

Exome and genome sequencing in clinical practice

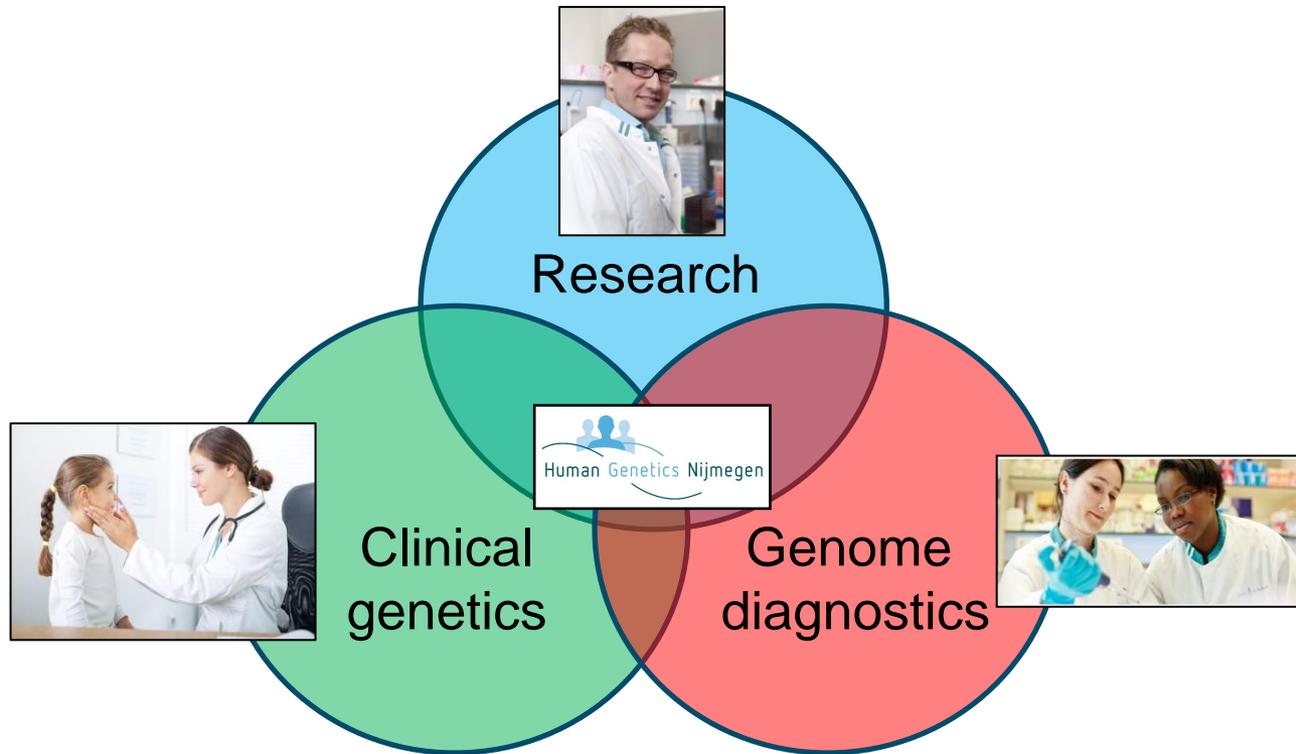
2nd IRDiRC conference - Shenzhen

Christian Gilissen Ph.D.
christian.gilissen@radboudumc.nl
07-11-2014

Human genetics Nijmegen



Human genetics Nijmegen



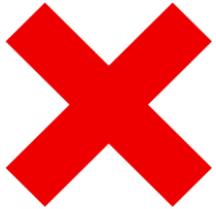
- Clinical genetics, diagnostics and research in one department
- Largest department of Human Genetics in the Netherlands

The challenge in diagnostics

- Single gene testing offered for more than 400 Mendelian diseases and more than 800 genes.



Good diagnostic yield for diseases with few causative genes (e.g. Noonan syndrome).



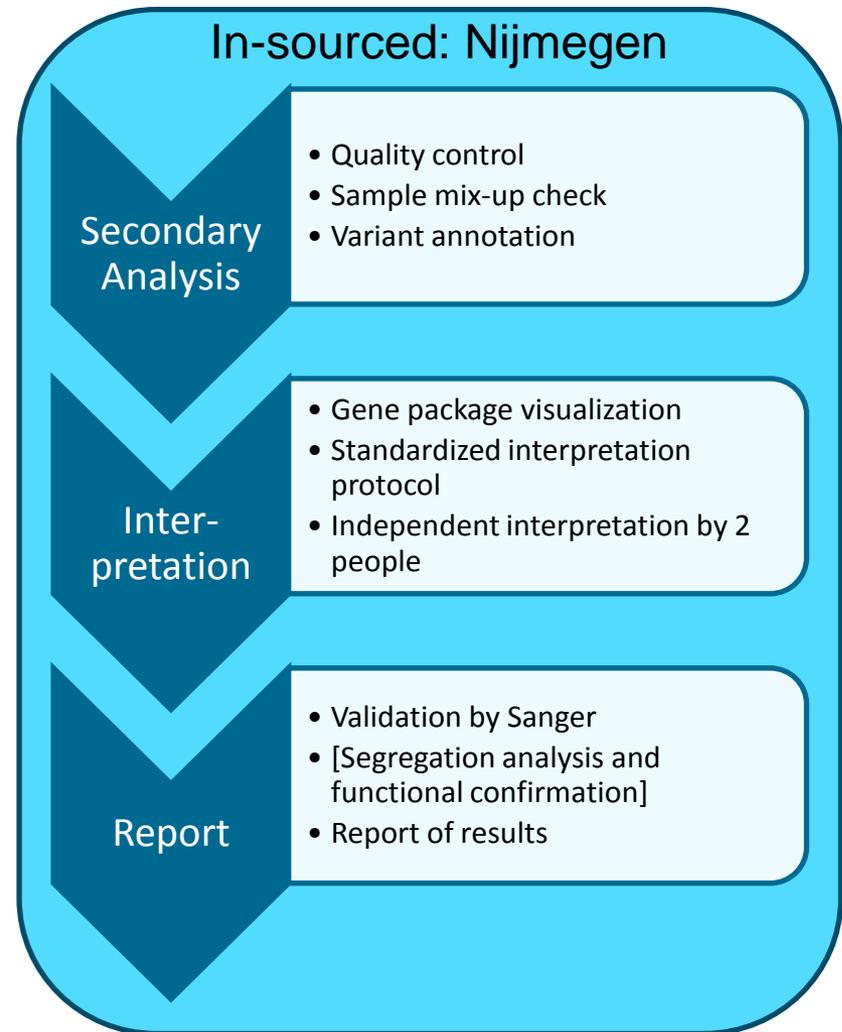
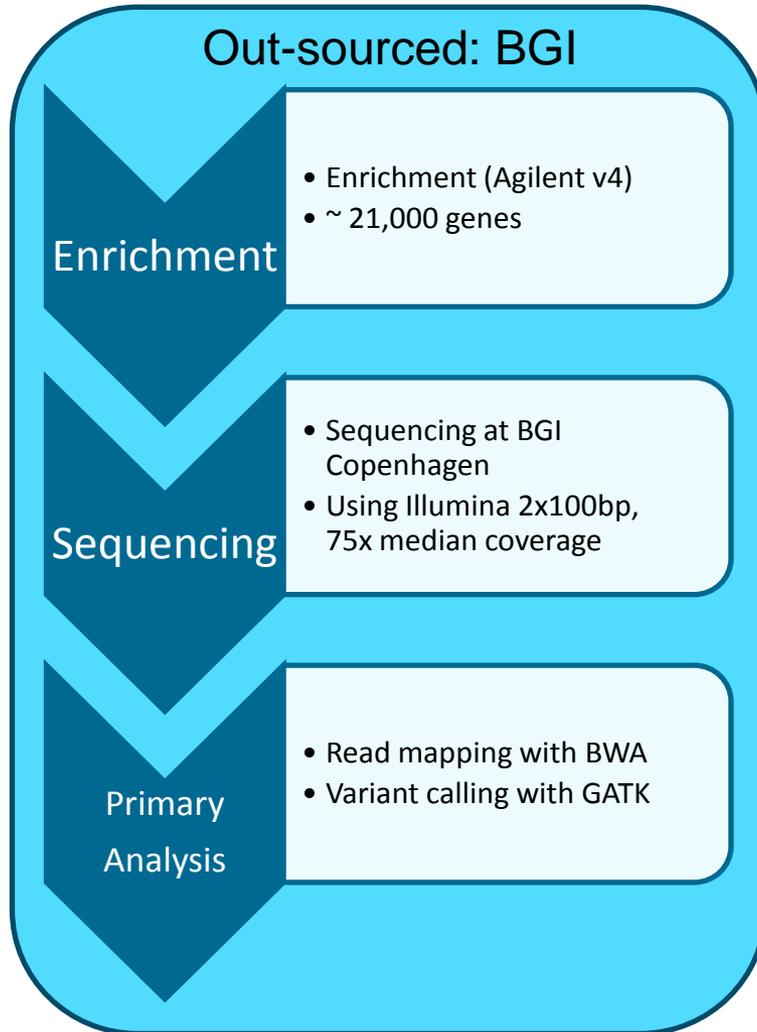
Poor diagnostic yield for genetically heterogeneous diseases. (e.g. Blindness, > 100 genes)

- **Solution:** Next generation sequencing?

Why exome sequencing?

Sanger	Targeted	Exome	Genome
<ul style="list-style-type: none">• Very accurate• Cheap per exon• High turn-around	<ul style="list-style-type: none">• Optimization possible• Low chance of incidental findings• “Easy” analysis• “Easy” interpretation	<ul style="list-style-type: none">• No bias for genes• Standardized workflow• Re-use of performed exomes to interpret new ones• Simple to add new genes	<ul style="list-style-type: none">• No bias in what you sequence• Little technical biases• Allows detection of SVs and SNVs in one experiment
<ul style="list-style-type: none">• Low diagnostic yield for genetically heterogeneous diseases	<ul style="list-style-type: none">• Design and re-design required• Different designs for different disorders• Sufficient patients required	<ul style="list-style-type: none">• Sequencing bias• No non-coding regions• Incidental findings	<ul style="list-style-type: none">• Data analysis bottleneck• Interpretation of non-coding variants• Expensive, time-consuming

Current exome sequencing workflow

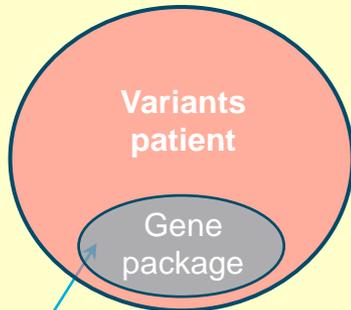


±350 exome samples per month!

