

Supplement to:

## Genome sequencing as a generic diagnostic strategy for rare disease

Gaby Schobers<sup>1,2</sup>, Ronny Derks<sup>1</sup>, Amber den Ouden<sup>1</sup>, Hilde Swinkels<sup>1</sup>, Jeroen van Reeuwijk<sup>1,2</sup>, Ermanno Bosgoed<sup>1</sup>, Dorien Lugtenberg<sup>1</sup>, Su Ming Sun<sup>3</sup>, Jordi Corominas Galbany<sup>1,2</sup>, Marjan Weiss<sup>1</sup>, Marinus J. Blok<sup>3</sup>, Richelle A.C.M. Olde Keizer<sup>1,2</sup>, Tom Hofste<sup>1</sup>, Debby Hellebrekers<sup>3</sup>, Nicole de Leeuw<sup>1</sup>, Alexander Stegmann<sup>3</sup>, Erik-Jan Kamsteeg<sup>1</sup>, Aimee D.C. Paulussen<sup>3</sup>, Marjolijn J.L. Ligtenberg<sup>1,2</sup>, Xiangqun Zheng Bradley<sup>4</sup>, John Peden<sup>4</sup>, Alejandra Gutierrez<sup>4</sup>, Adam Pullen<sup>4</sup>, Tom Payne<sup>4</sup>, Christian Gilissen<sup>1,2</sup>, Arthur van den Wijngaard<sup>3</sup>, Han G. Brunner<sup>1,2</sup>, Marcel Nelen<sup>1\*</sup>, Helger G. Yntema<sup>1,2\*</sup>, Lisenka E.L.M. Vissers<sup>1,2\*</sup>

<sup>1</sup> Department of Human Genetics, Radboudumc, Nijmegen, Netherlands

<sup>2</sup> Research Institute for Medical Innovation, Radboudumc, Nijmegen, Netherlands

<sup>3</sup> Department of Clinical Genetics, Maastricht University Medical Center, Maastricht, Netherlands

<sup>4</sup> Illumina Inc., Cambridge, United Kingdom

\*Contributed equally

<b>Figure S1</b> Scenario model to determine individuals eligible for GS-first strategy.....	2
<b>Figure S2</b> A cohort of 1000 cases with clinically relevant variants spanning the broad range of genome diagnostics.....	3
<b>Figure S3</b> The average output of 1000 genomes.....	4
<b>Figure S4</b> GS Technical validation by variant type and assessment of why variants were not identified.....	5
<b>Figure S5</b> Examples of comprehensive GS.....	6
<b>Figure S6</b> In silico coverage statistics at variant level and disease genes.....	7
<b>Figure S7</b> Schematic representation of referrals to Radboudumc and MUMC+ in 2022.....	8
<b>Figure S8</b> Schematic overview of assumptions made to evaluate the impact on diagnostic yield from transition to a generic GS approach.....	9
<b>Table S3</b> GS sensitivity: Overview of TPRs per workflow.....	10

