# Variant Analysis

**Bioinformatics for Systems and Synthetic Biology**

**Emidio Capriotti**

http://biofold.org/emidio

Department of Pharmacy and
Biotechnology (FaBiT)
University of Bologna

**Bio**molecules
**Fol**ding and
**Disease**

# Single Nucleotide Variants

Single Nucleotide Variants (SNVs)
is a DNA sequence variation occurring when a single nucleotide A, T, C, or G in the genome differs between members of the species.
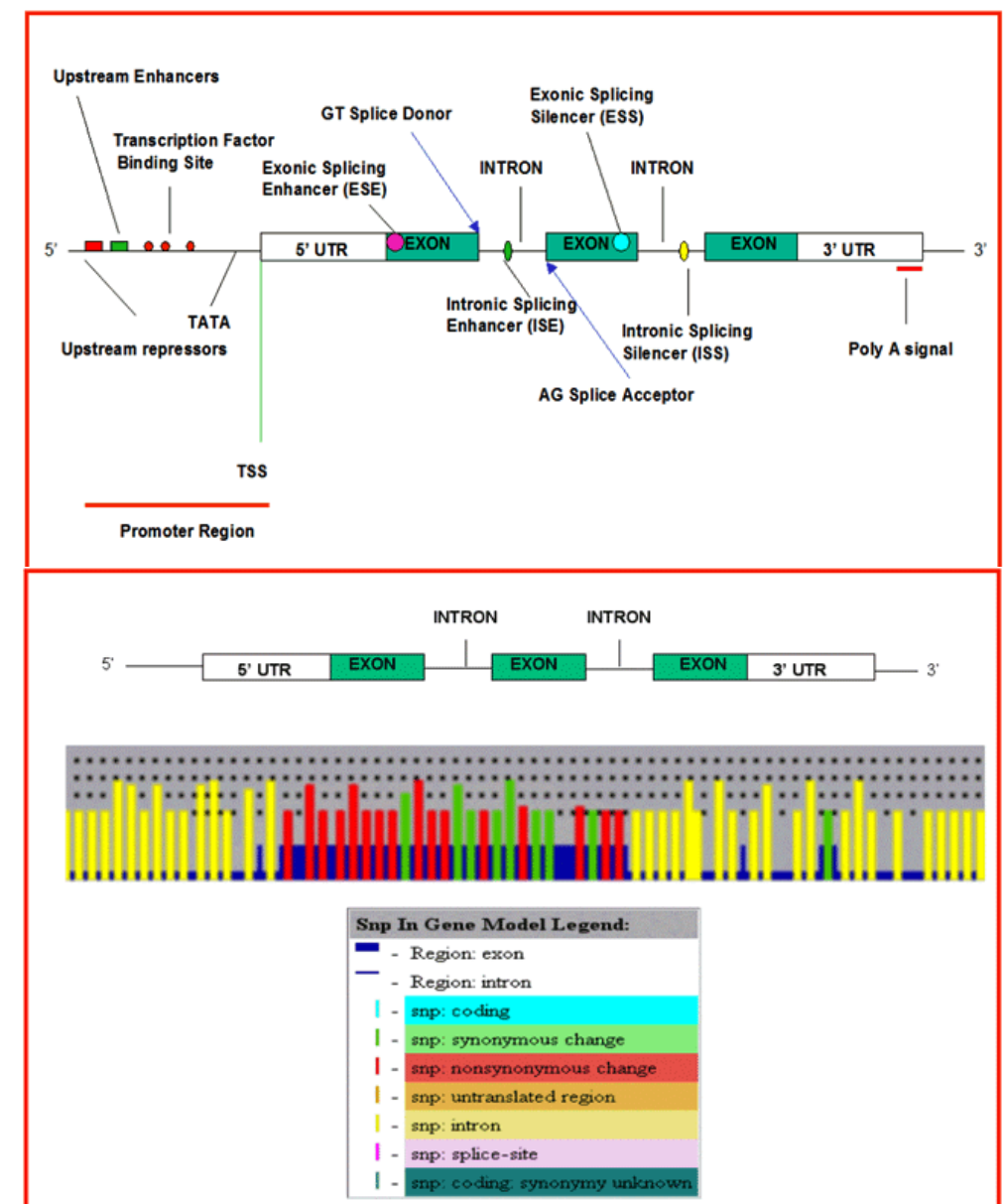It is used to refer to Polymorphisms when the population frequency is ≥ 1%

SNVs occur at any position and can be classified on the base of their locations.

Coding SNVs can be subdivided into two groups:

Synonymous: when single base substitutions do not cause a change in the resultant amino acid

Non-synonymous or Single Amino Acid Variants (SAVs): when single base substitutions cause a change in the resultant amino acid.

# Effects of variants

It is important to understand the functional effect of Single Nucleotide Polymorphisms (SNPs) that are very common type of variations, but also the impact rare variants which have allele frequencies below than 1%

Impact of coding variants
- Properties of amino acid residue substitution
- The evolutionary history of an amino acid position
- Sequence–function relationships
- Structure–function relationships

Impact of non-coding variants
- Transcription
- Pre-mRNA splicing
- MicroRNA binding
- Altering post-translational modification sites

*Cline and Karchin* (2011) Bioinformatics, 27; 441-448.