Diagnostic exome approaches

Gene package approach

Most genes known



Variants in known genes

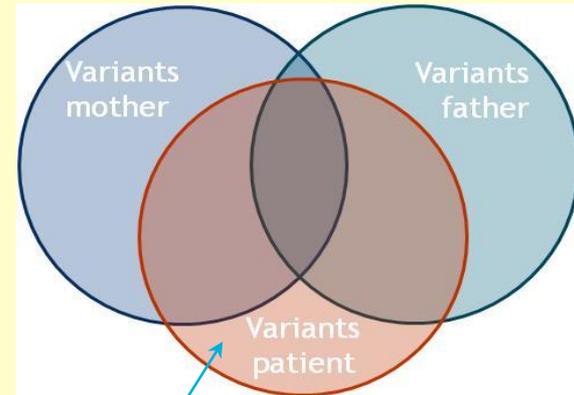
Pilot study: 50 exomes for 5 disorders

Neveling *et al.* Hum mut. 2013

- Sequence the exome of the patient
- Look only at known genes

Trio approach

Most genes unknown



De novo variants

Pilot study: 100 trios for intellectual disability

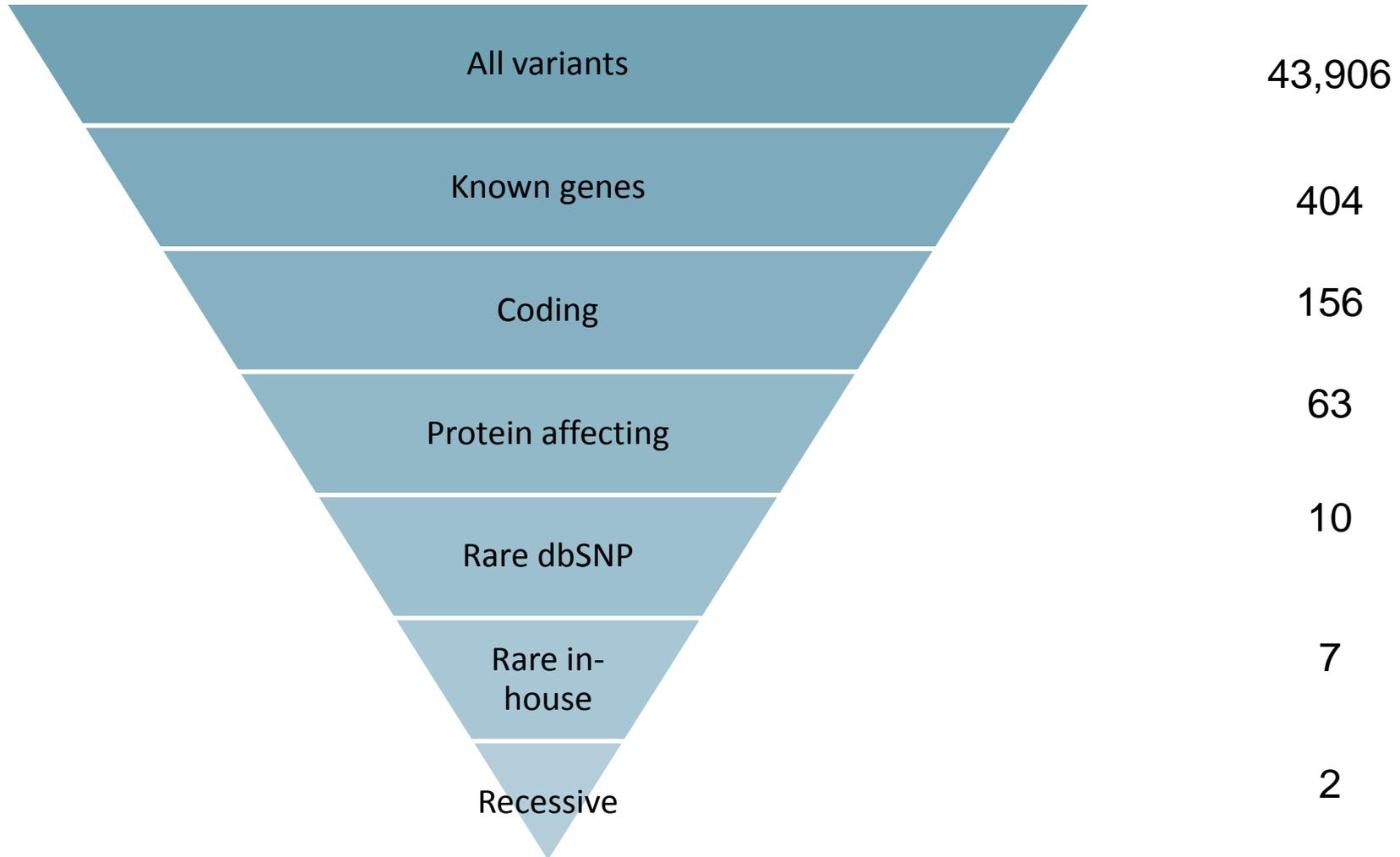
De Ligt *et al.* NEJM, 2012

- Sequence the exome of father, mother and patient.
- Look for *de novo* mutations

Challenges in clinical exome sequencing

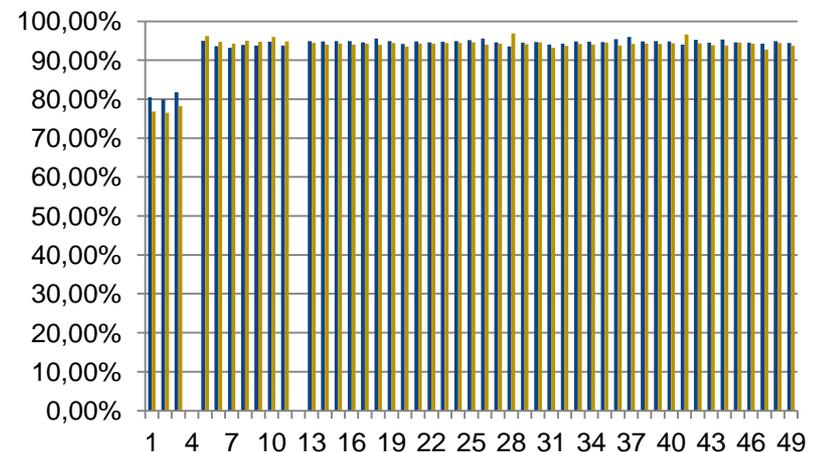
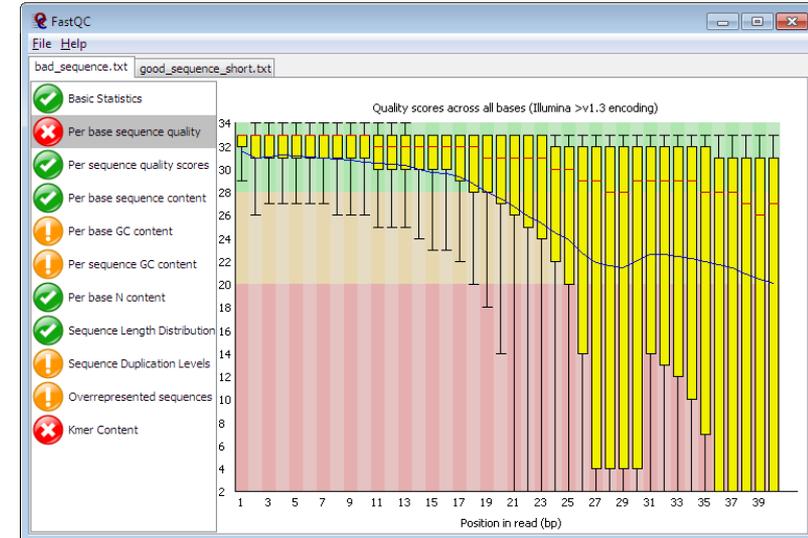
When it all works out...

- Prioritization of variants found by exome sequencing:

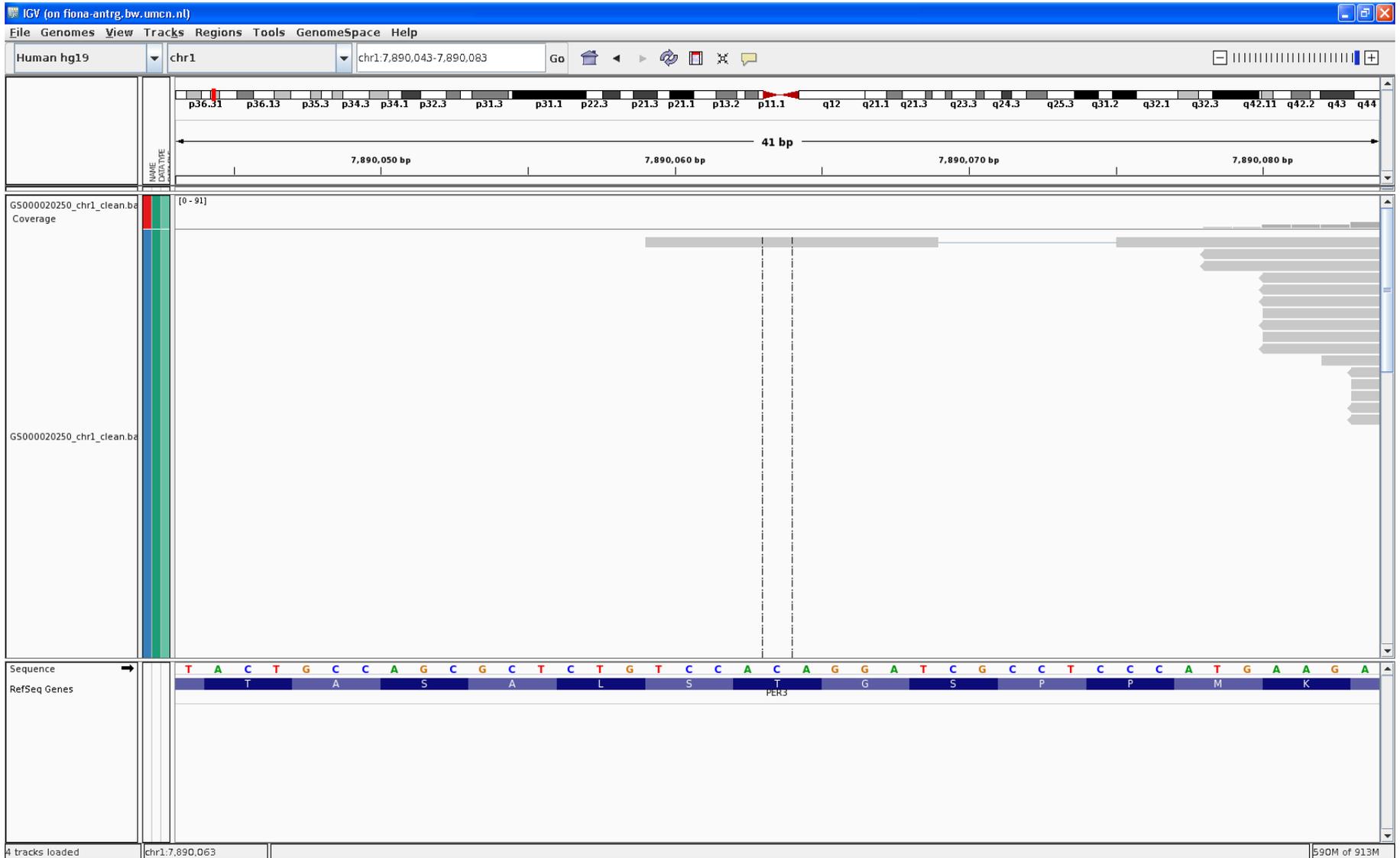


What can go wrong? Sample quality

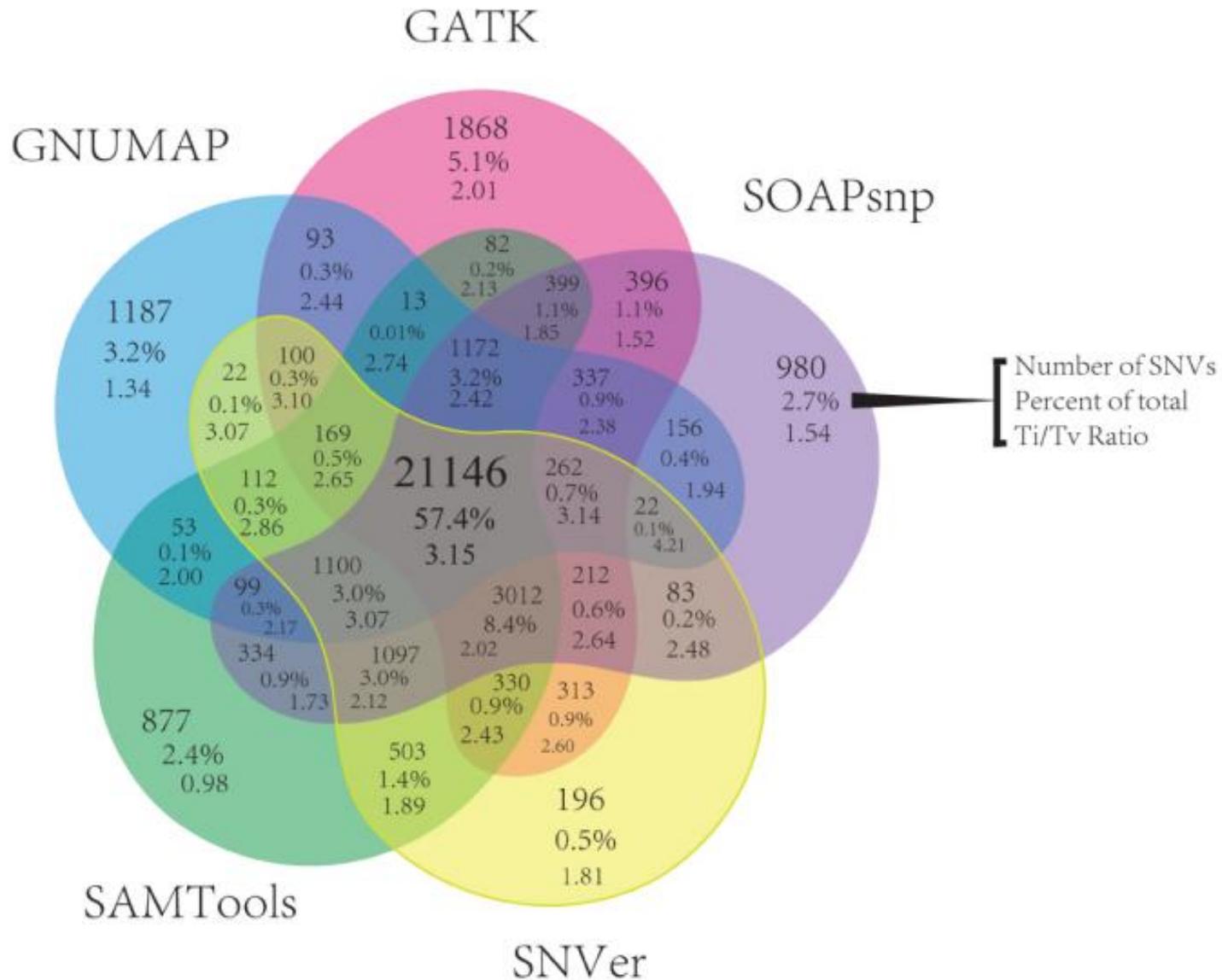
- **Raw sequence and mapping statistics**
 - FastQC tool
 - Bedtools – coverage statistics
 - Per gene / exon target coverage
- **Variant statistics:**
 - Overlap dbSNP
 - Number of truncating mutations
 - Tr/Ti ratio