### (Additional file 1)

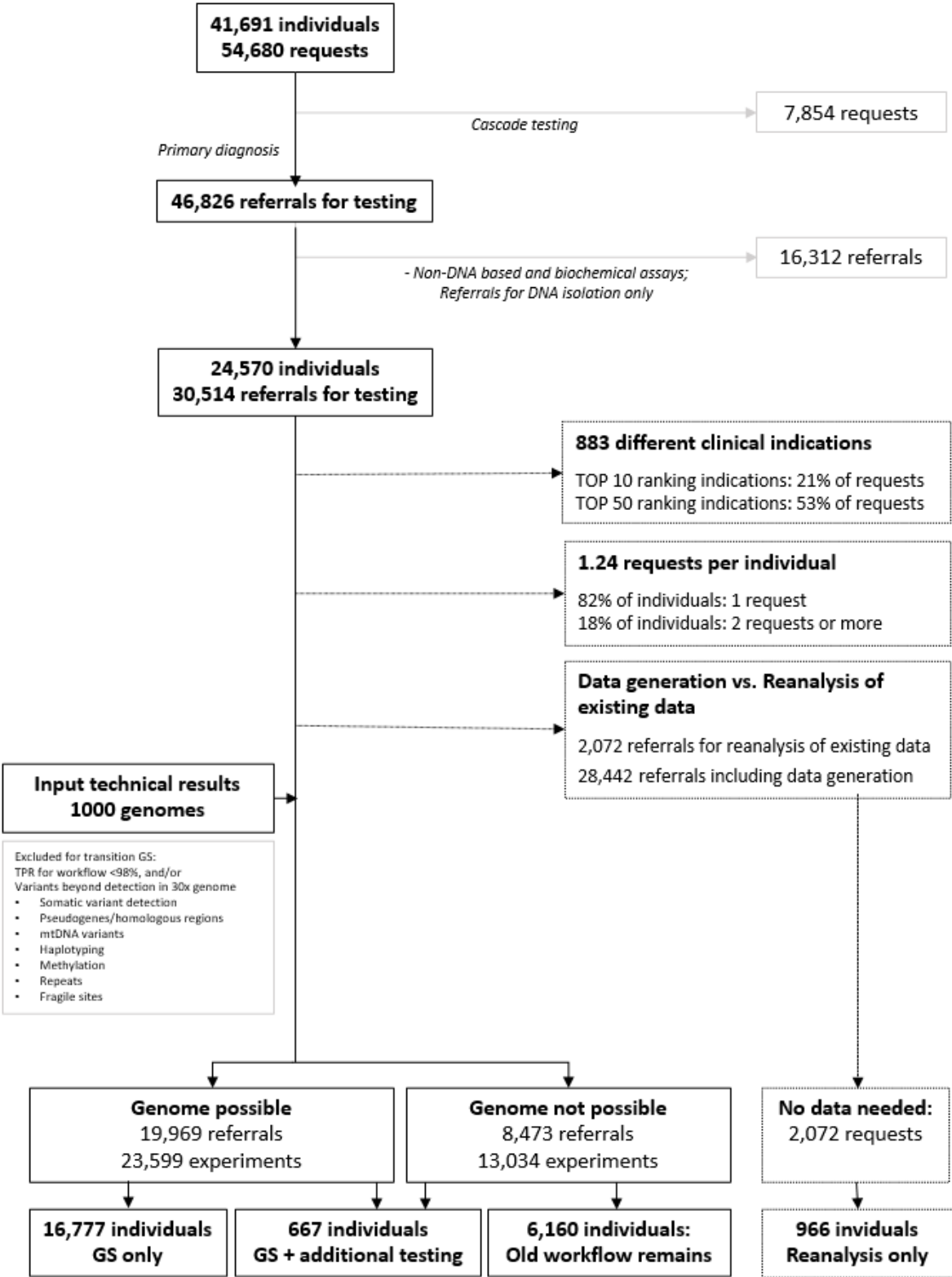
**Table S1** Online excel file providing overview of 1,000 individuals and workflows used

**Table S2** Online excel file providing 1,271 genetic variants in 1,000 individuals

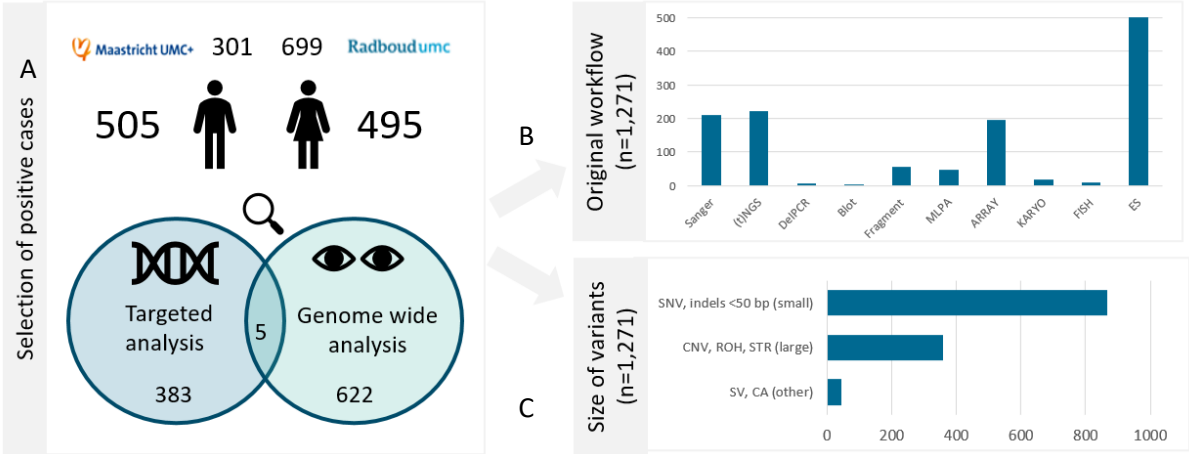
**Table S4** Online excel file providing coverage statistics for 58,393 variants for which previously (likely) pathogenic variants were described.

