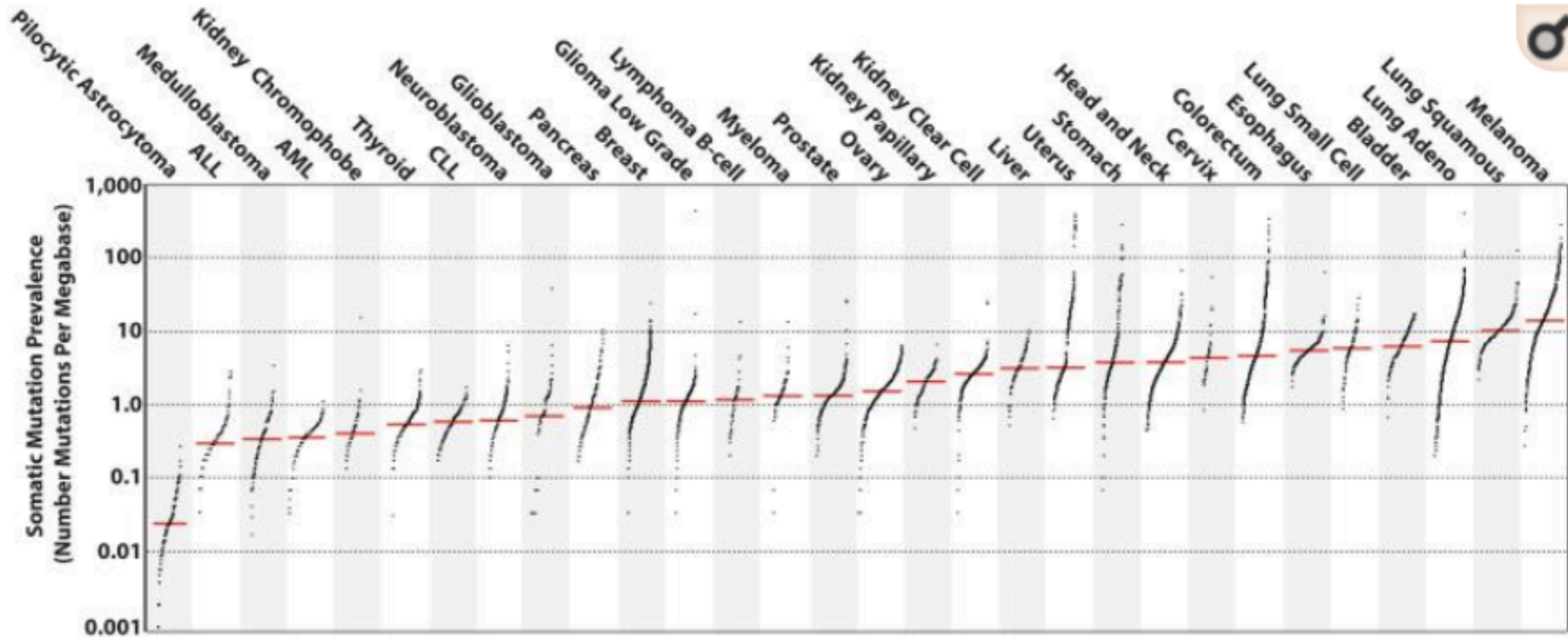**Transmitted** to offspring
**Fixate** in population (SNP)

**Somatic Mutations**
-Present only in **some** cells
-**Not transmitted** to offspring
-Do **not fixate** in population

**"Cancer is not a single disease, but rather 150+ different diseases."**

*Prof. Dr. Mariano Barbacid, former director of the Spanish National Cancer Research Center and discoverer of the first oncogene, RAS.*
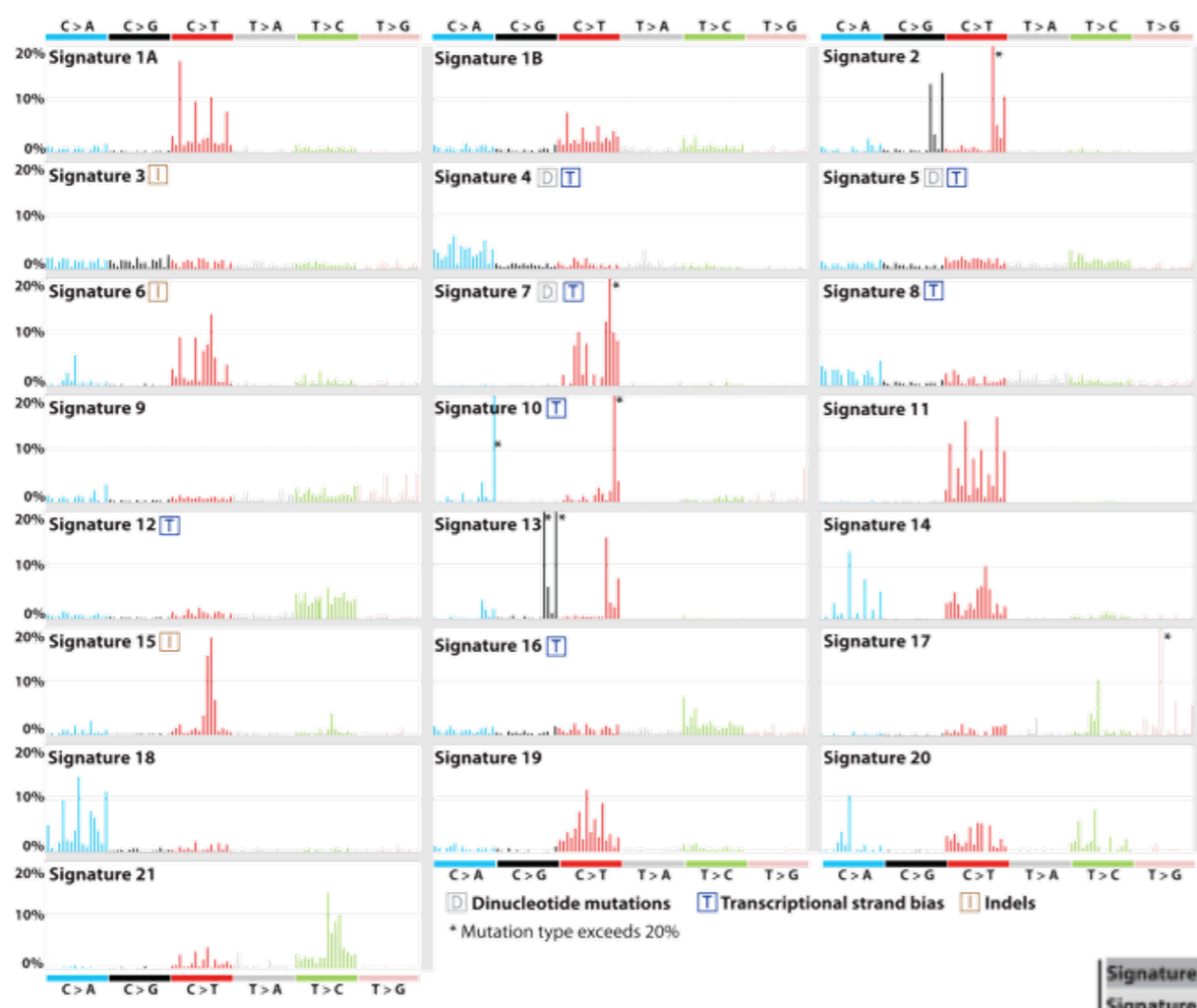
# Number of somatic mutations in different cancer types