What can go wrong? No sequence coverage



What can go wrong: uncalled variants



What can go wrong: incorrect interpretation

Variant Pos	rs ID	Alleles	EA Allele #	AA Allele #	All Allele #	Avg. Sample Read Depth	Genes	mRNA Accession #	GVS Function	Amino Acid
20:31021163	rs145699348	A/G	A=1/G=8599	A=0/G=4406	A=1/G=13005	81	ASXL1	NM_015338.5	missense	ILE,VAL
20:31021190	unknown	T/C	T=1/C=8599	T=0/C=4406	T=1/C=13005	90	ASXL1	NM_015338.5	missense	CYS,ARG
20:31021211	unknown	T/C	T=0/C=8600	T=2/C=4404	T=2/C=13004	92	ASXL1	NM_015338.5	stop-gained	stop,ARG
20:31021232	rs148964601	T/C	T=1/C=8599	T=0/C=4406	T=1/C=13005	90	ASXL1	NM_015338.5	missense	CYS,ARG
20:31021233	rs143719307	A/G	A=0/G=8600	A=1/G=4405	A=1/G=13005	90	ASXL1	NM_015338.5	missense	HIS,ARG
20:31021250	unknown	T/C	T=1/C=8599	T=0/C=4406	T=1/C=13005	87	ASXL1	NM_015338.5	stop-gained	stop,ARG
20:31021324	unknown	C/T	C=0/T=8600	C=1/T=4405	C=1/T=13005	98	ASXL1	NM_015338.5	coding-synonymous	none
20:31021332	unknown	G/C	G=1/C=8599	G=0/C=4406	G=1/C=13005	97	ASXL1	NM_015338.5	stop-gained	stop,SER
20:31021337	unknown	A/G	A=1/G=8599	A=0/G=4406	A=1/G=13005	97	ASXL1	NM_015338.5	missense	ILE,VAL
20:31021384	unknown	A/G	A=2/G=8598	A=0/G=4406	A=2/G=13004	102	ASXL1	NM_015338.5	coding-synonymous	none
20:31021389	unknown	A/G	A=1/G=8599	A=1/G=4405	A=2/G=13004	104	ASXL1	NM_015338.5	missense	ASN,SER
20:31021430	rs141346625	C/G	C=2/G=8598	C=17/G=4389	C=19/G=12987	107	ASXL1	NM_015338.5	missense	GLN,GLU
20:31021466	rs142172134	G/C	G=0/C=8600	G=10/C=4396	G=10/C=12996	106	ASXL1	NM_015338.5	missense	GLY,ARG
20:31021475	rs145913172	C/G	C=0/G=8600	C=2/G=4404	C=2/G=13004	109	ASXL1	NM_015338.5	missense	PRO,ALA
20:31021521	rs138971201	A/T	A=0/T=8600	A=14/T=4392	A=14/T=12992	131	ASXL1	NM_015338.5	missense	ASN,ILE
20:31021544	unknown	A/G	A=0/G=8600	A=1/G=4405	A=1/G=13005	140	ASXL1	NM_015338.5	missense	MET,VAL



Variation Color Code:
splice or nonsense or frameshift
missense
coding-synonymous
coding
utr
codingComplex

Bohring-Opitz syndrome is often fatal in early childhood.

<http://evs.gs.washington.edu/EVS/>

Better diagnostics by
whole exome sequencing?

-

Pilot for the gene panel approach

Pilot study – gene package approach

- **250 exomes:** 50 exomes for 5 genetically heterogeneous diseases
- Gene package design:
 - Only known genes are allowed, no candidate disease genes
 - Gene lists must be up-to-date and is updated every ~3 months
 - Created by team of experts from clinic, diagnostic and research division

	Number of genes (Sept. 2011)
Blindness	144
Deafness	98
Early onset colorectal cancer	115
Mitochondrial disorders	207
Movement disorders	152



How to do 400 samples per month?

The image shows a screenshot of the Variant Interface software. The main window displays a table of variants with columns for Chromosome, Start, End, reference, mutation, reads, variant, % var., Abbe., Score, Novel., Uniq., De novo assessment, SNP id, SNP, Causative - Freque., Causative, NonCausative - Frequency, NonCausative - Projects, and Gene name. The table lists various variants across chromosomes 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, X, and Y.

Overlaid on the table is the text "Variants and annotation".

Below the table is a "Meta data" panel with the following information:

Info	Data
mapping_folder	/data/results/MR/DNA1124224Z
Run date	2012-05-18
analysis_date	110612
ANALYSIS_OUTPUT_FOLDER	diag_out_110612
project_name	DNA1124224Z
project_dir_name	MR
diffs_file	/data/results/MR/DNA1124224Z/MR...
indel_file	/data/results/MR/DNA1124224Z/MR...
bam_file	/data/results/MR/DNA1124224Z/MR...
genome_build	hg19
Regions file	/data/raid0/references/hg19/SureSel...
dbSNP file	/data/raid0/pipe/lenp135
Pipeline version	110612
RefGene version	110612
Uniprot version	0125
Omic version	0125
bioscope_version	2.1
flow_cell	FC1
machine	5500ml_2
run	D18
XSQ file check	6
AANVRAGER	Mw. dr. T. Kleefstra
GEPLANDE_EINDDATUM	20-11-2012
ONDE_STATUS	G
ONDERZOEKSTYPE	D
ONDE_INDICATIES	mr
ALLE_FAMILIENRS	07-0236
GEBORTE DATUM	05-10-2006
GESLACHT	M
ONDE_FRACTIENUMMERS	DNA09-01550.DNA11-22681.DNA11...
FRACTIENUMMER	DNA11-24224Z
PyroSNP Concordance	1.0
Total mapped reads in regions	97087491

Overlaid on the meta data panel is the text "Quality control".

On the right side of the interface is a "Filters" panel with the following options:

- Open exome (OE)
- Less 1% in DBE (M)
- Disease variant (D)
- Truncating variants (S)
- De novo (N)
- Exonic/canonical SS (E)
- Causative (CA)
- Conserved (C)
- Non-Synonymous (NS)
- Dominant (AD)
- Less 5% in dbSNP (P)
- All Splice Sites (SS)
- Recessive (AR)

Below the filters are "Predefined views" (Disease variants (No filters)), a "Search" box, and "Hidden columns".

At the bottom right, the following statistics are shown:

- Total variants: 43092
- Disease variants: 1046
- Shown variants: 1046
- Highlighted variants: 0

Overlaid on the bottom right is the text "Filtering".

How to do 400 samples per month?

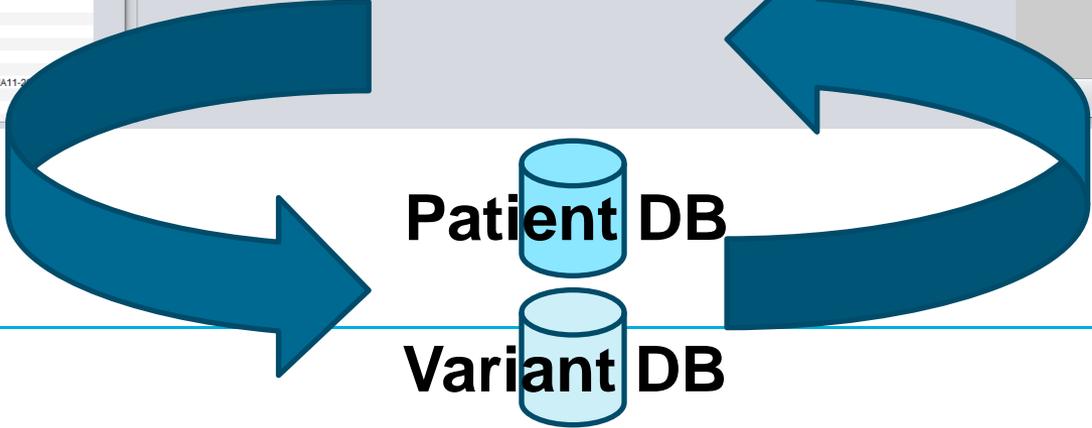
The screenshot shows the Variant Interface software with the following components:

- Variant data table:** A table with columns for Chromosome, Start, End, reference, mutation, reads, variant, % var., Abbe., Score, Novel., Uniq., De novo assessment, SNP id, SNP, Causative - Freq., Causative, NonCausative - Frequency, NonCausative - Projects, and Gene name. It lists various variants across chromosomes 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24.
- Meta data table:** A table with columns for Info and Data, listing project details like mapping_folder, Run date, analysis date, project_name, diff_file, bam_file, genome_build, Regions file, dsNP file, Pipeline version, RefGene version, Uniprot version, Omim version, bioscope_version, flow_cell, machine, run, XSQ file check, AANVRAGER, GEPLANDE_EINDDATUM, ONDE_STATUS, ONDERZOEKSTYPE, ONDE_INDICATIES, ALLE_FAMILIENRS, GEBORTE DATUM, GESLACHT, ONDE_FRACTIENUMMERS, FRACTIENUMMER, and PyroSNP Concordance.
- Filters panel:** A panel on the right with checkboxes for various filters such as Open exome (OE), Less 1% in MDB (M), Disease variant (D), Truncating variants (S), De novo (N), Exonic/canonical SS (E), Causative (CA), Conserved (C), Non-Synonymous (NS), Dominant (AD), Less 5% in dbSNP (P), All Splice Sites (SS), and Recessive (AR). It also includes checkboxes for Low % variation reads (QP) and Low number variation reads (QV).
- Predefined views:** A dropdown menu showing 'Disease variants (No filters)'. A search box and 'Clear search box' button are also present.
- Hidden columns:** A dropdown menu for selecting hidden columns.
- Status:** Total variants: 43092, Disease variants: 1046, Shown variants: 1046, Highlighted variants: 0.
- Buttons:** 'Filter', 'Search', 'Clear search box', 'Uncheck and reload'.