**Table S5** Online excel file providing coverage statistics for 4,266 disease-associated genes

Figure S1: Scenario model to determine individuals eligible for GS-first strategy



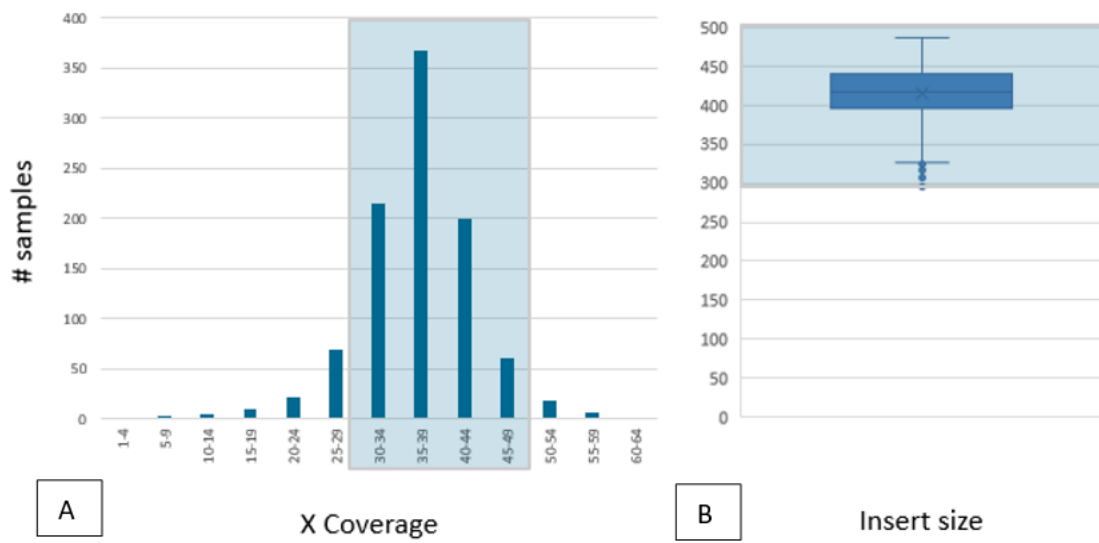
**Figure S2: A cohort of 1000 cases with clinically relevant variants spanning the broad range of genome diagnostics.**



**A** The 1000 genomes cohort consisted of 505 males and 495 females, who were genetically diagnosed in the Radboudumc or Maastricht UMC+ in 2018. The assays that were performed to find these diagnoses were either targeting specific variants and single (or a small set of) genes or complete gene panels or chromosomes were analyzed based on the patient’s phenotype. **B** In these cases, a total of 1,271 variants were identified, requiring >10 different workflows to diagnose them. **C** The variants were grouped in small (<50 bp), large (50 bp and up), and other variants (SVs and CA).

Abbreviations: targeted next generation sequencing ((t)NGS), deletion polymerase chain reaction (DeI/PCR), multiplex ligation-dependent probe amplification (MLPA), fluorescent in situ hybridisation (FISH), exome sequencing (ES), single nucleotide variants (SNV), short tandem repeat expansions (STRs), regions of homozygosity (ROH), copy number variants (CNV), structural variants (SV), chromosome anomalies (CA)

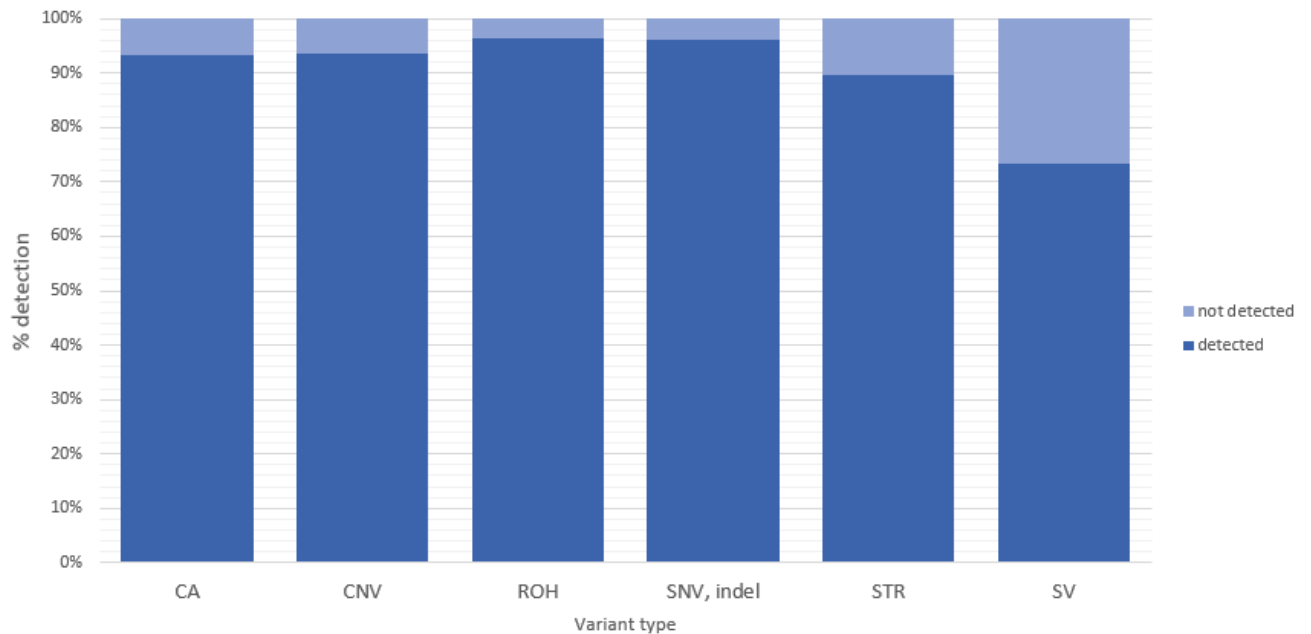
Figure S3: The average output of 1000 genomes.



