

# NGS in practice

---

*where things go wrong  
( practical experiences )*

Your last  
mistake  
is your  
best  
teacher

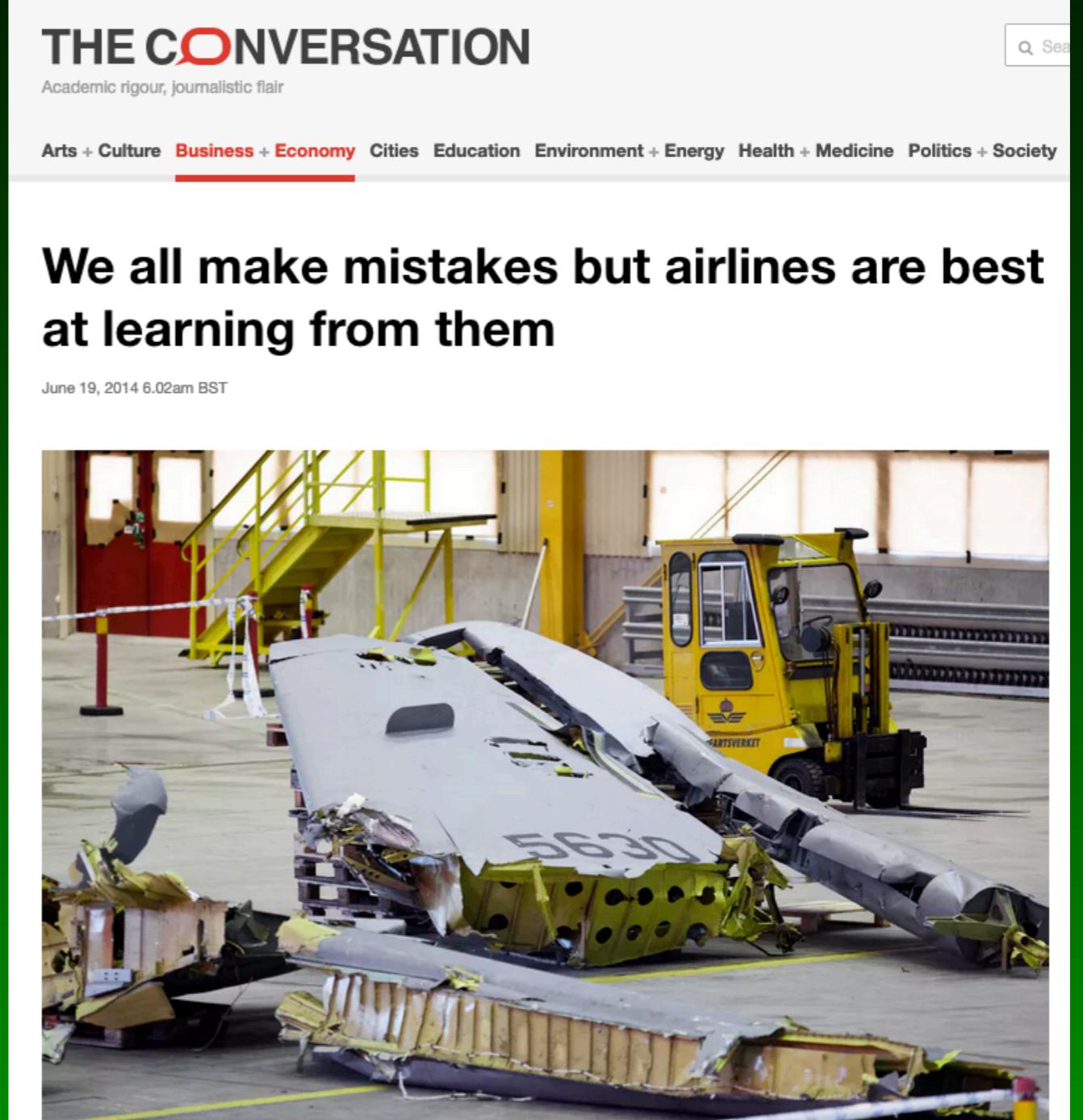


*i.c.w. Anna Benet-Pages  
Gijs Santen  
Christian Gilissen*

*Johan den Dunnen*



# Accept mistakes



# Accept mistakes

---

- report incidents and near incidents (!)
- perform extensive analysis
- do not focus on blame
- recommend changes

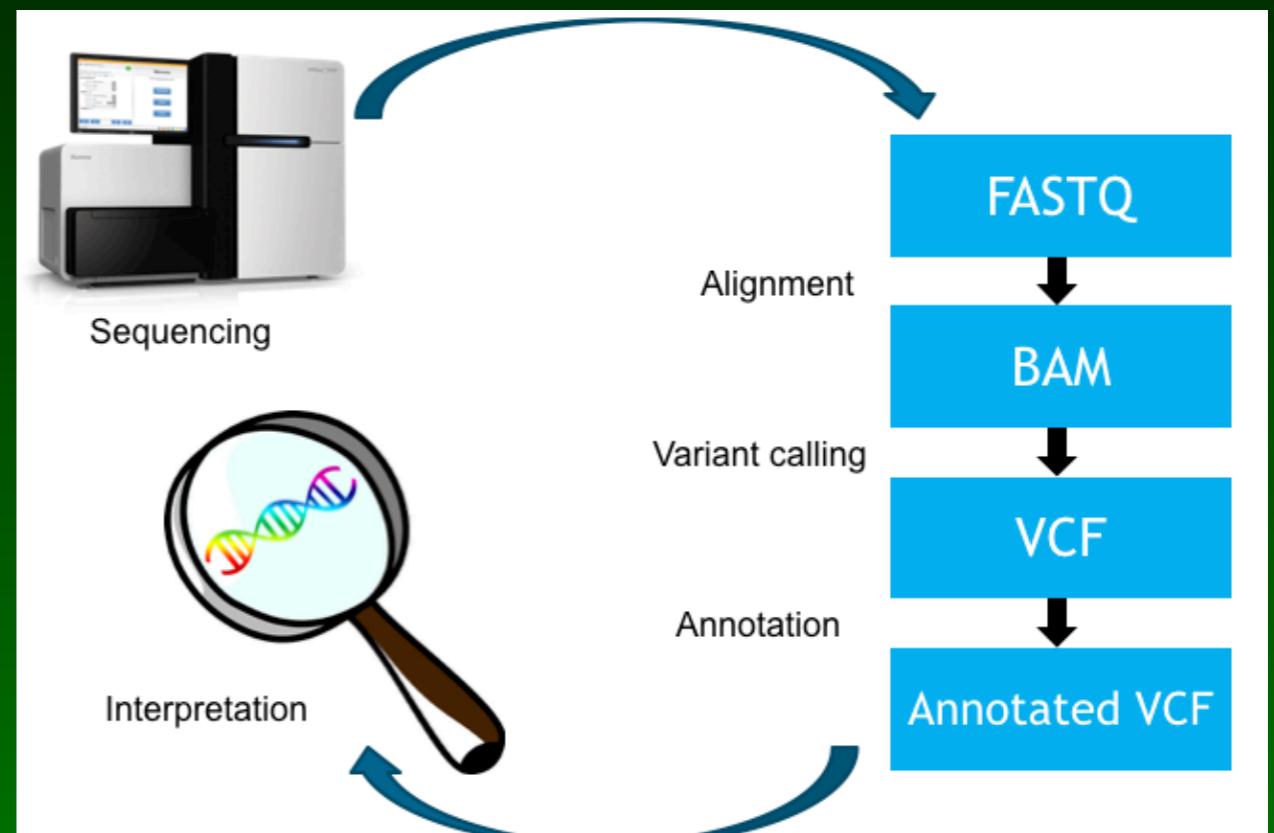
Anything  
that can  
go wrong...  
will go  
**—wrong—**

*unfortunately*



# NGS workflow

- sampling
- sample preparation
- sequencing
- mapping reads
- call variants
- annotate variants
- variant interpretation



*image from Christian Gilissen*

# Our first ever WES

---

- capture exome X-chromosome only
- X-linked disease  
*TOD (Terminal Osseus Dysplasia)*
- used pipeline with standard settings
  - threshold for calling variants***  
*number of reads before variant is called (10 or more)*
  - range of accepted allele frequency***

***No possible causative variant remaining***

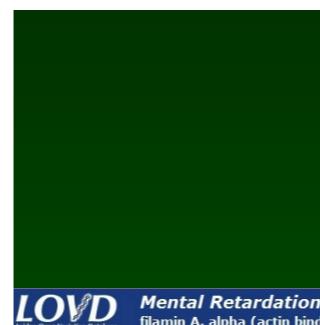
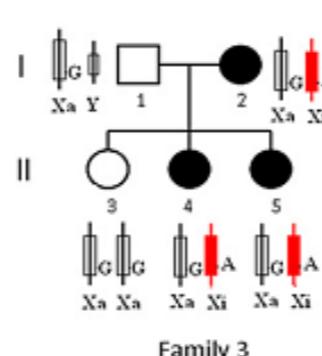
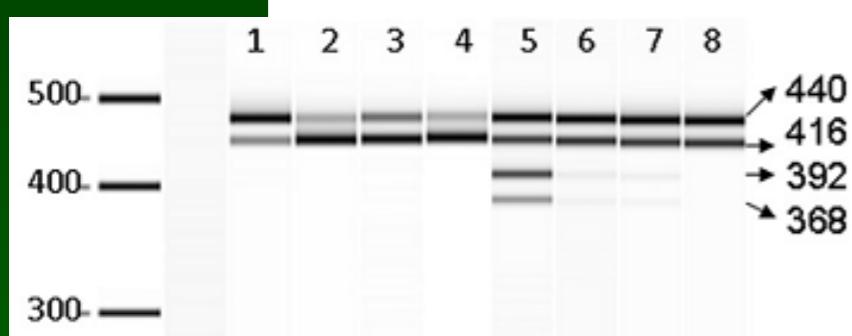
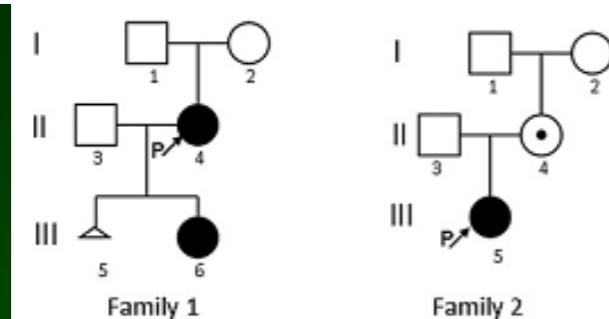
# X-linked TOD

## REPORT

The American Journal of Human Genetics 87, 146–153, July 9, 2010

### Terminal Osseous Dysplasia Is Caused by a Single Recurrent Mutation in the *FLNA* Gene

Yu Sun,<sup>1,11</sup> Rowida Almomani,<sup>1,11</sup> Emmelien Aten,<sup>1</sup> Jacopo Celli,<sup>1</sup> Jaap van der Heijden,<sup>1</sup> Hanka Venselaar,<sup>2</sup> Stephen P. Robertson,<sup>3</sup> Anna Baroncini,<sup>4</sup> Brunella Franco,<sup>5,6</sup> Lina Basel-Vanagaite,<sup>7</sup> Emiko Horii,<sup>8</sup> Ricardo Drut,<sup>9</sup> Yavuz Ariyurek,<sup>1,10</sup> Johan T. den Dunnen,<sup>1,10</sup> and Martijn H. Breuning<sup>1,\*</sup>



Detailed view of the LOVD database entry for the FLNA gene. General information includes: Gene name: filamin A, alpha (actin binding protein 280); Gene symbol: FLNA; Chromosome Location: Xq28; Curator: Johan den Dunnen; Database reference for citations: Sun et al. 2010, Am J Hum Genet. 87: 146–153; Date of creation: March 06, 2009; Last update: October 24, 2010; Version: FLNA101024; Add sequence variant: Submit a sequence variant; First time submitters: Register here; Reference sequence: coding DNA reference sequence for describing sequence variants FLNA\_NG\_011506.1.gb; Total number of unique DNA variants reported: 84; Total number of individuals with variant(s): 289; Total number of variants reported: 309.

**threshold to call variants set to 10 or more reads  
we missed 9 reads, all variant  
on X-chromosome (male) coverage is half**

# NGS workflow

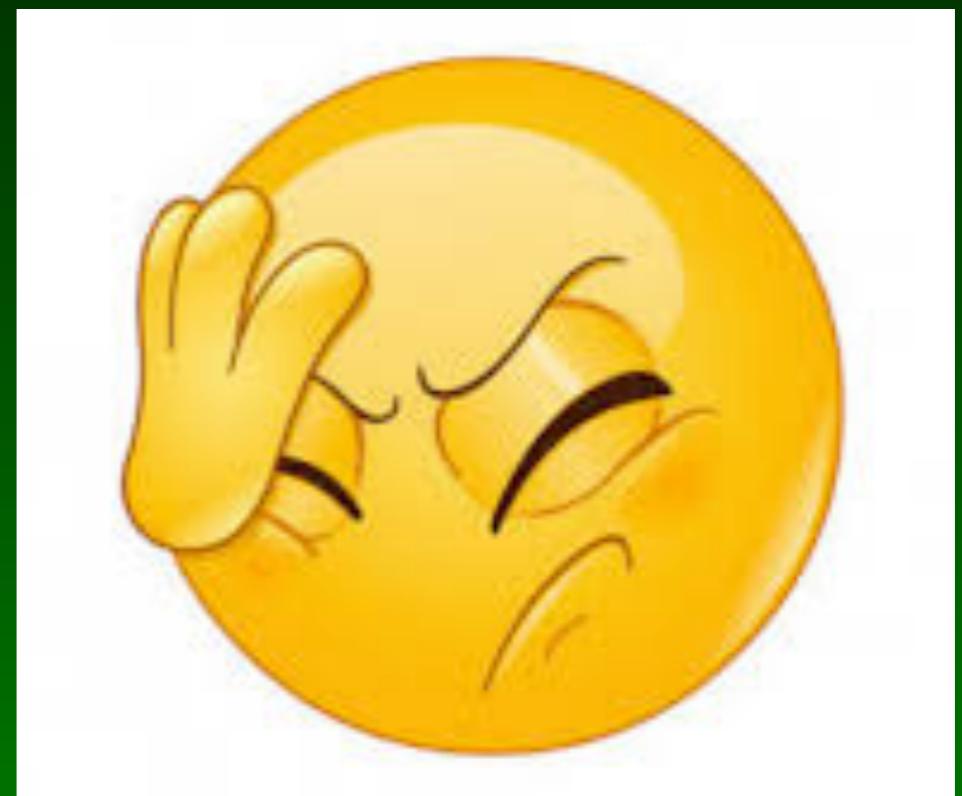
---

- what can go wrong ??