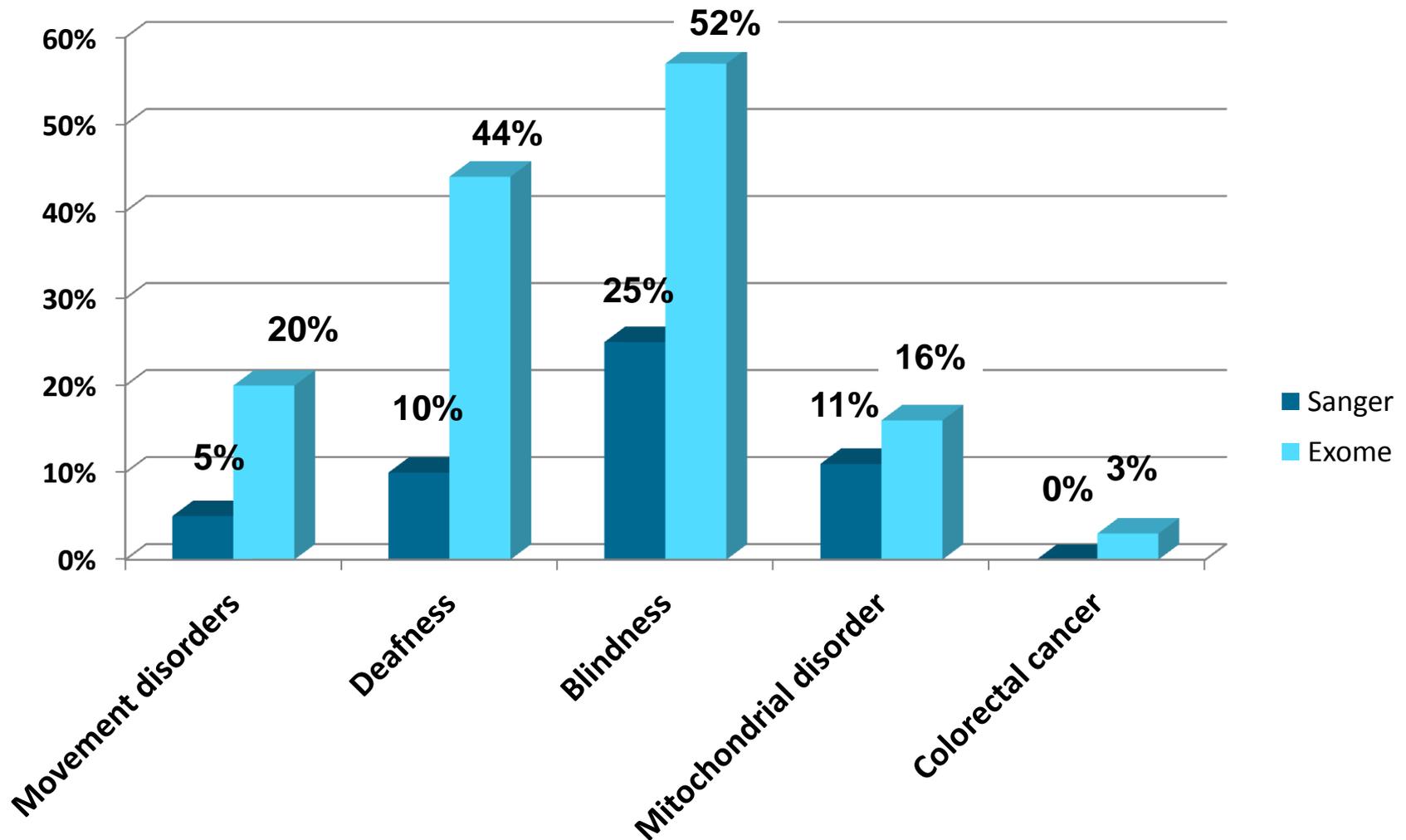
Variants and annotation

Quality control

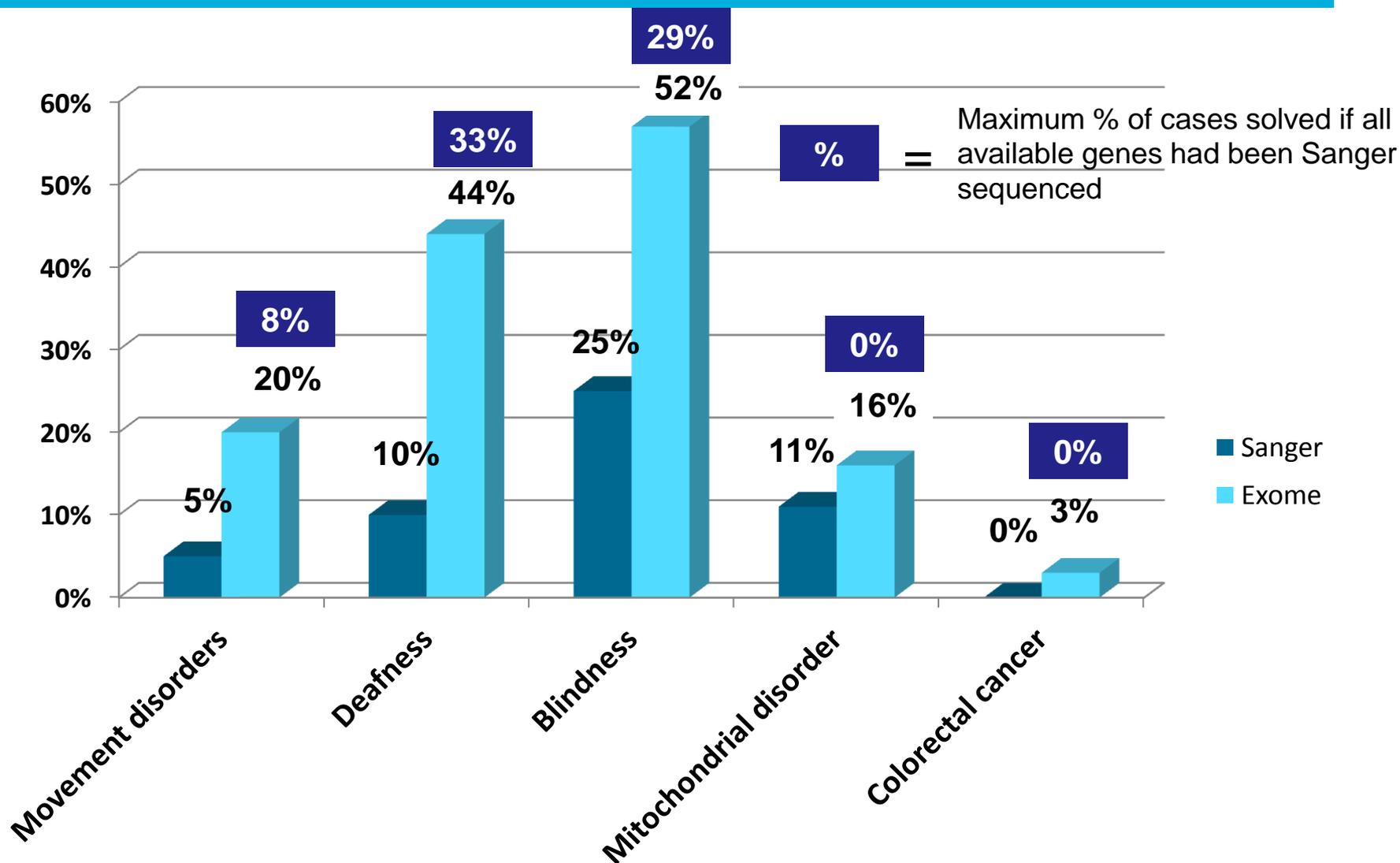
Filtering



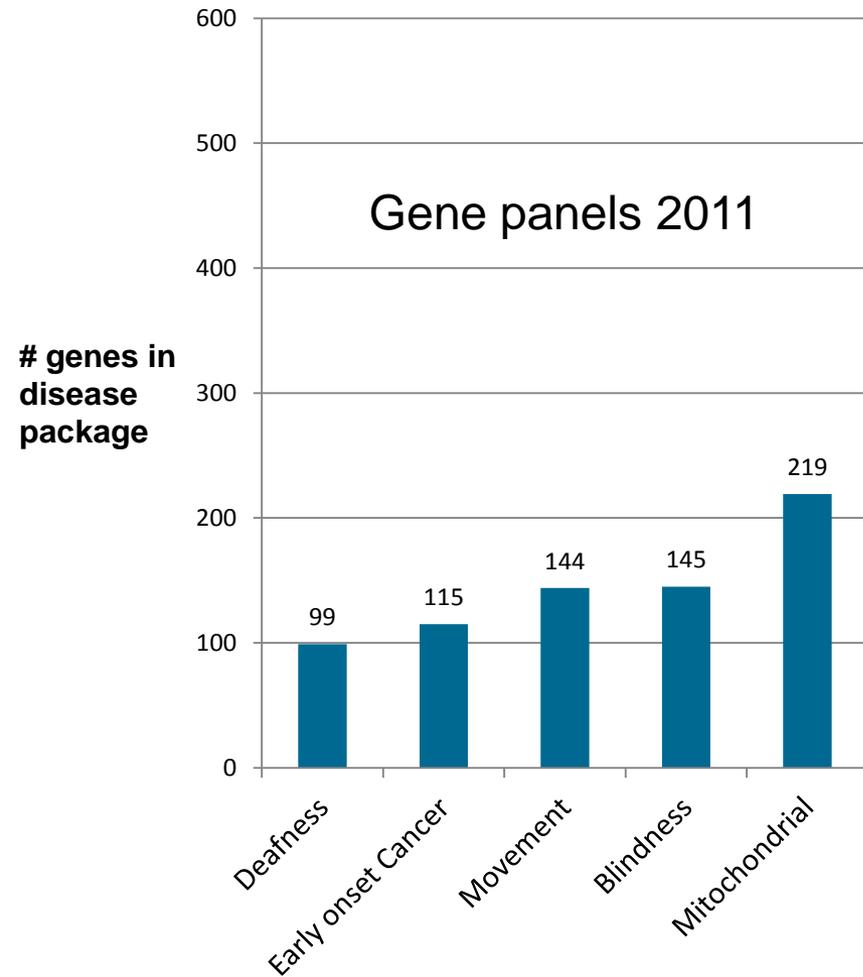
Diagnostic yield from exome sequencing



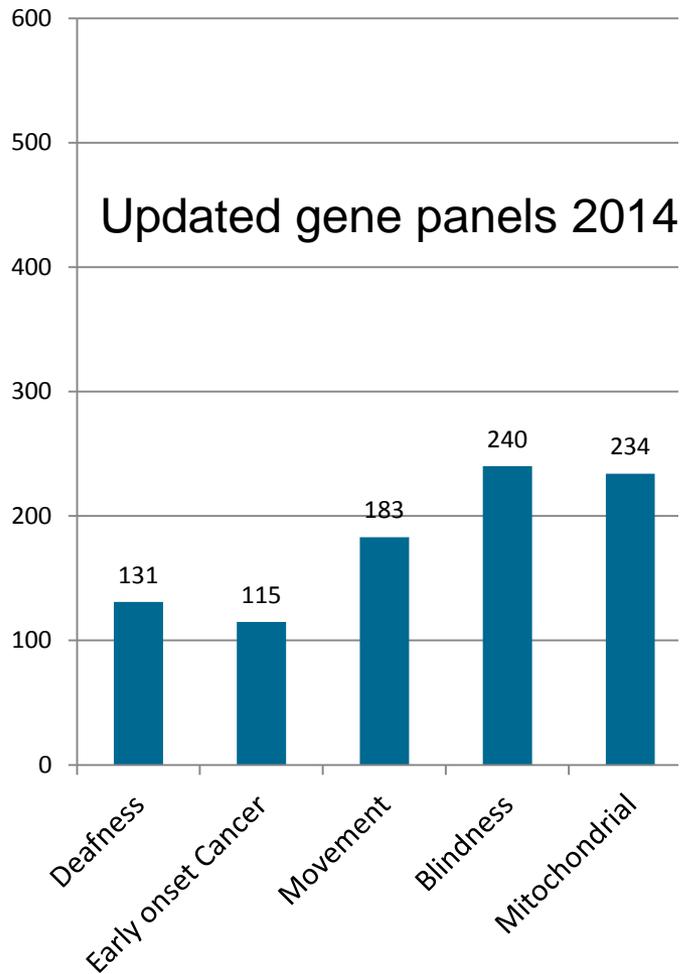
Diagnostic yield from exome sequencing



Current exome panels

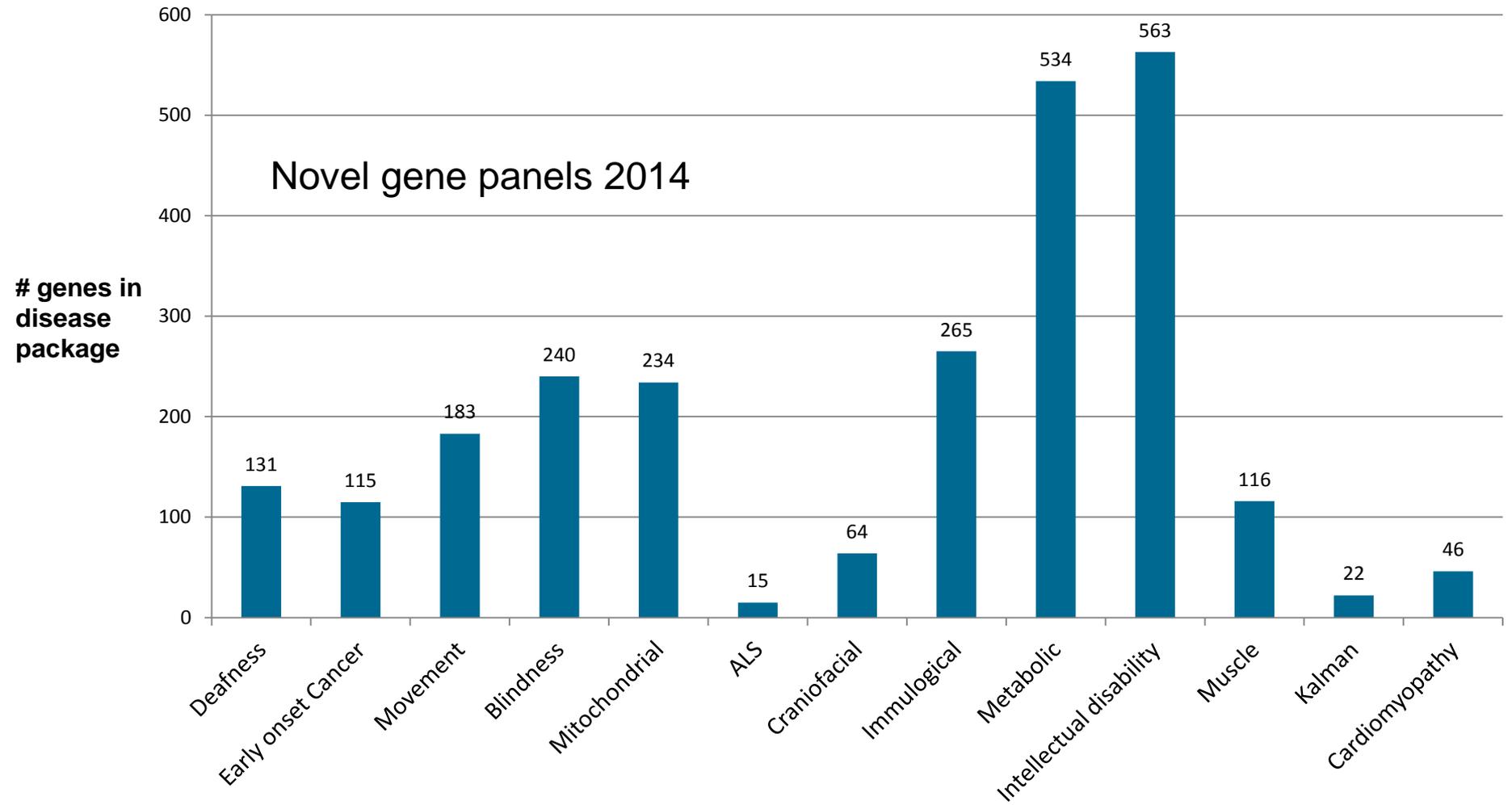


Current exome panels



Exome sequencing can be cost-efficient compared to Sanger when sequencing 3 genes or more..

Current exome panels

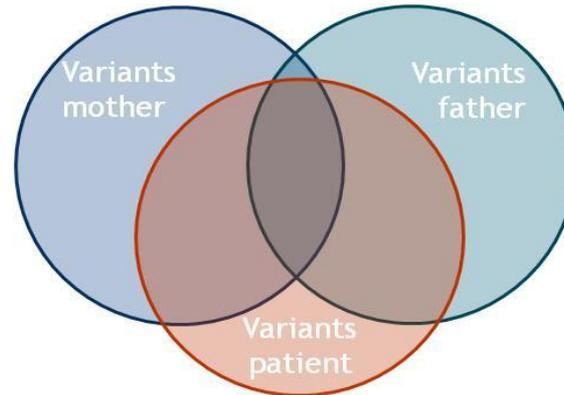


Better diagnostics by
whole exome sequencing?