*Nature et al., 2013*

# Mutation signatures in different cancer types



E.g.: Signature 4

- Smoking induced mutations

- Lung cancers

*Alexandrov, Nature et al., 2013*

# Why study cancer genomes?

**For the researcher**:

- Identify recurrent mutations that represent druggable targets
- Identify specific mutations or patterns that predict benefit from specific drugs
- Study the evolutionary process -- mutation, selection
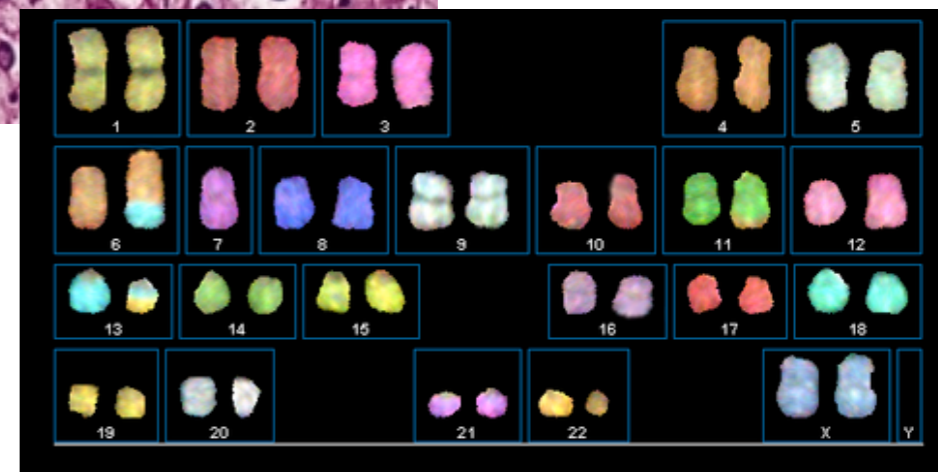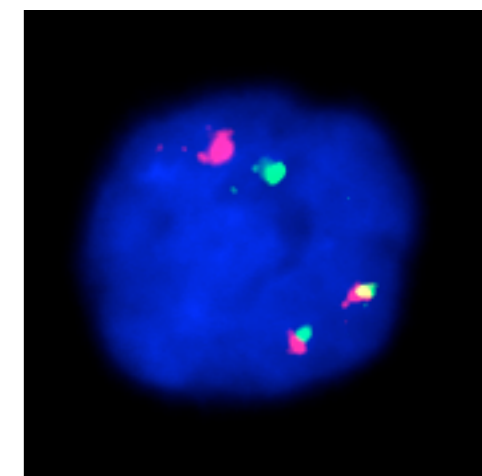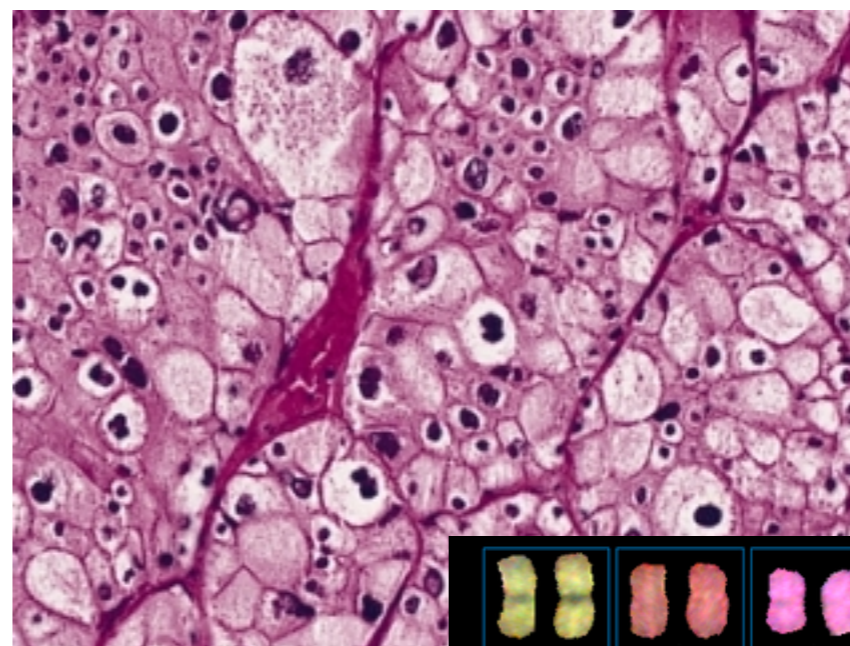
**For the cancer patient**:

Identify "actionable" mutations - inform treatment decisions
Aid in diagnosis

# Characterising a tumour specimen

**Measured in individual cells:**
- Cellular/tissue morphology
- Protein expression
- Gene copy number (FISH)
- Karyotype
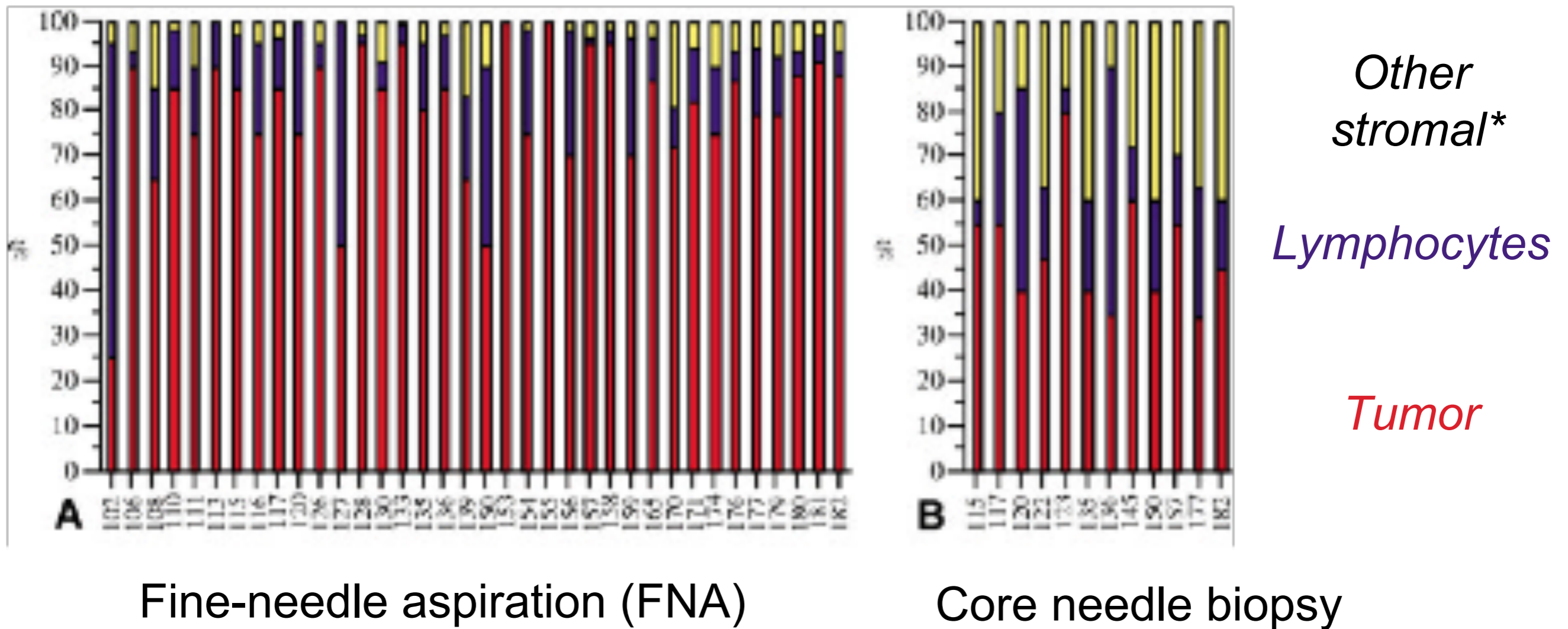
**Measured in bulk tissue (usually):**
- Gene expression (from microarray)
- Copy number profile (from SNP array)
- WES/WGS/RNA-seq

Bulk tissue includes non-tumor cells!

```
@HS21_6684:1:1306:6031:9563#14
CTTCCGATCTGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTT
+HS21_6684:1:1306:6031:9563#14
9((A=CA;2FDEEE>E=IIIIIGGGIIHEF?CEHFDIGIIGGGEGGHHHHIHFBBBBGHEIHHFDHDDFDD?@@
```