*...suggestions please*



*image from www.caracaschronicles.com*



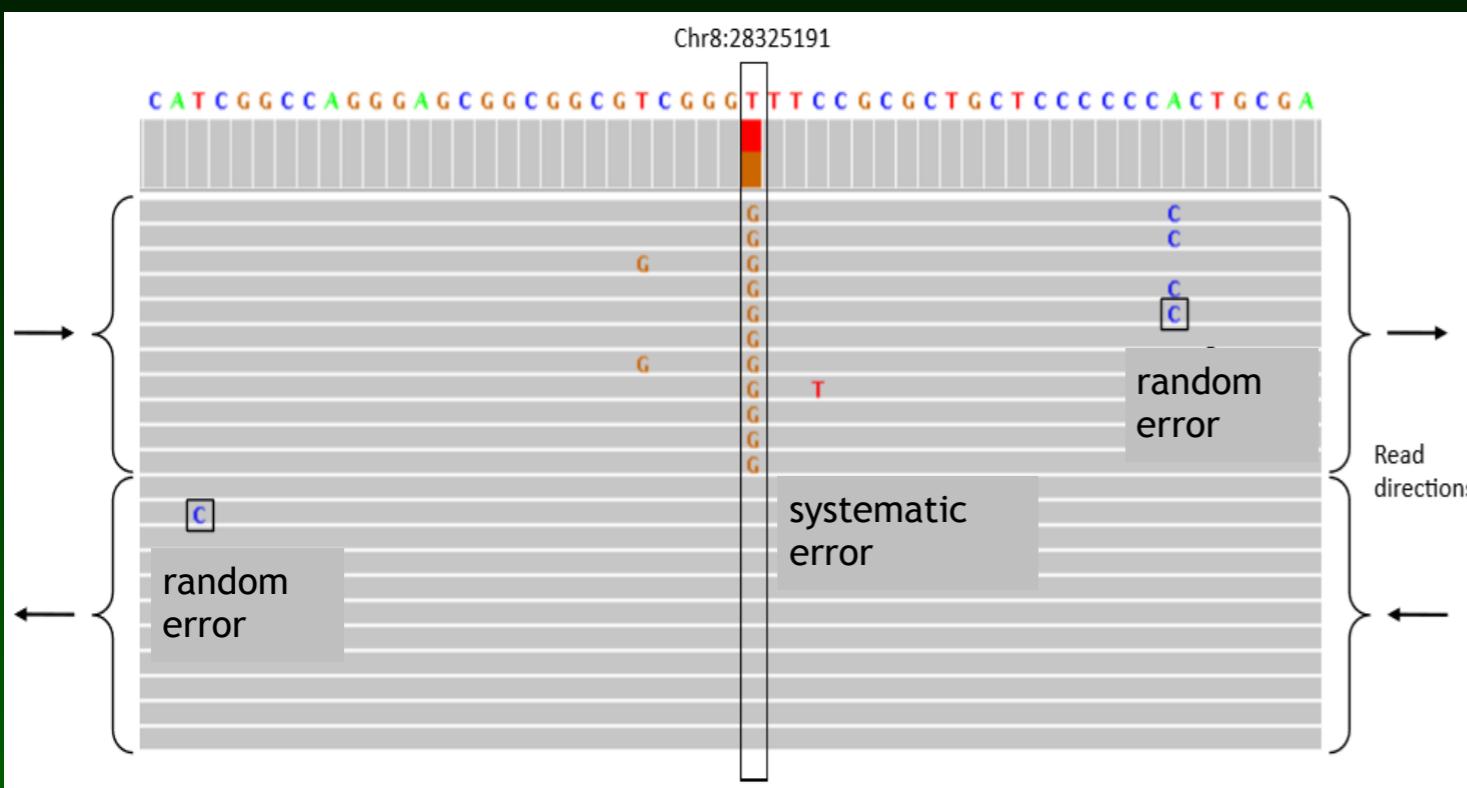
*image from coinspice.io*

# Sampling

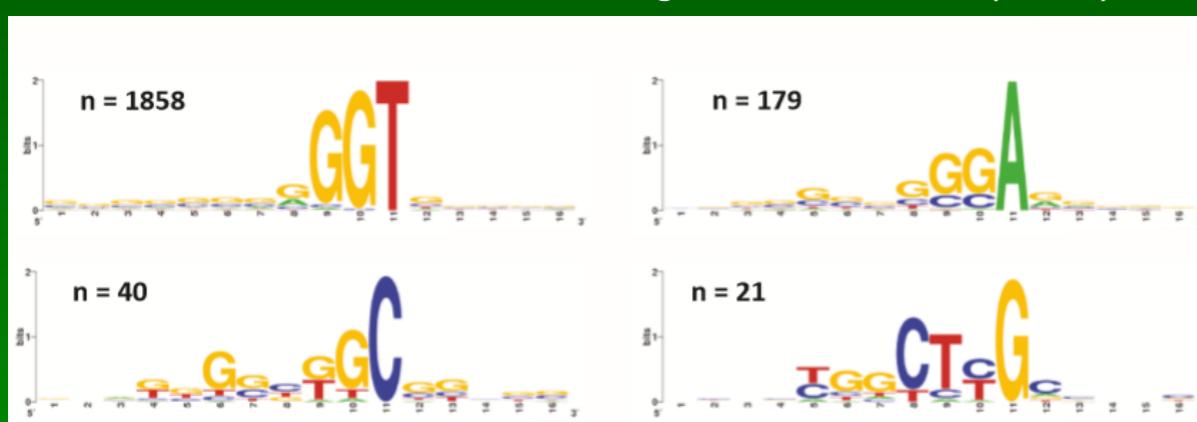
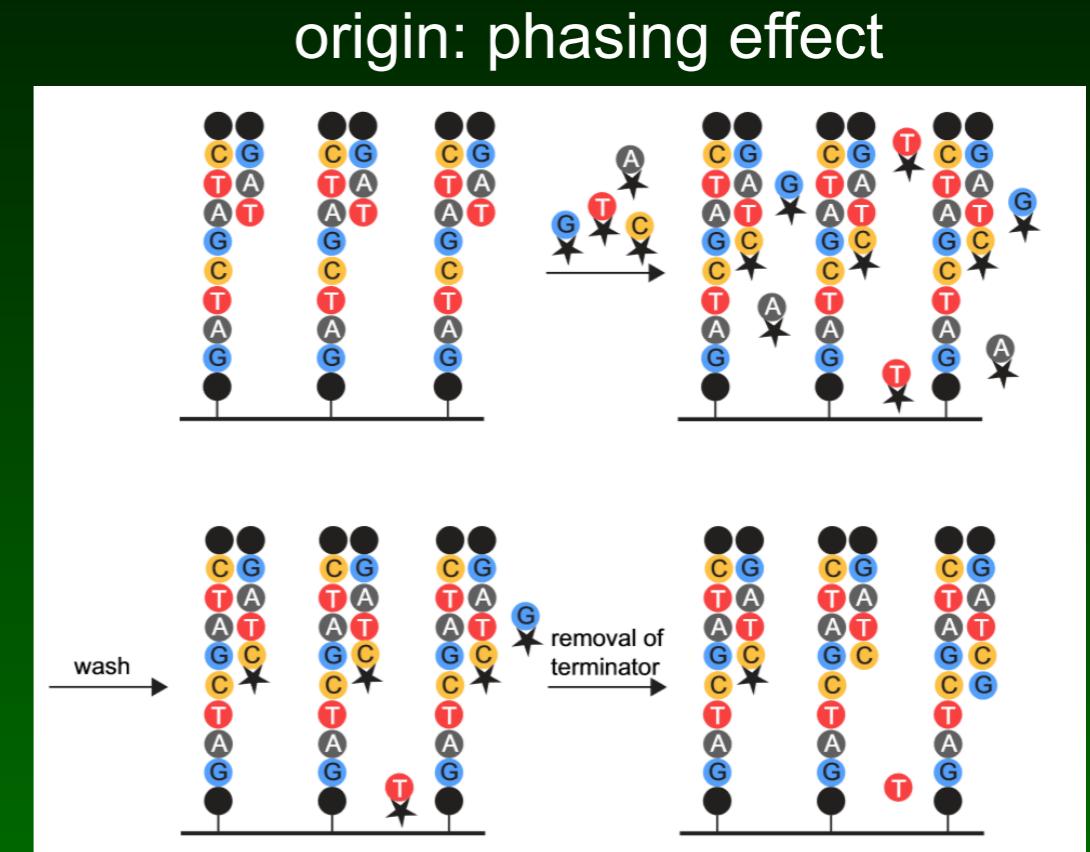
---

- sampling
- sample preparation
- sequencing
  - samples swaps will occur  
*at every step*
  - build in checks  
*female / male*  
*SNV-assay + SNV calls from NGS*  
*process every sample twice*  
*check Mendelian inheritance*

# Systematic errors



images from Pfeifer (2018)



sequence around systematic errors

**mono-nucleotide  
stretch problem**

# QC

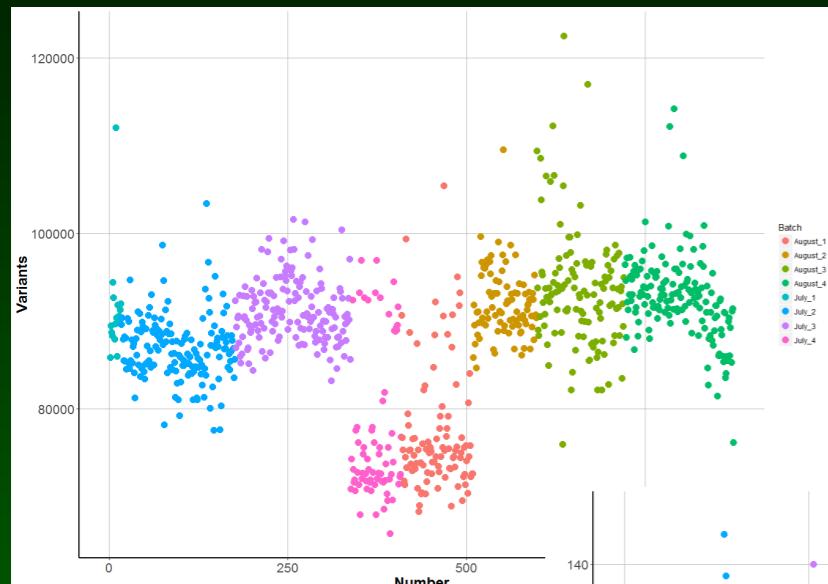
---

- options ?

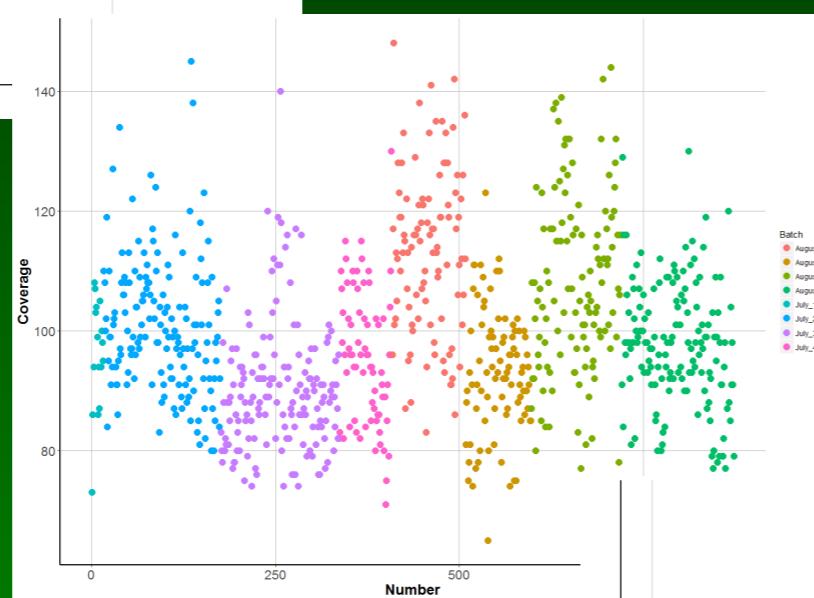


# QC

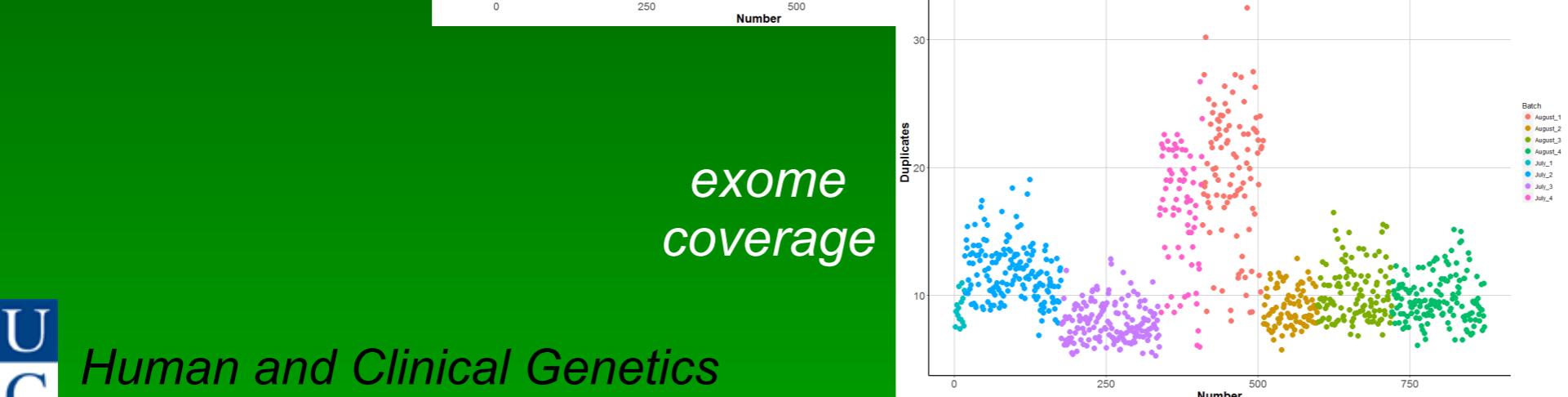
*variants called per exome*



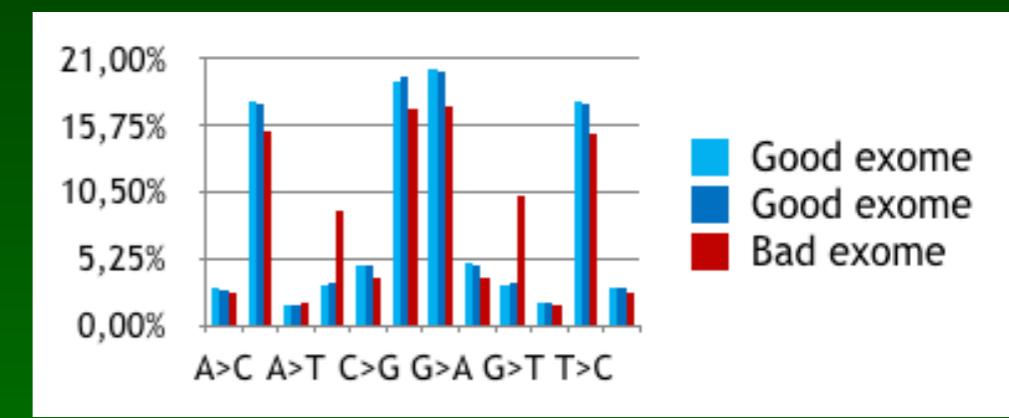
*number of  
duplicate reads*



*exome  
coverage*



**% mapped reads  
% reads on target  
transversion / transition rate  
contamination other organisms  
compare with all previous analysis**



© images from  
Christian Gilissen

# Updated pipeline



New pipeline / component:  
first re-run standard controls and  
compare results

## Contents

1. What are the GATK Best Practices?
2. Analysis phases
3. Experimental designs
4. Workflow scripts provided as reference implementations
5. Scope and limitations
6. What is *not* GATK Best Practices?
7. Beware legacy scripts

Best Practices

Introduction to the GATK Best Practices

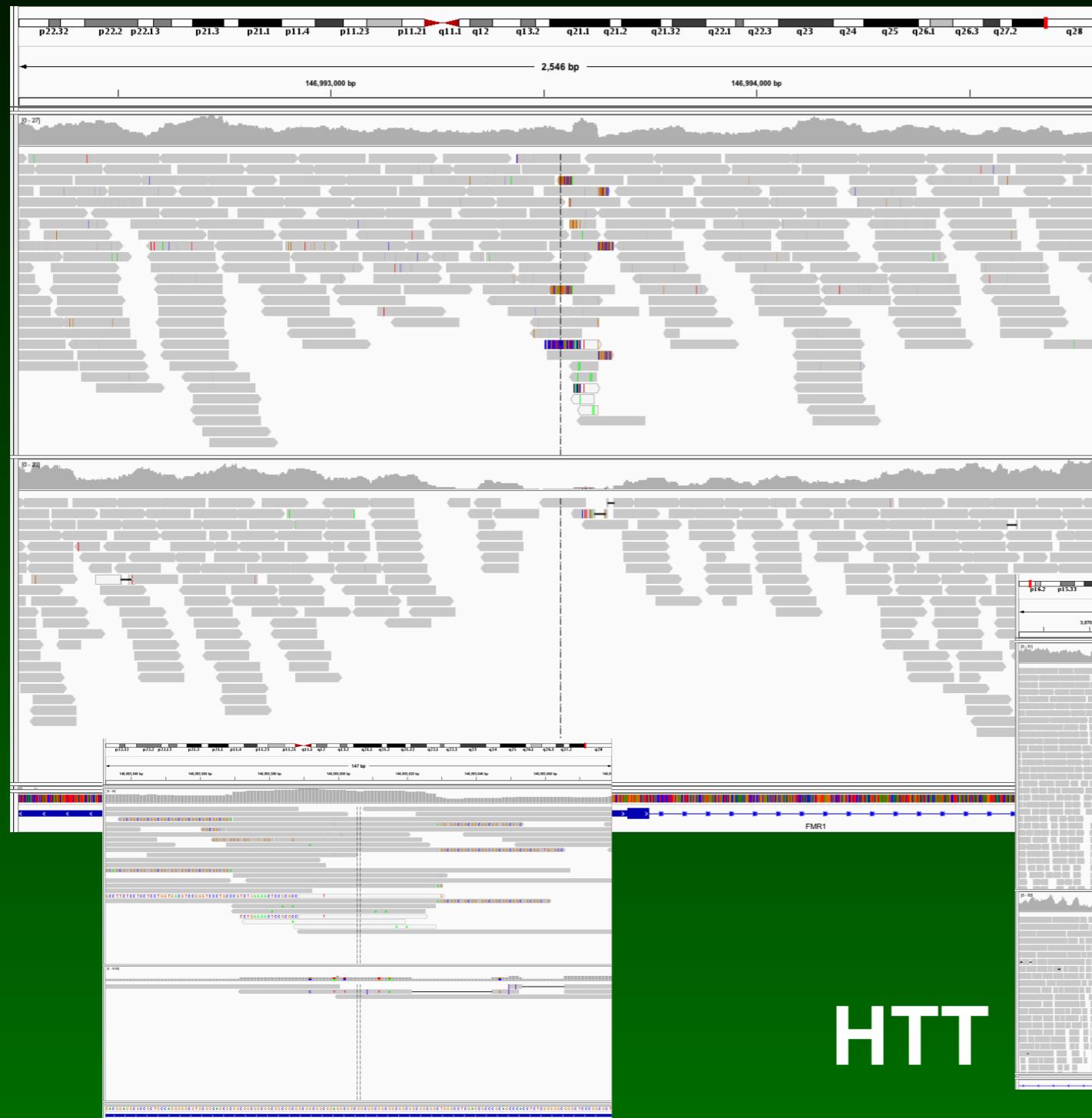
Best Practices Workflows | Created 2018-01-09 | Last updated 2018-01-09