# 1000 Genomes

The 1000 Genomes Project aims to create the largest public catalogue of human variations and genotype data. Last version released the genotype of ~2,500 individuals.
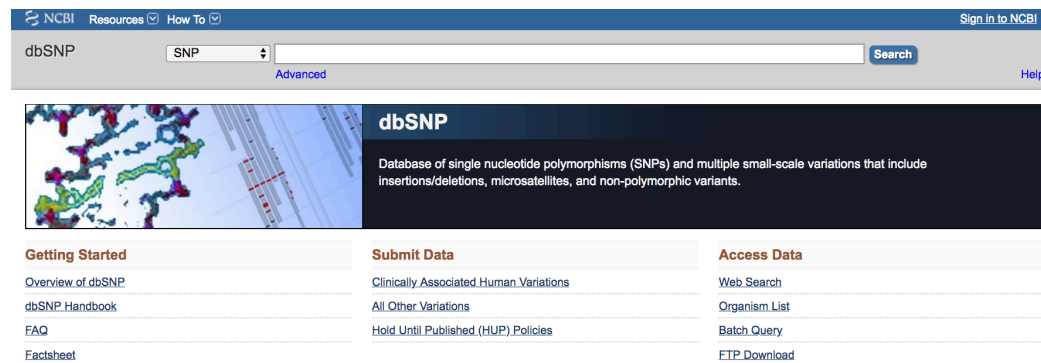
**Table 1 | Variants discovered by project, type, population and novelty**

**a** Summary of project data including combined exon populations

| Statistic | Low coverage | | | | Trios | | | Exon (total) | Union across projects |
|---|---|---|---|---|---|---|---|---|---|
| | CEU | YRI | CHB+JPT | Total | CEU | YRI | Total | | |
| Samples | 60 | 59 | 60 | 179 | 3 | 3 | 6 | 697 | 742 |
| Total raw bases (Gb) | 1,402 | 874 | 596 | 2,872 | 560 | 615 | 1,175 | 845 | 4,892 |
| Total mapped bases (Gb) | 817 | 596 | 468 | 1,881 | 369 | 342 | 711 | 56 | 2,648 |
| Mean mapped depth ($\times$) | 4.62 | 3.42 | 2.65 | 3.56 | 43.14 | 40.05 | 41.60 | 55.92 | NA |
| Bases accessed (% of genome) | 2.43 Gb (86%) | 2.39 Gb (85%) | 2.41 Gb (85%) | 2.42 Gb (86.0%) | 2.26 Gb (79%) | 2.21 Gb (78%) | 2.24 Gb (79%) | 1.4 Mb | NA |
| No. of SNPs (% novel) | 7,943,827 (33%) | 10,938,130 (47%) | 6,273,441 (28%) | 14,894,361 (54%) | 3,646,764 (11%) | 4,502,439 (23%) | 5,907,699 (24%) | 12,758 (70%) | 15,275,256 (55%) |
| Mean variant SNP sites per individual | 2,918,623 | 3,335,795 | 2,810,573 | 3,019,909 | 2,741,276 | 3,261,036 | 3,001,156 | 763 | NA |
| No. of indels (% novel) | 728,075 (39%) | 941,567 (52%) | 666,639 (39%) | 1,330,158 (57%) | 411,611 (25%) | 502,462 (37%) | 682,148 (38%) | 96 (74%) | 1,480,877 (57%) |
| Mean variant indel sites per individual | 354,767 | 383,200 | 347,400 | 361,669 | 322,078 | 382,869 | 352,474 | 3 | NA |
| No. of deletions (% novel) | ND | ND | ND | 15,893 (60%) | 6,593 (41%) | 8,129 (50%) | 11,248 (51%) | ND | 22,025 (61%) |
| No. of genotyped deletions (% novel) | ND | ND | ND | 10,742 (57%) | ND | ND | 6,317 (48%) | ND | 13,826 (58%) |
| No. of duplications (% novel) | 259 (90%) | 320 (90%) | 280 (91%) | 407 (89%) | 187 (93%) | 192 (91%) | 256 (92%) | ND | 501 (89%) |
| No. of mobile element insertions (% novel) | 3,202 (79%) | 3,105 (84%) | 1,952 (76%) | 4,775 (86%) | 1,397 (68%) | 1,846 (78%) | 2,531 (78%) | ND | 5,370 (87%) |
| No. of novel sequence insertions (% novel) | ND | ND | ND | ND | 111 (96%) | 66 (86%) | 174 (93%) | ND | 174 (93%) |

*1000 Genomes Project Consortium* (2010). Nature. 467: 1061-1073.