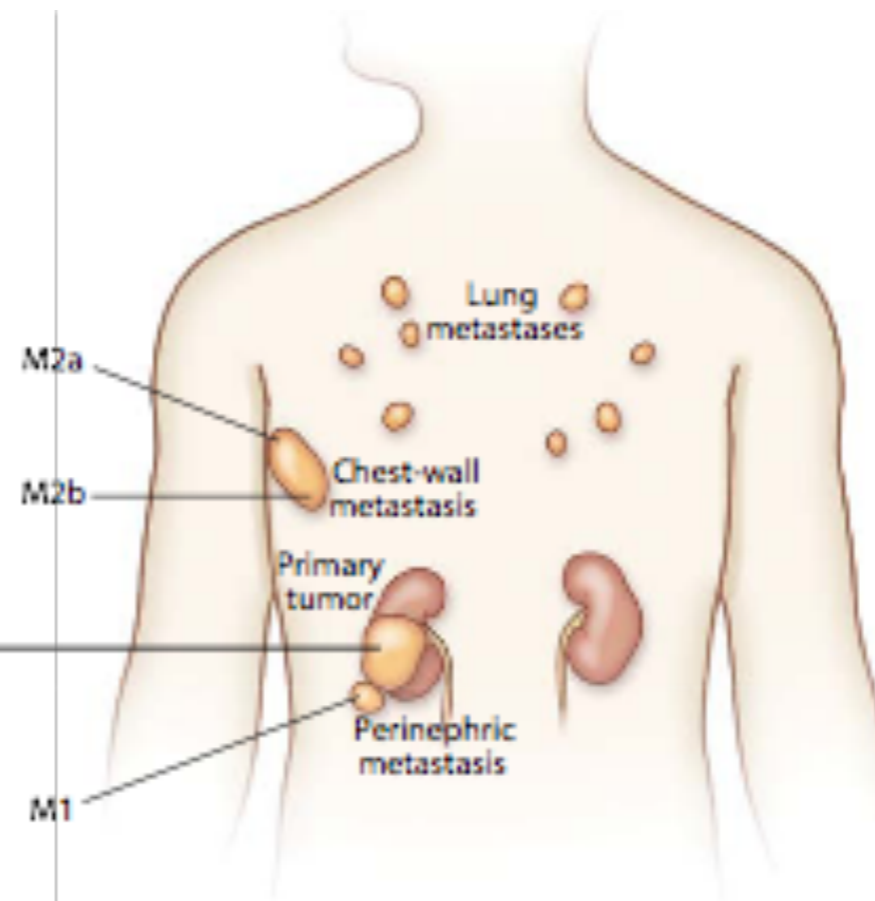# Caution 1: Tumour specimens are not just tumour cells

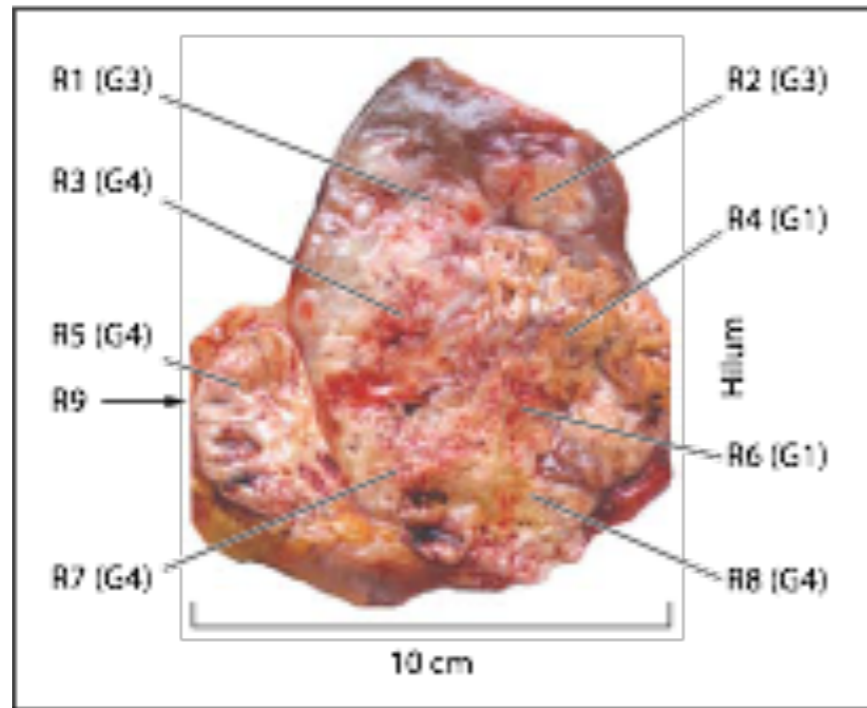Histological characterization of a set of tumor specimens



*Other stromal\**

*Lymphocytes*

*Tumor*

Fine-needle aspiration (FNA)