# Uneven coverage - CG%

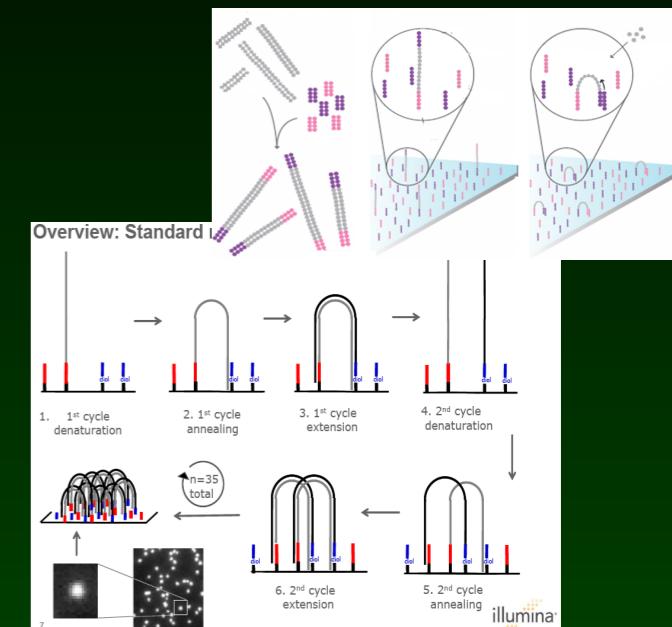
FMR1

without  
PCR

with  
PCR

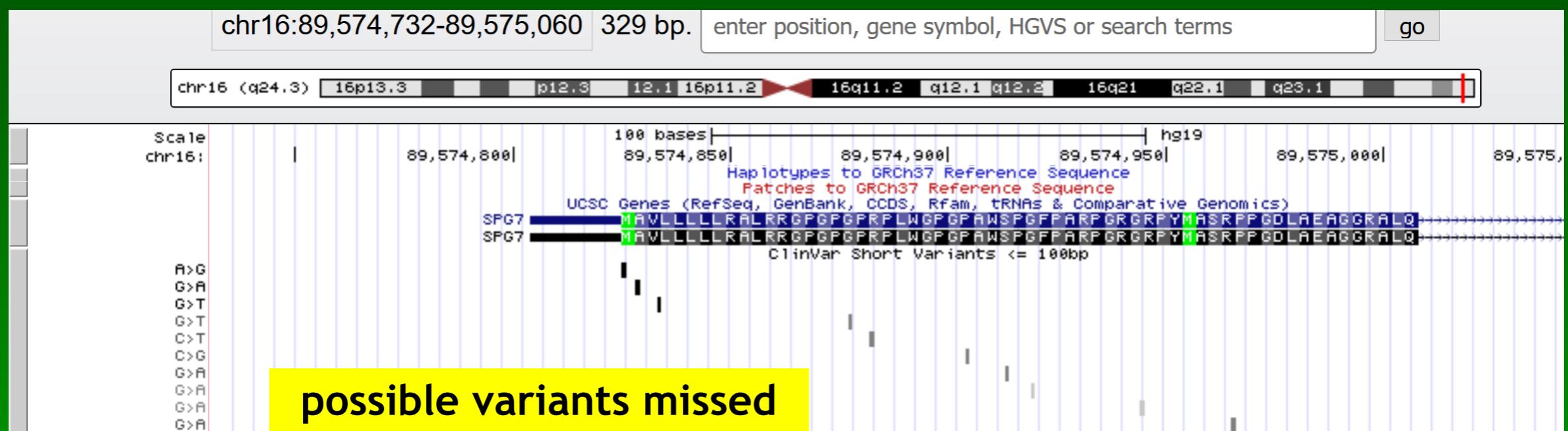
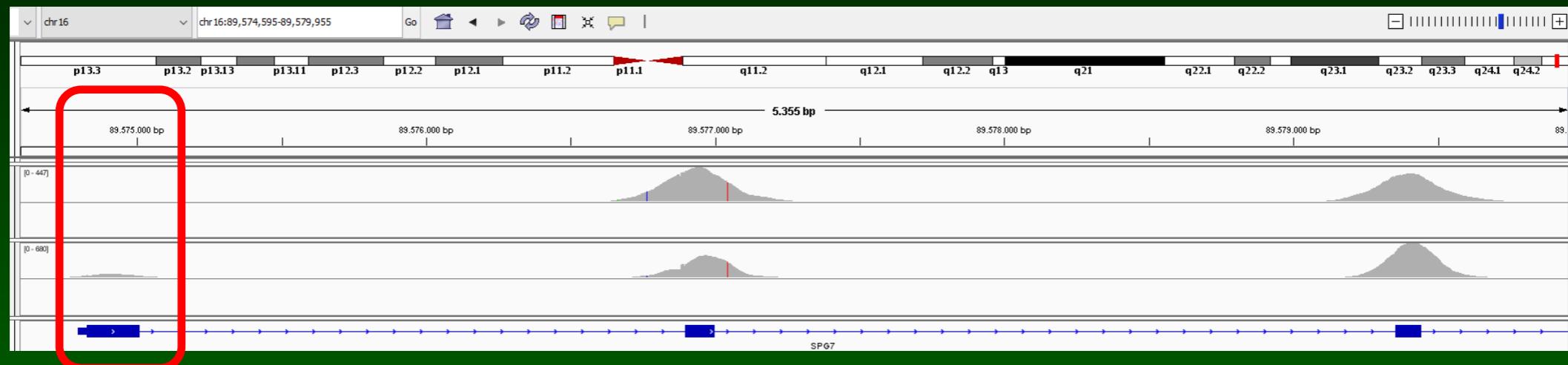


HTT



# Uneven coverage

*sample preparation  
gene panel / exome capture*



# Different tools

CNV calling – Tool performance

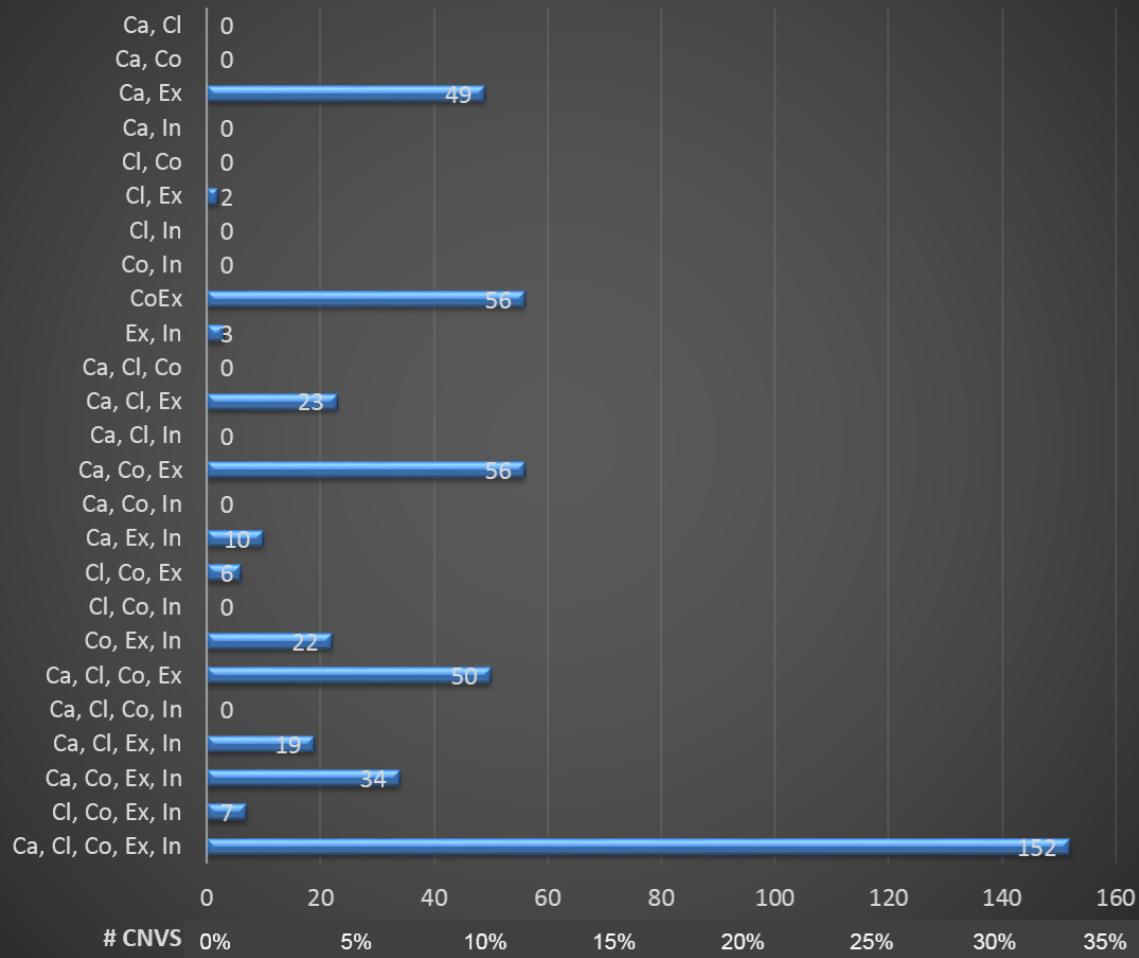


image from Anna Benet-Pages

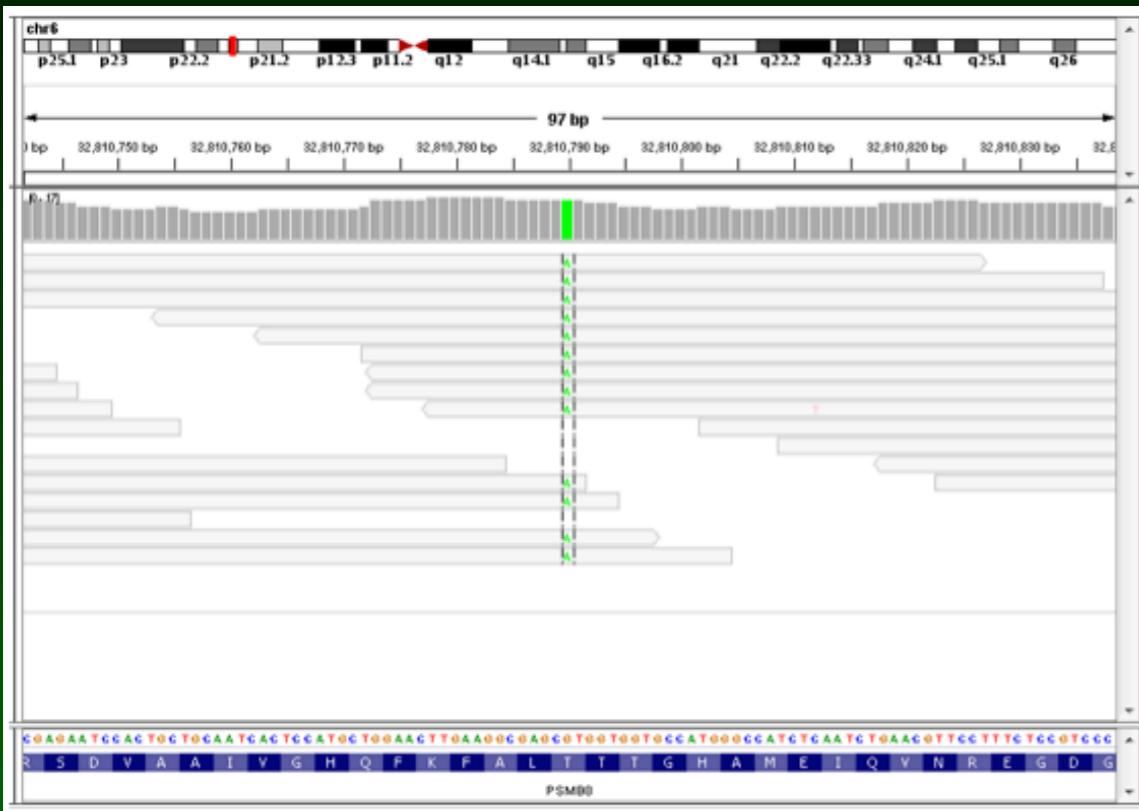
*...different results*

Tool	Mean Calls per Sample
Canoes	~ 1.4
Clamms	~ 5.1
Codex	~ 2.6
ExomeDepth	~ 7.7
Inhouse	~ 4.4

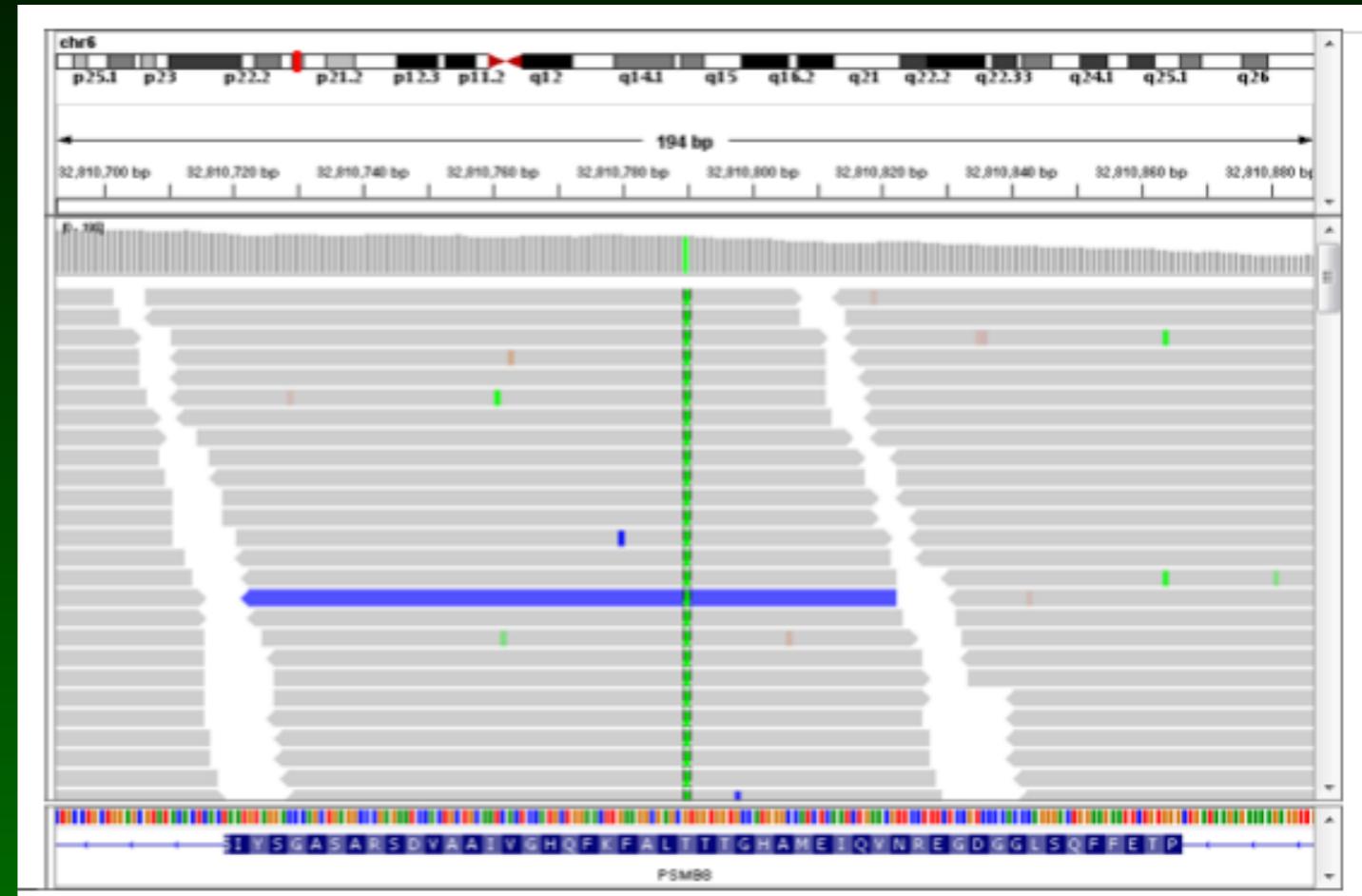
*esp. for CNV detection given coverage  
is variable*

# Mapping reads

images from Christian Gilissen



mapped to reference genome  
+ unmapped contigs  
+ alternate haplotypes



mapped to reference genome  
+ unmapped contigs  
(no alternate haplotypes)

# Mapping

*true situation*

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
+  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG

*0% variant*

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
+  
GATTGGGTAGAGCGATGG  
GATTGGGTAGAGCGATGG

*50% variant*

( e.g. globin genes )

# Non-unique<sup>2</sup>

---

*goes to dust bin*

*GATTGGGCAGAGCGATGG*

*GATTGGGCAGAGCGATGG*

*no coverage*

*no coverage*

*"deletion"*

*"deletion"*

# Non-unique

---

*map to first position*

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
**GATTGGGTAGAGCGATGG**  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
**GATTGGGTAGAGCGATGG**  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG

*25% variant*

*( no data )*

# Non-unique

*map to all positions*

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG

*25% variant*

*25% variant*

*variant missed / found twice*

# Non-unique

*probabalistic mapping*

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG

GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG

*25% variant*

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG

*25% variant*

*or*

GATTGGGCAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGTAGAGCGATGG  
GATTGGGCAGAGCGATGG

*50% variant*

*0% variant*

GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG  
GATTGGGCAGAGCGATGG

*0% variant*

*or*

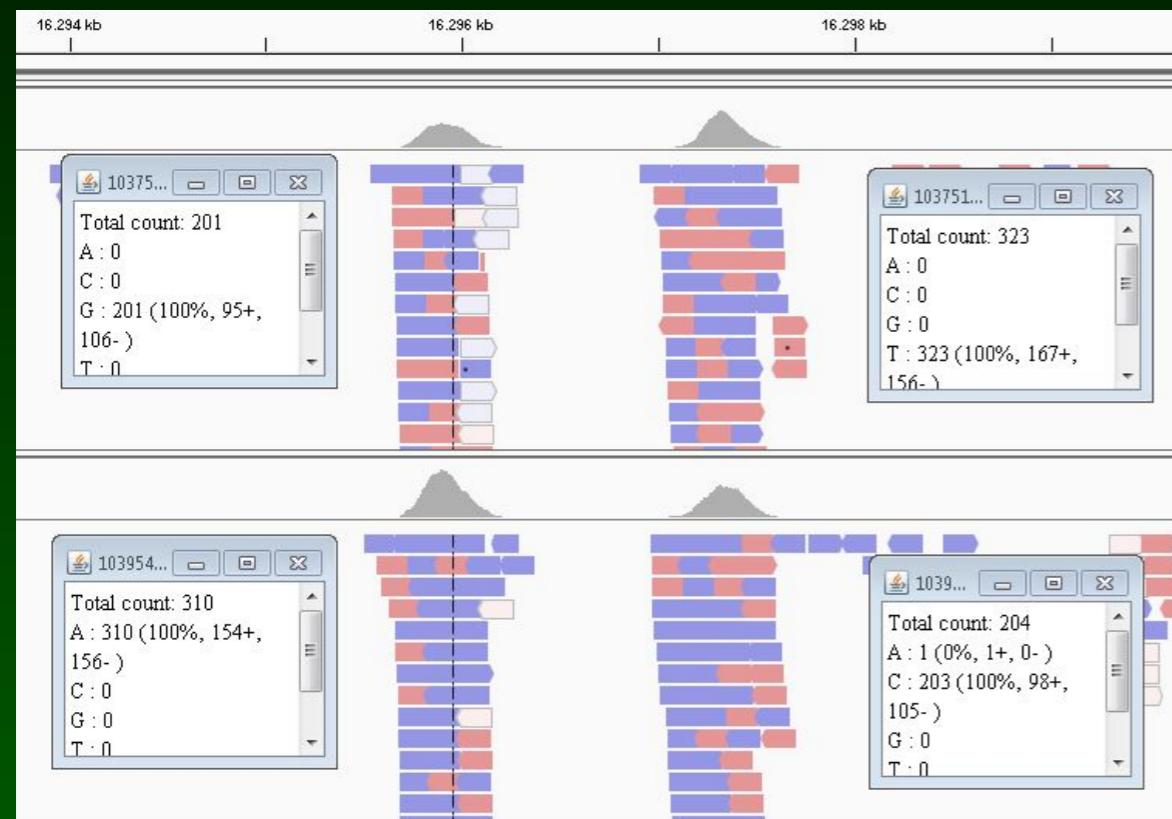
*50% variant*



# Gene conversion

*Pseudoxanthoma elasticum:* c.4182del (p.Lys1394Asnfs\*9)

exon 9 coverage 0.5



ABCC6P1



image from Anna Benet-Pages

**gene conversion causes reads to map to ABCC6P1**

- no variant in ABCC6*
- no variant pseudogene*
- coverage differs !!*

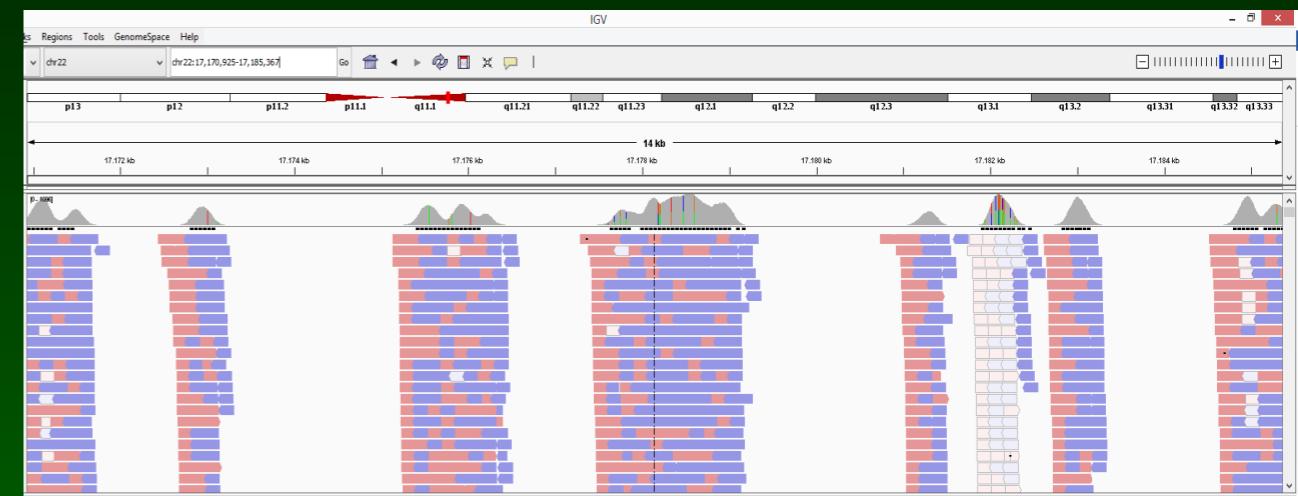
# Homologous sequences

VWF gene causes "*von Willebrand Syndrome*"

mechanism: SNVs (~90%), CNVs (~10%), 6-335 bp gene conversions



VWF - von Willebrand factor gene  
chr12:6,118,707-6,137,253



VWFP1 - von Willebrand factor pseudogene 1  
chr22:17,170,925-17,185,367

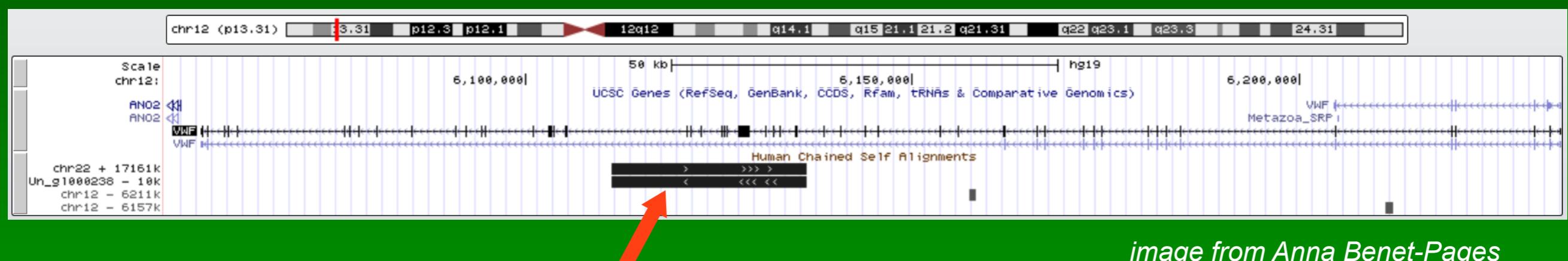


image from Anna Benet-Pages

# Calling variants

whole exome sequencing

call variants in targeted regions

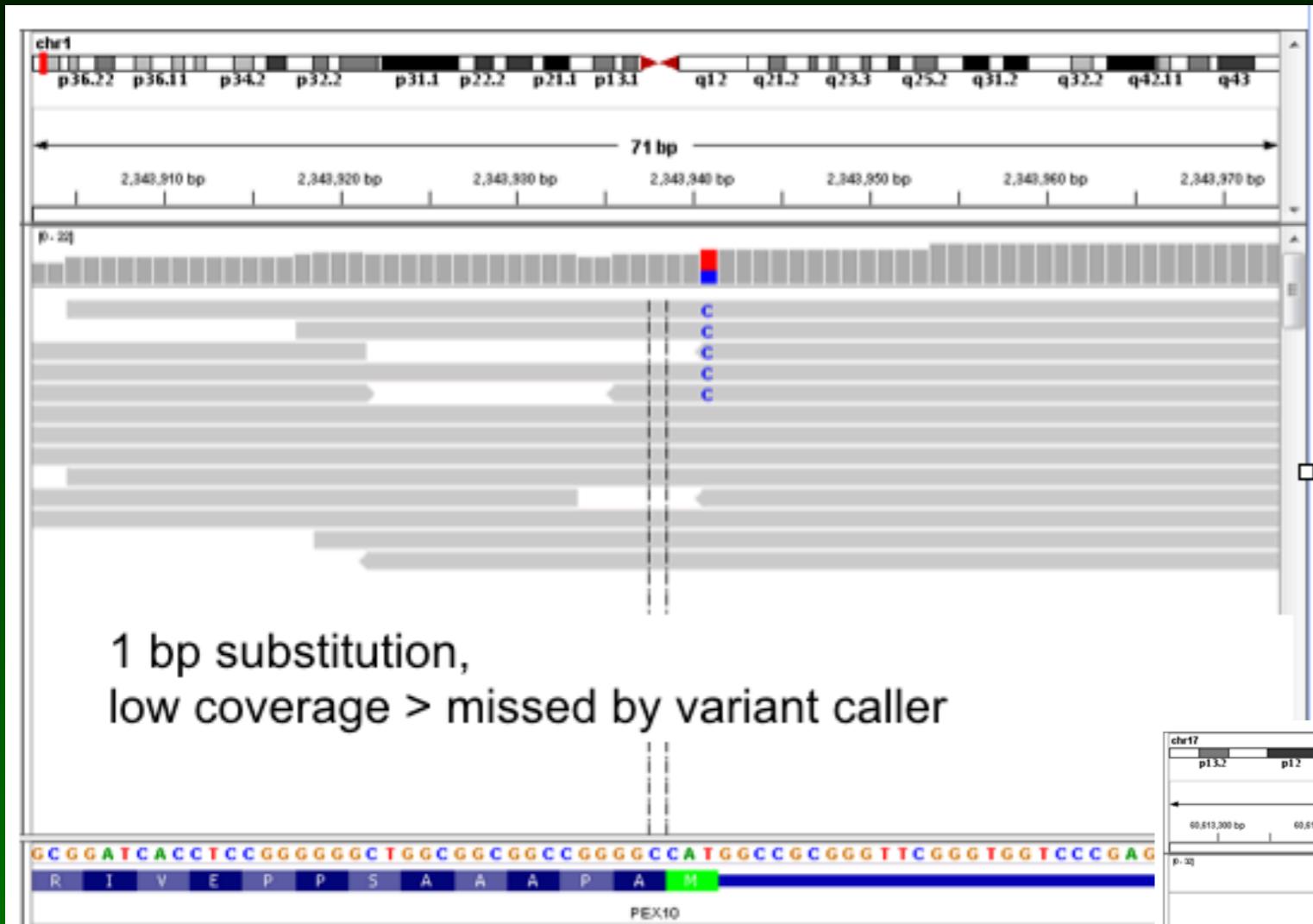


image from Christian Gilissen

probe = targeted but you capture more

call variants in targeted regions +/- ... nucleotides

# Calling thresholds



1 bp substitution,  
low coverage > missed by variant caller

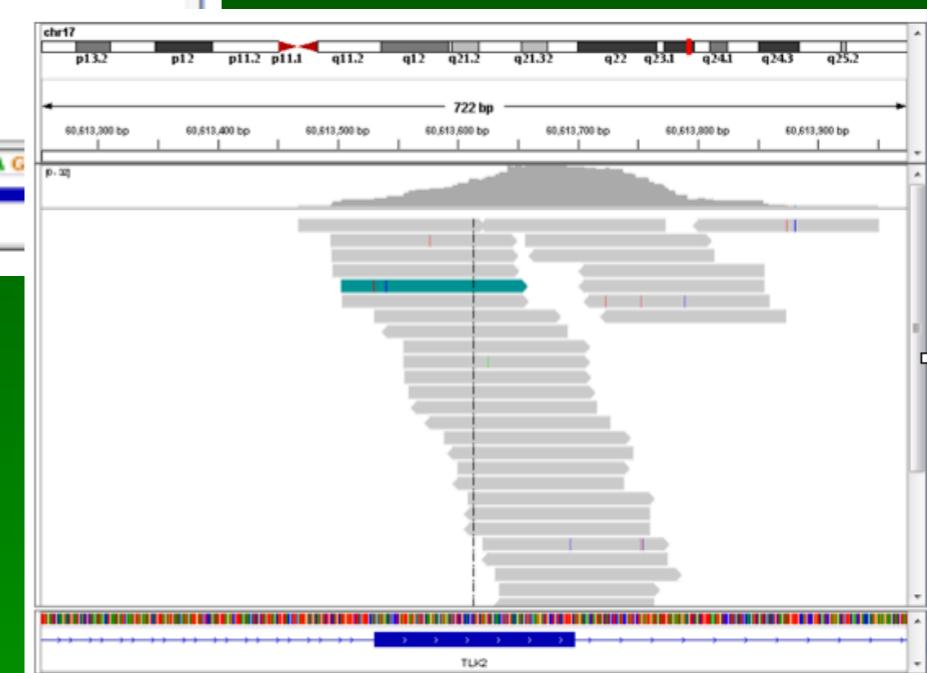
image from Christian Gilissen

Thresholds are essential

too low = many false positive  
too high = miss true variants

a delicate balance

NOTE: X/Y coverage in males



# De novo variants

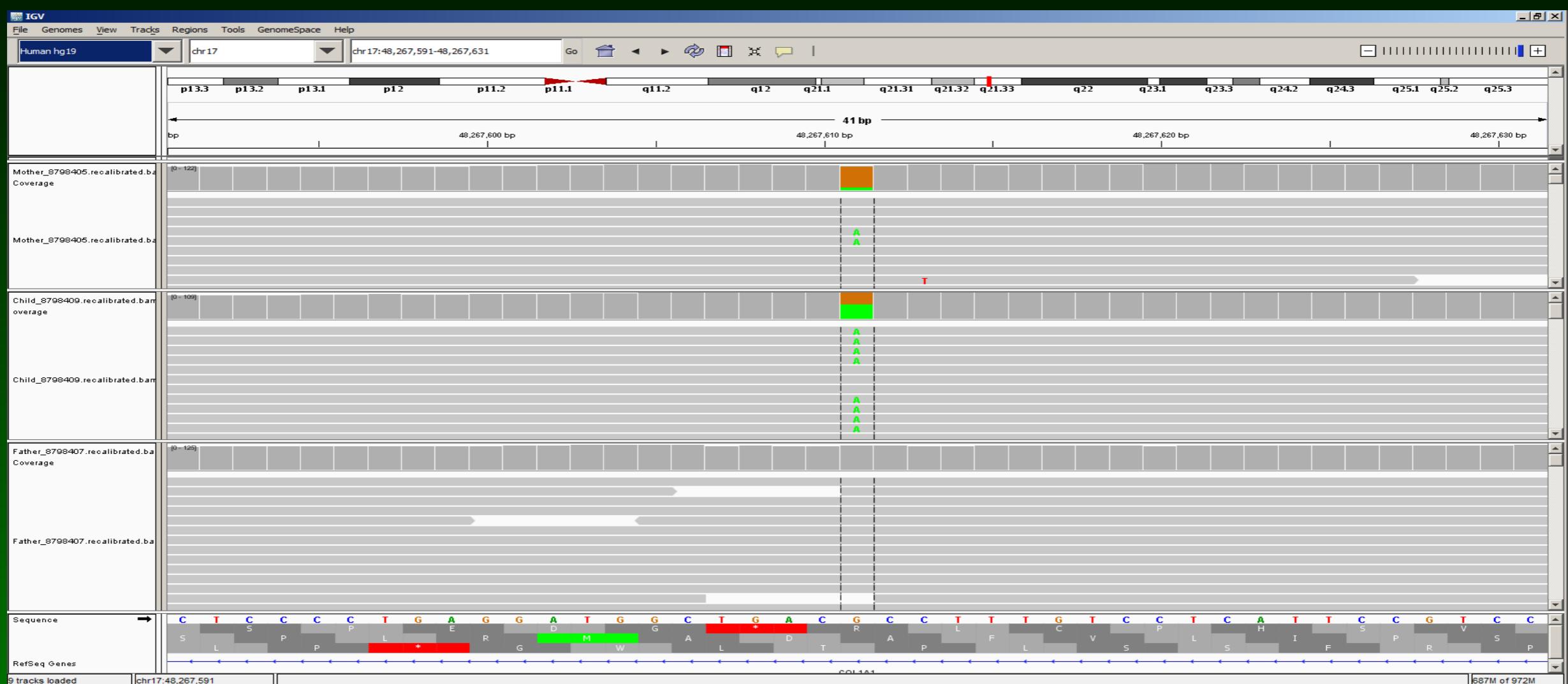


image from Gjjs Santen

variant missed as de novo

why?      no variant read allowed in parent

# Reference genome

---

*Which genome build do you use?*

*From which year is it ?*

## Human Assembly

- ✓ Dec. 2013 (GRCh38/hg38)
- Feb. 2009 (GRCh37/hg19)
- Mar. 2006 (NCBI36/hg18)
- May 2004 (NCBI35/hg17)
- July 2003 (NCBI34/hg16)

