

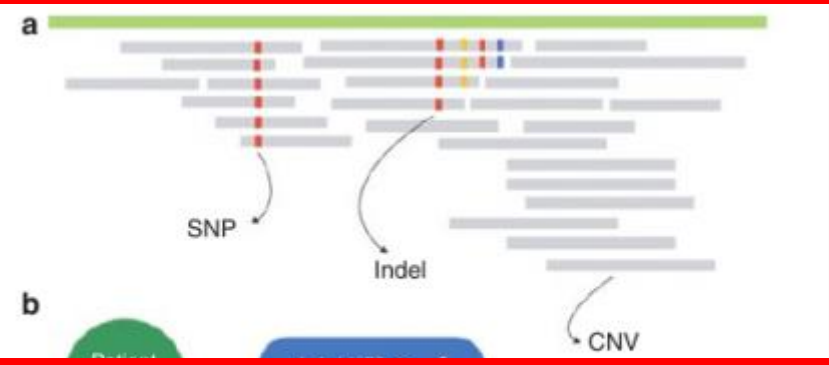


# Variant calling and analysis

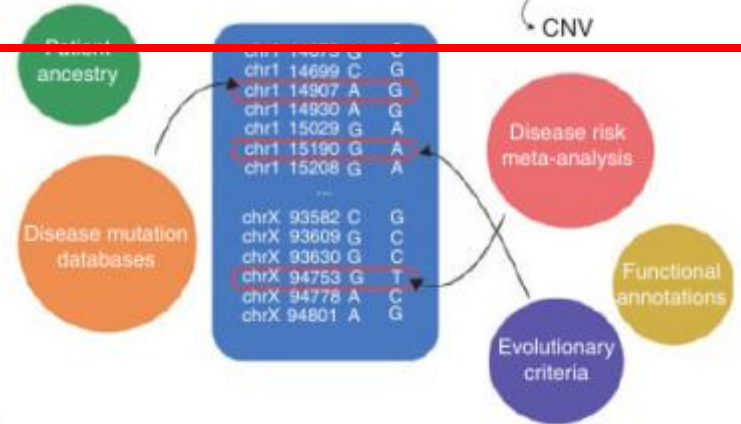
## From reads to diagnosis

# Overview

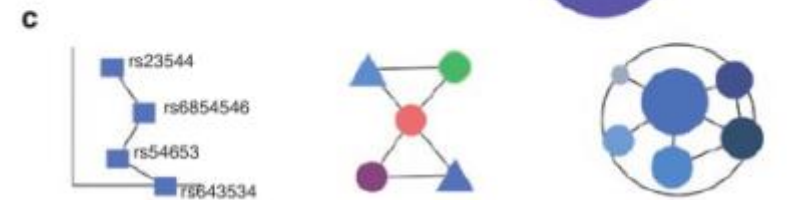
Variant  
calling



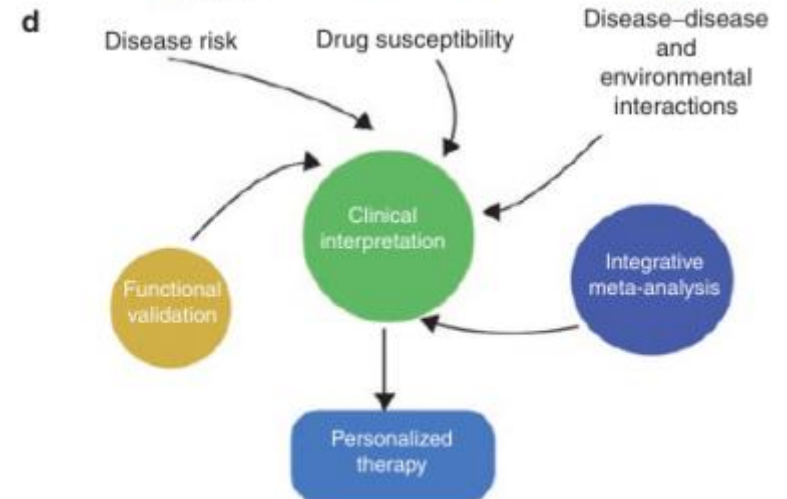
Variant  
filtering



Variant  
relations

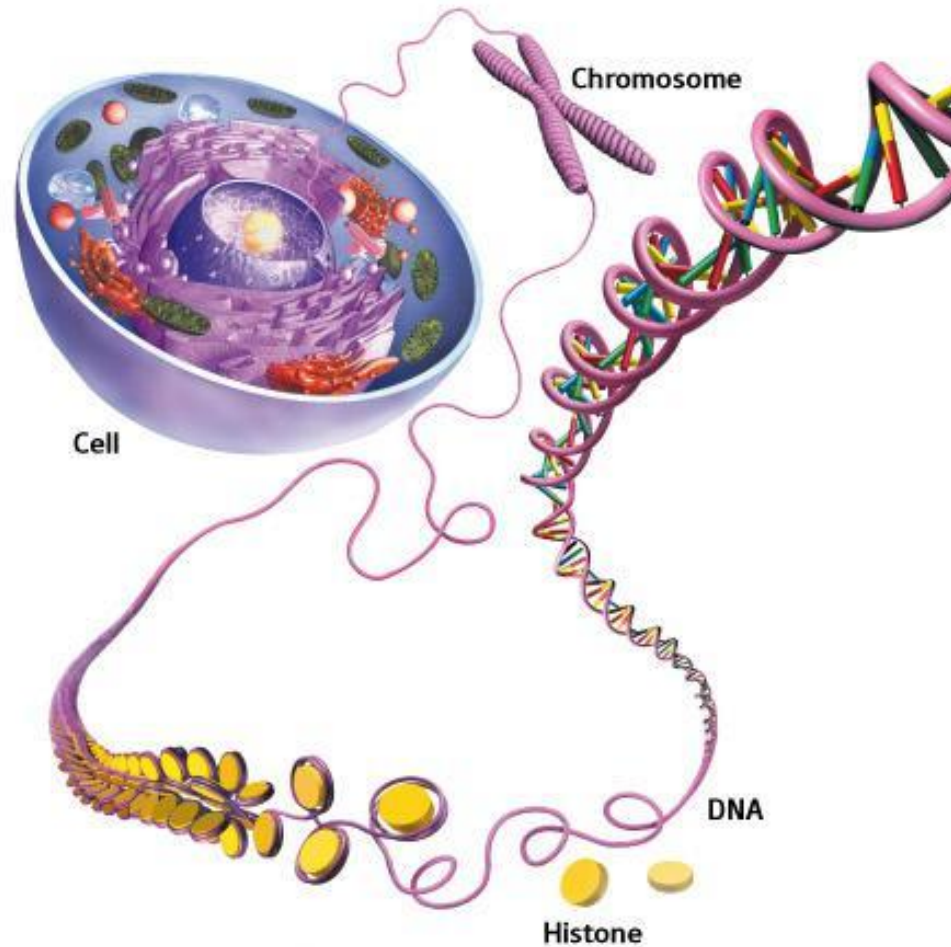


Variant  
visualization



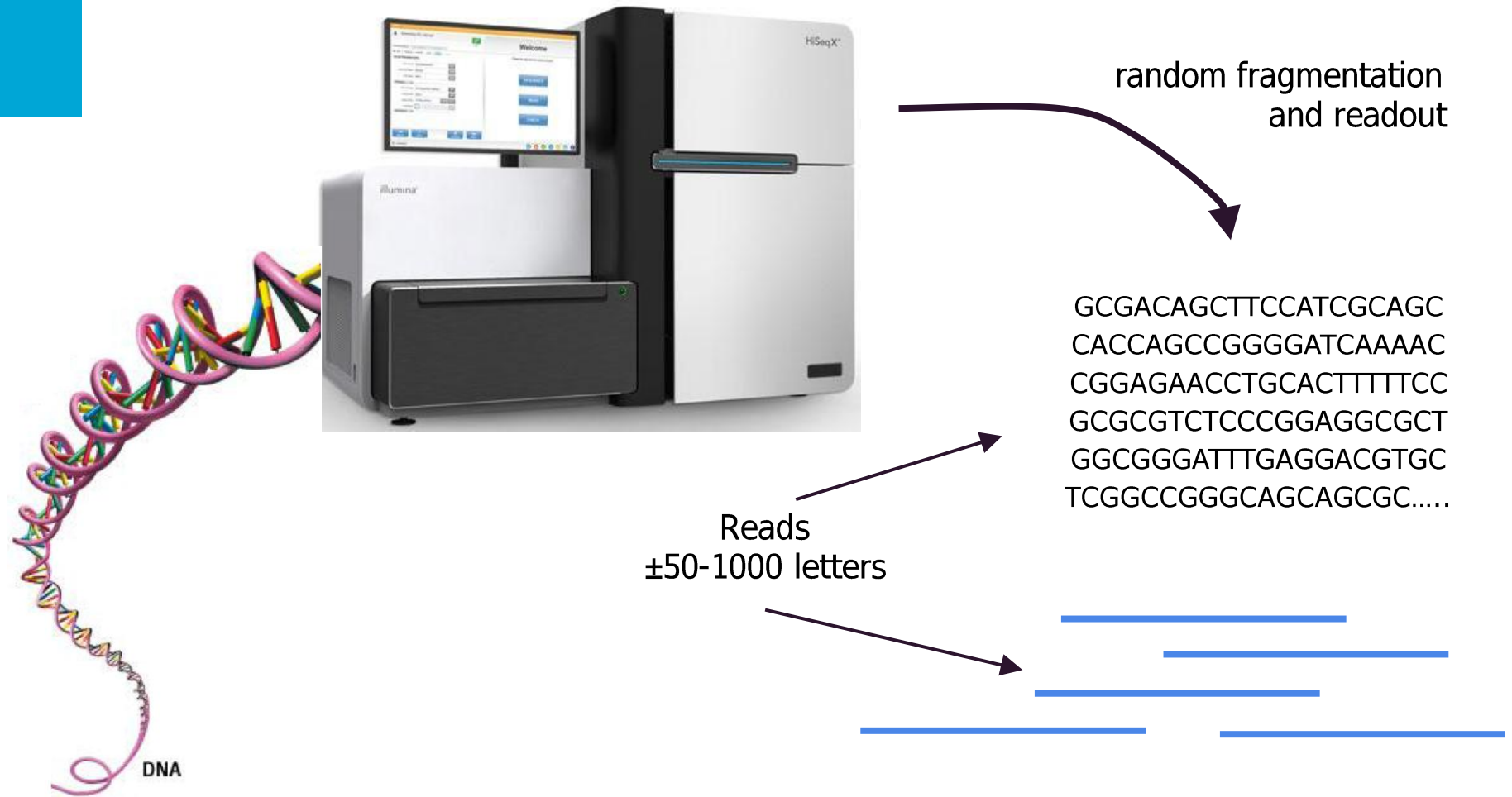
# Reading the genome

DNA in living cells



# Sequencing

As means of reading out the genome





# Variant calling

Requires a reference genome



Read depth: ~60  
~180 billion letters!



CCCTGCGCCGCGTGCGCGACAGCTTCCATCGCAGCCTG  
CTGTGGATAGGACACCAGCCGGGGATCAAAAC  
CCGCCTGACGGCGCGGGAGAACCTGCACTTTTTCCACC  
CCGGCGACGGCGCGCGTCTCCCGGAGGCGCTG  
GCGCAGGCCGGGCTGGCGGGATTTGAGGACGTGCCGG  
TCGCTCAGCTCTCGGCCGGGCAGCAGCGCCGGG.....