Core needle biopsy

\* Other stromal = fibroblasts, endothelials, histocytes, adipocytes

W. F. Symmans *et al.*, *Cancer* 97:2960 (2003)

Primary renal cell carcinoma

Gerlinger et al. (2012) *NEJM*
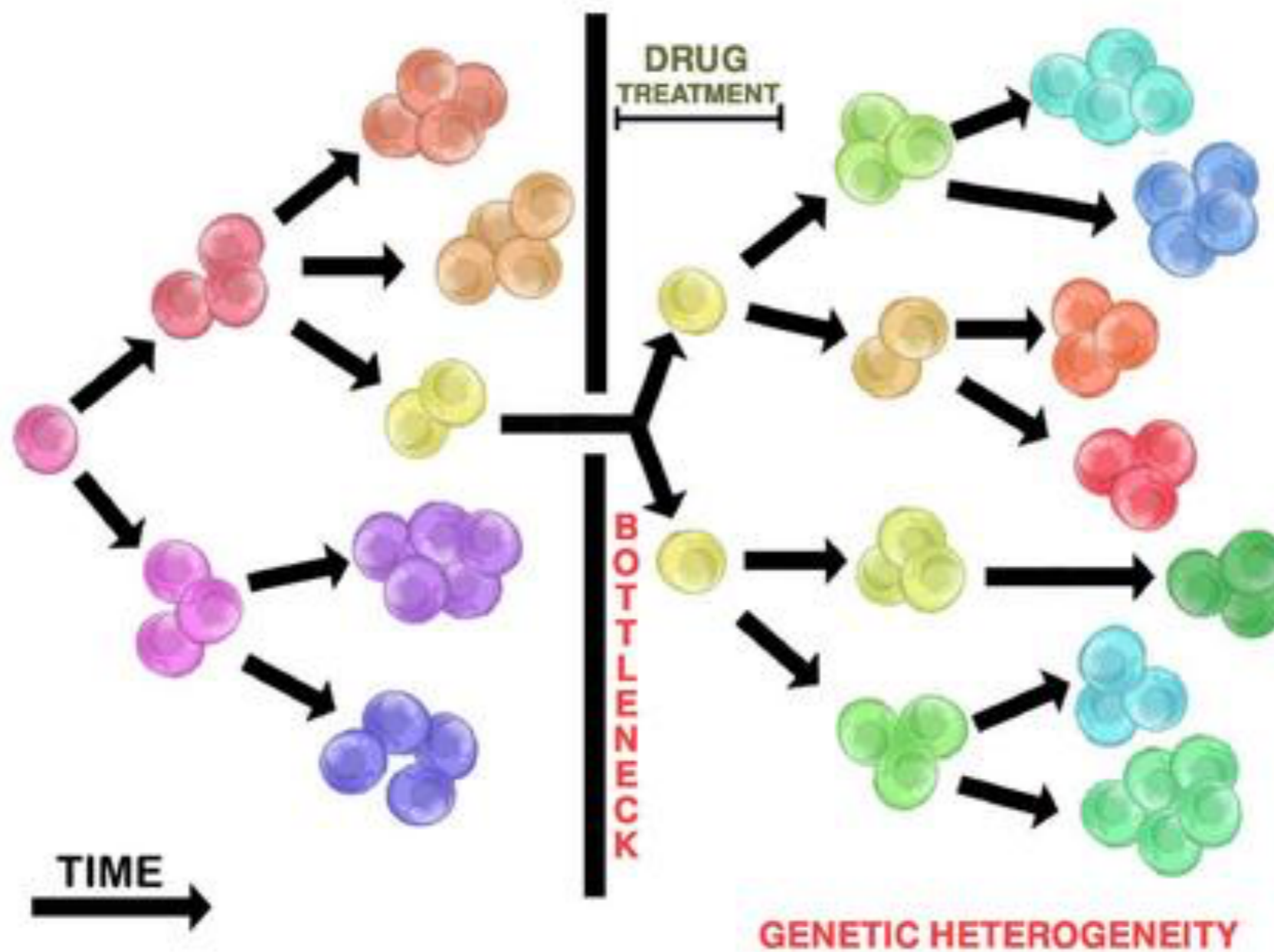
LINEAR EVOLUTION

TIME

BRANCHED EVOLUTION

TIME

- Tumour heterogeneity ➔ differential subclonal evolution
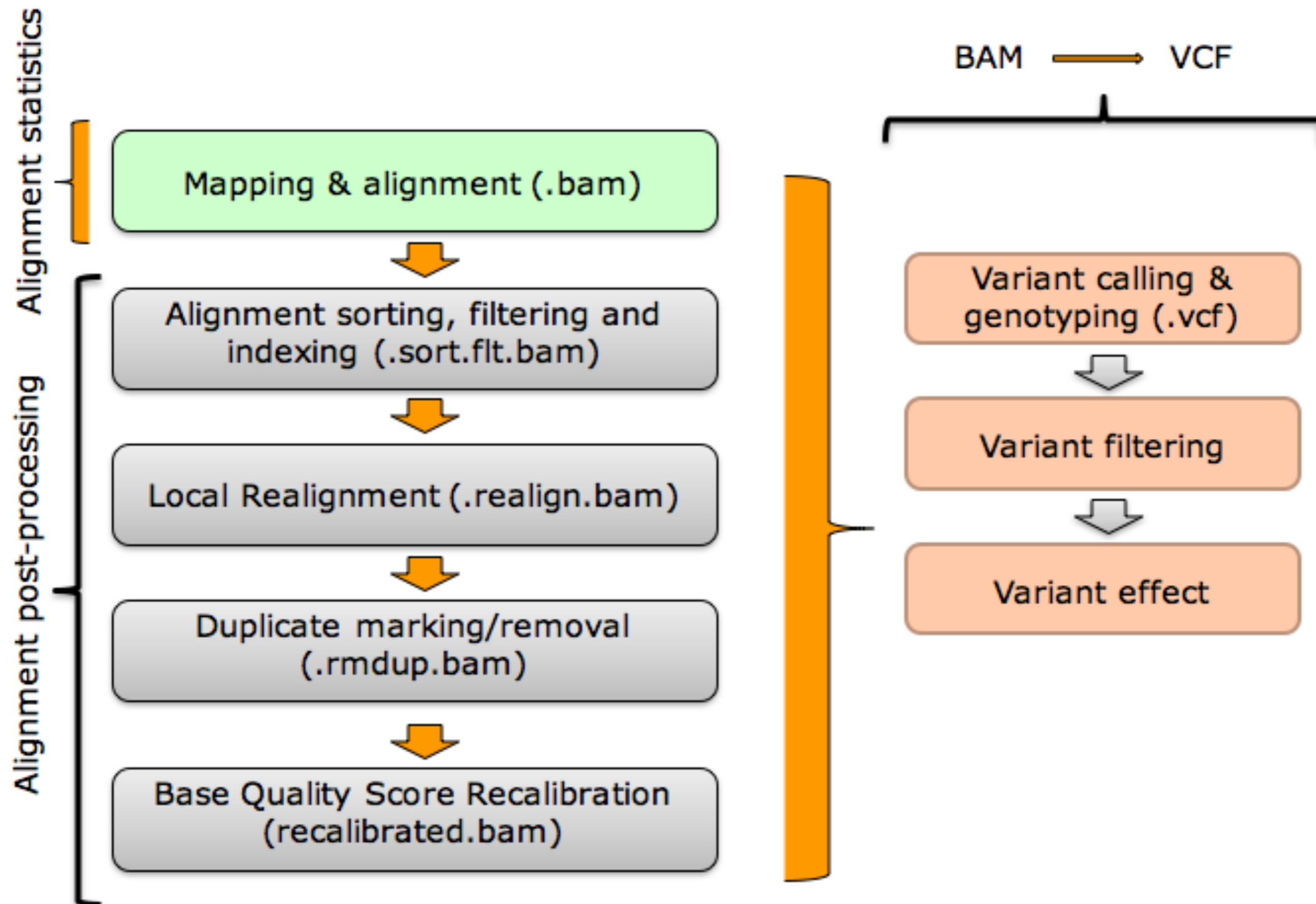
- Clones accumulate different mutations as they diverge.