*Is the genome complete ?*

# XYLT1

- known disease gene  
*Desbuquois dysplasia type 2 (DBQD2)*  
*many patients without deleterious variants*
- detailed analysis negative  
*Sanger sequencing, WES, WGS*
- XYL1 methylated in patients  
*allele not expressed*

GGC Repeat Expansion and Exon 1 Methylation  
of XYL1 Is a Common Pathogenic Variant  
in Baratela-Scott Syndrome

American Journal of Human Genetics 104, 35–44

Amy J. LaCroix,<sup>1,11</sup> Deborah Stabley,<sup>2,11</sup> Rebecca Sahraoui,<sup>2,3</sup> Margaret P. Adam,<sup>1,4</sup> Michele Mehaffey,<sup>1</sup> Kelly Kernal,<sup>1</sup> Candace T. Myers,<sup>5</sup> Carrie Fagerstrom,<sup>6</sup> George Anadiotis,<sup>6</sup> Yassmine M. Akkari,<sup>6</sup> Katherine M. Robbins,<sup>2</sup> Karen W. Gripp,<sup>2</sup> Wagner A.R. Baratela,<sup>7,8</sup> Michael B. Bober,<sup>7</sup> Angela L. Duker,<sup>7</sup> Dan Doherty,<sup>1,4</sup> Jennifer C. Dempsey,<sup>1</sup> Daniel G. Miller,<sup>1</sup> Martin Kircher,<sup>9</sup> Michael J. Bamshad,<sup>1,4</sup> Deborah A. Nickerson,<sup>4,9</sup> University of Washington Center for Mendelian Genomics,  
Heather C. Mefford,<sup>1,4,12,\*</sup> and Katia Sol-Church<sup>2,10,12,\*</sup>

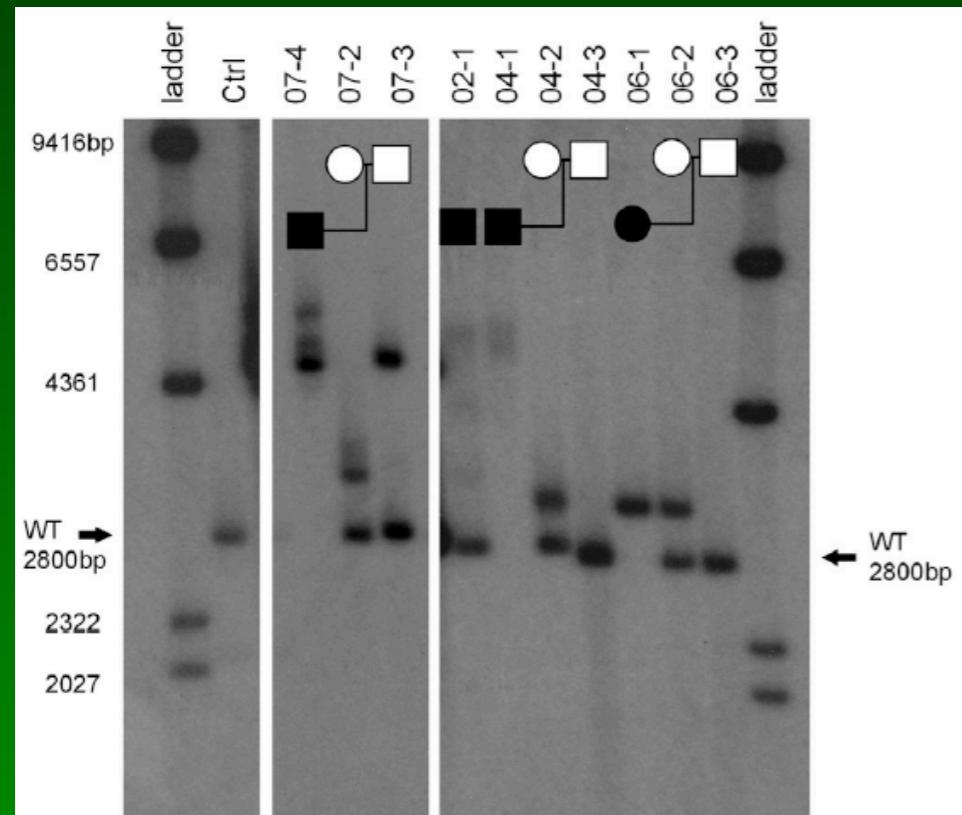


# XYLT1

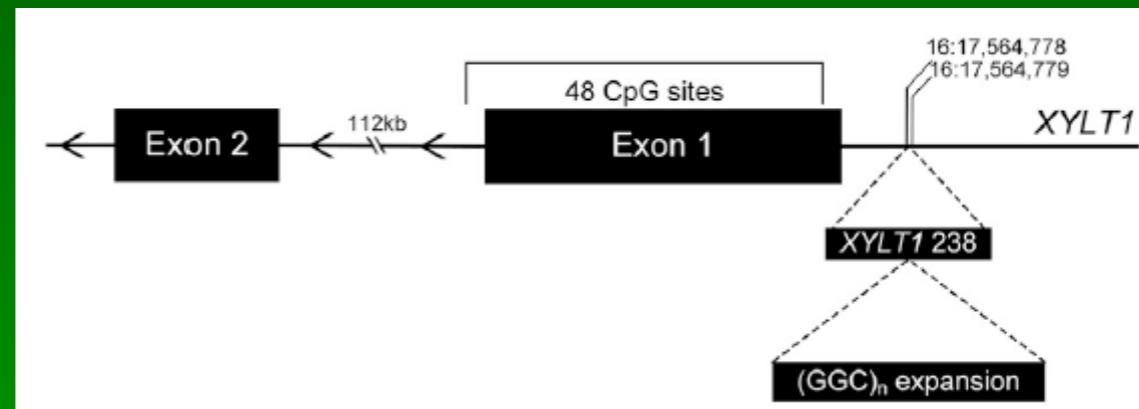
## GGC Repeat Expansion and Exon 1 Methylation of XYL<sub>T</sub>1 Is a Common Pathogenic Variant in Baratela-Scott Syndrome

American Journal of Human Genetics 104, 35–44

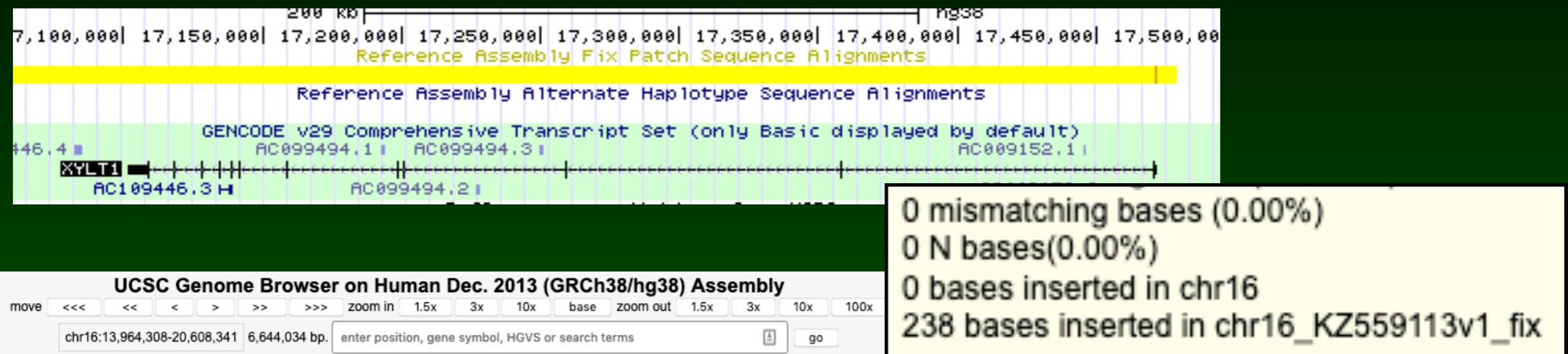
Amy J. LaCroix,<sup>1,11</sup> Deborah Stabley,<sup>2,11</sup> Rebecca Sahraoui,<sup>2,3</sup> Margaret P. Adam,<sup>1,4</sup> Michele Mehaffey,<sup>1</sup> Kelly Kerman,<sup>1</sup> Candace T. Myers,<sup>5</sup> Carrie Fagerstrom,<sup>6</sup> George Anadiotis,<sup>6</sup> Yassmine M. Akkari,<sup>6</sup> Katherine M. Robbins,<sup>2</sup> Karen W. Gripp,<sup>2</sup> Wagner A.R. Baratela,<sup>7,8</sup> Michael B. Bo Dan Doherty,<sup>1,4</sup> Jennifer C. Dempsey,<sup>1</sup> Daniel G. Miller,<sup>1</sup> Martin Kircher,<sup>9</sup> Mich Deborah A. Nickerson,<sup>4,9</sup> University of Washington Center for Mendelian Geno Heather C. Mefford,<sup>1,4,12,\*</sup> and Katia Sol-Church<sup>2,10,12,\*</sup>



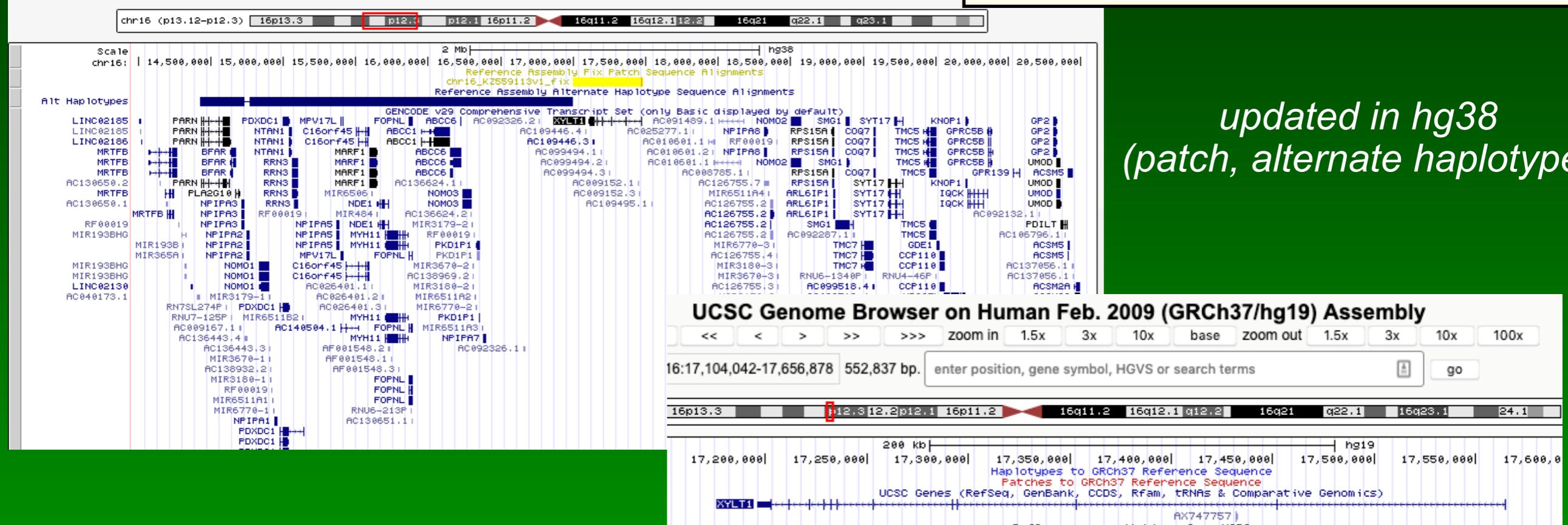
sulfate and dermatan sulfate proteoglycans.<sup>9</sup> The promoter region of XYL<sub>T</sub>1 was recently characterized and found to have 238 bp of sequence (XYLT1-238) that is not present in the reference genome (GRCh37/hg19, GRCh38/hg38) but appears to be fixed in the human population (100/100 individuals tested)<sup>10</sup> and is evolutionarily conserved in chimpanzee and mouse. The 238 bp region also contains a variable GGC repeat with a range of 9–20 repeats in 100 individuals tested; repeat length within this range did not



# XYLT1 hg38



*updated in hg38  
(patch, alternate haplotype)*



# Dutch genome

**Go•NL**  
GENOMEoftheNETHERLANDS

Home About us Access to the data Resources The GoNL team Wiki

Ultra-sharp genetic group portrait of the Dutch

Posted on July 15, 2012 by r

SNP calling complete

B B M R I • N L

230 x One child  
10 x Dizygotic twins  
10 x Monozygotic twins  
770 Individuals

Next Gen Sequencing – 12x coverage Illumina HiSeq 2000 platform

Genotyping Minimal 2 array platforms/sample ImmunoChip + others

used as reference ?

contains sequence not found in reference genome

The diagram illustrates the GoNL project's data collection process. It starts with four blue boxes labeled 'BioBank Amsterdam', 'BioBank Groningen', 'BioBank Leiden', and 'BioBank Rotterdam'. Arrows point from these boxes to a central orange box containing the text: '230 x One child', '10 x Dizygotic twins', '10 x Monozygotic twins', and '770 Individuals'. From this central box, two arrows point down to green boxes: 'Next Gen Sequencing – 12x coverage Illumina HiSeq 2000 platform' and 'Genotyping Minimal 2 array platforms/sample ImmunoChip + others'. To the right of this central area is a white box titled 'News' with the text 'SNP calling complete'. Above the 'News' box is a small graphic showing colored dots (green, pink, red, purple, blue) arranged in a grid pattern. To the right of the 'News' box is another white box containing the letters 'B B M R I • N L' above a grid of colored dots (green, pink, red, purple, blue). Below the 'News' box is a large green area with the text 'used as reference ?' and 'contains sequence not found in reference genome'.

# Repeat expansion

- disease gene mapped for many years  
*extensive families available*  
*...no deleterious variants*  
*Sanger sequencing, WES, WGS, RNA-seq*



# Reference genome

---

*Which genome build do you use?*

*From which year is it ?*

- Human Assembly**
- ✓ Dec. 2013 (GRCh38/hg38)
- Feb. 2009 (GRCh37/hg19)
- Mar. 2006 (NCBI36/hg18)
- May 2004 (NCBI35/hg17)
- July 2003 (NCBI34/hg16)

*Is the genome complete ?*

*Is the transcript annotation complete ?*

# Gene annotation

call variants in genes / transcripts

RefSeq

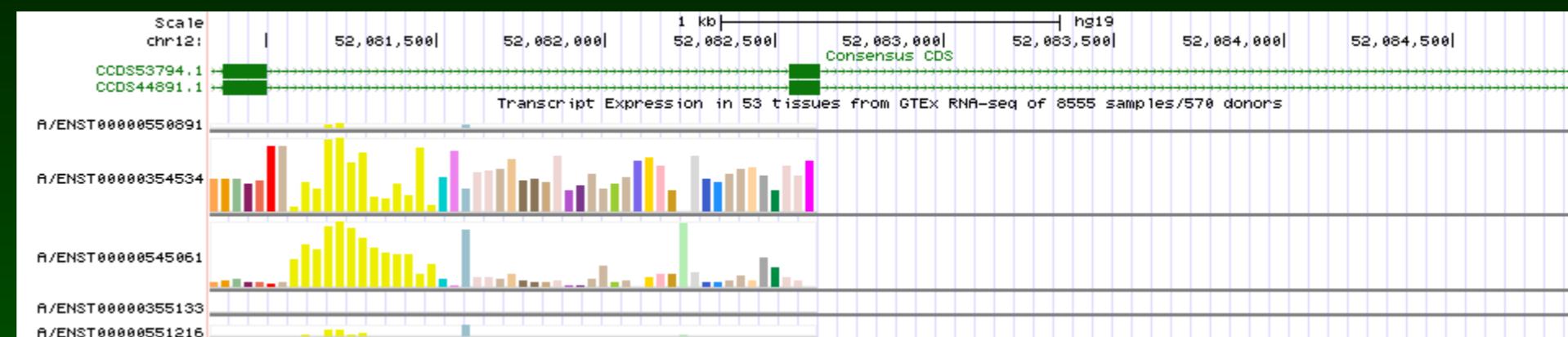
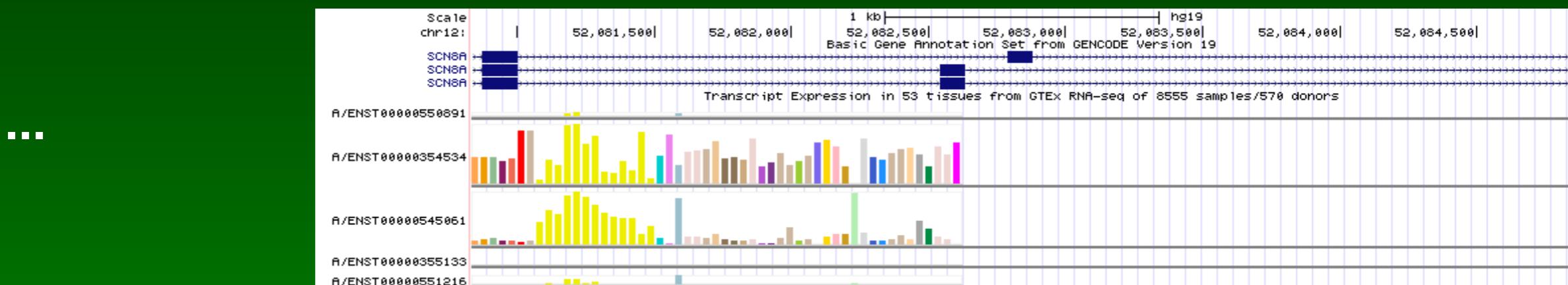


image from Christian Gilissen



calling variants linked to all transcripts gives an enormous list

select transcripts, but be careful

# Frequently variant



Shyr (2014) BMC Medical Genomics

not likely candidates:

- alignment issues (e.g. MUC16)
- size of the gene (e.g. TTN)

# Why cause missed?

- family  
*parents consanguineous*

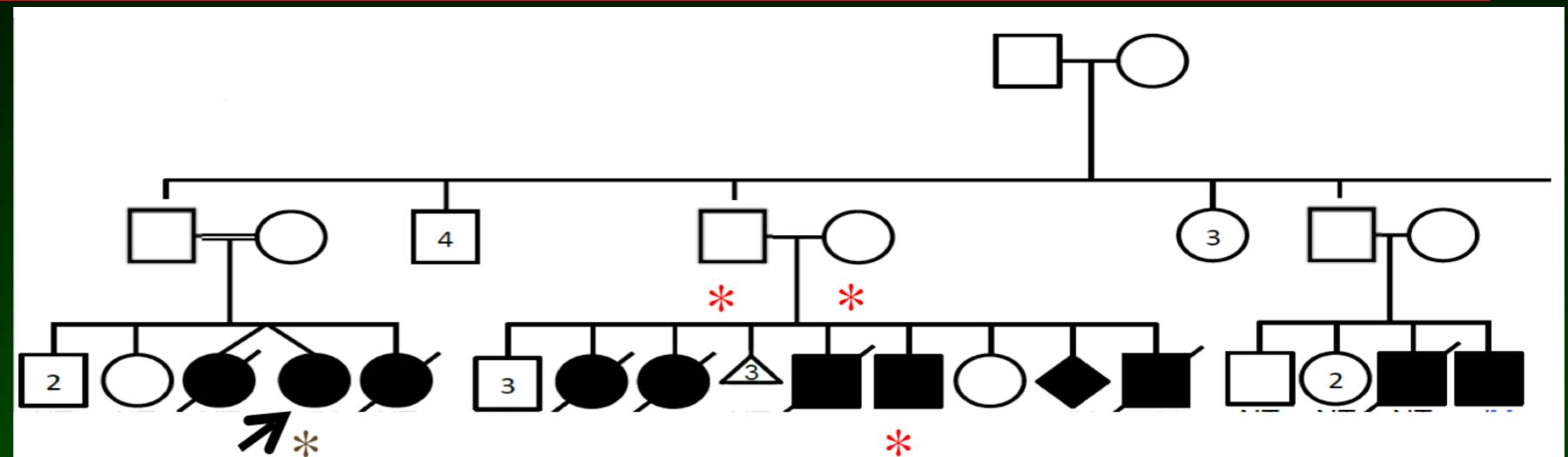


image from Gijs Santen

- index case  
*decreased fetal movement (HP:0001558), severe neonatal hypotonia (HP:0001319), severe feeding difficulties (HP:0011968)*
- exome analysis  
*no causative variants*

# Imprinting disorders

RESEARCH ARTICLE

Open Access

## Paraganglioma and pheochromocytoma upon maternal transmission of *SDHD* mutations

Jean-Pierre Bayley<sup>1\*†</sup>, Rogier A Oldenburg<sup>6†</sup>, Jennifer Nuk<sup>8</sup>, Attje S Hoekstra<sup>1</sup>, Conny A van der Meer<sup>6</sup>, Esther Korpershoek<sup>7</sup>, Barbara McGillivray<sup>8</sup>, Eleonora PM Corssmit<sup>5</sup>, Winand NM Dinjens<sup>7</sup>, Ronald R de Krijger<sup>7</sup>, Peter Devilee<sup>1,3</sup>, Jeroen C Jansen<sup>4</sup> and Frederik J Hes<sup>2</sup>

BMC Medical Genetics 2014, 15:111

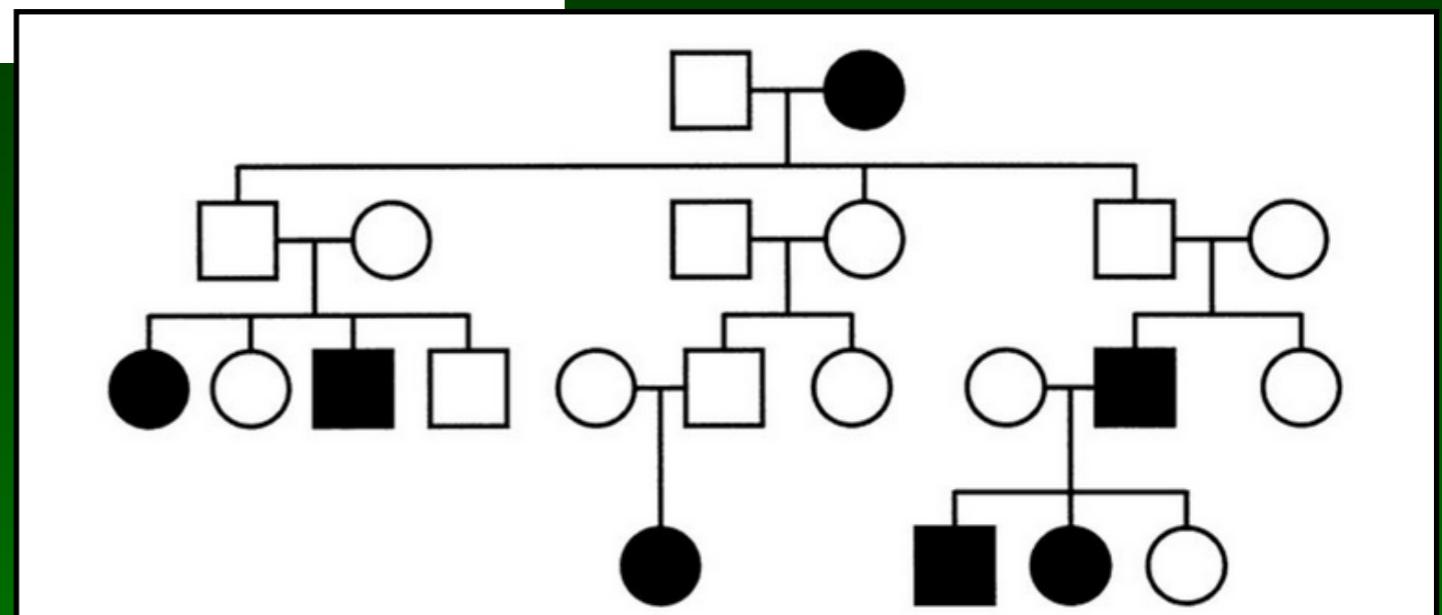


image from Anna Benet-Pages

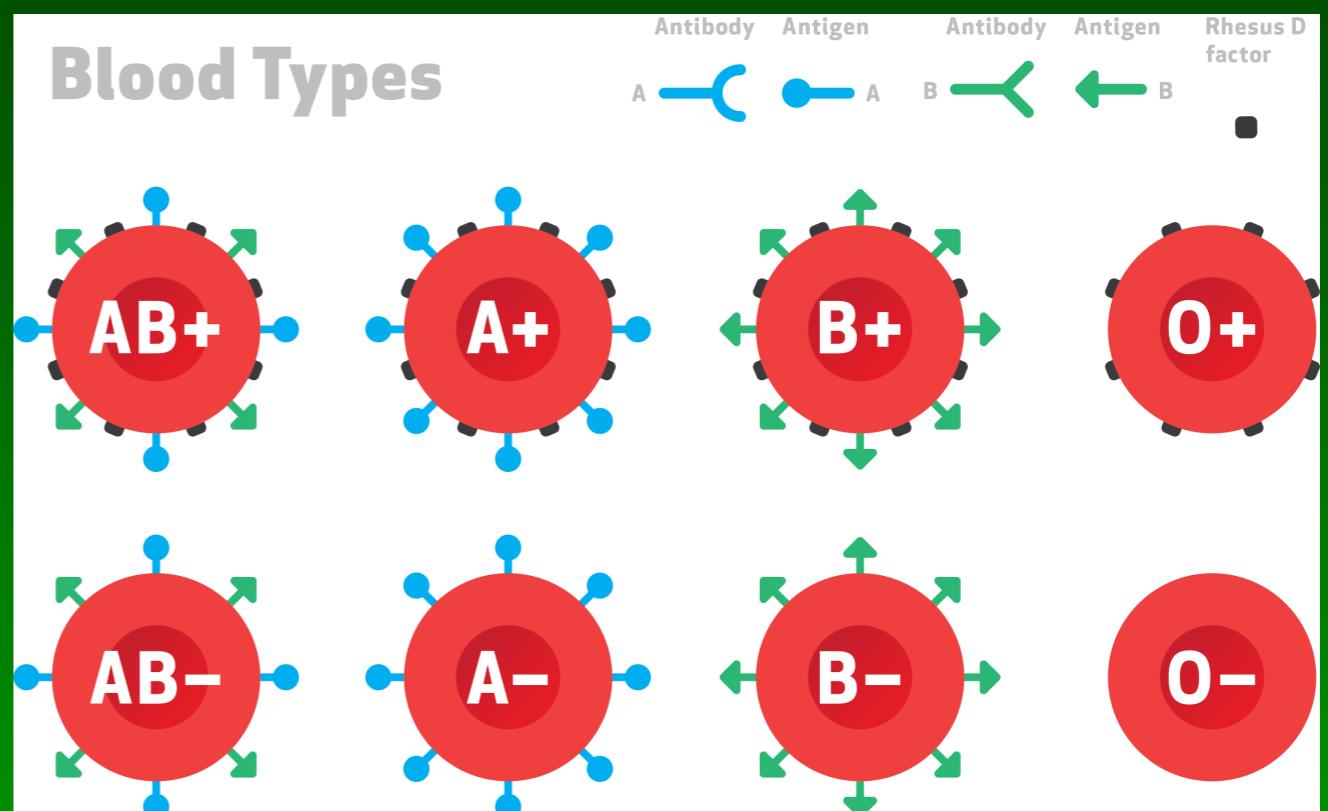
variant inherited from father:      **affected**  
variant inherited from mother:      **healthy**

# Reference phenotype

- phenotype reference genome?
- ABO blood group gene  
*genome has no variant called  
ABO blood group type is A, B or O ??*

***blood group O***

*so not the “normal” gene*



*image from www.blood.ca*

# Reference genome

- not always the major allele
- ref transcript/ref genome may differ  
*sometimes with length differences*

Reference SNP (refSNP) Cluster Report: rs8176719      **\*\*Clinical Channel\*\***

RefSNP	Allele	HGVS Names
Organism: human ( <i>Homo sapiens</i> )	Variation Class: DIV: deletion/insertion variation	CM000671.2:g.133257521_133257522insC
Molecule Type: Genomic	RefSNP Alleles: -/G (REV)	NC_000009.11:g.136132908_136132909insC
Created/Updated in build: 117/151	Allele Origin:	NC_000009.12:g.133257521_133257522insC
Map to Genome Build: 108/Weight 1	Ancestral Allele: A	NG_006669.1:g.20145_20147insG
Validation Status:	Variation Viewer:	NM_020469.2:c.260_262insG
Citation: PubMed LitVar <sup>NEW</sup>	Clinical Significance: NA	NP_065202.2:p.Val87_Thr88=fs
Association: NHGRI GWAS	MAF/MinorAlleleCount: C=0.3764/45143 (ExAC) C=0.3438/1722 (1000 Genomes) C=0.3400/4154 (GO-ESP) C=0.3456/43394 (TOPMED)	XP_005276905.1:p.Thr87Aspfs XP_005276906.1:p.Thr69Aspfs

strange variant descriptions

how about g. to c. descriptions ?

Alleles      Location      Evidence status      HGVS names

Synonyms      Original source      About this variant

PharmGKB PA166170086

Variants (including SNPs and indels) imported from dbSNP (release 150) | [View in dbSNP](#)

This variant overlaps [4 transcripts](#), has [2550 sample genotypes](#), is associated with [2 phenotypes](#) and is mentioned in [54 citations](#).

# Simple sequences

*MSH2 : NM\_000251:c.942+3A>T*

*most frequent pathogenic MSH2 variant in HNPCC / Lynch Syndrome*

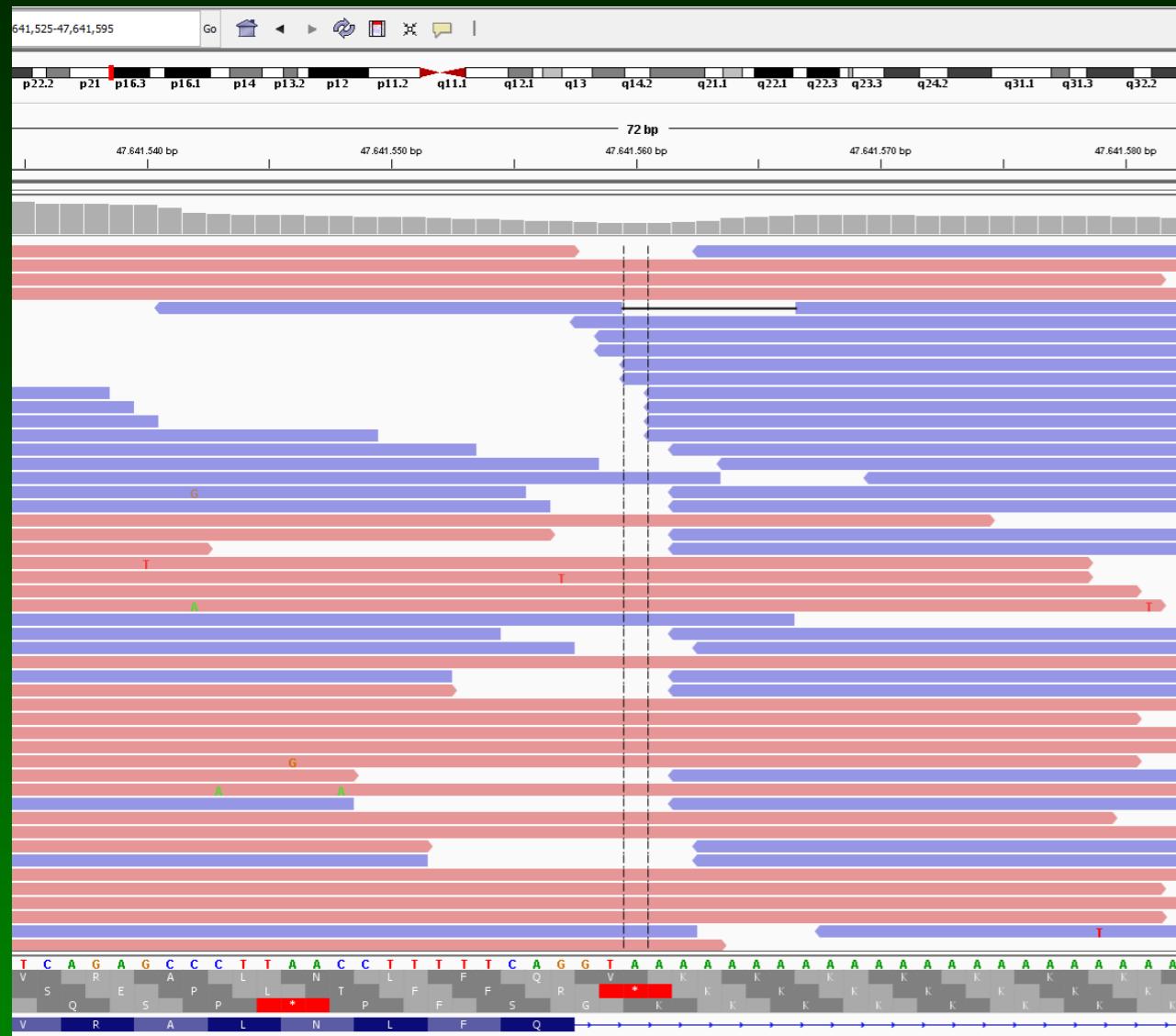
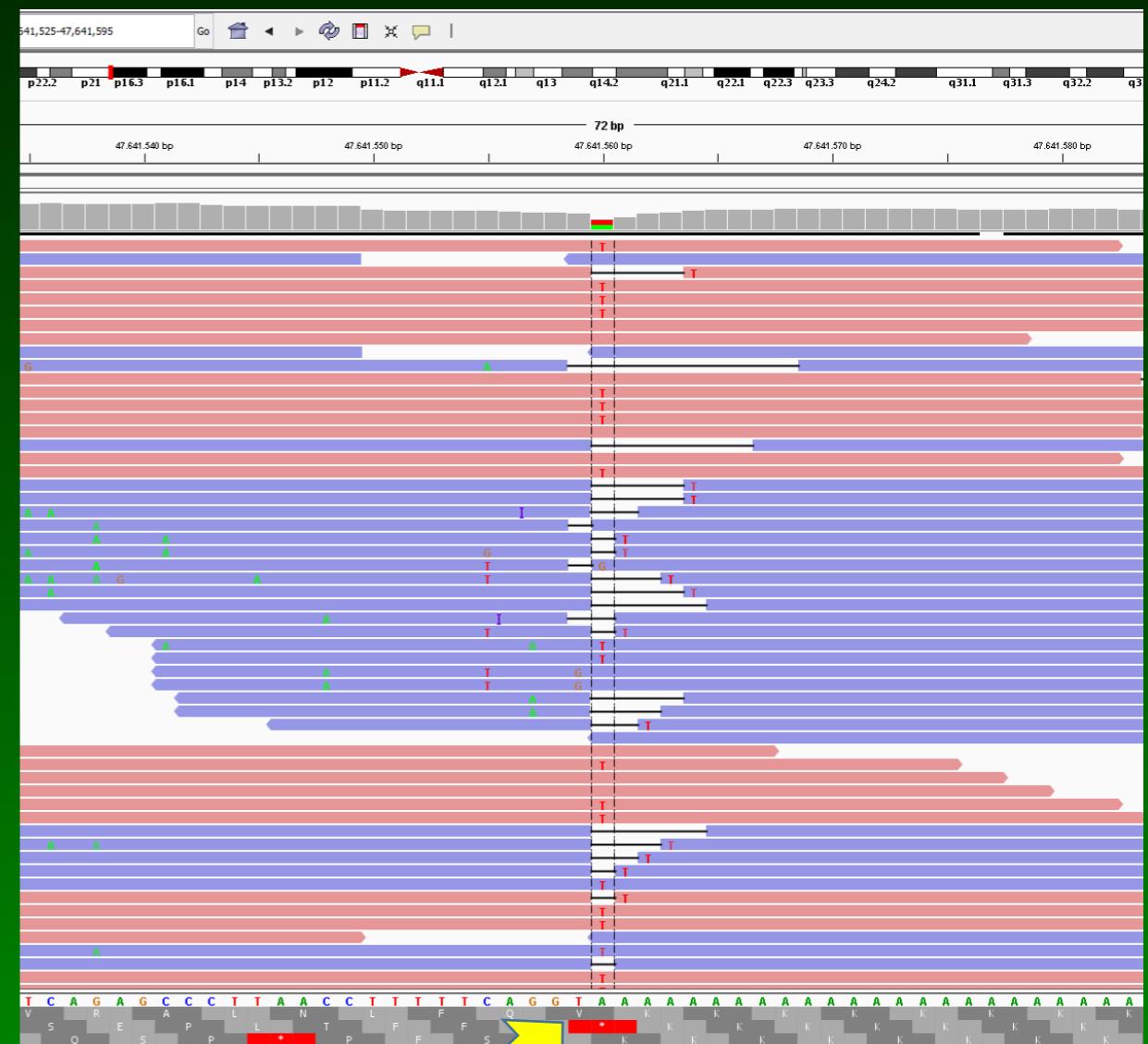