# SNVs and SAVs databases

dbSNP @ NCBI



http://www.ncbi.nlm.nih.gov/snp

Single Nucleotide Variants

***Homo sapiens***     **904,623,795**

SwissVar @ ExPASy



http://www.expasy.ch/swissvar/

Single Amino acid Variants

| | |
|---|---|
| *Homo sapiens* | 83,996 |
| *Disease* | *32,930* |
| *Polymorphisms* | *39,938* |

*Jan 2026*

# Variant Call Format

The final result of the variant calling procedure is a VCF file.

```
##fileformat=VCFv4.1
##tcgaversion=1.1
##reference=<ID=hg19,source=.>
##phasing=none
##geneAnno=none
##INFO=<ID=VT,Number=1,Type=String,Description="Variant type, can be SNP, INS or DEL">
##INFO=<ID=VLS,Number=1,Type=Integer,Description="Final validation status relative to non-adjacent Normal, ......">
##FILTER=<ID=CA,Description="Fail Carnac (Tumor and normal coverage, tumor variant count, mapping quality, ......">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read depth at this position in the sample">
##FORMAT=<ID=AD,Number=.,Type=Integer,Description="Depth of reads supporting alleles 0/1/2/3...">
##FORMAT=<ID=BQ,Number=.,Type=Integer,Description="Average base quality for reads supporting alleles">
##FORMAT=<ID=SS,Number=1,Type=Integer,Description="Variant status relative to non-adjacent Normal,0=wildtype, ......">
##FORMAT=<ID=SSC,Number=1,Type=Integer,Description="Somatic score between 0 and 255">
##FORMAT=<ID=MQ60,Number=1,Type=Integer,Description="Number of reads (mapping quality=60) supporting variant">
#CHROM    POS      ID    REF   ALT   QUAL   FILTER   INFO          FORMAT                NORMAL                  PRIMARY
1         10048    .     C     CCT   .      CA       VT=INS;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:66:.,0:.:0:.:0    0/1:32:.,2:.:2:.:0
1         10078    .     CT    C     .      CA       VT=DEL;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:25:.,0:.:0:.:0    0/1:13:.,2:.:2:.:0
1         10177    .     A     AC    .      CA       VT=INS;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:57:.,0:.:0:.:0    0/1:22:.,2:.:2:.:0
. . . . . .
. . . . . .
1         900505   .     G     C     .      PASS     VT=SNP;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/1:188:.,89:26:1:.:81  0/1:210:.,113:24:1:.:100
. . . . . .
. . . . .
1         1991007  .     G     T     .      PASS     VT=SNP;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:222:.,1:2:0:.:1   0/1:88:.,41:25:2:50:34
. . . . . .
```

# File Content

The file contains information about single nucleotide variants and indels of single or multiple samples.

For each variant the number of supporting reads for reference and alternative alleale

The original VCF does not contain any information functional effect of the variants.

# Main data sources

Single genetic variants are collected in different databases:

- dbSNP - variation from all species. http://www.ncbi.nlm.nih.gov/SNP/
- EVS - specific for human.  http://evs.gs.washington.edu/EVS/
- ClinVar - Variants and human health.  http://www.ncbi.nlm.nih.gov/clinvar/
- Cosmic - Somatic mutation in cancer. http://cancer.sanger.ac.uk/

This information is important for variant calling but useless for capturing the complexity of genotype/phenotype relationship. The VCF more informative because we can analyze co-occurring events. The major sources are:

- 1000 Genomes: WGS data of individuals http://www.1000genomes.org/
- TCGA: Cancer Genomes https://tcga-data.nci.nih.gov/

# All in One

Exomizer is a variant analysis tools that tests presence of variants associated to specific phenotypes



http://www.sanger.ac.uk/resources/software/exomiser/submit/

# The complexity of cancer

Cancer is **complex disorder** characterized by high level of mutation rate.

Mutations can be classified in germline and somatic whether they are inherited from parents or the result of error in DNA replication.

Another classification is between driver and passenger mutations whether they provide selective advantage with respect to normal cells increasing their proliferation rate or not.

# Hallmarks of cancer

The six hallmarks of cancer - distinctive and complementary capabilities that enable tumor growth and metastatic dissemination.

# Oncogene vs Suppressor

Oncogenes have highly recurrent mutations, Tumor suppressors have sparse variants.



*Vogelstein et al.* Science 2013, **339**:1546

# Main challenges

Computational methods for cancer genome interpretation have been developed to address the following issues:

- Detection of recurrent somatic mutations and cancer driver genes;

- Prediction of driver variants and their functional impact;

- Estimate the impact of multiple variants at network and pathway level;

- Differentiate subclonal populations and their variation pattern.