23 chromosomes

Reference  
sequence



# Variant calling

Requires a reference sequence



- Coordinate system needed
  - for annotations (where are the genes)
  - for reporting variants
- Currently build 38 (GRCh38)
  - 20<sup>th</sup> version (Hg20)
  - 350 gaps left
  - Based mostly on a few anonymous donors (from Buffalo, NY)
- Many still use Hg19 or Hg18

# Mapping read to a reference (1)

allows for detecting variation

Perfect match!

GCAGCCACCAGCCGGGGATCAAA



...GCGACAGCTTCCATCGCAGCCACCAGCCGGGGATCAAAACCGGAGAACCTGCACTTTTTTCC...

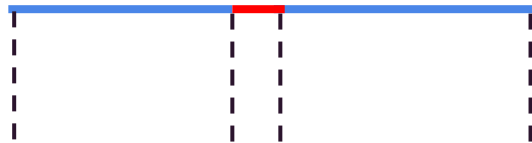
# Mapping read to a reference (3)

allows for detecting variation

SNP = Single Nucleotide  
Polymorphism



GCAGCCACCA T CCGGGGATCAAA

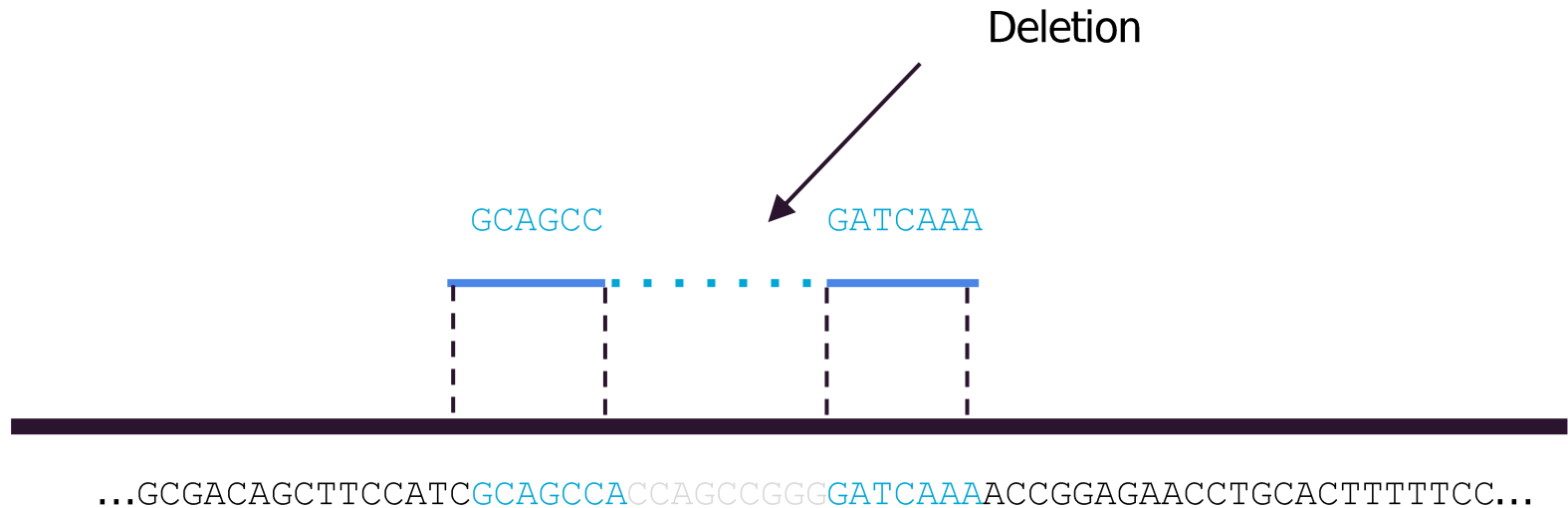


...GCGACAGCTTCCATCGCAGCCACCAGCCGGGGATCAAAACCGGAGAACCTGCACTTTTTTCC...