-

Pilot for the trio approach

Pilot study – *de novo* approach



- **100 patients + 200 parents!**
 - Severe intellectual disability (IQ<50)
 - No etiological or syndromic diagnosis
 - Negative family history
- Patients have reached the end stage of conventional strategies
 - Targeted gene tests negative
 - Genomic array profile negative

Yield in 100 ID patients

Positive diagnosis	June 2012	June 2013
All mutations	16	29
<i>De novo</i> mutations	13	28
Autosomal dominant	10	23
X-linked	2	4
Autosomal recessive	1	1
Inherited mutations	3	1
X-linked	3	1
Autosomal recessive	0	0
Candidates	19	11

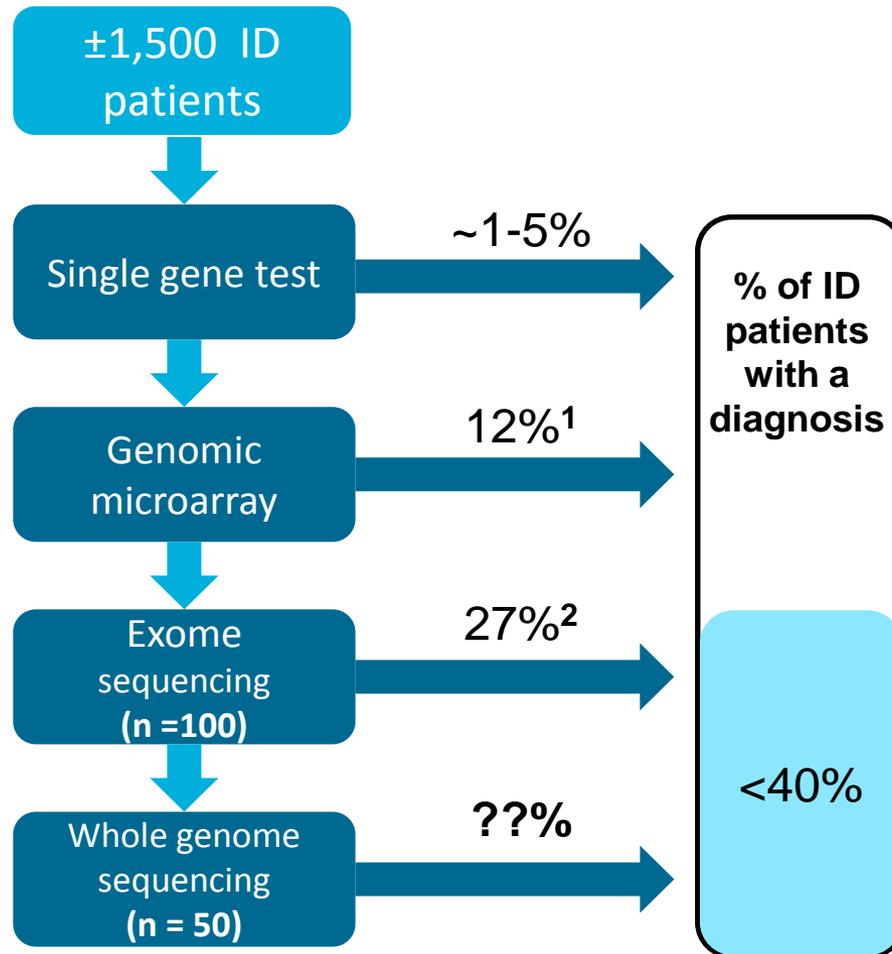
Yield of ~30% in patients with severe ID

Better diagnostics by
whole genome sequencing?

-

pilot study

Diagnosis in patients with severe intellectual disability (ID)

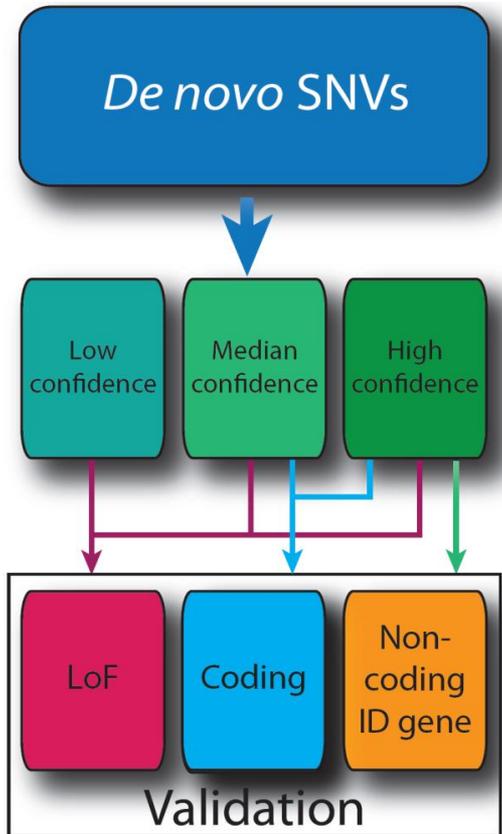


- Whole genome sequencing
- 50 trios at 80x coverage
- A *de novo* approach

¹ Vulto-van Silfhout, A. T. *et al* Hum mut 2013

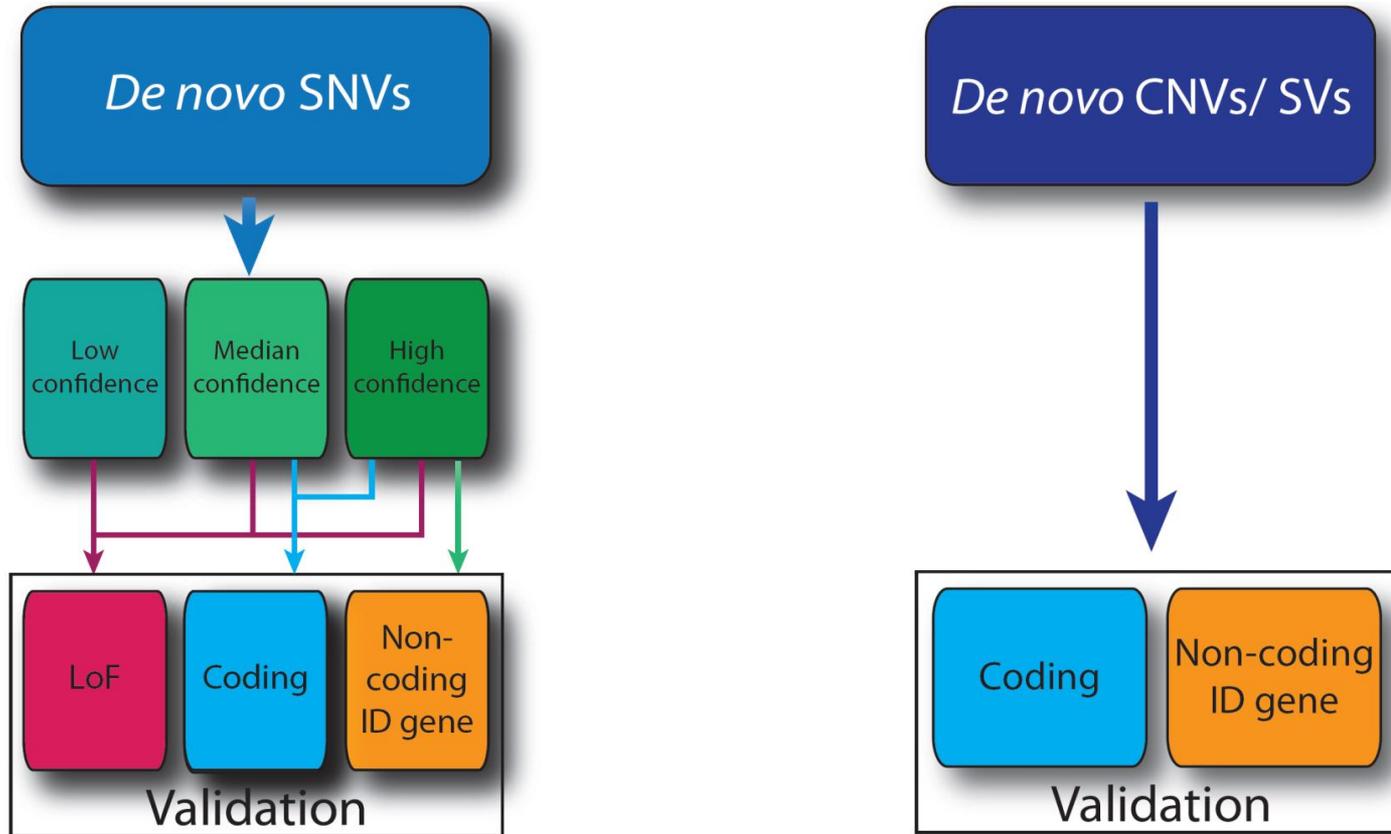
² de Ligt, J. *et al*. NEJM 2012

Can we identify *de novo* mutations?



- Validation rate: **38%**
(80% for high confidence!)
- Coding *de novo* SNVs: **84**
- Comparison to WES SNVs: **12%**

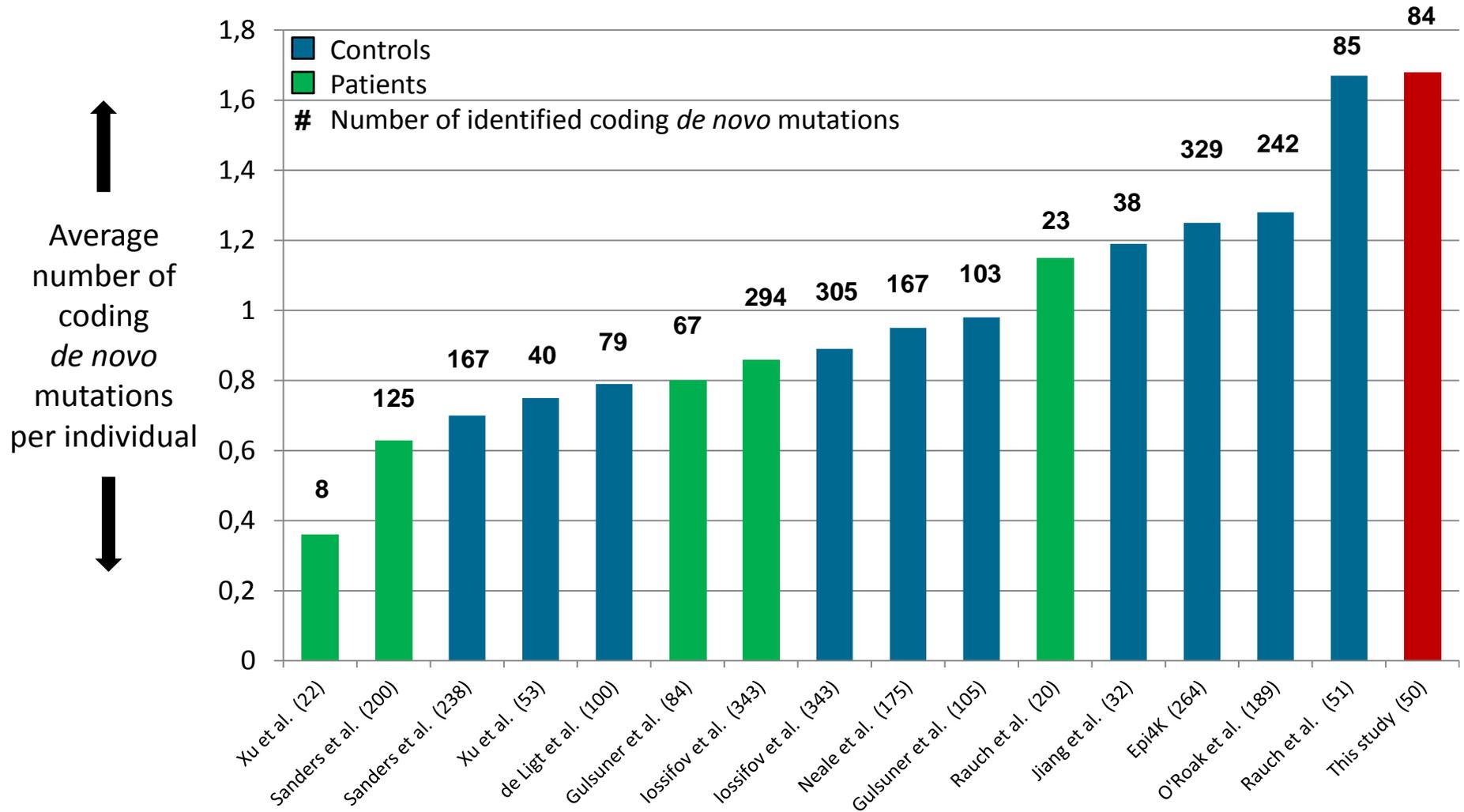
Can we identify *de novo* mutations?



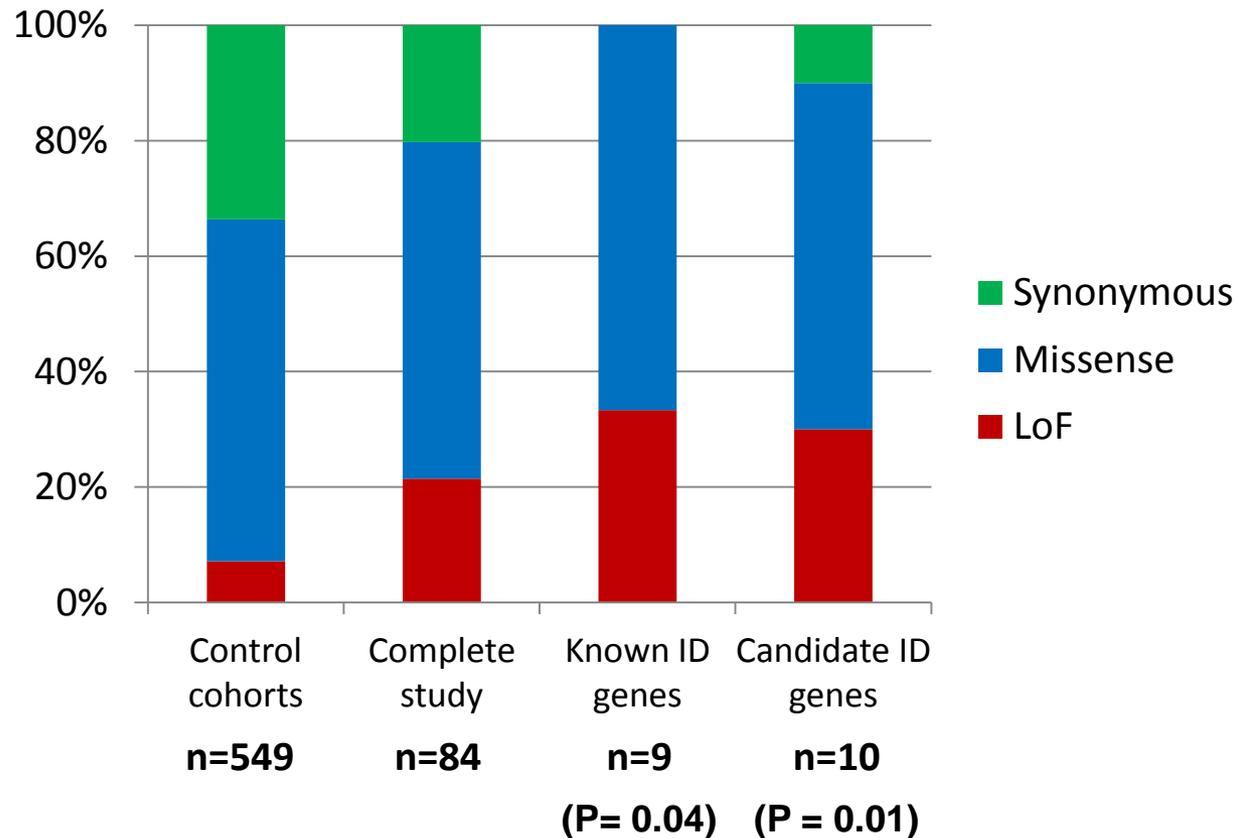
- Validation rate: **38%**
(80% for high confidence!)
- Coding *de novo* SNVs: **84**
- Comparison to WES SNVs: **12%**

- Validation rate: **82%**
- Coding *de novo* SVs: **9**

Did we find more coding *de novo* mutations than other studies?