*Raphael et al.* Genome Medicine 2014, **6**:5

# How data looks like?

Variant Calling File (VCF) with germline and somatic variants

```
##fileformat=VCFv4.1
##tcgaversion=1.1
##reference=<ID=hg19,source=.>
##phasing=none
##geneAnno=none
##INFO=<ID=VT,Number=1,Type=String,Description="Variant type, can be SNP, INS or DEL">
##INFO=<ID=VLS,Number=1,Type=Integer,Description="Final validation status relative to non-adjacent Normal, ......">
##FILTER=<ID=CA,Description="Fail Carnac (Tumor and normal coverage, tumor variant count, mapping quality, ......">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read depth at this position in the sample">
##FORMAT=<ID=AD,Number=.,Type=Integer,Description="Depth of reads supporting alleles 0/1/2/3...">
##FORMAT=<ID=BQ,Number=.,Type=Integer,Description="Average base quality for reads supporting alleles">
##FORMAT=<ID=SS,Number=1,Type=Integer,Description="Variant status relative to non-adjacent Normal,0=wildtype, ......">
##FORMAT=<ID=SSC,Number=1,Type=Integer,Description="Somatic score between 0 and 255">
##FORMAT=<ID=MQ60,Number=1,Type=Integer,Description="Number of reads (mapping quality=60) supporting variant">
#CHROM    POS       ID     REF    ALT    QUAL    FILTER    INFO          FORMAT                NORMAL                      PRIMARY
1         10048     .      C      CCT    .       CA        VT=INS;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:66:.,0:.:0:.:0         0/1:32:.,2:.:2:.:0
1         10078     .      CT     C      .       CA        VT=DEL;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:25:.,0:.:0:.:0         0/1:13:.,2:.:2:.:0
1         10177     .      A      AC     .       CA        VT=INS;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:57:.,0:.:0:.:0         0/1:22:.,2:.:2:.:0
. . . . . .
1         900505    .      G      C      .       PASS      VT=SNP;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/1:188:.,89:26:1:.:81     0/1:210:.,113:24:1:.:100
. . . . . .
1         1991007   .      G      T      .       PASS      VT=SNP;VLS=5  GT:DP:AD:BQ:SS:SSC:MQ60  0/0:222:.,1:2:0:.:1         0/1:88:.,41:25:2:50:34
. . . . . .
```

# The TCGA data

The Cancer Genome Atlas Consortium

TCGA data (https://portal.gdc.cancer.gov/)
- 91 cancer projects (~50,270 cases)
- BAM files available

# The ICGC AR%GO

The International Cancer Genome Consortium

ICGC (https://platform.icgc-argo.org/)

- 5,528 Donors
- 429.22 TB data
- 77.4 million simple somatic mutations.