# Treatment can re-shape tumour heterogeneity

How to identify somatic mutations in a tumor

# Variant calling pipeline



## Recommended workflow[1]

**Alignment statistics**

**Alignment post-processing**

Mapping & alignment (.bam)

Alignment sorting, filtering and indexing (.sort.flt.bam)

Local Realignment (.realign.bam)

Duplicate marking/removal (.rmdup.bam)

Base Quality Score Recalibration (recalibrated.bam)

BAM → VCF

Variant calling & genotyping (.vcf)

Variant filtering

Variant effect

# Matched samples for variant calling



(Healthy tissue in tumour fraction)

Tumour

Blood (healthy)

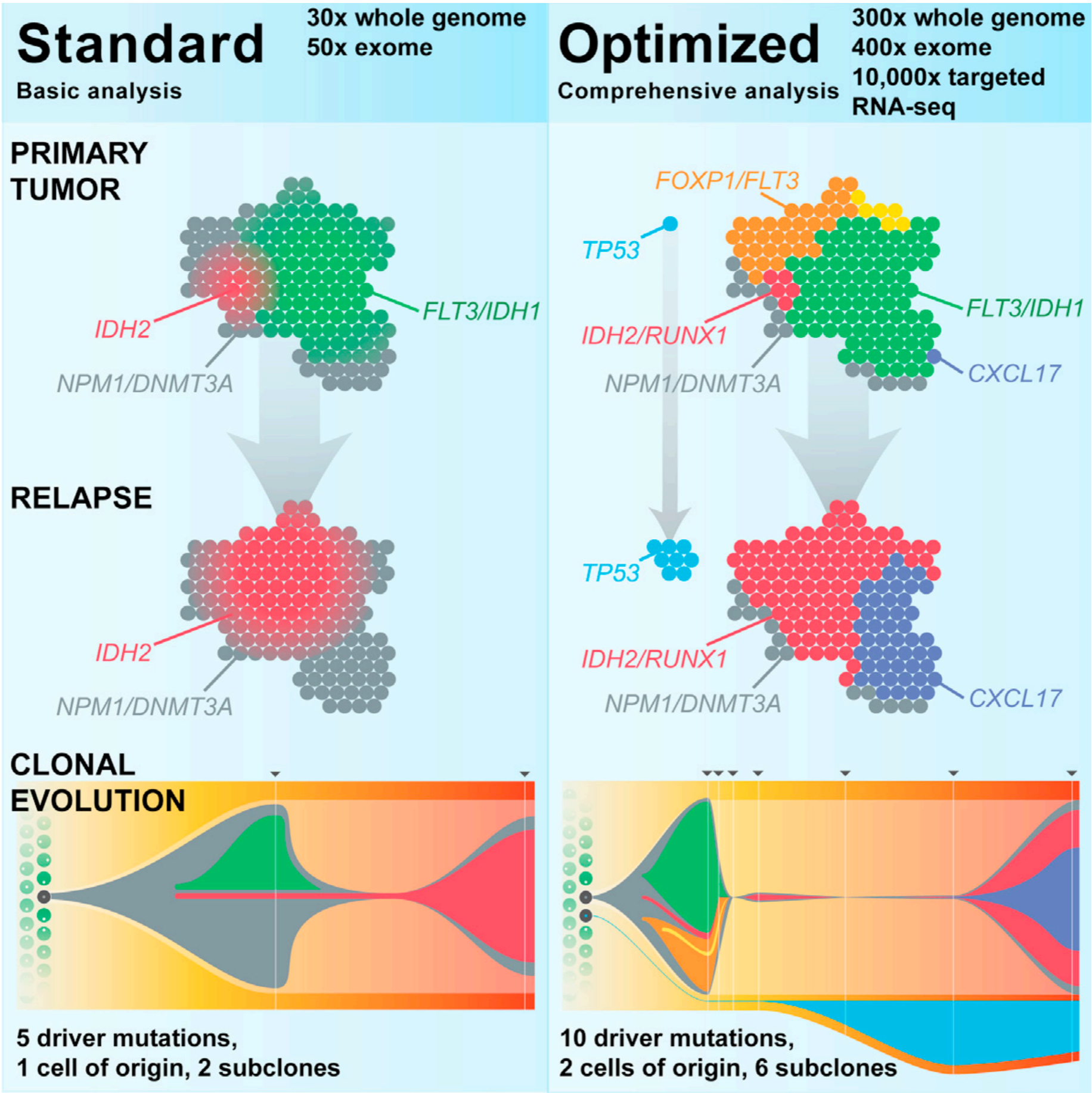germline          somatic          germline

# Somatic mutation calling vs."regular" variant calling

1. We are interested in somatic mutations:  differences between the **tumor genome and the normal genome** (NOT the reference genome).