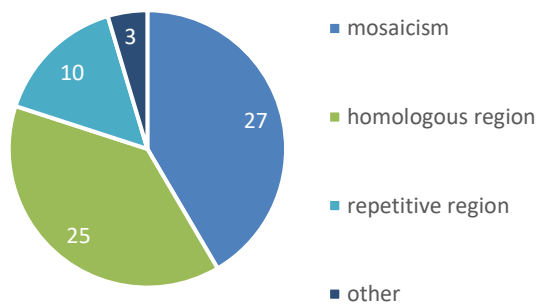
**A** As multiple observations per base are needed to come to a reliable base call, the recommended sequencing depth for genome sequencing is 30x to 50x. **B** Insert sizes are also important for the sequencing. For efficient sequencing, small insert sizes (risk of overlapping paired sequences) as well as larger fragments (decrease of cluster efficiency) must be avoided. We therefore aimed for a 300-500bp range for our 2x150bp paired-end sequencing. In this project we reached an average sequencing depth of 37x and an insert size of around 400-450bp.

**Figure S4: GS Technical validation by variant type and assessment of why variants were not identified**

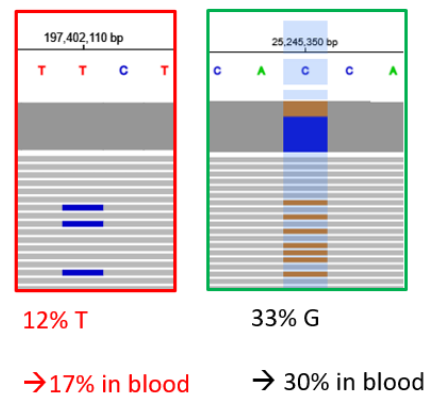
**A**



**B**



**C**



**A** In total, 94.9 % (1,206/1,271) of all variants were detected with GS. Small variants (<50bp) were detected in 96.1% (833/867), large variants (123 bp - 72.8Mb) in 93.3% (334/359), and other variants in 86.7% (39/45). The total list of variants and whether they were present in the GS data ('detected' vs. 'not detected') can be found in **Supplementary Table S2**. **B** In the 5% undetected variants (N=65), we identified common themes that are attributable to short-read 30x GS and downstream analysis. Undetected variants were mostly found in mosaic cases (n=27, 2.4-20%), homologous regions (n=25), i.e. pseudogenes or paralogues genes, or likewise in repetitive regions (n=10), i.e. repeats, telomeres or centromeres, and 3 others. **C** A mosaic variants in the *SF3B1* gene (Chr2(GRCh38):g.197402110T>C), which was originally detected with a targeted NGS approach in 17% of the blood sample, was present in 6/50 (12%) of the reads and not present in the VCF file of the GS data. A mosaic variant in the *KRAS* gene (Chr12(GRCh38):g.25245350C>G), originally detected with a targeted NGS approach in 30% of the blood sample, was present in 15/46 (33%) of the reads and in the VCF file of the GS data.

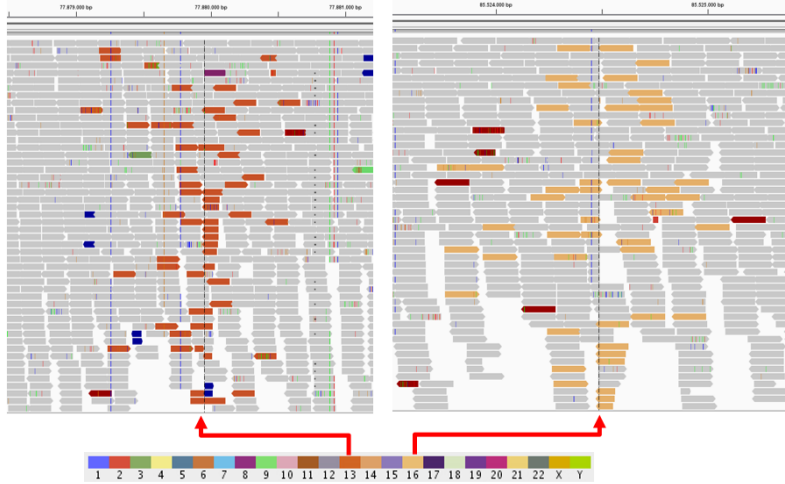
**Figure S5: Examples of comprehensive GS**