# Somatic Mutations

Number of somatic mutations per sample vary significantly across cancer types



Number of Somatic Mutations in Donor's Exomes Across Cancer Projects

Donor Distribution

24,077 Donors across 84 Projects

Top 20 Mutated Cancer Genes with High Functional Impact SSMs
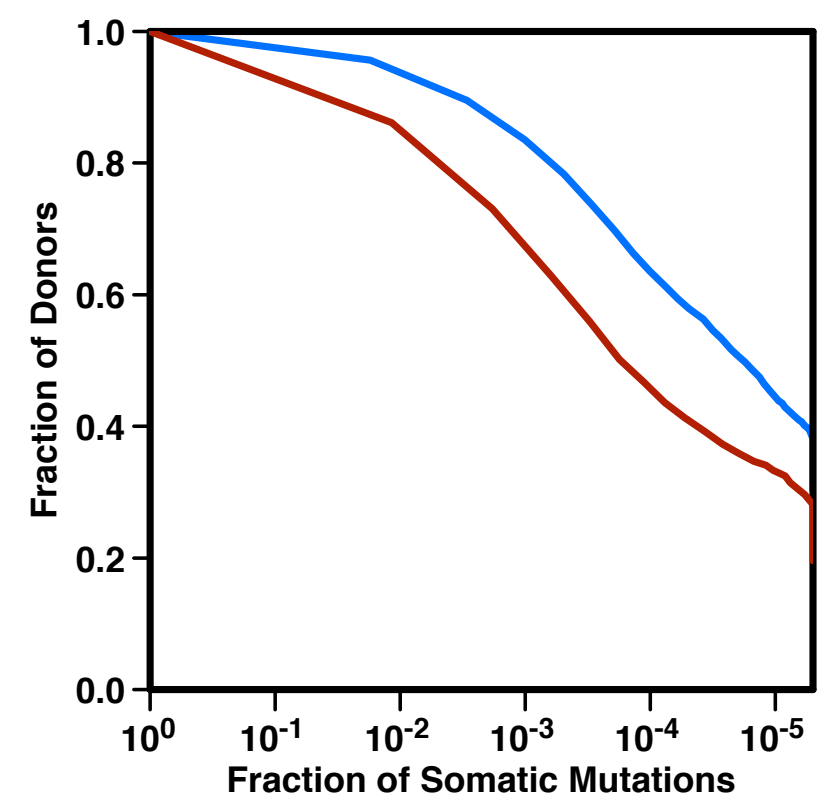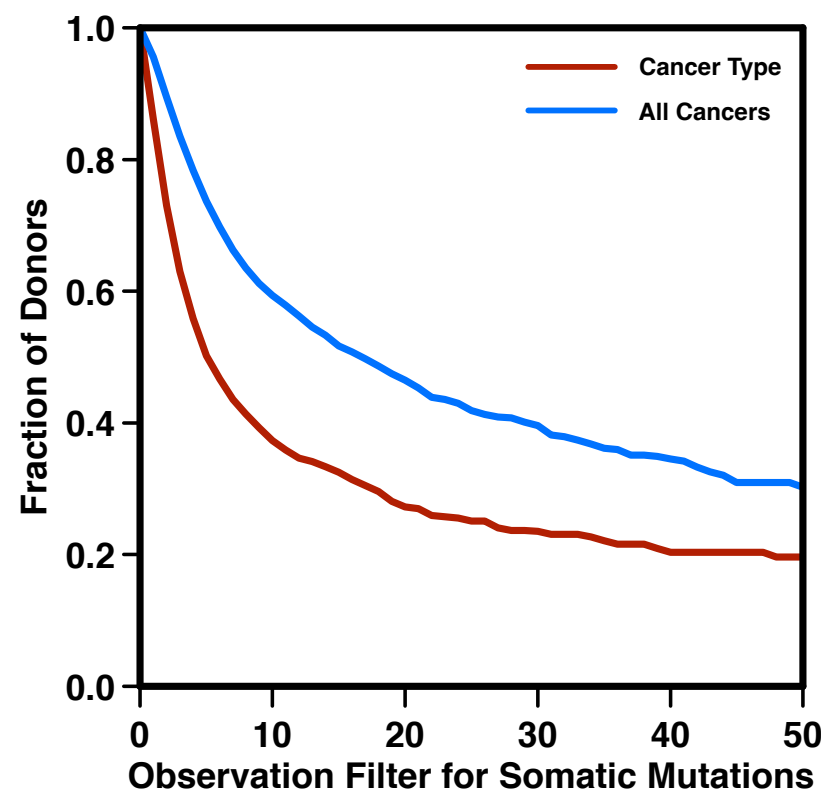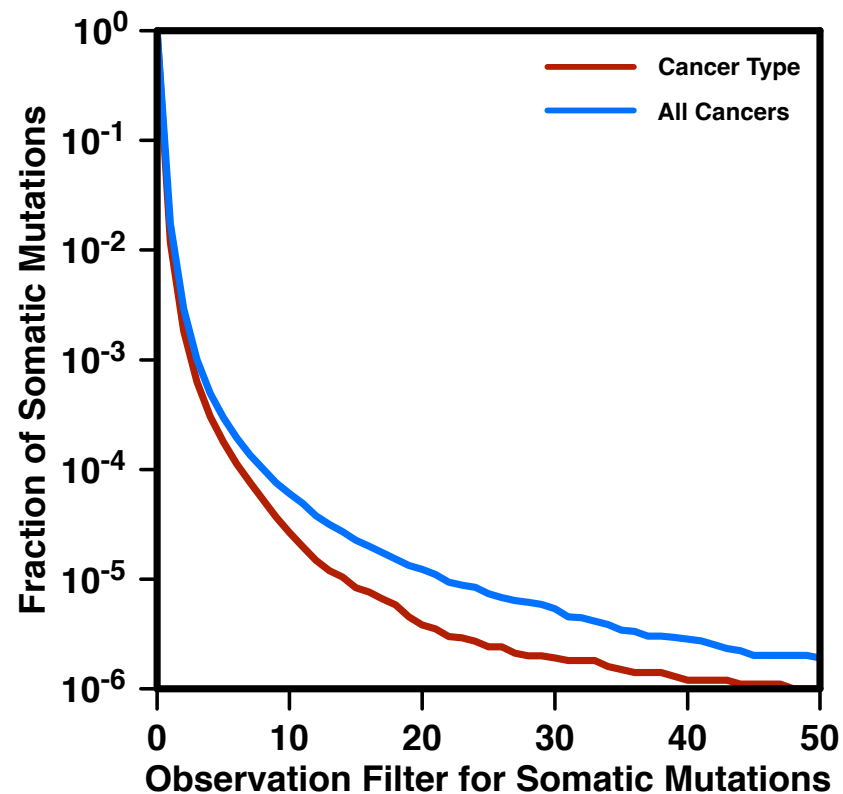
15,285 Unique SSM-Tested Donors