2. The tumor data represents a **mixture** of reads from tumor cells and from normal cells, so we need **deeper sequencing** and **more sensitive analysis** to detect variants.

3. Tumors are often heterogeneous, and relevant mutations may be present at low allelic frequency.  So we need **even deeper** sequencing.

Also: we are often interested in copy number changes, translocations, and clonal architecture

**Standard** — 30x whole genome, 50x exome — Basic analysis

**Optimized** — 300x whole genome, 400x exome, 10,000x targeted, RNA-seq — Comprehensive analysis

PRIMARY TUMOR

Standard: FLT3/IDH1, IDH2, NPM1/DNMT3A

Optimized: FOXP1/FLT3, TP53, IDH2/RUNX1, NPM1/DNMT3A, FLT3/IDH1, CXCL17

RELAPSE

Standard: IDH2, NPM1/DNMT3A

Optimized: TP53, IDH2/RUNX1, NPM1/DNMT3A, CXCL17

CLONAL EVOLUTION

Standard: 5 driver mutations, 1 cell of origin, 2 subclones

Optimized: 10 driver mutations, 2 cells of origin, 6 subclones

M. Griffith *et al.*, *Cell Systems* (2015)

# Cancer gene panel amplicon sequencing

TruSeq Amplicon - Cancer Panel Gene List

| | | | | |
|---|---|---|---|---|
| ABL1 | EGFR | GNAS | MLH1 | RET |
| AKT1 | ERBB2 | HNF1A | MPL | SMAD4 |
| ALK | ERBB4 | HRAS | NOTCH1 | SMARCB1 |
| APC | FBXW7 | IDH1 | NPM1 | SMO |
| ATM | FGFR1 | JAK2 | NRAS | SRC |
| BRAF | FGFR2 | JAK3 | PDGFRA | STK11 |
| CDH1 | FGFR3 | KDR | PIK3CA | TP53 |
| CDKN2A | FLT3 | KIT | PTEN | VHL |
| CSF1R | GNA11 | KRAS | PTPN11 | |
| CTNNB1 | GNAQ | MET | RB1 | |

# Understanding variation in –omics times

**Traditionally**

1 Mutation
=
1 Disease



Phenotype
Function
Mechanism

Lots of hard work

**Now (High Throughput Sequencing, NGS)**

*X* Mutations
In
Y Patients
And
*Z* Conditions



Prediction of
Pathogenicity /
Unfeasible
Prioritization

http://www.ensembl.org/Homo_sapiens/Info/Index

# Ensembl Variant Effect Predictor (I)

# Ensembl Variant Effect Predictor (II)

## Options

**Transcript database to use:**
- ◉ Ensembl transcripts
- ○ RefSeq and other transcripts

**Get regulatory region consequences (human and mouse only):** ☑

**Type of consequences to display:** `Sequence Ontology terms ▼`

**Check for existing co-located variants:** `Yes ▼`

**Get 1000 Genomes global allele frequency for existing variants:** ☑

**Return results for variants in coding regions only:** ☐

**Show HGNC identifier for genes where available:** ☐

**Show Ensembl protein identifiers where available:** ☐

**Show HGVS identifiers for variants where available:** `No ▼`

## Missense SNP predictions (human only)

**SIFT predictions:** `Prediction and score ▼`

**PolyPhen predictions:** `Prediction and score ▼`

## Frequency filtering of existing variants (human only)

**Filter variants by frequency:** ☐

**NB:** Enabling frequency filtering may be slow for large datasets. The default options will filter out common variants found by the 1000 Genomes project.

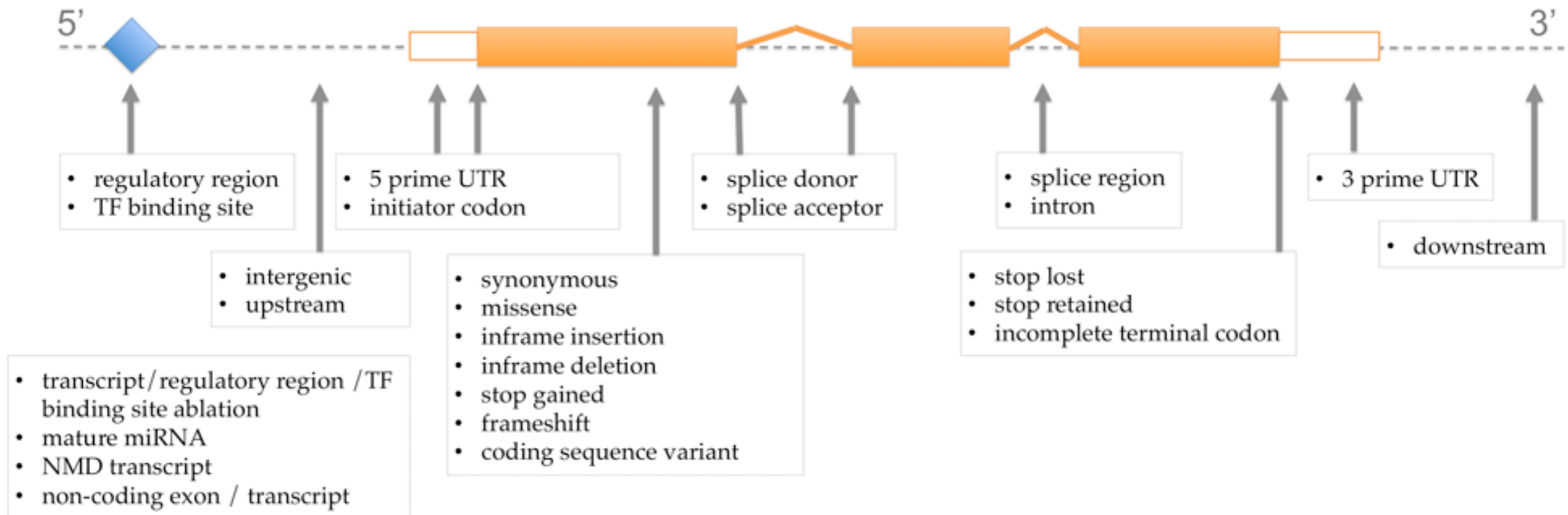**Filter:** `Exclude ▼` `variants with MAF greater than ▼` `0.01` `in 1000 genomes (1KG) combined population ▼`