**A**

**P1-B3 46XX,t(13;16)(q22;q22)**

chr16:77,879,000-77,881,000

chr13:85,524,000-85,525,000

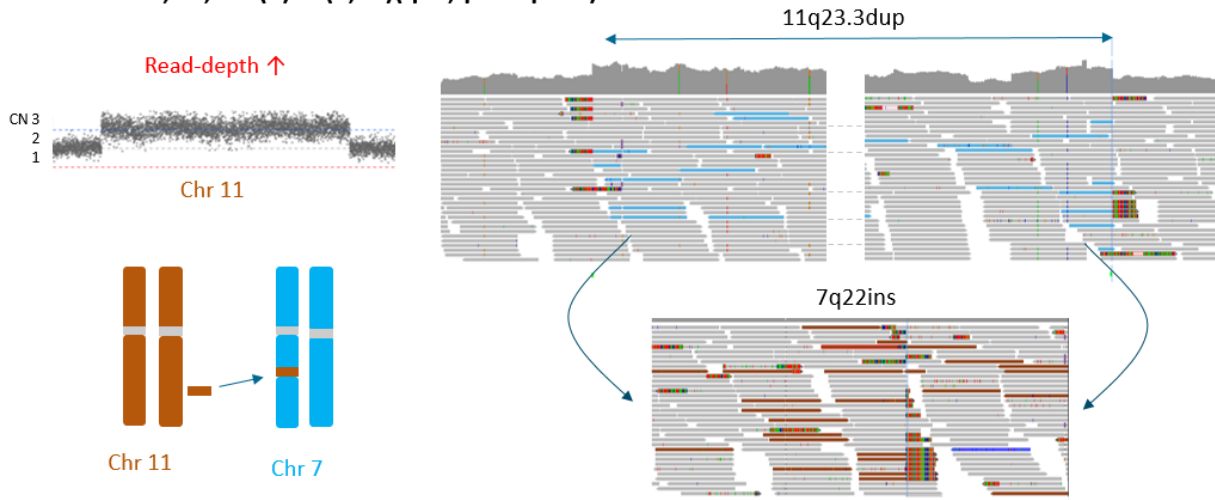


VCF

```
chr13 85524488[chr16:77879956]
chr16 77879956[chr13:85524488]
```

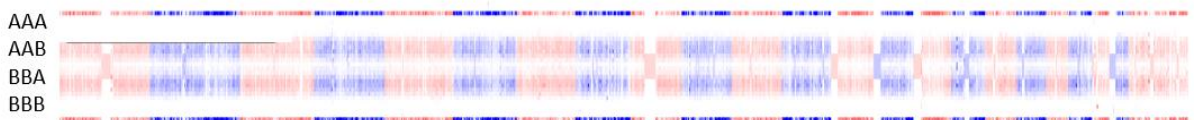
**B**

**P1-H6 46, XY, der(7)ins(7;11)(q22;q23.3q23.3)**



**C**

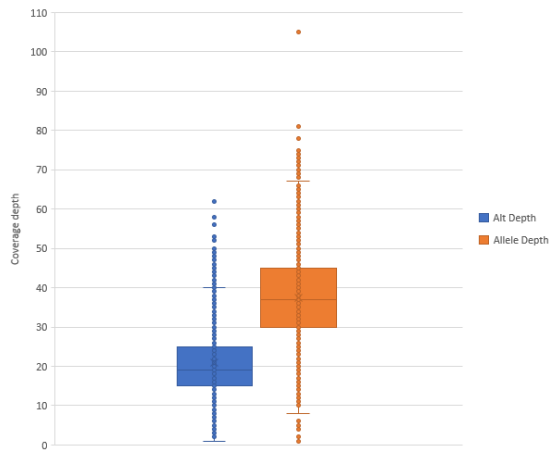
**P4-F4 69, XXX**



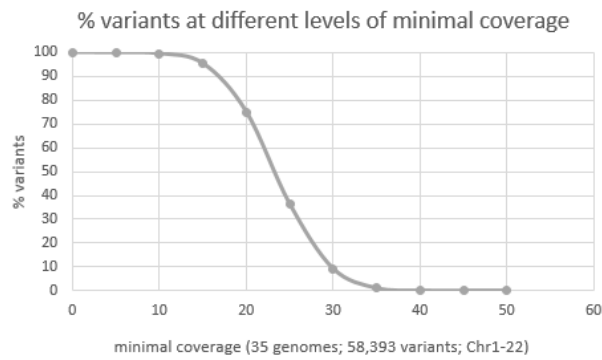
**A** Based on visual inspection and targeted manual search of the variant calling format (VCF) file, we could identify the previously detected translocation between chromosomes 13 and 16. **B** Likewise, we detected a copy-number gain on chromosome 11, which translocated to chromosome 7. In the diagnostic trajectory this derivative chromosome was detected with a targeted FISH analysis performed subsequent to an array analysis, in which only the gain was identified. **C** GS B-allele frequency plots can identify triploidies.

**Figure S6: *In silico* coverage statistics at variant level and disease genes**

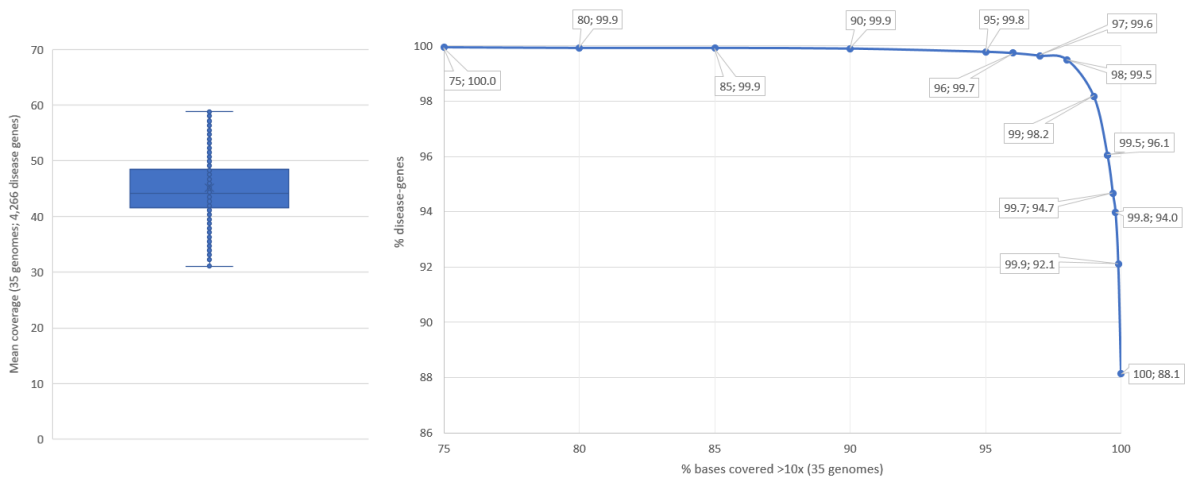
**A: Coverage statistics for 794 detected SNVs from the 1000 Genomes**