image from Anna Benet-Pages

A



A>T

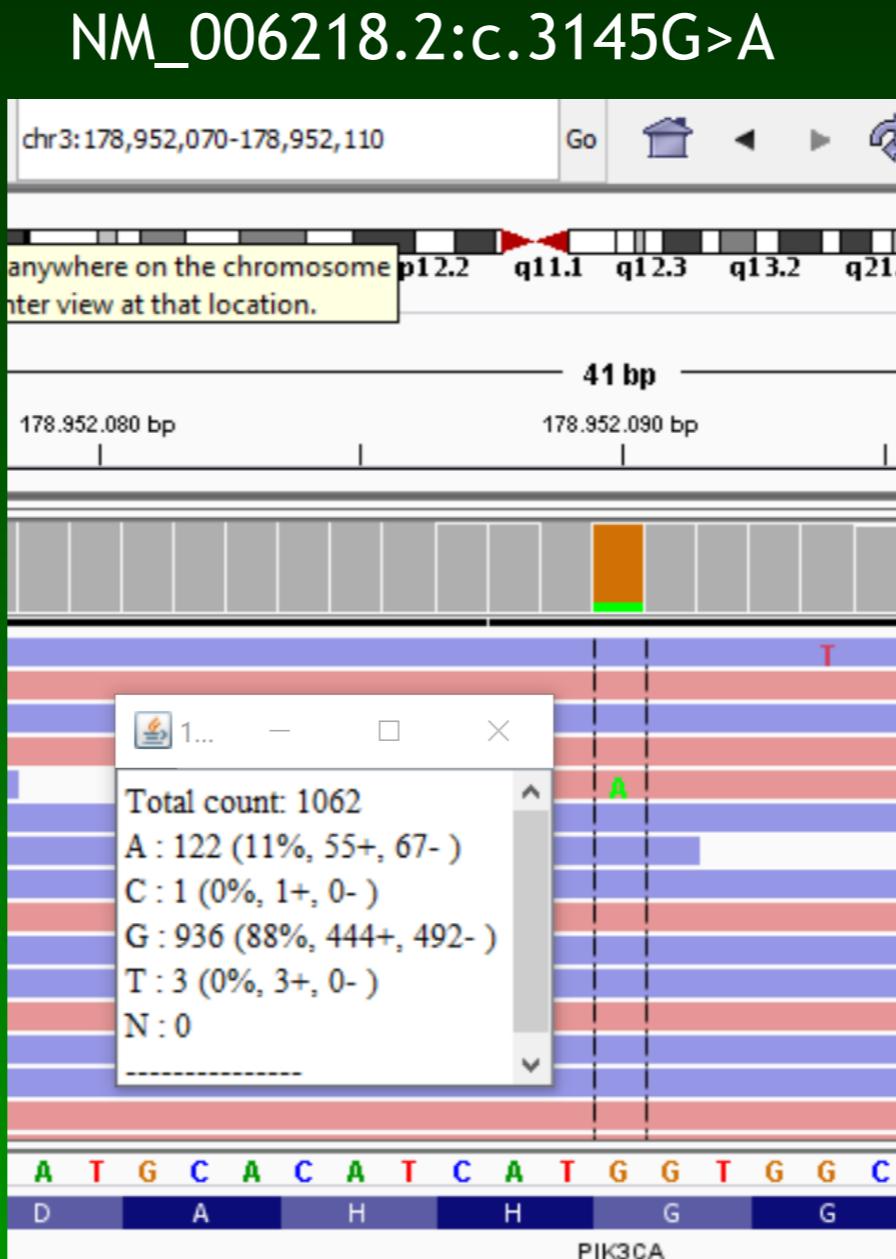
# Somatic variants

Somatic Mosaic Activating Mutations in *PIK3CA* Cause CLOVES Syndrome

Kyle C. Kurek,<sup>1</sup> Valerie L. Luks,<sup>2</sup> Ugur M. Ayturk,<sup>2,9</sup> Ahmad I. Alomari,<sup>3,6</sup> Steven J. Fishman,<sup>4,6</sup> Samantha A. Spencer,<sup>2,6</sup> John B. Mulliken,<sup>5,6</sup> Margot E. Bowen,<sup>2,9</sup> Guilherme L. Yamamoto,<sup>7</sup> Harry P.W. Kozakewich,<sup>1,6</sup> and Matthew L. Warman<sup>2,6,8,9,\*</sup>

only detected with specific  
LOW-FREQUENCY pipeline

*...and when you look for  
the right variant feature  
in the right gene/disease*



NM\_006218.2:c.2740G>A

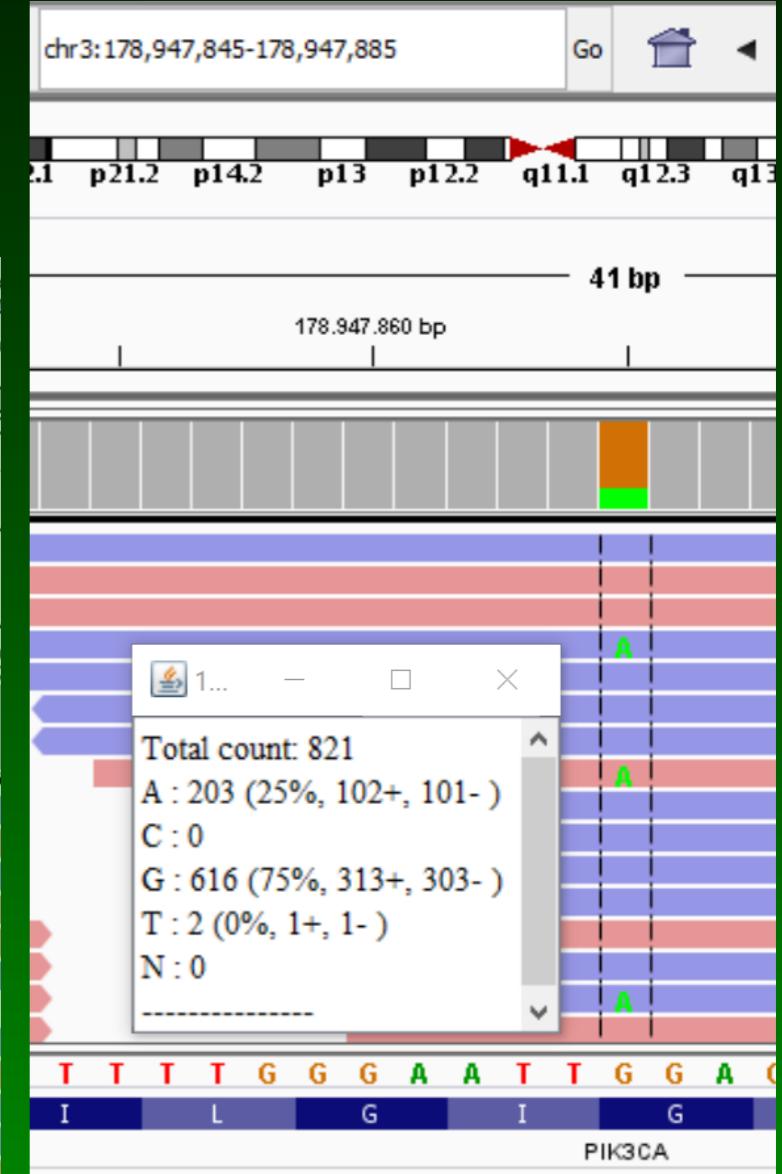


image from Anna Benet-Pages

# Variant calling

RASA1



image from Anna Benet-Pages

called: c.1016\_1017insTTA p.(Gly340\*) and c.1017G>C p.(=)

correct: c.1017delinsTTAC p.(Val39\_Gly340insTyr)

**be careful when two variant affect one codon**

# Variant calling

## BRCA2

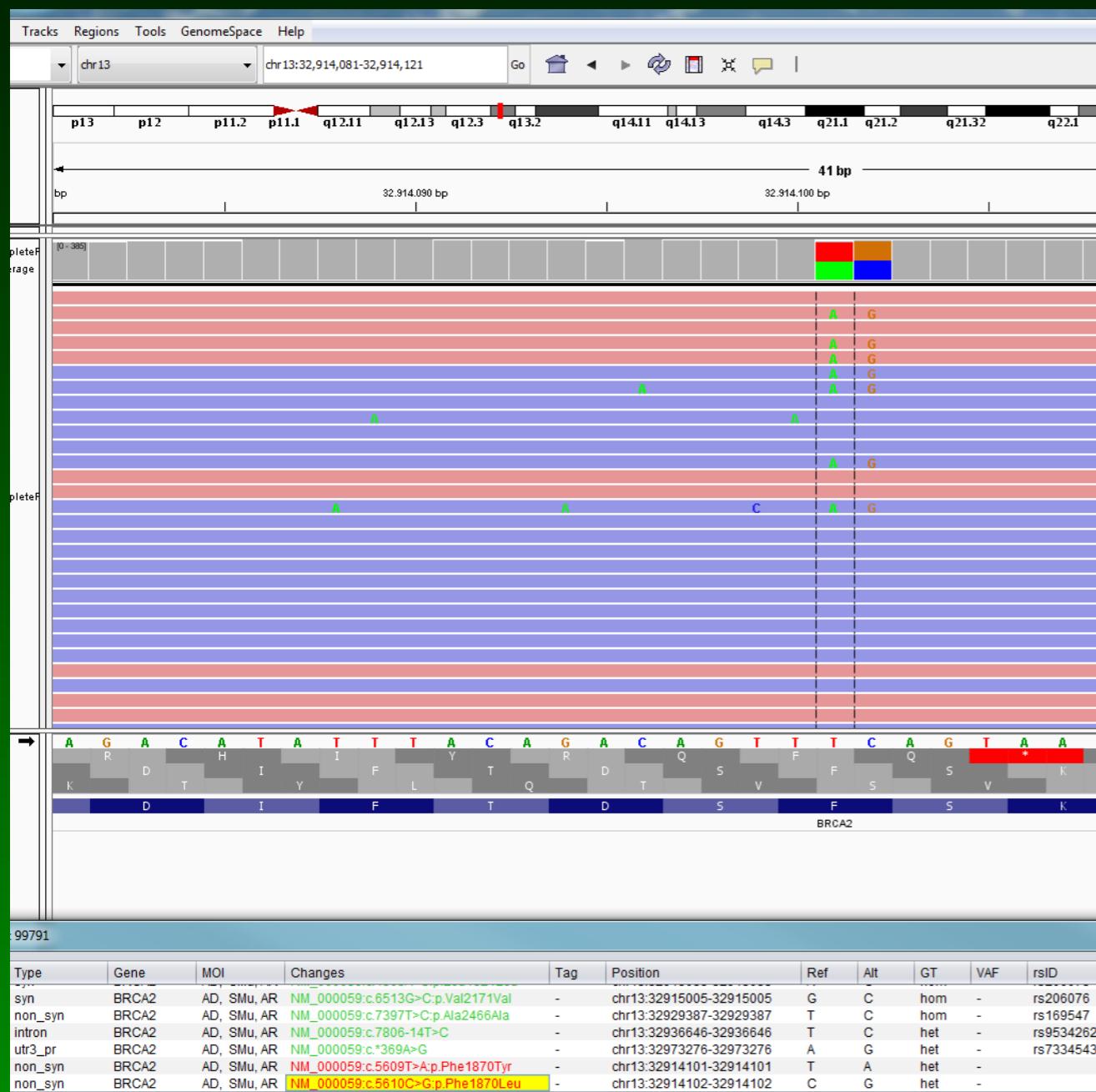


image from Anna Benet-Pages

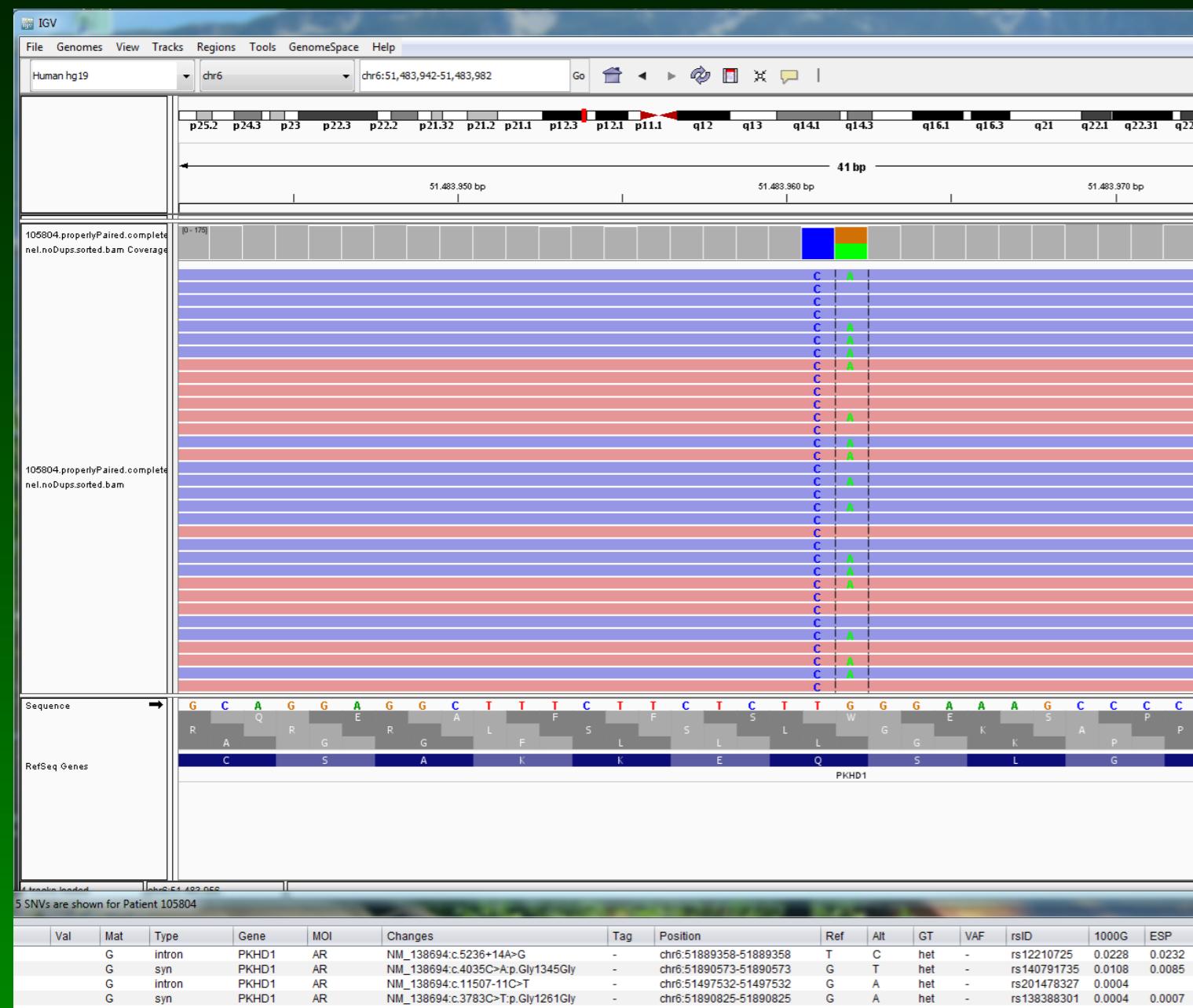
called: c.5609T>A p.(Phe1870Tyr)  
and c.5610C>G p.(Phe1870Leu)

correct: c.5609\_5610delinsAG  
p.(Phe1870\*)

How does your  
pipeline handle such  
cases ??