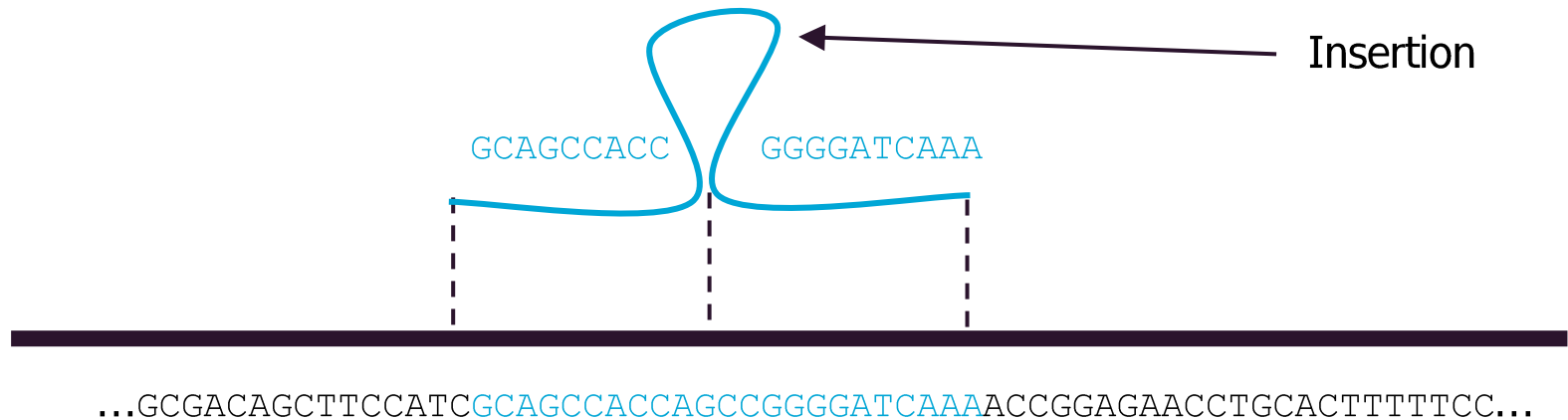
# Mapping read to a reference (4)

allows for detecting variation



# Mapping read to a reference (5)

allows for detecting variation



# Types of variants

Single nucleotide variant

ATTGGCCTTAACCC**C**CCGATTATCAGGAT  
ATTGGCCTTAACCC**T**CCGATTATCAGGAT

Insertion–deletion variant

ATTGGCCTTAACCC**GAT**CCGATTATCAGGAT  
ATTGGCCTTAACCC**---**CCGATTATCAGGAT

Block substitution

ATTGGCCTTAAC**CCCC**GATTATCAGGAT  
ATTGGCCTTAAC**AGTG**GATTATCAGGAT

Inversion variant

ATTGGCCTT**AACCCCCG**ATTATCAGGAT  
ATTGGCCTT**CGGGGGTT**ATTATCAGGAT

Copy number variant

ATT**GGCCTTAGGCCTTA**ACCCCCGATTATCAGGAT  
ATT**GGCCTTA**-----ACCTCCGATTATCAGGAT

Structural variants

# True variant or a sequencing error?

Statistics helps!

Reference: GCGACAGCTTCAATCGCAGCCACCAGCCGGGGATCAAAACCGGAGAACCTGCACTTTTTTCC  
Sample: GCGACAGCTTCGATCGCAGCCACCAGCCGGGGATCAAAACCGGAGAACCTGCACTTTTTTCC  
Read 1 ACAGCTTCGATCGCAGCCACCAG  
Read 2 CTTCGATCGCAGCCACCAGCCGGG  
Read 3 CGATCGCAGCCACCAGCCGGGGATC  
Read 4 CCAGCCGGGGATCAAAACCGGA  
Read 5 GGGGATCAAAACCGGAGAACAT  
Read 6 TCAAAACCGGAGAACCTGCACTTTT  
Read 7 AGAACCTGCACTTTTTTCC

# Variant Call Format (VCF) file

- Reports changes w.r.t. to reference genome
- Contains usually millions of variants
- Describes for each variant:
  - Reference sequence (e.g. 'A')
  - Alternate sequence (e.g. 'G')
  - Genotype (homozygote/heterozygote for variant)
  - Quality score
  - and more...



# Variant Call Format (VCF) file

- The first 25 lines of a VCF file of a trio (father, mother, child):

```
##fileformat=VCFv4.0
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele, ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot_data/technical/reference/ancestral_alignments/README"
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=HM2,Number=0,Type=Flag,Description="HapMap2 membership">
##INFO=<ID=HM3,Number=0,Type=Flag,Description="HapMap3 membership">
##reference=human_b36_both.fasta
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##INFO=<ID=GP,Number=1,Type=String,Description="GRCh37 position(s)">
##INFO=<ID=BN,Number=1,Type=Integer,Description="First dbSNP build #">
##INFO=<ID=NR,Number=0,Type=Flag,Description="No dbSNP 132 map weight=1 rs number assigned to position">
##INFO=<ID=OR,Number=1,Type=String,Description="Previous rs number">
##INFO=<ID=MP,Number=0,Type=Flag,Description="Maps to multiple positions on GRCh37">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA12891 NA12892 NA12878
1 52066 rs28402963 T C . PASS AA=C;DP=84;GP=1:62203;BN=125 GT:GQ:DP 1/0:44:23 1/0:43:20 1/0:70:36
1 695745 rs72631875 G A . PASS AA=. ;DP=124;GP=1:705882;BN=130 GT:GQ:DP 1/0:100:34 0/0:62:20 1/0:100:56
1 742429 rs3094315 G A . PASS AA=g;DP=132;HM2;GP=1:752566;BN=103 GT:GQ:DP 1/1:100:38 1/1:59:30 1/1:100:44
1 742584 rs3131972 A G . PASS AA=a;DP=160;HM3;GP=1:752721;BN=103 GT:GQ:DP 1/1:100:50 1/1:100:33 1/1:100:60
1 744366 rs3115859 G A . PASS AA=g;DP=127;GP=1:754503;BN=103 GT:GQ:DP 1/1:80:31 1/1:100:34 1/1:100:45
1 746243 rs3131963 T A . PASS AA=t;DP=105;GP=1:756380;BN=103 GT:GQ:DP 1/1:52:29 1/1:43:24 1/1:100:56
1 746775 rs6699990 A G . PASS AA=N;DP=120;GP=1:756912;BN=116 GT:GQ:DP 0/1:100:34 0/0:89:30 0/0:100:46
1 747503 rs3115853 G A . PASS AA=- ;DP=113;GP=1:757640;BN=103 GT:GQ:DP 1/1:100:37 1/1:61:25 1/1:100:27
1 747597 rs4951929 C T . PASS AA=c;DP=129;GP=1:757734;BN=111 GT:GQ:DP 1/1:100:36 1/1:86:28 1/1:100:56
1 747799 rs4951862 C A . PASS AA=c;DP=97;GP=1:757936;BN=111 GT:GQ:DP 1/1:67:26 1/1:44:13 1/1:100:62
```

↑                      ↑                      ↑  
Father                  Mother                  Daughter