# Driver vs Passenger
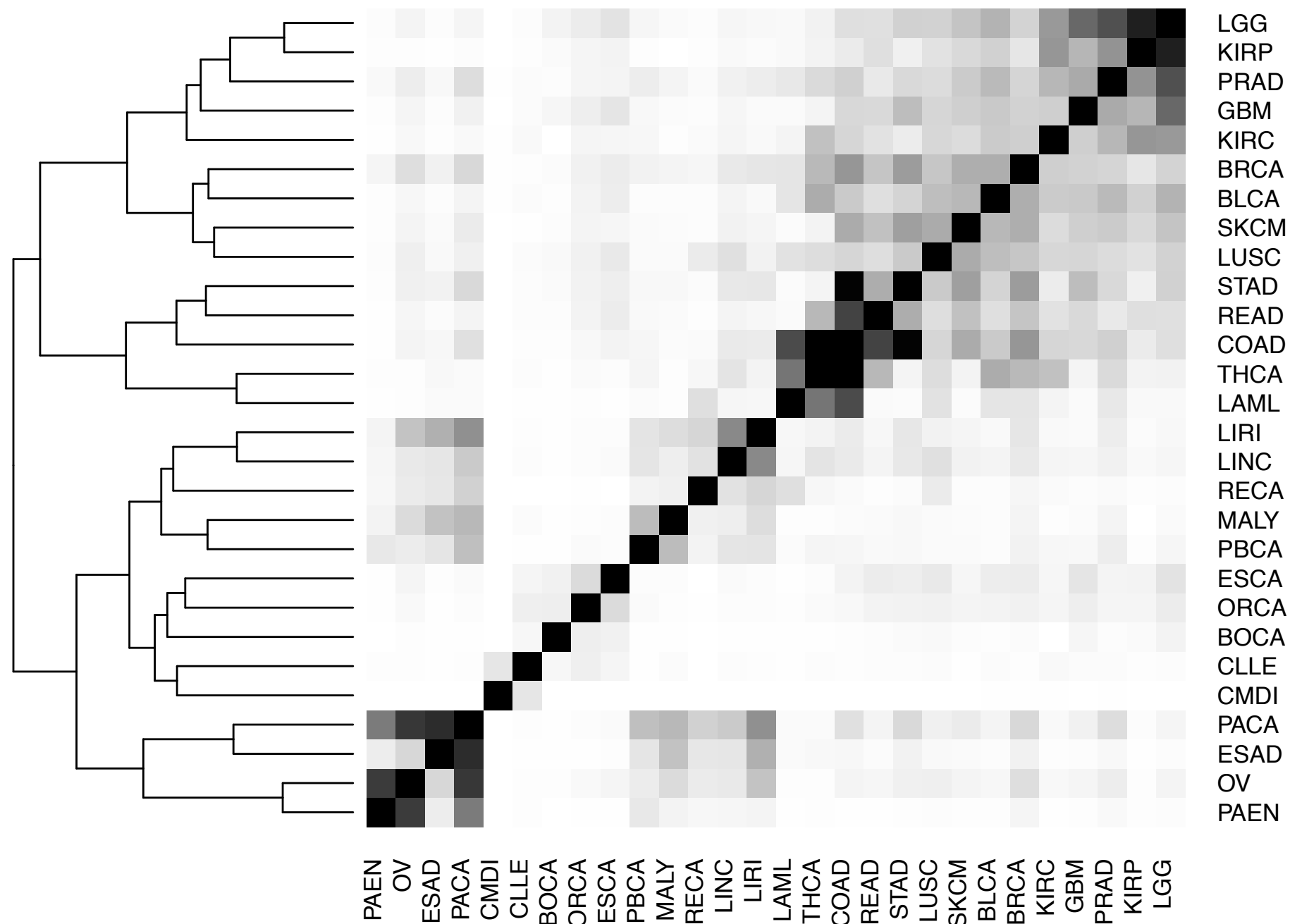
Number of recurrent mutations decrease exponentially.
On average a small fraction of variants a present in the majority of the samples.

Selecting mutations that are repeated at least twice we filter out ~98% mutations
and are still able to recover ~96% of the patients
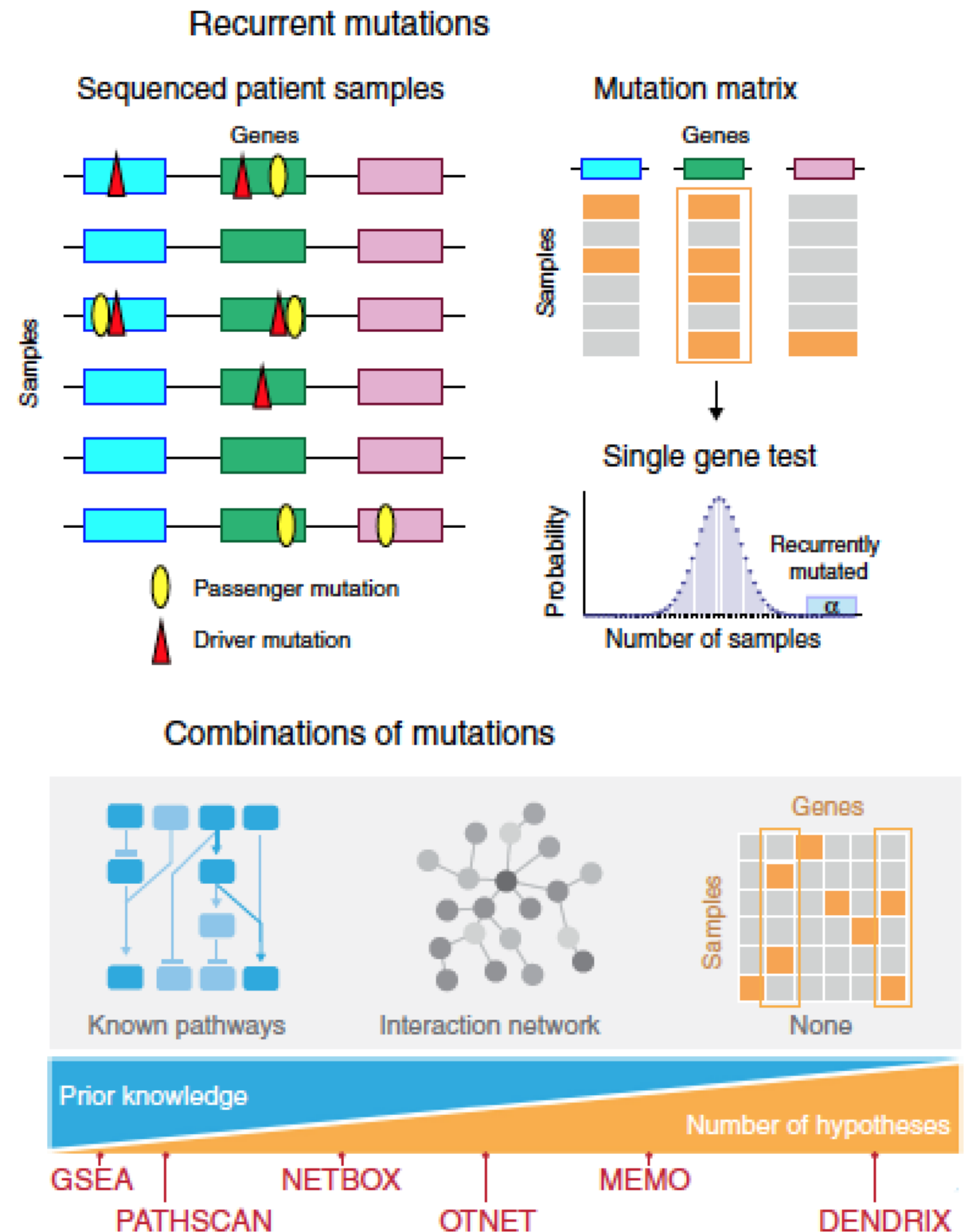