**B: Coverage statistics of 58,393 ClinVar and VKGL variants**

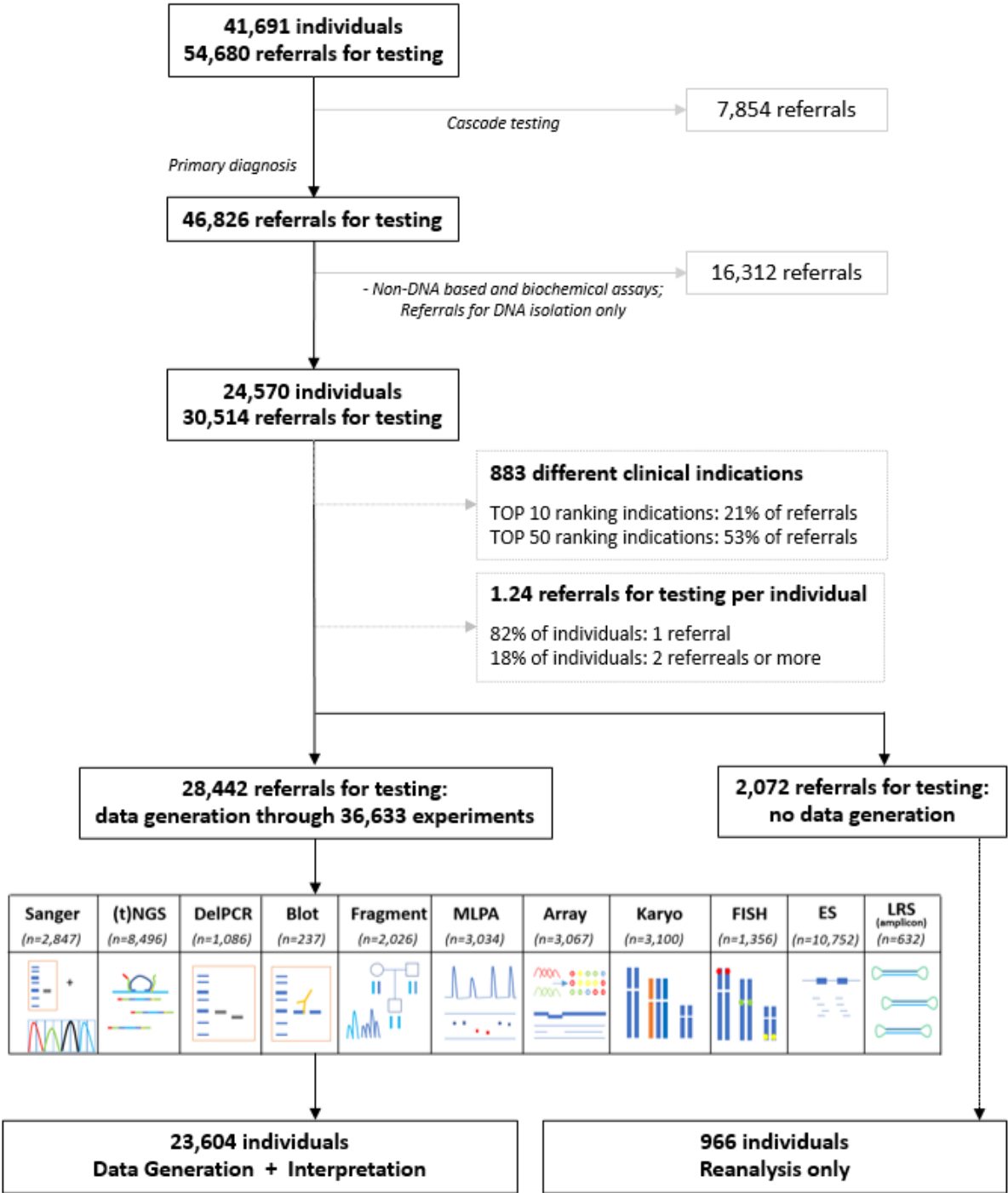


**C-D: Coverage statistics of 4,266 disease genes**



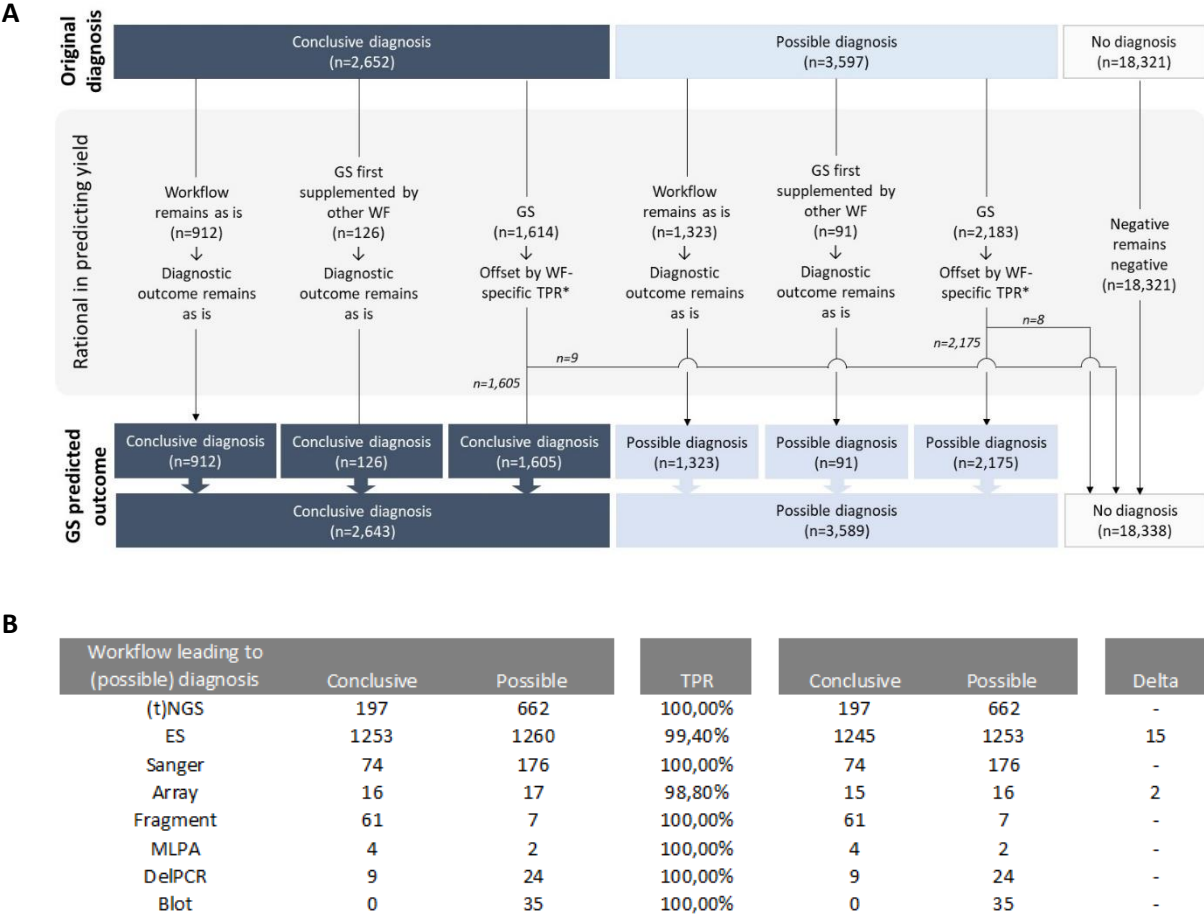
**A** Coverage data of 794 detected SNVs in our cohort, where allele depth ranged from 1-105, and (variant) alternative allele depth ranged from 1-62, with a 13-100% variant range. **B** Sequence depth at genomic positions that are known to harbor (likely) pathogenic variation and **C** Mean coverages for all coding positions of genes with well-established rare disease associations were calculated from 35 randomly selected genomes. **D** The fraction of genes versus the percentage of bases of the gene with  $\geq 10x$  coverage.