# Variant Call Format (VCF) file

```
##fileformat=VCFv4.0
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele, ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot_data/technical/reference/ancestral_alignments/README">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=HM2,Number=0,Type=Flag,Description="HapMap2 membership">
##INFO=<ID=HM3,Number=0,Type=Flag,Description="HapMap3 membership">
##reference=human_b36_both.fasta
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##INFO=<ID=GP,Number=1,Type=String,Description="GRCh37 position(s)">
##INFO=<ID=BN,Number=1,Type=Integer,Description="First dbSNP build #">
##INFO=<ID=NR,Number=0,Type=Flag,Description="No dbSNP 132 map weight=1 rs number assigned to position">
##INFO=<ID=OR,Number=1,Type=String,Description="Previous rs number">
##INFO=<ID=MP,Number=0,Type=Flag,Description="Maps to multiple positions on GRCh37">
```

```
#CHROM POS ID REF
1 52066 rs28402963
1 695745 rs72631875
1 742429 rs3094315
1 742584 rs3131972
1 744366 rs3115859
1 746243 rs3131963
1 746775 rs6699990
1 747503 rs3115853
1 747597 rs4951929
1 747799 rs4951862
```

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
1	52066	rs28402963	T	C	.	PASS	
1	695745	rs72631875	G	A	.	PASS	
1	742429	rs3094315	G	A	.	PASS	
1	742584	rs3131972	A	G	.	PASS	
1	744366	rs3115859	G	A	.	PASS	
1	746243	rs3131963	T	A	.	PASS	
1	746775	rs6699990	A	G	.	PASS	
1	747503	rs3115853	G	A	.	PASS	
1	747597	rs4951929	C	T	.	PASS	
1	747799	rs4951862	C	A	.	PASS	

```
6
|1:100:44
|1:100:60
5
6
6
7
6
2
```

# dbSNP

- Variant database
  - Assign identifiers to variants
    - Remain constant across different builds of reference genome
  - Population frequencies

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
1	52066	rs28402963		T	C	.	PASS
1	695745	rs72631875		G	A	.	PASS
1	742429	rs3094315		G	A	.	PASS
1	742584	rs3131972		A	G	.	PASS
1	744366	rs3115859		G	A	.	PASS
1	746243	rs3131963		T	A	.	PASS
1	746775	rs6699990		A	G	.	PASS
1	747503	rs3115853		G	A	.	PASS
1	747597	rs4951929		C	T	.	PASS
1	747799	rs4951862		C	A	.	PASS

## Reference SNP (refSNP) Cluster Report: rs3131972

RefSNP	Allele	HGVS Names
Organism: human ( <a href="#">Homo sapiens</a> )	SNV: single nucleotide variation	NC_000001.10:g.752721A>G
Molecule Type: Genomic	<a href="#">Variation Class</a> :	NT_004350.19:g.231353A>G
Created/Updated in build: 103/138	RefSNP Alleles: C/T ( <a href="#">REV</a> )	XR_108280.1:n.-30A>G
Map to Genome Build: <a href="#">37.5/Weight 1</a>	Allele Origin:	
<a href="#">Validation Status</a> :	Ancestral Allele: T	
	Clinical Channel: unknown	
	Clinical Significance: NA	
	<a href="#">MAF/MinorAlleleCount</a> : A=0.321/700	
	MAF Source: 1000 Genomes	

SNP Details are organized in the following sections:

[GeneView](#) [Map](#) [Submission](#) [Fasta](#) [Resource](#) [Diversity](#) [Validation](#)

## Integrated Maps (Hint: click on 'Chr Pos' or 'Contig Pos' column value to see variation in NCBI sequence viewer)

Assembly	Genome Build	Chr	Chr Pos	Contig	Contig Pos	SNP to Chr	Contig allele	Contig to Chr	Neighbor SNP	Map Method
GRCh37.p10	104.0	<a href="#">1</a>	<a href="#">752721</a>	<a href="#">NT_004350.19</a>	<a href="#">231353</a>	Rev	A	Fwd	<a href="#">view</a>	remap
NCBI36	36.3	<a href="#">1</a>	<a href="#">742584</a>	<a href="#">NT_004350.18</a>	<a href="#">231353</a>	Rev	A	Fwd	<a href="#">view</a>	blast
HuRef	36.3	<a href="#">1</a>	<a href="#">21185</a>	<a href="#">NW_001838585.1</a>	<a href="#">1938</a>	Rev	G	Fwd	<a href="#">view</a>	blast
HuRef	104.0	<a href="#">1</a>	<a href="#">21185</a>	<a href="#">NW_001838585.1</a>	<a href="#">1938</a>	Rev	G	Fwd	<a href="#">view</a>	remap
CHM1_1.0	104.0	<a href="#">1</a>	<a href="#">740191</a>	<a href="#">NW_004077988.1</a>	<a href="#">227600</a>	Rev	G	Fwd	<a href="#">view</a>	remap

## GeneView

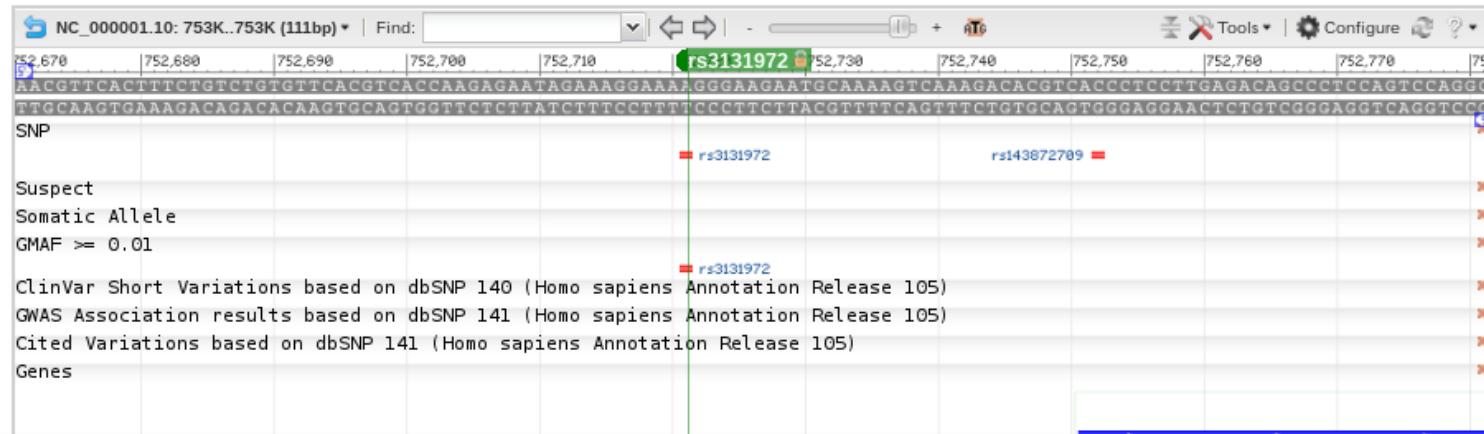
GeneView via analysis of contig annotation: [FAM87B](#) family with sequence similarity 87, member B

View more variation on this gene (click to hide).