*Tian R, Basu M, Capriotti E. BMC Genomics 2015*

# Recurrent variations

Recurrent mutations that are found in more samples than would be expected by chance are good candidates for driver mutations.

To identify such recurrent mutations, a statistical test is performed which usually collapses all the non-synonymous mutations in a gene.

Identification of recurrent mutations in predefined groups such as pathways and protein-protein interaction networks and de novo identification of combinations, without relying on a priori definition.

# The main idea

**Genes implicated in cancer** should have **high mutation rate**

In comparison to normal, tumor cells should have higher occurrence of functional mutations in genes involved in the insurgence and progression of the disease.
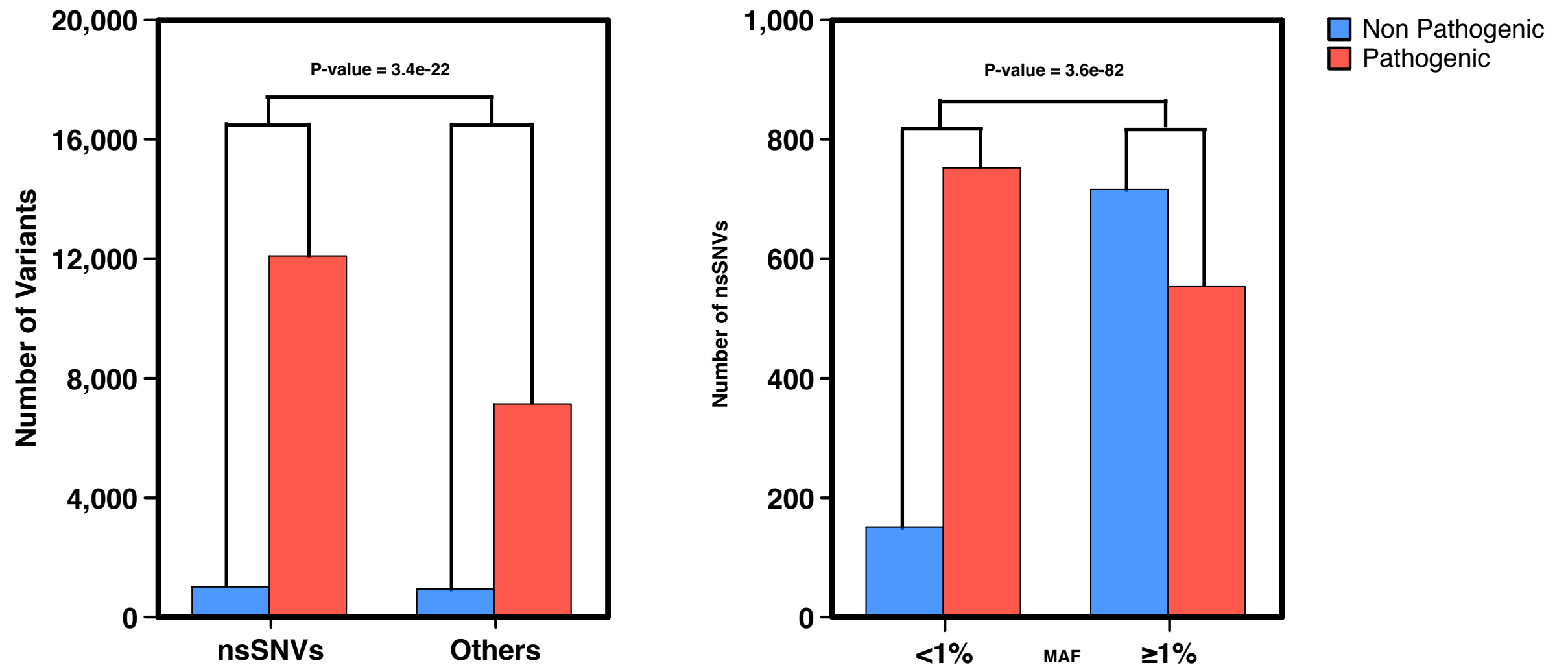
**Problem:**

How can we select mutations with functional impact?