# Variant calling

PKHD1



*image from Anna Benet-Pages*

**standard pipelines filter out  
frequent variants  
(population based)**

# c.12142C>T p.(Gln4048\*) truncating variant

c.12141A>G p.(=)  
5% frequency

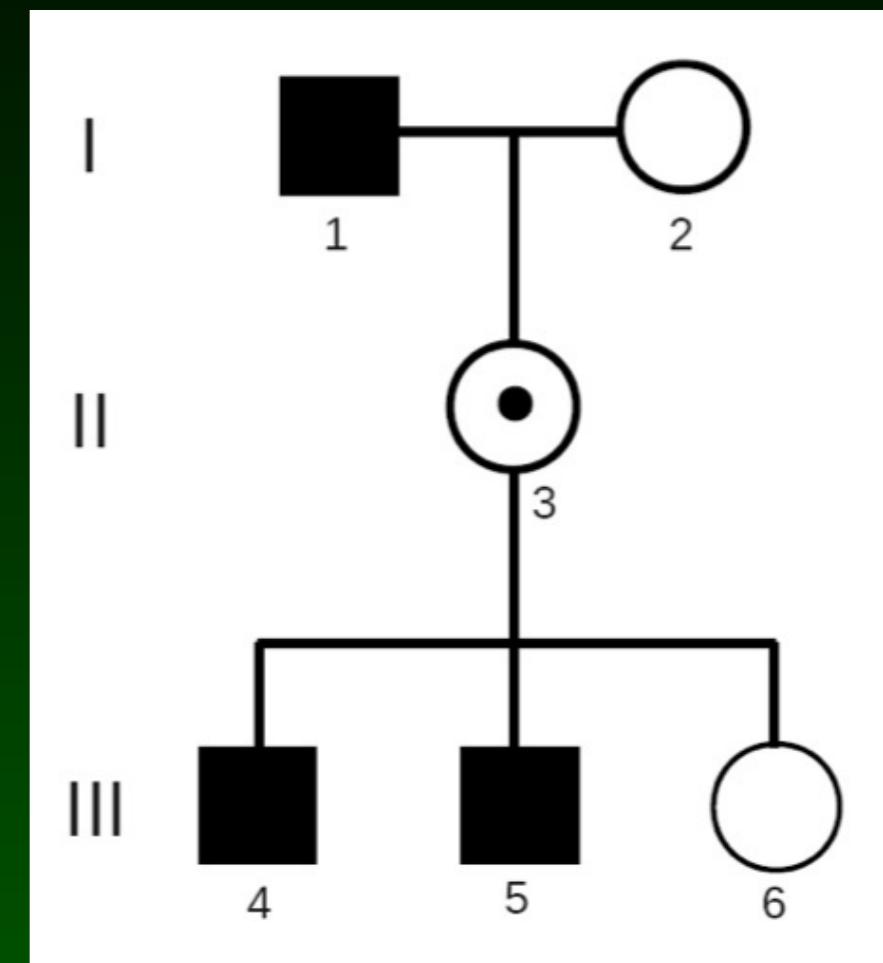
**combination automatically  
filtered out**

# Exome focus

- Aarskog-Scott syndrome  
*FGD1 gene screened*  
*> no variants*
- whole exome capture  
*no obvious variants*  
*change thresholds filtering*



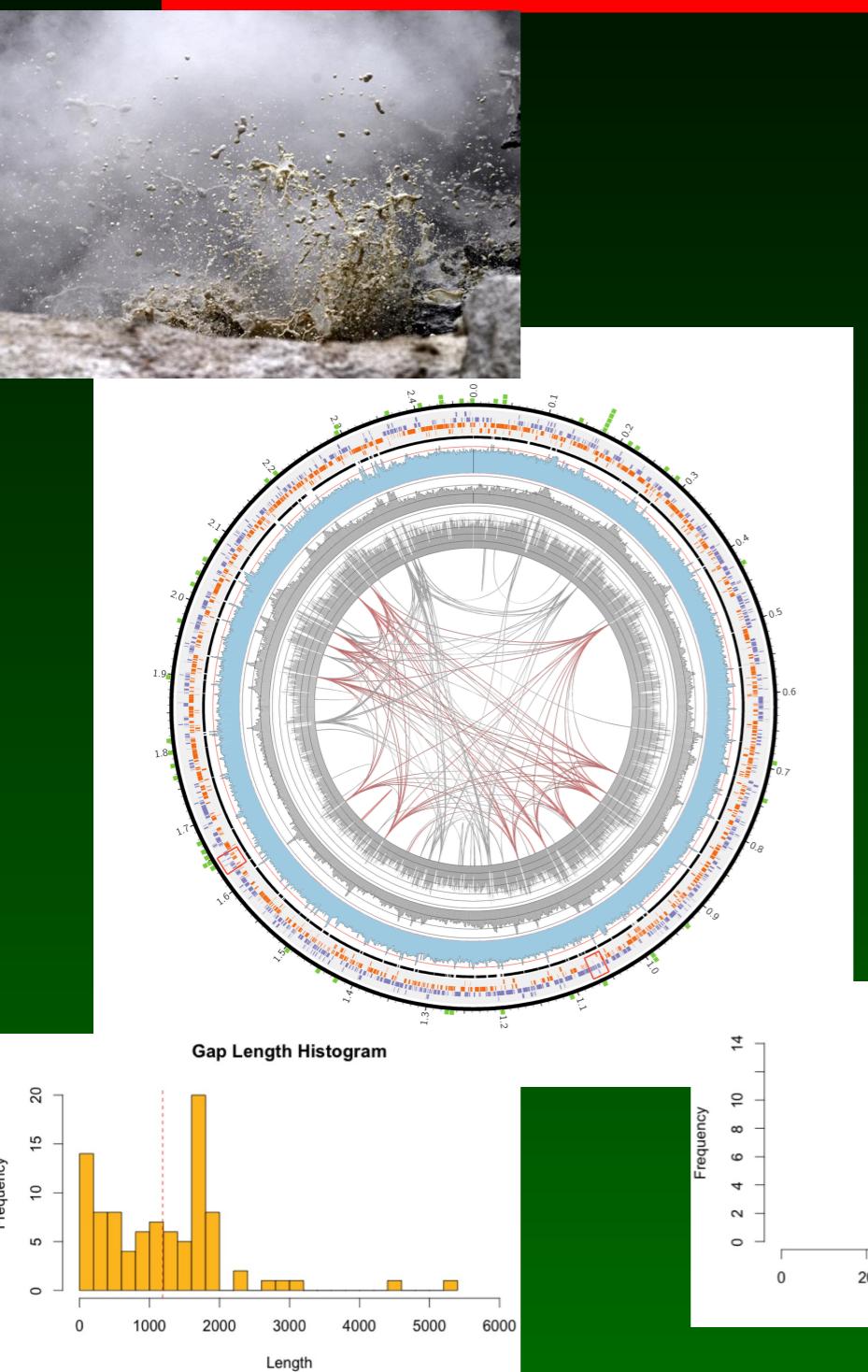
*intron -35delA variant*



©Yu Sun  
Emmelien Aten

# Bacterial genome

- sequence across GC-rich regions
- resolved repetitive regions



Anvar et al. BMC Genomics 2014, 15:914  
<http://www.biomedcentral.com/1471-2164/15/914>

**RESEARCH ARTICLE** **Open Access**

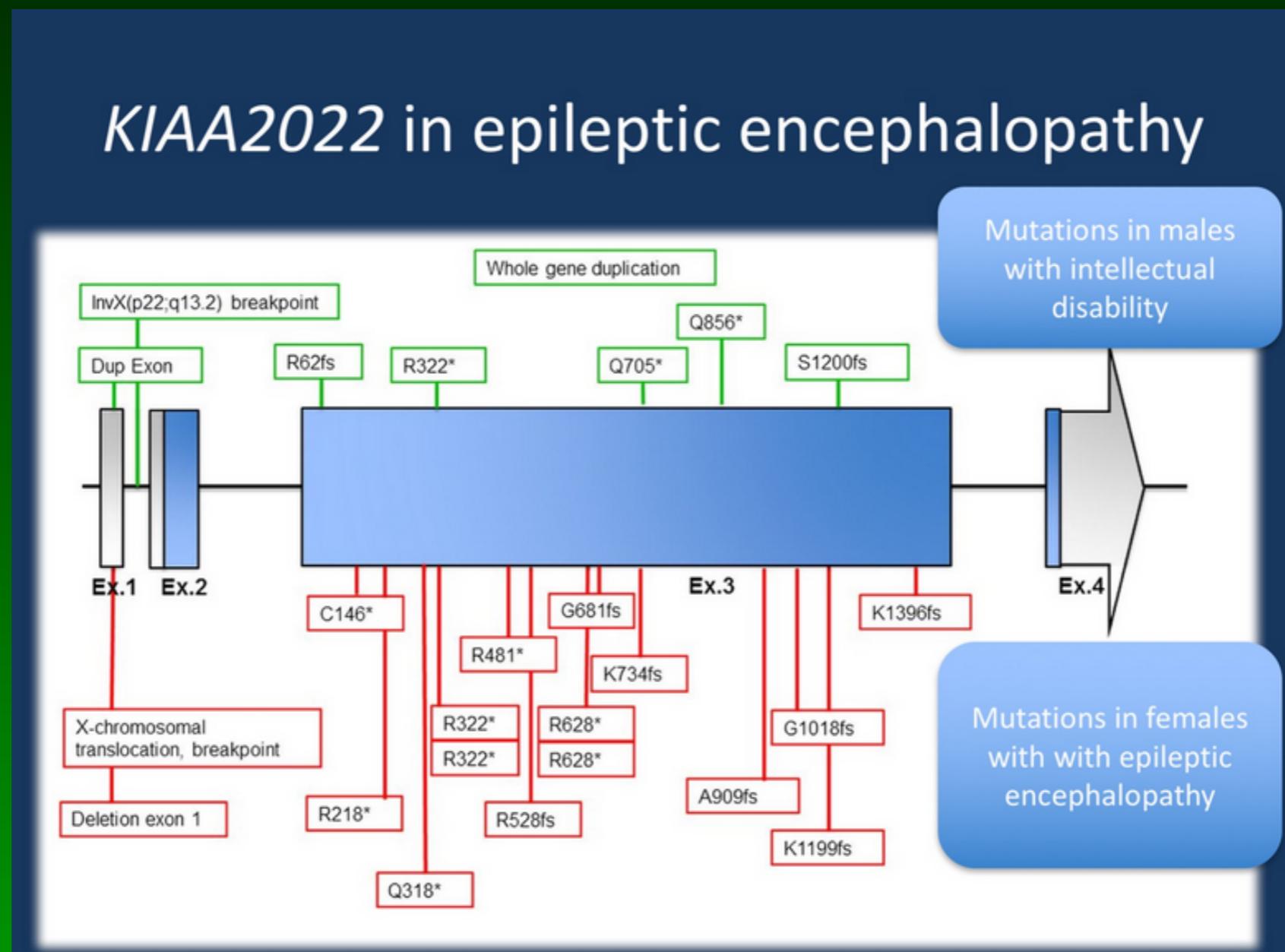
The genomic landscape of the verrucomicrobial methanotroph *Methylacidiphilum fumariolicum* SolV

Seyed Yahya Anvar<sup>1,2\*</sup>, Jeroen Frank<sup>2</sup>, Arjan Pol<sup>3</sup>, Arnoud Schmitz<sup>2</sup>, Ken Kraaijeveld<sup>2,4</sup>, Johan T den Dunnen<sup>1,2,5</sup> and Huub JM Op den Camp<sup>3\*</sup>

**best result using long reads only**  
*filter reads 8kb or larger*  
**missed 5 kb plasmid**

# Special cases

KIAA2022 - X-linked  
Females (epileptic encephalopathy) / males (intellectual disability)



*image from Anna Benet-Pages*

# Special cases

PCDH19: X-linked  
... but ONLY heterozygous females and mosaic males are affected

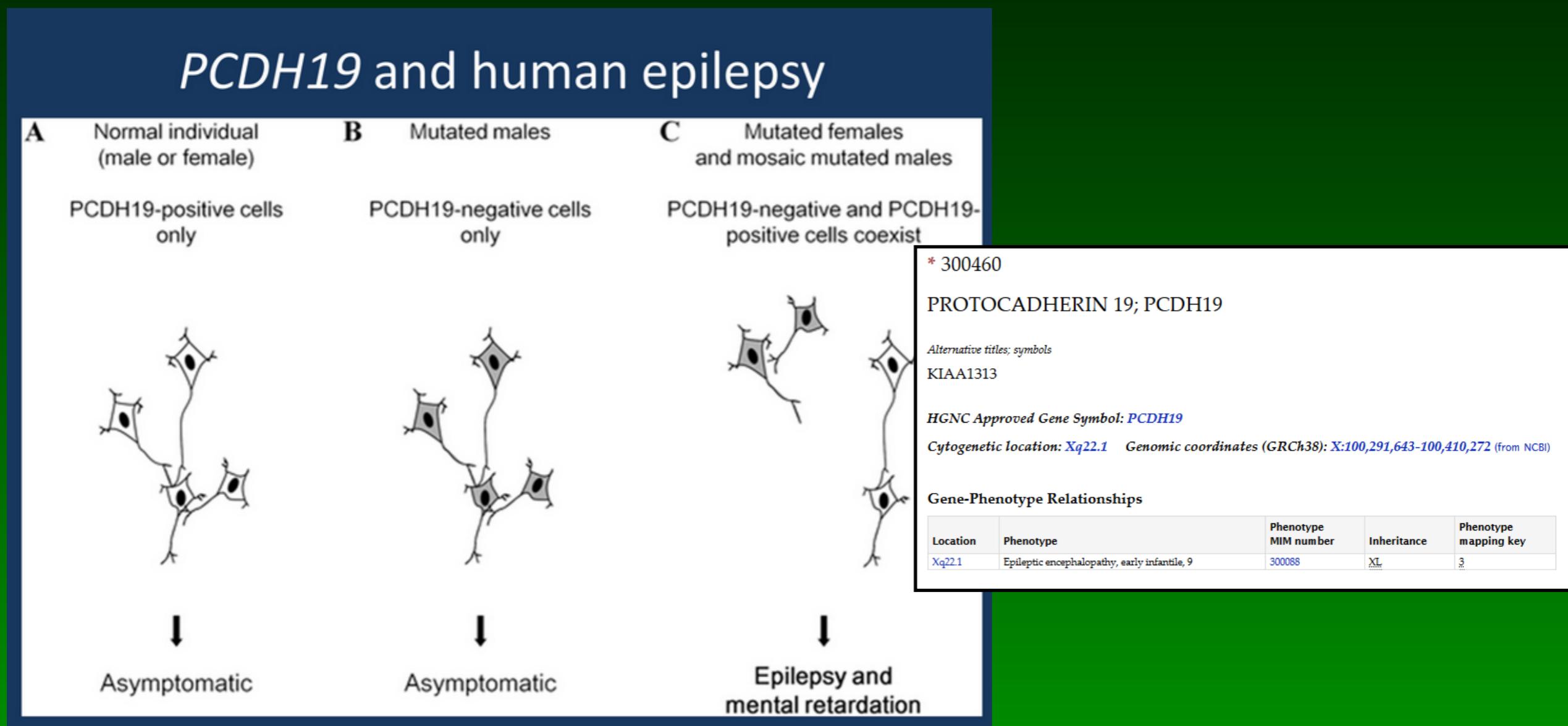


image from Anna Benet-Pages

# Murphy's law

---

*...also applies in NGS*



Your last  
mistake  
is your  
best  
teacher

# Acknowledgement

---

*Presentation prepared by:*  
***Johan den Dunnen***

*Human Genetics & Clinical Genetics  
Leiden University Medical Center  
Leiden, Nederland*



*date: April 2019*

*with contributions from Anna Benet-Pages,  
Christian Gilissen and Gijs Santen*