☒ Clinical Source: ☐ in gene region ☒ cSNP ☐ has frequency ☐ double hit 

## Primary Assembly Mapping

Assembly		SNP to Chr	Chr	Chr position	Contig	Contig position	Allele
GRCh37.p10		Rev	1	752721	NT_004350.19	231353	A
Gene Model(s)							
Function	mRNA				Protein		
	SNP to mRNA	Accession	Position	Allele change	Accession	Position	Residue change
nearGene-5	NA	XR_108280.1	NA	NA ⇒ NA	NA	NA	NA



ss#	Sample Ascertainment				Genotype Detail				Alleles	
	Population	Individual Group	Chrom. Sample Cnt.	Source	C/C	C/T	T/T	HWP	C	T
<a href="#">ss118438193 YRI</a>			2	IG	1.000				1.000	
<a href="#">ss138899069 ENSEMBL_Venter</a>			2	IG	1.000				1.000	
<a href="#">ss162980826 YRI</a>		Sub-Saharan African	2	IG	1.000				1.000	
<a href="#">ss163702698 CEU</a>		European	2	IG	1.000				1.000	
<a href="#">ss165981005 PGP</a>			2	IG	1.000				0.500	0.500
<a href="#">ss197885385 BUSHMAN_POP2</a>			2	IG	1.000				0.500	0.500
<a href="#">BANTU</a>			2	IG	1.000				0.500	0.500
<a href="#">ss218190360 pilot_1_YRI_low_coverage_panel</a>			118	AF					0.305	0.695
<a href="#">ss230395425 pilot_1_CEU_low_coverage_panel</a>			120	AF					0.858	0.142
<a href="#">ss238114952 pilot_1_CHB+JPT_low_coverage_panel</a>			120	AF					0.742	0.258
<a href="#">ss78643137 HapMap-CEU</a>		European	226	IG	0.681	0.301	0.018	0.439	0.832	0.168
<a href="#">HapMap-HCB</a>		Asian	86	IG	0.488	0.488	0.023	0.150	0.733	0.267
<a href="#">HapMap-JPT</a>		Asian	172	IG	0.547	0.407	0.047	0.439	0.750	0.250
<a href="#">HapMap-YRI</a>		Sub-Saharan African	226	IG	0.044	0.345	0.611	1.000	0.217	0.783
<a href="#">HAPMAP-ASW</a>		African-American	98	IG	0.143	0.531	0.327	0.527	0.408	0.592
<a href="#">HAPMAP-CHB</a>		Asian	82	IG	0.683	0.317		0.527	0.841	0.159
<a href="#">HAPMAP-CHD</a>		Chinese-Americans	170	IG	0.694	0.271	0.035	0.752	0.829	0.171
<a href="#">HAPMAP-GIH</a>		Indian-Americans	176	IG	0.557	0.420	0.023	0.100	0.767	0.233
<a href="#">HAPMAP-LWK</a>		Luhya, Kenya	180	IG	0.122	0.456	0.422	1.000	0.350	0.650
<a href="#">HAPMAP-MEX</a>		Mexican	100	IG	0.580	0.360	0.060	1.000	0.760	0.240
<a href="#">HAPMAP-MKK</a>		Maasai, Kenya	286	IG	0.105	0.559	0.336	0.050	0.385	0.615
<a href="#">HAPMAP-TSI</a>		Toscani, Italia	176	IG	0.739	0.239	0.023	1.000	0.858	0.142
<a href="#">ss97913182 J. Craig Venter</a>			2	IG	1.000				1.000	

Summary	Average Het.+/- std err:	Individual Count	Founders Count	Individual Overlap	Genotype Conflict
	0.436+/-0.167	1213	1008	2	0

# Genome of the Netherlands

- GoNL: 769 persons genotyped

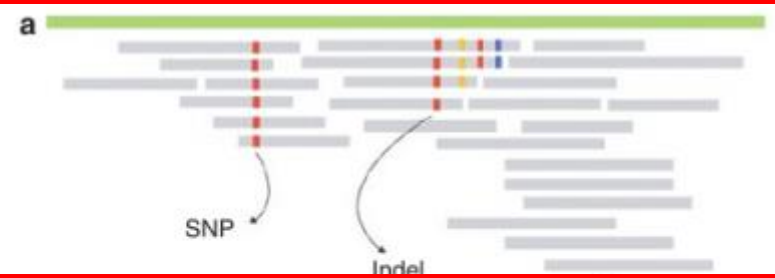
```
1 752721 rs3131972 A G . PASS AC=792;AF=0.794;AN=998;DB;set=SNP
```



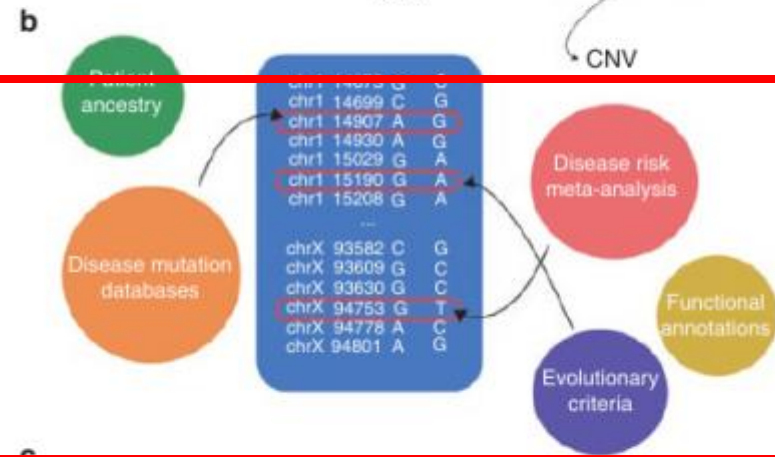


# Overview

Variant  
calling



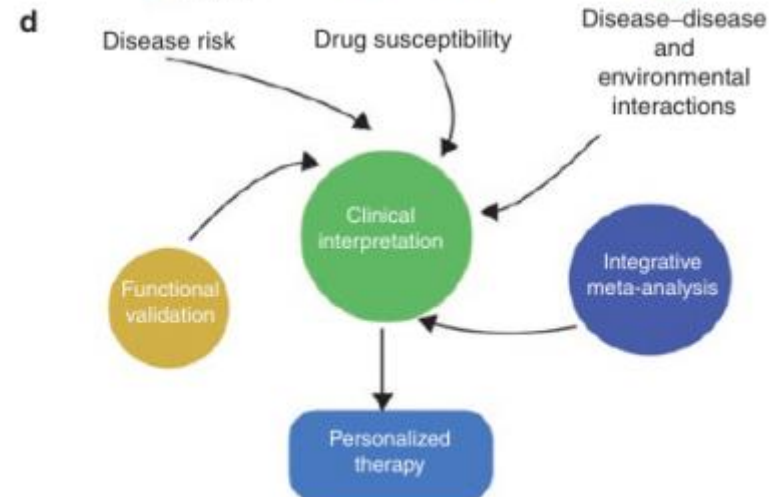
Variant  
filtering



Variant  
relations



Variant  
visualization



# Millions of SNPs....


- Europeans have about 3 million SNPs w.r.t. reference genome.
- Which variant is the cause for a disease?


# Method 1: Look for known 'disease mutations'

- Databases: OMIM, HGMD, **ClinVar**
- Clinvar classifies mutations as 'untested', (likely) 'benign', (likely) 'pathogenic'

Reference SNP (refSNP) Cluster Report: rs80187739

**\*\* With pathogenic,probable-pathogenic,untested allele \*\***

RefSNP
Organism: human ( <a href="#">Homo sapiens</a> )
Molecule Type: Genomic
Created/Updated in build: 131/138
Map to Genome Build: <a href="#">37.5/Weight 1</a>
<a href="#">Validation Status:</a> 

Allele	
<a href="#">Variation Class:</a>	SNV: single nucleotide variation
RefSNP Alleles:	A/C/G/T ( <b>REV</b> )
Allele Origin:	
Ancestral Allele:	G
Clinical Channel:	 <b>VarView</b>
Clinical Significance:	<b>With pathogenic,probable-pathogenic,untested allele</b> <a href="#">[ClinVar]</a>
MAF/MinorAlleleCount:	NA
MAF Source:	

**Gene Model(s)**

NC\_000017.10: 41M..41M (101bp) Find:

41,219,580 41,219,590 41,219,600 41,219,610 41,219,620 rs80187739 41,219,630 41,219,640 41,219,650

GTGGTTTATGTCAGCAGATGCAAGGTATTCTGTAAAGGTTCTTGGTATACCTGTTTTCATAACAACATGAGTAGTCTCTTTCAGCACCAAAATACGTCGTCTACGTTCCATAAGACATTTCCAAGAACCATATGGACAAAAGTATTGTTGTACTCATCAGAGAAGTCT