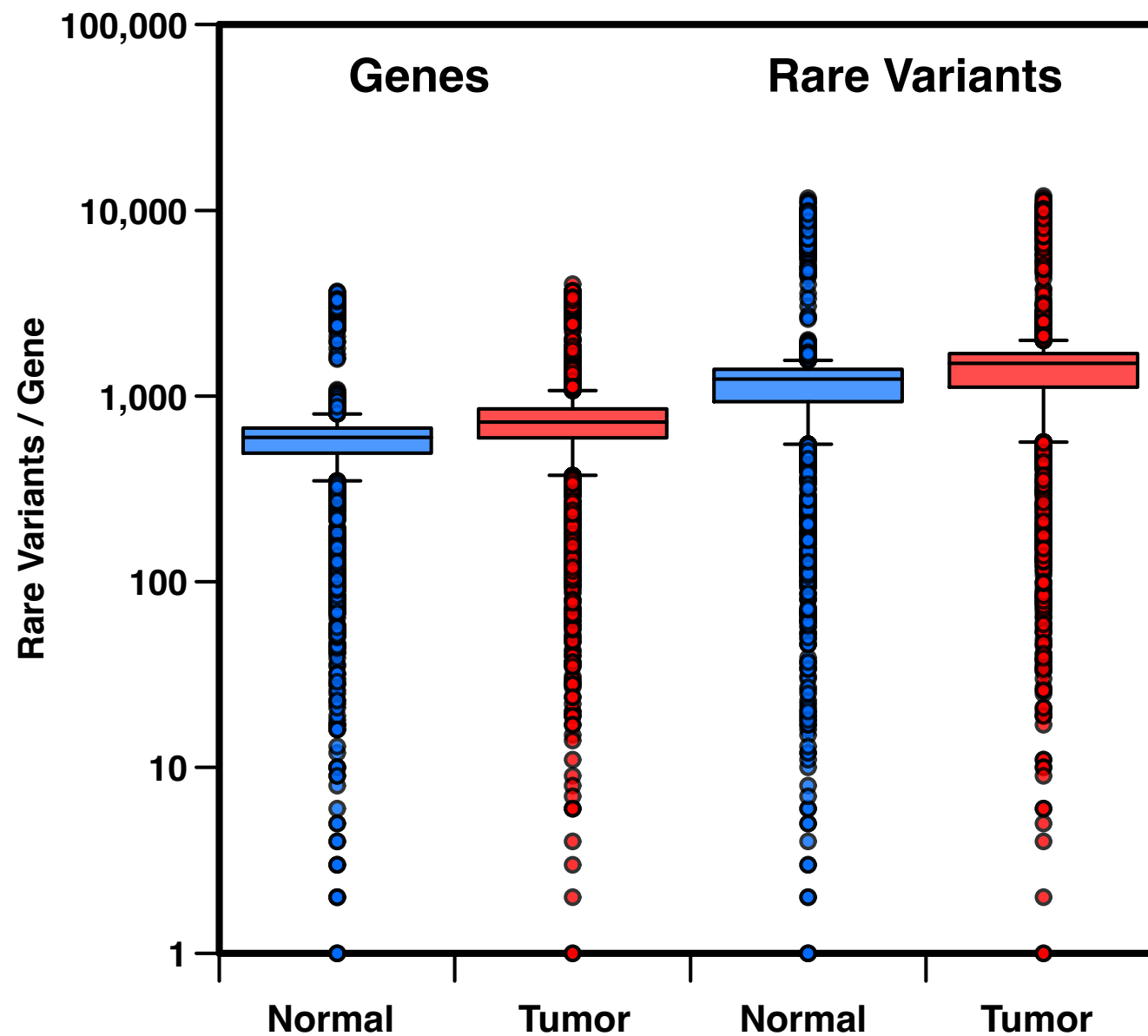| | |
|---|---|
| Average number of variants | ~3,000,000 |
| Average exome variants | ~23,000 |
| Average nonsynonymous single nucleotide variants | ~10,000 |
| Average rare (MAF≤0.5%) nonsynonymous single nucleotide variants | ~300 |

*The 1000 Genomes Project* (2010). Nature. 467; 1062-1073.

# Variants and MAF

Rare variants are more likely to be associated to disease than high frequency variants



*Tian R, Basu M, Capriotti E* (2014). Bioinformatics. 30: i572-i578

# Rate Variants and Genes

On average tumor samples (COAD) have ~150 more rare missense variants and mutated genes

# Mutation rates

The analysis of **1000 Genomes, The Cancer Genome Atlas (TCGA)** normal and tumor samples shows an **increasing number of genes with rare nonsynonymous SNVs**.

| Cohort | %Genes PDR≤0.05 | %Genes PDR>0.05 |
|---|---|---|
| 1000 Genomes | 95% | 5% |
| TCGA Normal | 92% | 8% |
| TCGA Tumor | 82% | 18% |

Tumor = Colon Adenocarcinoma
PDR = Gene Putative Defective Rate
     Fraction of samples in which a gene has ≥1 nonsynonymous variant with MAF≤0.5%