Are these *de novo* SNVs related to ID?

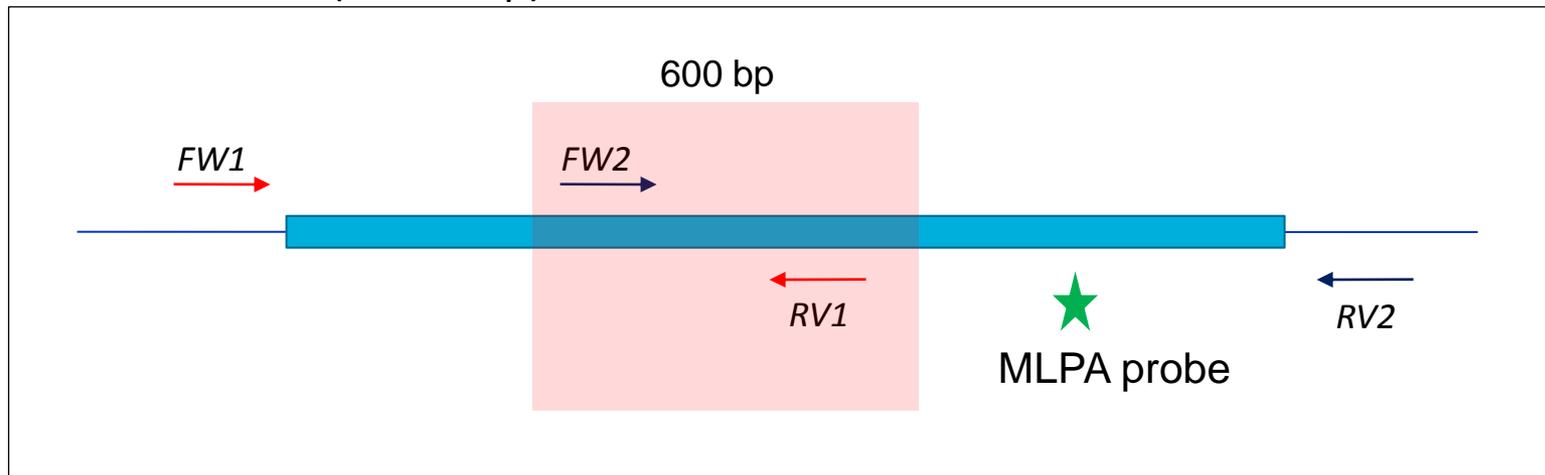


- Significantly more *de novo* mutations in known/candidate ID genes
- Significantly more loss-of-function mutations ($P=4.8^{-06}$, $P=0.02$)

De novo structural variants, example 1

- A patient with the clinical suspicion of Rett syndrome
- *MECP2* gene tested by Sanger sequencing but no mutations identified

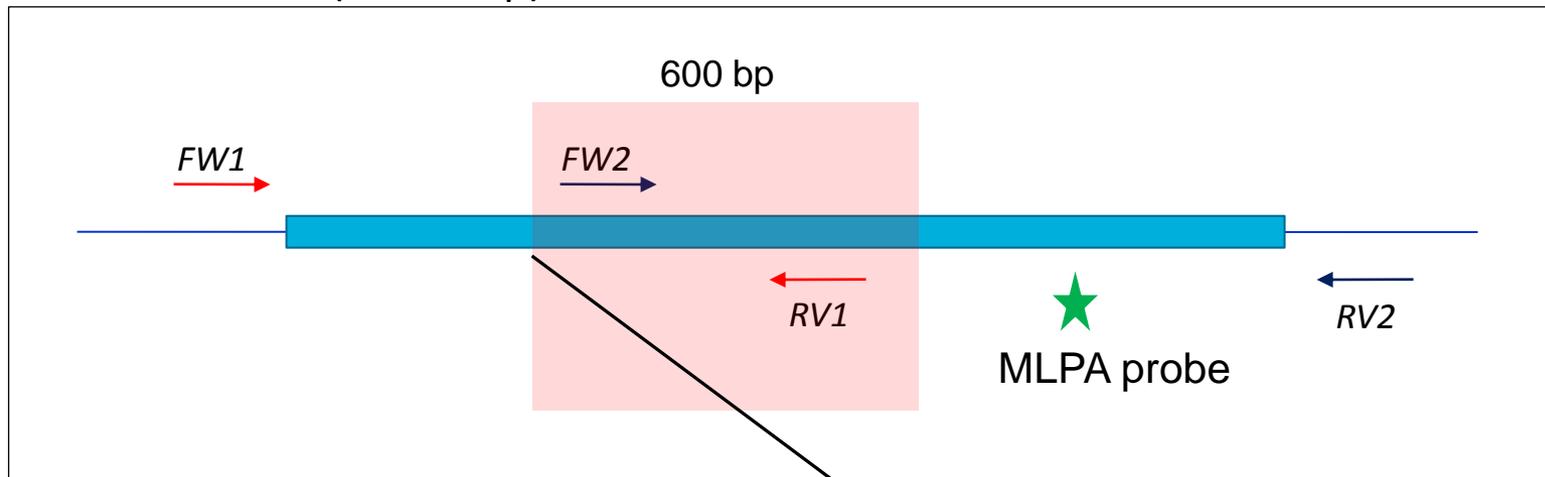
MECP2, exon 4 (~1000 bp)



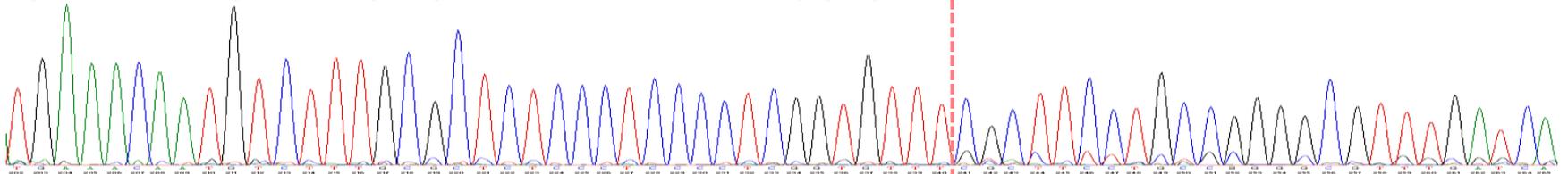
De novo structural variants, example 1

- A patient with the clinical suspicion of Rett syndrome
- *MECP2* gene tested by Sanger sequencing but no mutations identified

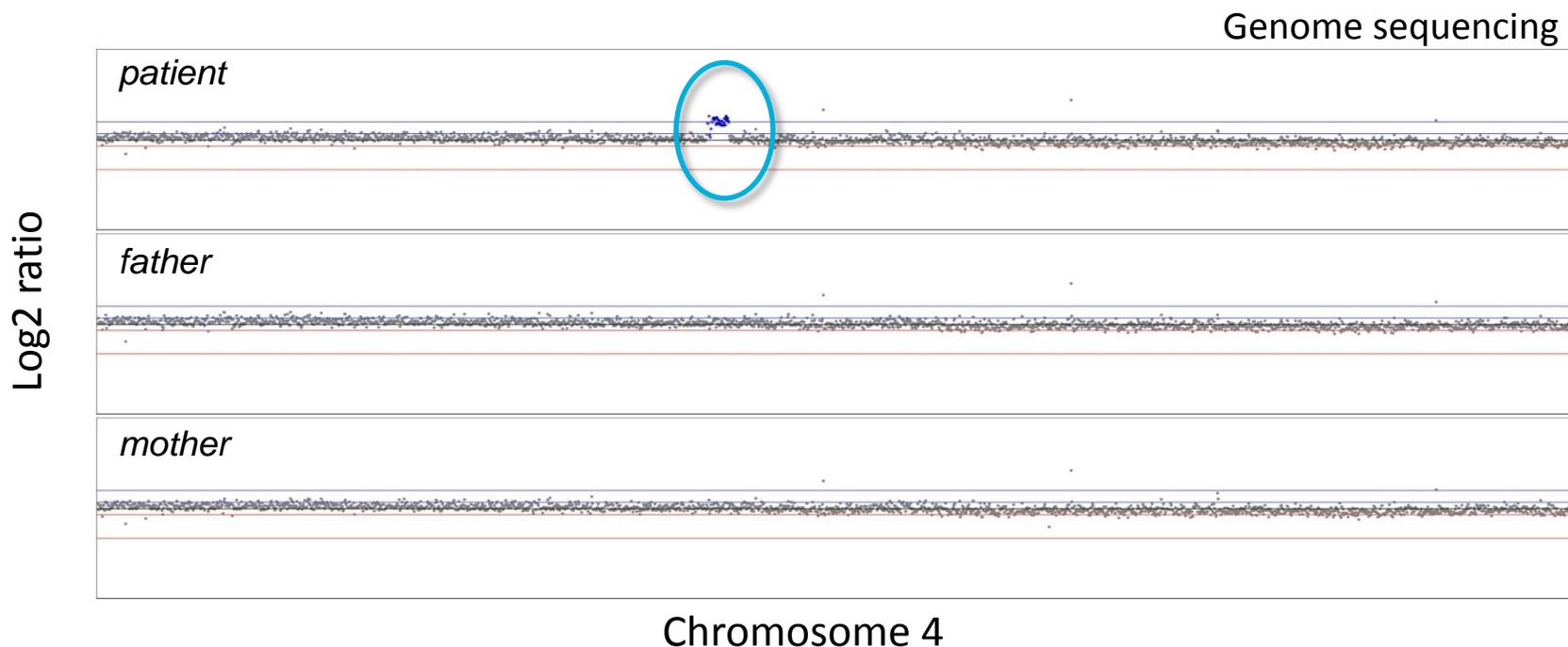
MECP2, exon 4 (~1000 bp)



t g a a c a a t g t c t t t g c g c t c t c c c t c c c c t c g g t g t t t c g c t t c c t g c c g g g g c g t t t g a t c a