SNP

Suspect

Somatic Allele

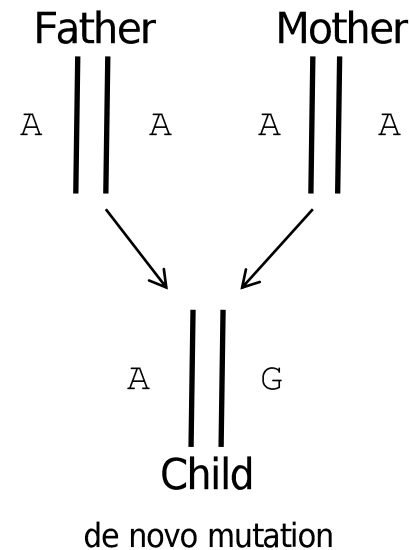
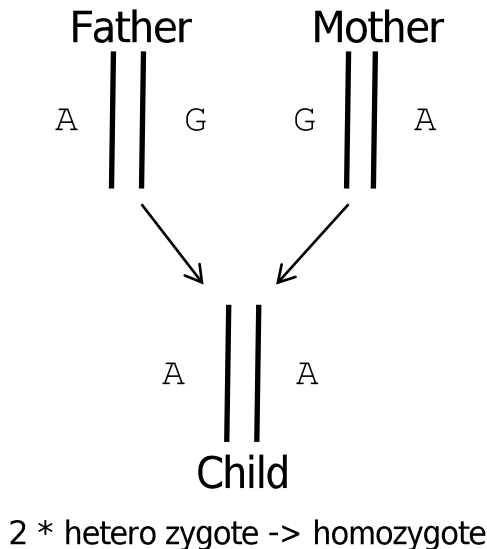
GMAF >= 0.01

ClinVar Short Variations based on dbSNP 140 (Homo sapiens Annotation Release 105)

# Method 2: use of family



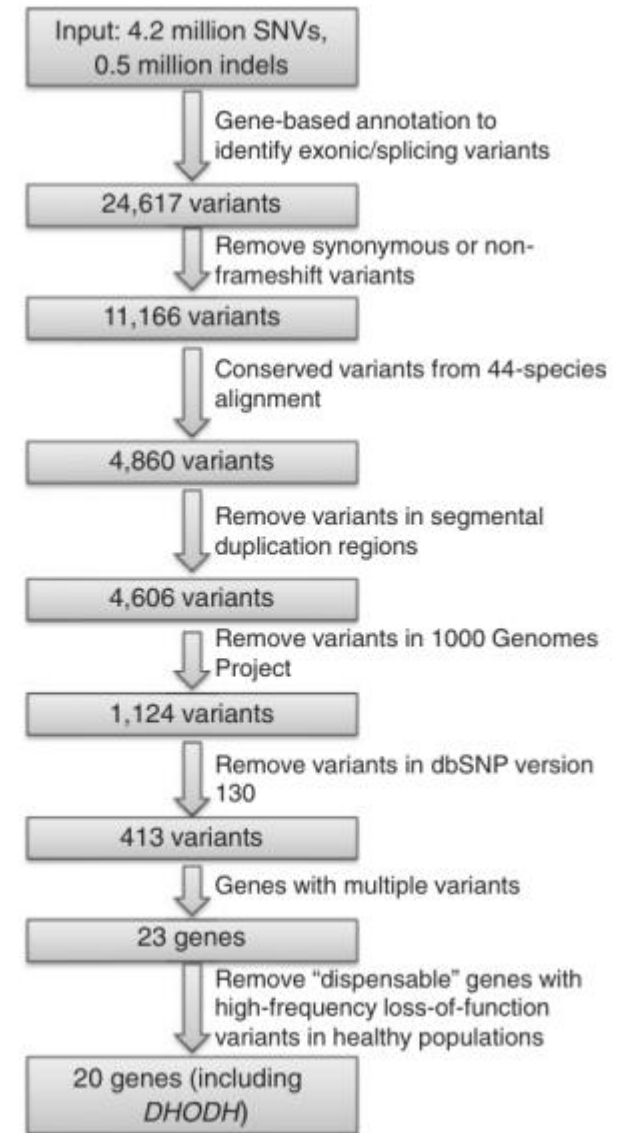
- Trio: Father, Mother, Child
- If child has disease, and not the father or mother, then:



- Possible to go to extended family

# Method 3: Filtering

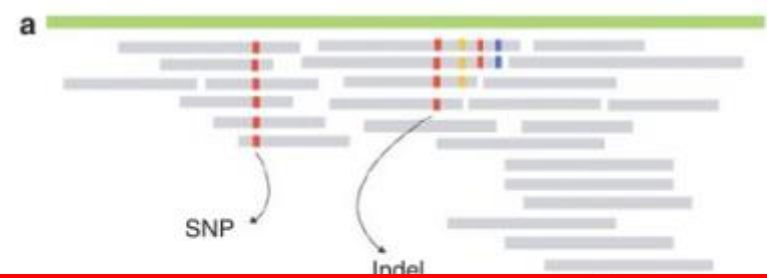
- Is variant:
  - In a gene?
  - Does it cause potential damage?
  - Is sequence conserved in other species?
  - Is it a common, well known variant?
- ANNOVAR pipeline



