

# Supplementary Materials

1		
2		
3	Title .....	2
4	Authors.....	2
5	Author Contributions.....	2
6	Supplementary Methods .....	3
7	CNV Detection.....	3
8	Alignment and definition of capture regions of interest .....	3
9	ClinCNV Workflow .....	3
10	Conifer Workflow.....	5
11	ExomeDepth Workflow.....	6
12	CNV Classes .....	7
13	CNV Visualisation.....	8
14	Diagnostic interpretation.....	8
15	ERN EURO-NMD.....	9
16	ERN GENTURIS.....	10
17	ERN ITHACA .....	10
18	ERN RND .....	11
19	References.....	12
20	Supplementary Figures .....	13
21	Supplementary Figure 1.....	13
22	Supplementary Figure 2.....	14
23	Supplementary Figure 3.....	15
24	Supplementary Figure 4.....	16
25	Supplementary Tables (see Excel file).....	16
26		

27 **Title**

28 *Comprehensive reanalysis for CNVs in ES data from unsolved rare disease cases*  
29 *results in new diagnoses*

30 **Authors**

31 German Demidov<sup>1,2\*‡</sup>, Burcu Yaldiz<sup>3,4\*</sup>, José Garcia-Pelaez<sup>5,6,7\*</sup>, Elke de Boer<sup>3,8,9\*</sup>, Nika  
32 Schuermans<sup>10\*</sup>, Liedewei Van de Vondel<sup>11,12\*</sup>, Ida Paramonov<sup>13</sup>, Lennart F. Johansson<sup>14</sup>,  
33 Francesco Musacchia<sup>15,16</sup>, Elisa Benetti<sup>17</sup>, Gemma Bullich<sup>13</sup>, Karolis Sablauskas<sup>3,18</sup>, Sergi  
34 Beltran<sup>13,19,20</sup>, Christian Gilissen<sup>3</sup>, Alexander Hoischen<sup>3,21,22</sup>, Stephan Ossowski<sup>1,2</sup>, Richarda  
35 de Voer<sup>3,23</sup>, Katja Lohmann<sup>24</sup>, Carla Oliveira<sup>5,6,7</sup>, Ana Topf<sup>25</sup>, Lisenka E.L.M. Vissers<sup>3,8</sup>, the  
36 Solve-RD Consortia, Steven Laurie<sup>13\*\*</sup>

37

38 **Author Contributions**

39 \* These authors contributed equally to this work

40 ‡ Corresponding authors

41 The authors declare no conflicts of interest.

42

43

44

## 45 Supplementary Methods

### 46 CNV Detection

#### 47 Alignment and definition of capture regions of interest

48 Sequencing data was submitted in BAM, CRAM, or FastQ format. Where data was submitted  
49 in BAM or CRAM format it was reconverted to FastQs at read-group level prior to being  
50 realigned to the hs37d5 human genome reference version, as used in phase 2 of the 1000  
51 genomes project<sup>1</sup> with BWA-MEM<sup>2</sup> (v0.7.8-r455). As GC-rich enrichment targets are known to  
52 amplify poorly, resulting in unreliable CNV calling<sup>3</sup>, the GC-content for each target in each  
53 enrichment kit was calculated and any targets in which the GC-content was >80% were  
54 removed from the corresponding target BED file prior to CNV calling. This resulted in the  
55 removal of <0.5% of target regions per kit. Ensembl version 75 was used for gene and  
56 transcript definition.

#### 57 ClinCNV Workflow

58 Analysis was performed separately for experiments generated by different exome enrichment  
59 kits. Initially, ClinCNV calculates the average read coverage of targeted regions of the  
60 enrichment kit divided into 120bp windows. As the first step of preprocessing, coverage is  
61 corrected for GC-content and library size for each sample individually. Following  
62 normalisation, systematically poorly covered regions (i.e. where 90% of samples had a  
63 normalised coverage <0.3) were excluded, followed by the application of variance stabilisation  
64 of read counts (square root transformation). To ameliorate the potential impact of batch effects  
65 on coverage calculation, samples were further clustered based on their global coverage  
66 profiles. In generating these clusters, target regions in the top and bottom quintiles for variance  
67 were excluded to minimise the potential impact of polymorphic regions on cluster generation,  
68 and coverage profiles smoothed using the rolling median. Uniform manifold approximation and