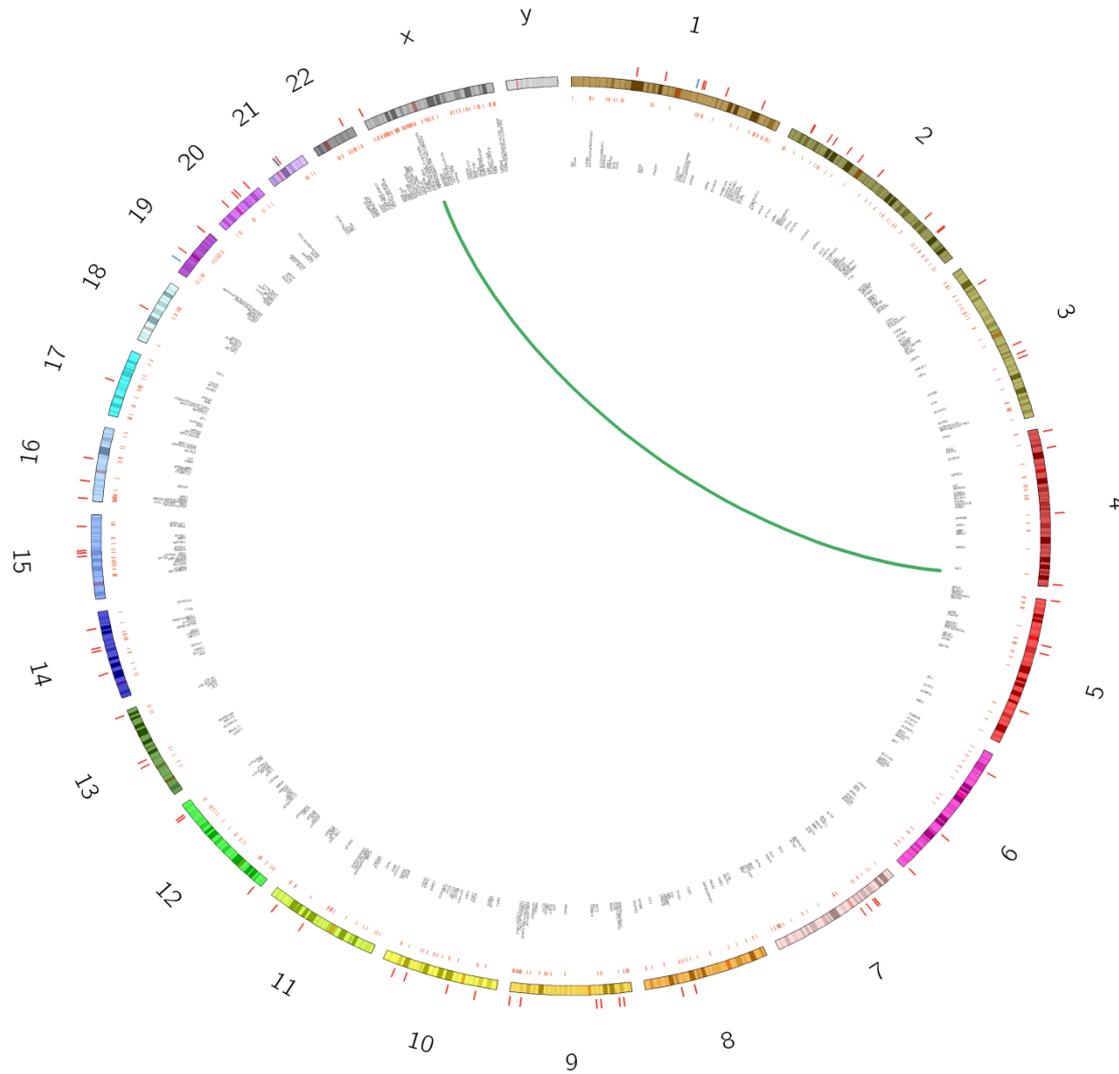


Example 2: A *de novo* duplication on chromosome 4

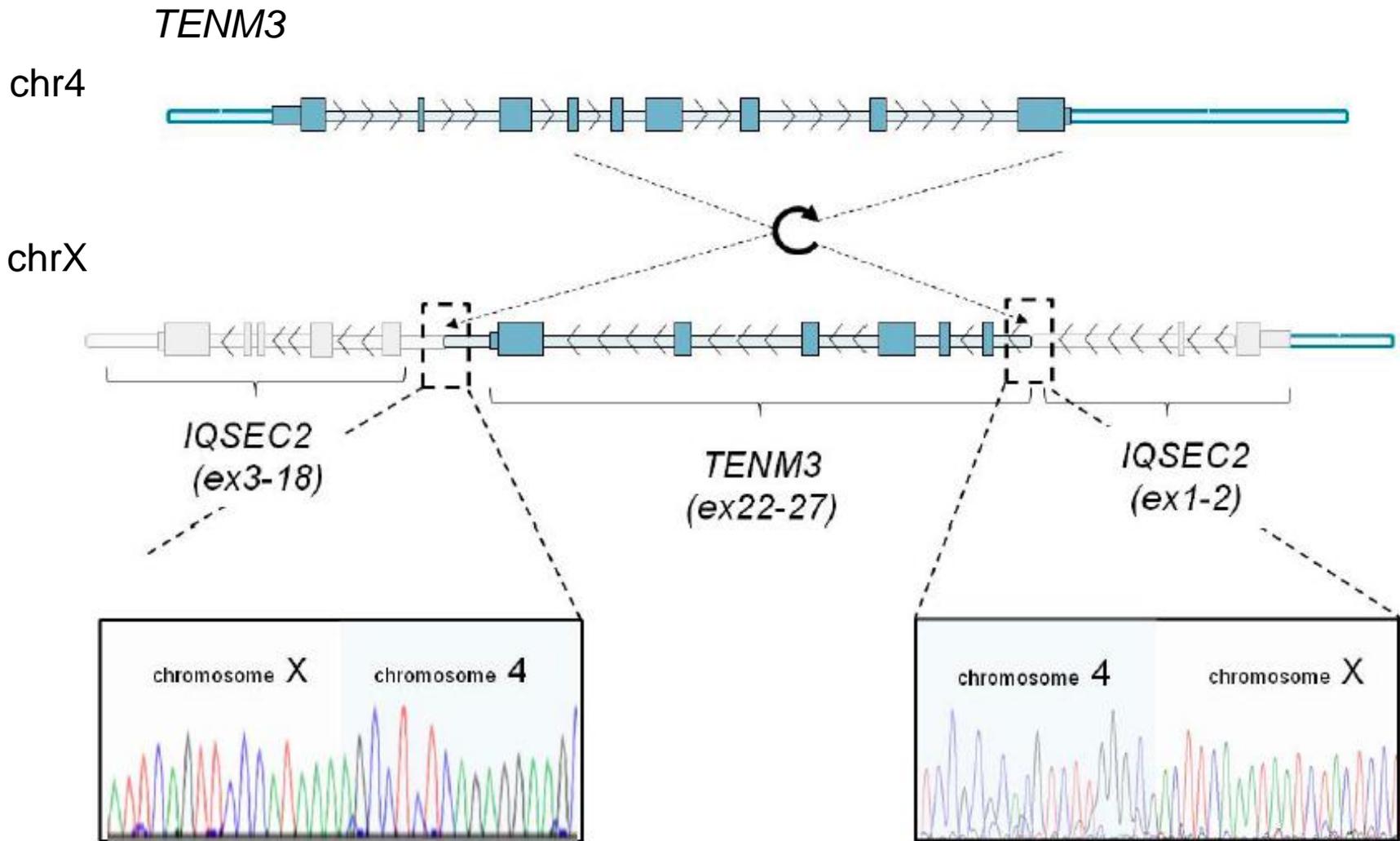


- ~60 kb *de novo* duplication on chromosome 4
- Affecting the last 6 exons of the *TENM3* gene
- *TENM3* is associated with coloboma, and microphthalmia

Duplication from chr4 to chrX



Duplication from chr4 to chrX



Only possible with genome sequencing!

Diagnostic yield genome sequencing

Highly likely diagnosis¹

SNV	SV
<i>TBR1 (2x)</i>	<i>SHANK3</i>
<i>WDR45</i>	<i>VPS13B</i> *
<i>SMC1A</i>	<i>MECP2</i>
<i>SPTAN1</i>	<i>IQSEC2</i>
<i>RAI1</i>	<i>STAG1</i>
<i>MED13L</i>	<i>SMC1A</i>
<i>SATB2</i>	16p11.2 microdel. syndrome
<i>PPP2R5D</i>	Multiple genes
<i>KCNA1</i>	
<i>SCN2A</i>	
<i>POGZ</i>	
<i>KANSL2</i>	

21/50 cases diagnosed: 42%

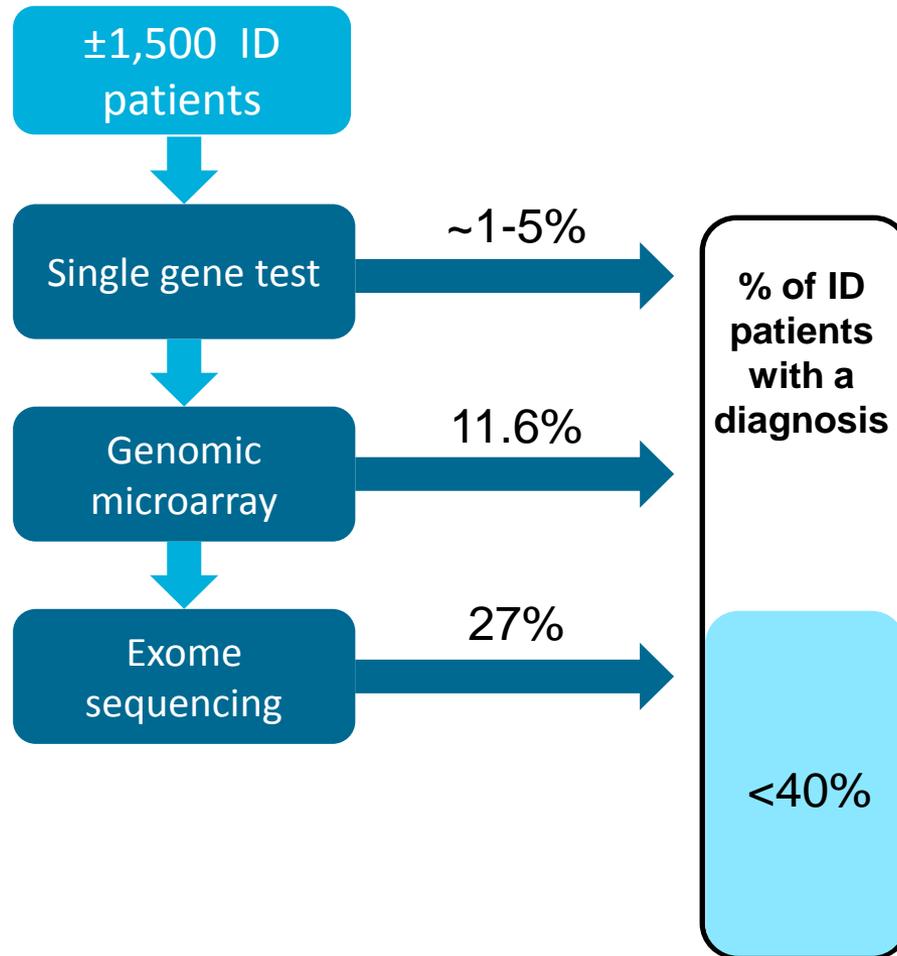
* Recessive variants. Known genes in bold.

Diagnostic yield genome sequencing

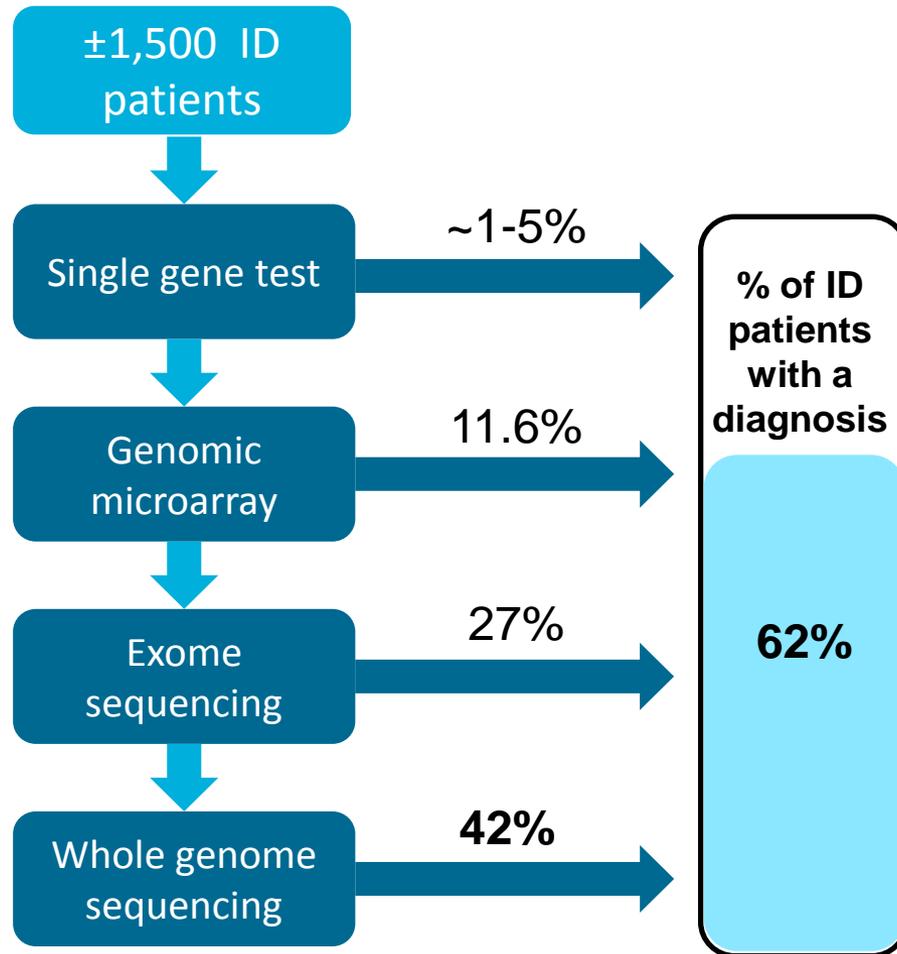
Highly likely diagnosis ¹		Possible diagnosis
SNV	SV	SNV
<i>TBR1 (2x)</i>	<i>SHANK3</i>	<i>NGFR</i>
<i>WDR45</i>	<i>VPS13B</i> *	<i>GFPT2</i>
<i>SMC1A</i>	<i>MECP2</i>	<i>WWP2</i>
<i>SPTAN1</i>	<i>IQSEC2</i>	<i>ASUN</i>
<i>RAI1</i>	<i>STAG1</i>	<i>BRD3</i>
<i>MED13L</i>	<i>SMC1A</i>	<i>MAST1</i>
<i>SATB2</i>	16p11.2 microdel. syndrome	<i>APPL2</i>
<i>PPP2R5D</i>	Multiple genes	<i>NACC1</i>
<i>KCNA1</i>	21/50 cases diagnosed: 42%	
<i>SCN2A</i>		
<i>POGZ</i>		
<i>KANSL2</i>		

* Recessive variants. Known genes in bold.