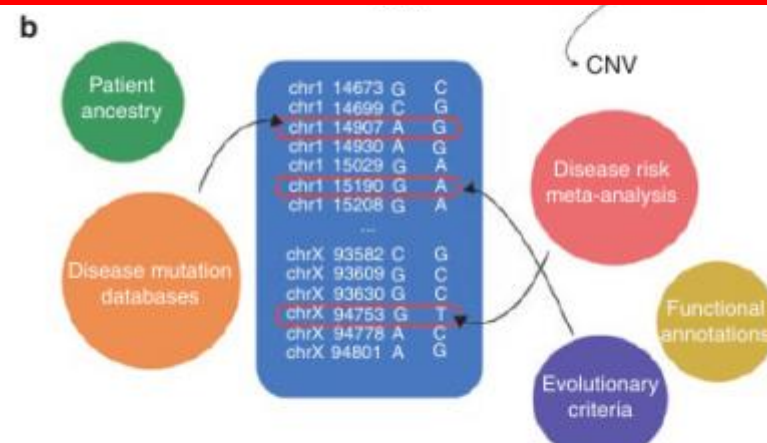


# Overview

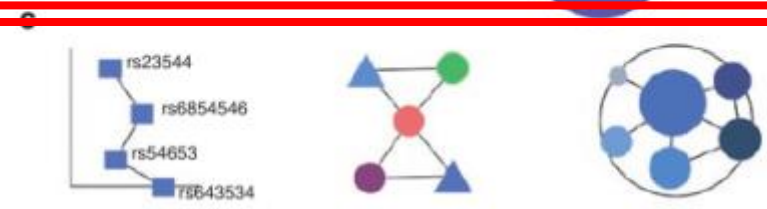
Variant  
calling



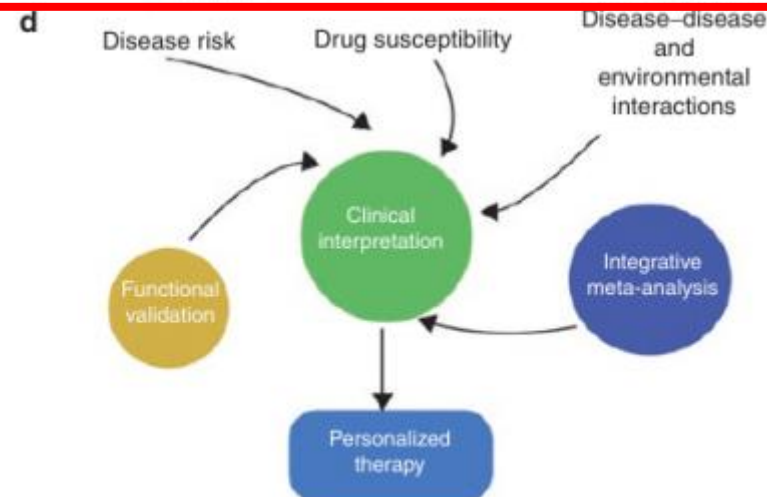
Variant  
filtering



Variant  
relations



Variant  
visualization



# Common versus rare variants

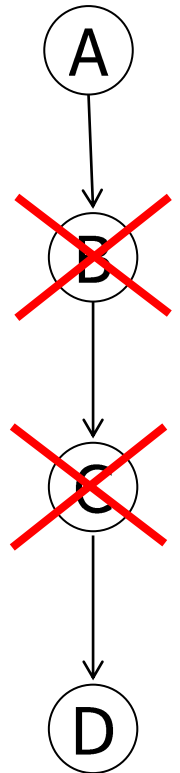
- Early days of GWAS: common variants found for e.g. diabetes
  - but could explain only part of observed heritability
- Hypothesis: most damage due to rare variants
- Problem: rare variants are rare --> not much known about them.

- Are there many other disease-associated mutations in the same gene?

- Burden test: have diseased patients more rare damaging variants in a certain gene than healthy people?

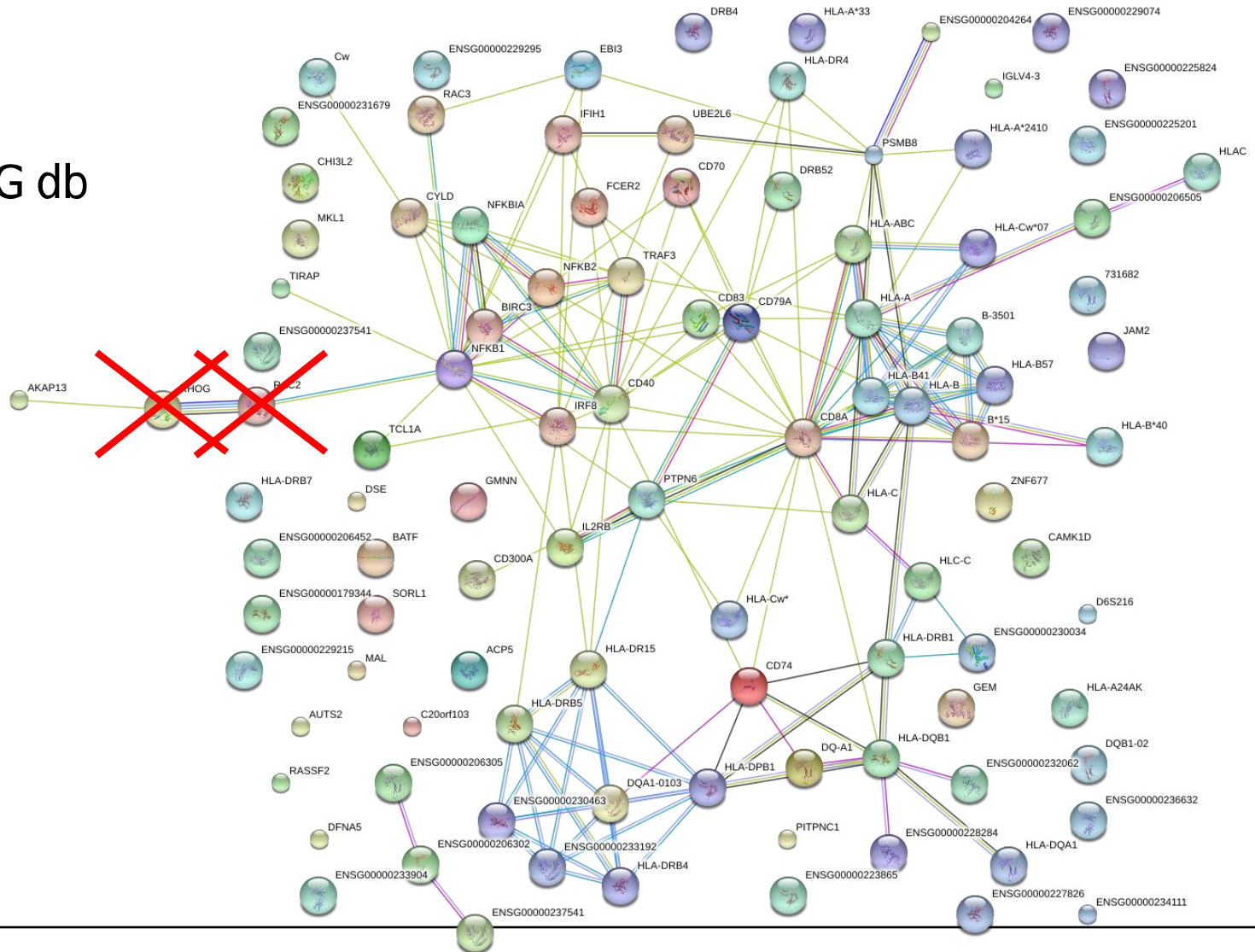
# Method 5: Nearby genes

- Gene relations based on common function
  - Similar activity
  - Similar functional annotations
  - Similar sequence
  - Proteins have physical interaction
  - Etc.
- Hypothesis: damage in genes with related functions can cause similar diseases.



# Method 5: Nearby genes

- STRING db







# Visualizations

- To assist interpretation of variant data

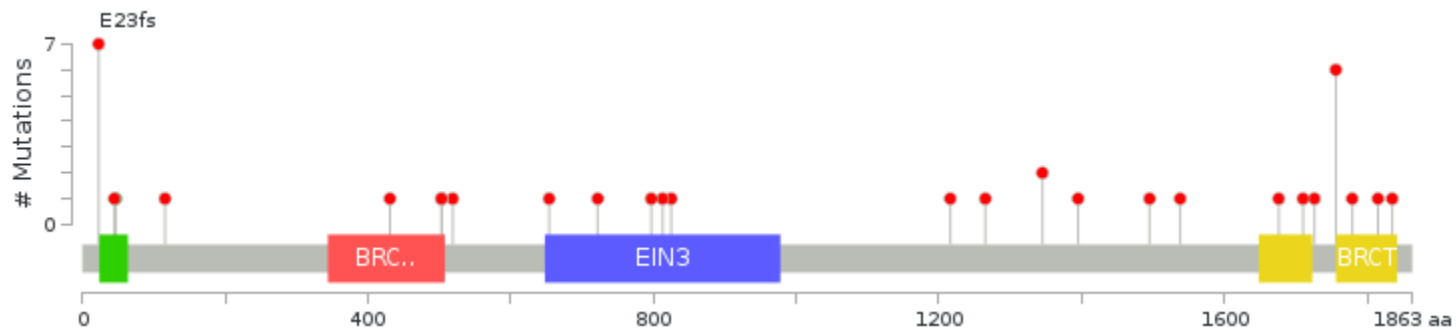
# Visualizations

Chromosome	Position	Ref/variant	Gene	Amino acid change
chr10	126691631	A/T	CTBP2	Thr626Ser
chr11	1016928	G/C	MUC6	Ser1958Thr
chr1	156202173	G/A	PMF1	Gln75Arg
chr12	58220841	T/C	CTDSP2	Asp98Asn
chr12	58220844	T/C	CTDSP2	Glu97Lys
chr14	20692643	T/C	OR11H6	Cys259Arg
chr17	45214648	C/G	CDC27	Gln595Glu
chr17	45214651	G/T	CDC27	Ile594Leu
chr17	45214654	T/C	CDC27	Ala593Thr
chr2	97877292	A/C	ANKRD36	Pro709Gln
chr3	75787876	T/G	ZNF717	Leu300Ile
chr4	190876242	A/G	FRG1	Gly123Glu
chr4	47901476	A/G	NFXL1	Pro246Leu
chr9	68455161	T/C	LOC100287354	Arg94Trp