Figure S7: Schematic representation of referrals to Radboudumc and MUMC+ in 2022





**Figure S8: Schematic overview of assumptions made to evaluate the impact on diagnostic yield from transition to a generic GS approach**



**A)** Based on clinical referrals being transferable to generic GS, the impact on diagnosis was evaluated for all 24,570. Top row shows original diagnosis per individual, where 'n' refers to number of individuals; \*Offset with workflow specific TPRs are provided in **B**. Assuming all negative diagnoses remain negative, this translates to a possible false negative diagnostic rate of 0.3% (17/6232).

**Table S3: GS sensitivity**

\*Excluded indications: Adenomatous polyposis coli, Chronic lymphocytic leukemia, PTEN Hamartoma tumor syndrome (diagnostic referrals that are under suspicion of harboring mosaic variants and/or added to include mosaic variants although not primarily aimed at germline testing); Excluded variants: mosaic variants <20%, variants in the *CYP21A2*, *SMN1*, *OTOA*, *STRC* or *OPSIN* genes.

TPR $\geq$ 98% indicated by grey marking

# variants Workflow	1) Technical validation				2) Technical validation + exclusion expected false negatives*			
	positive	false negative	total	TPR	positive	false negative	total	TPR
Sanger	197	12	209	94.3%	178	0	178	100.0%
(t)NGS	209	13	222	94.1%	193	0	193	100.0%
DelPCR	6	1	7	85.7%	4	0	4	100.0%
Blot	2	0	2	100.0%	2	0	2	100.0%
Fragment	51	6	57	89.5%	51	6	57	89.5%
MLPA	44	4	48	91.7%	33	0	33	100.0%
Array	183	14	197	92.9%	168	2	170	98.8%
Karyo	15	4	19	78.9%	15	1	16	93.8%
FISH	8	2	10	80.0%	8	2	10	80.0%
ES	491	9	500	98.2%	486	3	489	99.4%
<b>Total</b>	<b>1206</b>	<b>65</b>	<b>1271</b>	<b>94.9%</b>	<b>1138</b>	<b>14</b>	<b>1152</b>	<b>98.8%</b>
<b>Type variant</b>								
SNV, indels	827	34	861	96.1%	789	3	792	99.6%
STR	52	6	58	89.7%	52	6	58	89.7%
ROH	26	1	27	96.3%	24	0	24	100.0%
CNV	262	18	280	93.6%	239	2	241	99.2%
CA	28	2	30	93.3%	23	0	23	100.0%
SV	11	4	15	73.3%	11	3	14	78.6%
<b>Total</b>	<b>1206</b>	<b>65</b>	<b>1271</b>		<b>1138</b>	<b>14</b>	<b>1152</b>	

Abbreviations: targeted next generation sequencing ((t)NGS), deletion polymerase chain reaction (DelPCR) multiplex ligation-dependent probe amplification (MLPA), fluorescence in situ hybridisation (FISH), exome sequencing (ES), single nucleotide variants (SNV), short tandem repeat expansions (STRs), regions of homozygosity (ROH), copy number variants (CNV), chromosome anomalies (CA), structural variants (SV)