**Next >**

# Ensembl Variant Effect Predictor (Results)

**Variant Effect Predictor Results:**

Download text version

Show 10 entries | Show/hide columns | Filter

| Uploaded Variation | Location | Allele | Gene | Feature | Feature type | Consequence | Position in cDNA | Position in CDS | Position in protein | Amino acid change | Codon change | Co-located Variation | Extra |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1_881907_-/C | 1:881906-881907 | C | ENSG00000187634 | ENST00000466827 | Transcript | downstream_gene_variant | - | - | - | - | - | 1_- | DISTANCE=3724 |
| 5_140532_T/C | 5:140532 | C | ENSG00000249430 | ENST00000512035 | Transcript | downstream_gene_variant | - | - | - | - | - | rs12516846 | DISTANCE=554; GMAF=C:0.1534 |
| 5_140532_T/C | 5:140532 | C | ENSG00000199540 | ENST00000362670 | Transcript | downstream_gene_variant | - | - | - | - | - | rs12516846 | DISTANCE=3670; GMAF=C:0.1534 |
| 5_140532_T/C | 5:140532 | C | ENSG00000153404 | ENST00000283426 | Transcript | missense_variant | 160 | 110 | 37 | V/A | gTa/gCa | rs12516846 | PolyPhen=benign(0); SIFT=tolerated(1); GMAF=C:0.1534 |
| 5_140532_T/C | 5:140532 | C | ENSG00000153404 | ENST00000502646 | Transcript | upstream_gene_variant | - | - | - | - | - | rs12516846 | DISTANCE=149; GMAF=C:0.1534 |

Showing 11 to 15 of 15 entries    << < 1 2 > >>

# Predictors: SIFT

http://sift.jcvi.org/

**Predicting Deleterious Amino Acid Substitutions**

Pauline C. Ng and Steven Henikoff



- Based on the degree of conservation in a multiple sequence alignment (MSA)
- MSA generated from PSI-BLAST results (closely related sequences)
- Deleterious if SIFT ≤ 0.05

# Predictors: Polyphen-2

http://genetics.bwh.harvard.edu/pph2/

### A method and server for predicting damaging missense mutations

Ivan A. Adzhubei,[1,7] Steffen Schmidt,[2,7] Leonid Peshkin,[3,7] Vasily E. Ramensky,[4] Anna Gerasimova,[5] Peer Bork,[6] Alexey S. Kondrashov,[5] and Shamil R. Sunyaev[1]

MACHINE LEARNING
- Naïve Bayes Classifier

SEQUENCE BASED FEATURES
- Importance of site: DISULFID, CROSSLNK, BINDING, ACT_SITE, LIPID, METAL, SITE, MOD_RES, CARBOHYD, NON_STD…
- Importance of region: TRANSMEM, INTRAMEM, COMPBIAS, REPEAT, COILED, SIGNAL, PROPEP…
- PSIC conservation score

STRUCTURE BASED FEATURES
- Likeness to destroy hydrophobic core, electrostatic interactions, interactions with ligands, or other important features of proteins

# Predictors: Polyphen-2

# Automatic methods to predict the pathogenicity of mutations

**SNAP**

**SIFT**

**SNAP: predict effect of non-synonymous polymorphisms on function**

Yana Bromberg[1,2,4,*] and Burkhard Rost[1,2,3]

**Predicting Deleterious Amino Acid Substitutions**

Pauline C. Ng and Steven Henikoff

**SNPs&GO**

**Polyphen-2**

A method and server for predicting damaging missense mutations

Ivan A. Adzhubei,[1,7] Steffen Schmidt,[2,7] Leonid Peshkin,[3,7] Vasily E. Ramensky,[4] Anna Gerasimova,[5] Peer Bork,[6] Alexey S. Kondrashov,[5] and Shamil R. Sunyaev[1]

**Functional Annotations Improve the Predictive Score of Human Disease-Related Mutations in Proteins**

Remo Calabrese, Emidio Capriotti, Piero Fariselli, Pier Luigi Martelli, and Rita Casadio*

**PMUT**

**MutationAssessor**

**PMUT: a web-based tool for the annotation of pathological mutations on proteins**

Carles Ferrer-Costa[1], Josep Lluis Gelpí[1,2,*], Leire Zamakola[1,3], Ivan Parraga[1,3], Xavier de la Cruz[1,4] and Modesto Orozco[1,2,3,*]

**Predicting the functional impact of protein mutations: application to cancer genomics**

Boris Reva*, Yevgeniy Antipin* and Chris Sander*

**Torkamani**

**Accurate prediction of deleterious protein kinase polymorphisms**

Ali Torkamani[1] and Nicholas J. Schork[2,*]

Some of the (many) methods implemented during the last decade