Table 1: Non-synonymous variations found homozygously in the offspring and heterozygously in both parents

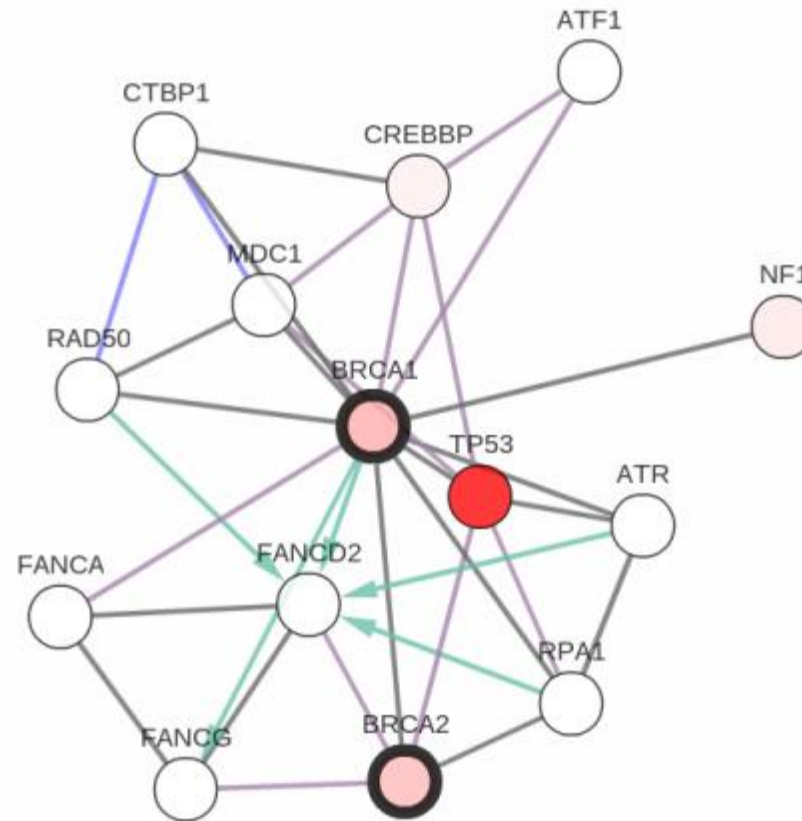
# Visualizations

- From cBioPortal:

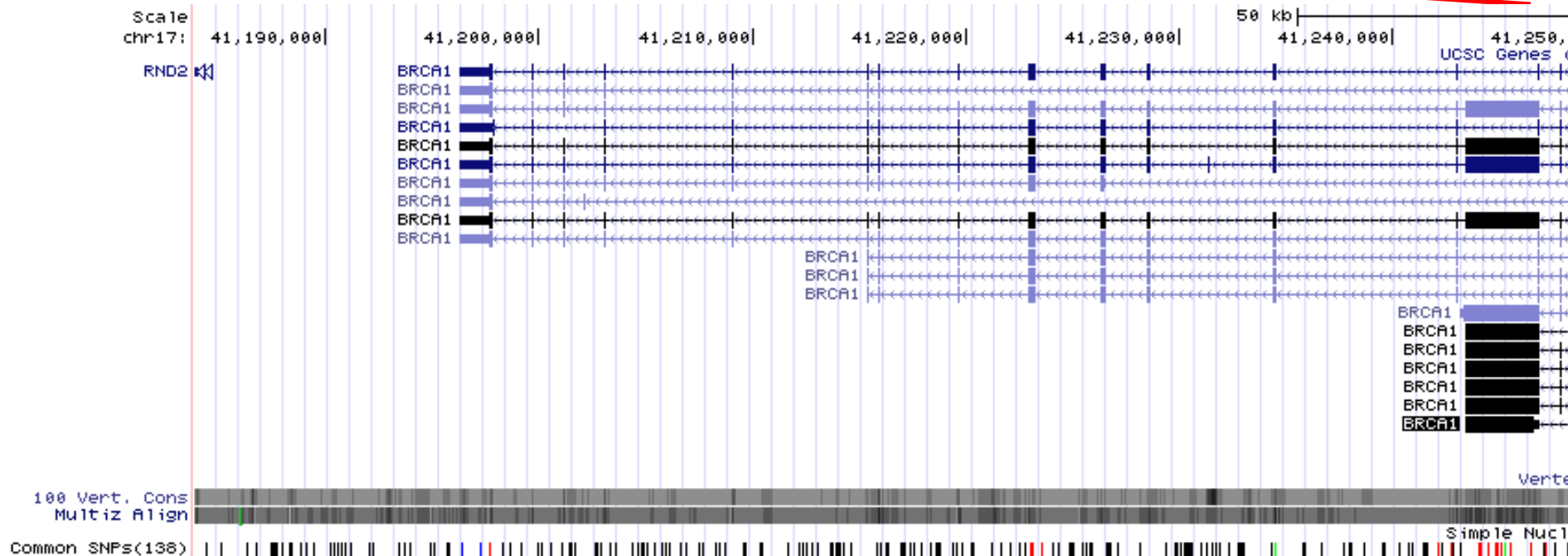
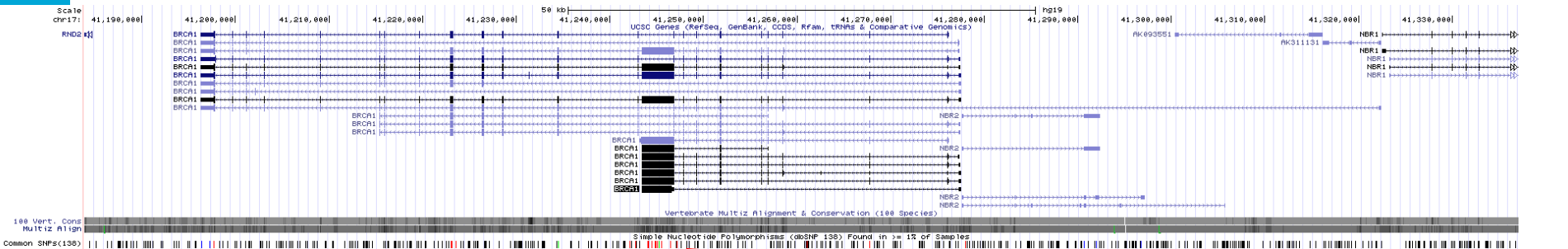


# Visualizations

- From cBioPortal:



# Visualizations







# The project

- Design an application which visualizes and assist in the interpretation of genomic variants
- We are specifically interested in an application that is able to visualize:
  - Trio data
  - Gene interactions



# The project

- Method 1 (known variants):
  - We will give access to dbSNP SQL database
- Method 2 (family relations):
  - We will supply trio variant calls in VCF format
- Method 3 (ANNOVAR):
  - We will supply files containing information useful for filtering
- Method 4 (SNP relations):
  - We will supply a file containing gene locations (GFF3 format)
- Method 5 (gene relations):
  - We will give access to STRING SQL database

# Goal

- Create a useful visualization application which
  - helps clinical geneticists
  - in the interpretation of variants
  - by making use of the aforementioned data sources

# Voor volgende keer (6 mei)

- Lees materiaal voor volgende keer (synthetische biologie)
- Maak samenvatting in en lever deze in voor het begin van het college.