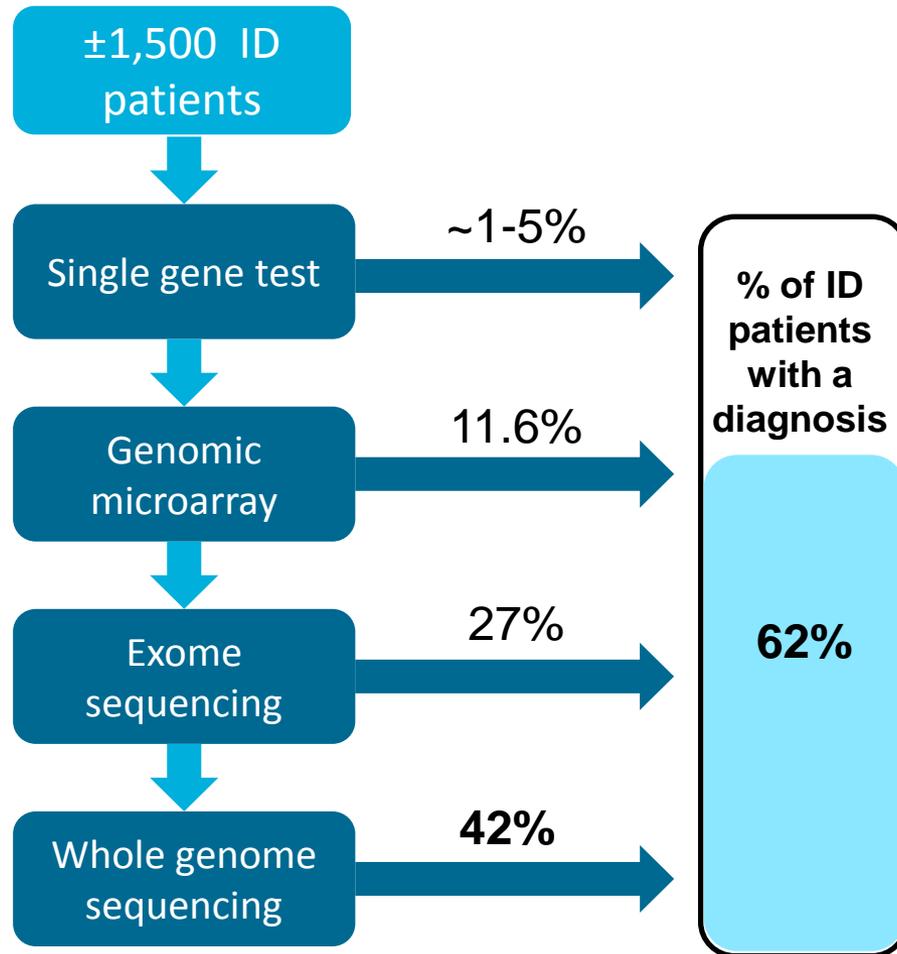
Diagnostic yield



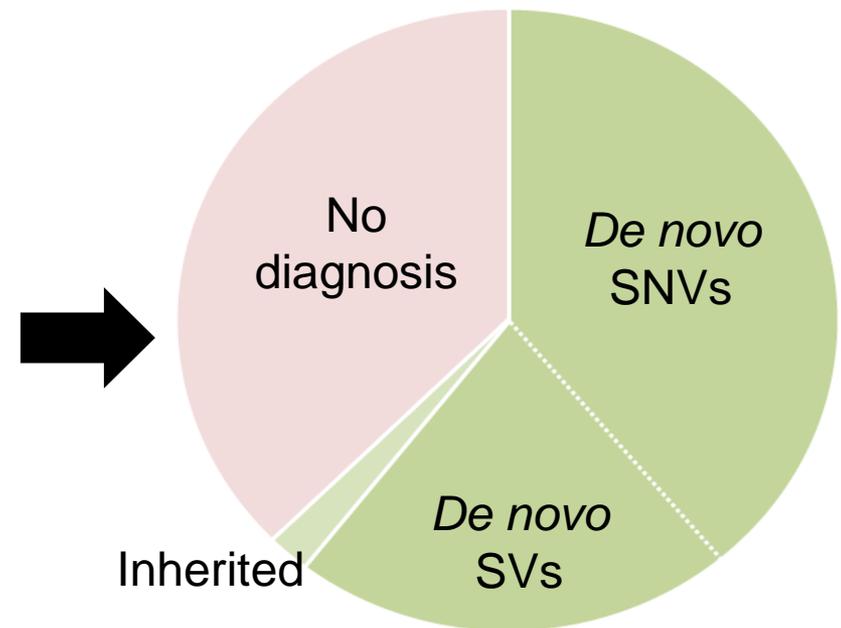
Diagnostic yield



Diagnostic yield



- **Majority is *de novo*!**



Gilissen *et al.* Nature 2014

Conclusions

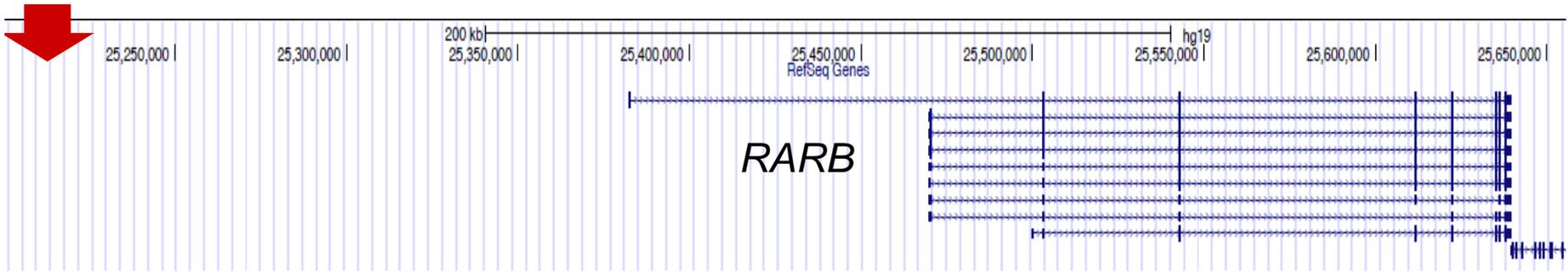
1. Exome sequencing results in a higher diagnostic yield for genetically heterogeneous diseases than Sanger-based approaches
2. Based on clinical diagnostic criteria we can provide a genetic diagnosis for the majority of severe ID cases by using WGS
1. *De novo* coding mutations are the major cause of severe ID

Discussion

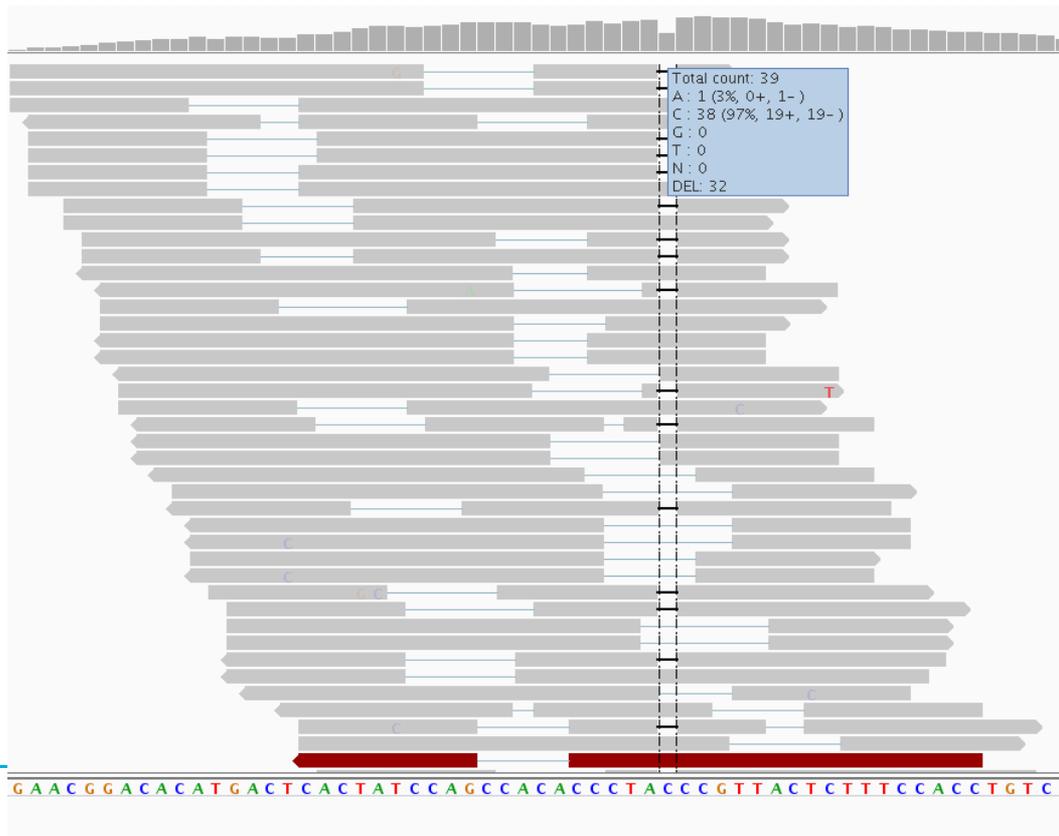
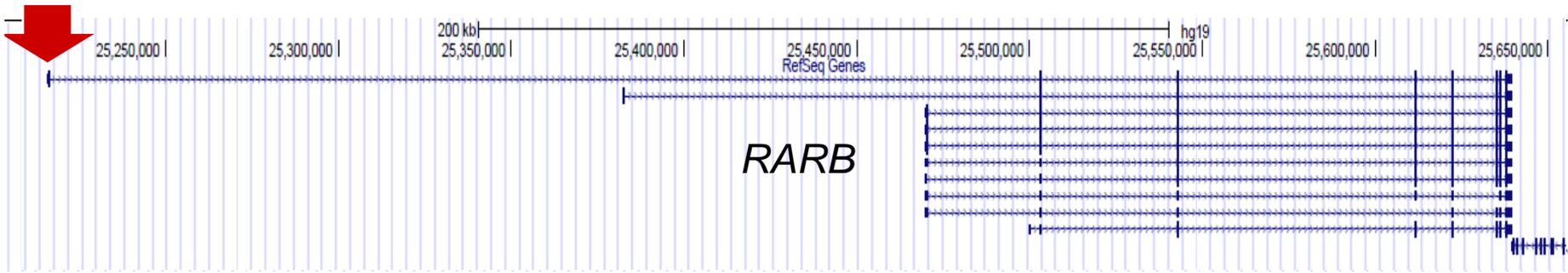
Other advantages of genome sequencing:

- Comprehensive set of variants: novel genes/exons annotations are still being added
- Non-coding variants: *de novo* mutations introducing novel splice sites
- Mosaic variants: variants occurring in sub-population of cells
- Phasing: identifying whether two variants are on the same allele

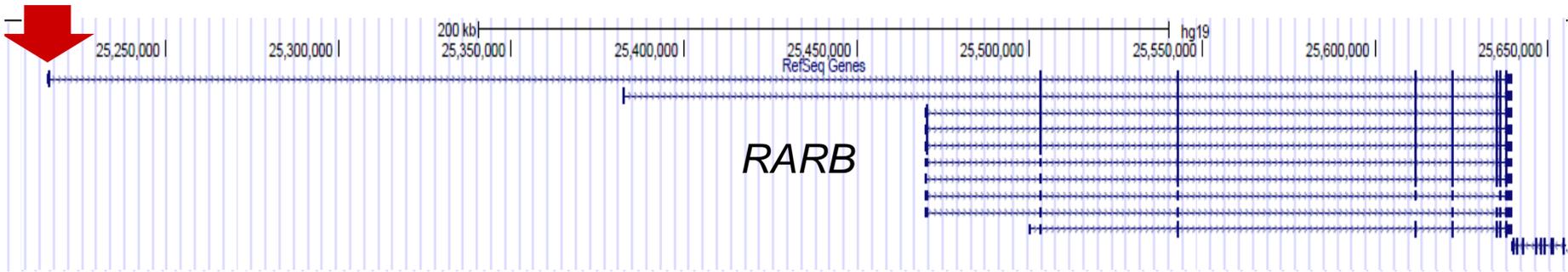
One more coding *de novo* variant?



One more coding *de novo* variant?



One more coding *de novo* variant?



- Mutations in the original *RARB* transcript are a known cause for Microphthalmia
- Validation confirmed: coding *de novo* mutation #85
- Expression analysis: this specific exon is expressed in brain / fetal brain.

Acknowledgments



Lisenka Vissers



Kornelia Neveling



Alex Hoischen



Jayne Hehir-Kwa



Joris Veltman



Marcel Nelen



Hans Scheffer



Han Brunner

Genome Diagnostics

Helger Yntema

Erik-Jan Kamsteeg

Lies Hoefsloot

Willy Nillesen

Marjolijn Ligtenberg

Arjen Mensenkamp

Dorien Lugtenberg

Rolph Pfundt

Genome Research

Djie Thung

Maartje van de Vorst

Rick de Reuver

Marisol del Rosario

Nienke Wieskamp

Petra de Vries

Michael Kwint

Irene Janssen

Marloes Steehouwer

Clinical genetics

Marjolein Willemsen

Tjitske Kleefstra

Ernie Bongers

David Koolen

Anneke Vulto-van

Silfthout

Wendy van Zelst-Stams



European
Research
Council



Diagnostic exome sequencing in 100 ID trios

Level of ID	Number of patients
IQ <30	62
IQ 30-50	38
IQ 50-70	0

Gender	
Male	47
Female	53

Age groups	
<10 yrs	37
10-20 yrs	41
>20 yrs	22

Sibship size	
1	12
2	47
3	36
4	1
≥5	2
unknown	2

Number of major congenital anomalies	
0	62
1	31
2	7
3	0