69 projection (UMAP)<sup>4</sup> was performed for the mapping of smoothed coverage profiles. Samples  
70 were clustered into subgroups with a minimum size of 15 using dbSCAN<sup>5</sup>. Finally, the coverage  
71 of each 120bp window was normalised using the median of coverages within the cluster.  
72 Different potential copy numbers are modelled using the theoretical expected value and  
73 estimated variance, and the log-likelihood of normalised coverage under different expected  
74 copy-number models is calculated for each window. Calling is performed analogously to  
75 Circular Binary Segmentation<sup>6</sup> using a Maximum Subarray Sum algorithm<sup>7</sup> *i.e.* the segment  
76 with the highest evidence supporting an alternative copy-number to that of the model is  
77 identified at each step of the segmentation, rather than the segment with the largest difference  
78 in mean.

79

80 Resulting CNV calls were filtered according to measures of within kit allele frequency of the  
81 CNV and the noisiness of the coverage at the CNV site, requiring a minimum log-likelihood  
82 ratio of 20 to be considered worthy of biological interpretation. A robust regression model is  
83 fitted, taking the 75% percentile rank of the per-chromosome number of CNVs as a response  
84 variable, and median read depth, enrichment kit, and predicted ancestry determined using  
85 SampleAncestry ([https://github.com/imgag/ngs-bits/blob/master/doc/tools/SampleAncestry](https://github.com/imgag/ngs-bits/blob/master/doc/tools/SampleAncestry/index.md)  
86 /index.md) as predictors. A sample was assessed as QC failed if the response variable was  
87 outwith the 99.5% prediction interval of the regression. The 75% percentile of the per-  
88 chromosome number of CNVs was chosen to overcome cases where long CNVs may have  
89 been segmented into many separate calls and thus an otherwise good sample could be falsely  
90 identified as QC failed if only the total number of CNV calls was used as a response. Where  
91 parents of a case were available (*i.e.* family trios), copy-number information from the parents  
92 was also provided to assist in interpretation, and to confirm if CNVs represented *de novo*  
93 events.

94

## 95 Conifer Workflow

96 Conifer<sup>8</sup> (<http://conifer.sourceforge.net/>) uses Singular Value Decomposition (SVD) to identify  
97 rare CNVs from exome sequencing data. Samples with similar read lengths were analysed in  
98 the same batch, and sex-specific sample pools were created for generating accurate X-  
99 Chromosome calls. RPKM (Reads Per Kilobase per Million mapped reads) values were  
100 calculated independently by enrichment kit for all corresponding targets. Following SVD to  
101 identify biases in coverage introduced by batch effects, 3 to 15 components were removed  
102 from each group based on manual inspection of the inflection points of scree plots generated  
103 by the program.

104

105 Within each analysis batch, if all experiments had less than 30 calls, the results were  
106 considered ready for further filtering. On the contrary where any experiment in a batch had  
107 more than 30 calls, then if the median number of calls per experiment in the batch was less  
108 than 10, any experiment with more than 30 calls was discarded as failing QC, and the results  
109 from the remaining experiments were considered ready for filtering. However, if the median  
110 number of calls within the batch was more than 10 per experiment then the SVD value was  
111 increased, and the batch analysis rerun, until either all experiments had less than 30 calls, or  
112 the median number of calls was less than 10, at which point any experiment with more than  
113 30 calls was discarded as described above. CNVs with an SVD-ZRPKM value greater than  
114 1.75 or less than -1.75 were considered as *bona fide* duplication or deletion calls, respectively,  
115 worthy of biological interpretation. Conifer does not provide any guidance as to the exact copy  
116 number identified at a particular locus and provides no further indicators of the quality of a  
117 detected event other than the SVD-RPKM metric.

118

## 119 ExomeDepth Workflow

120 ExomeDepth<sup>9</sup>, applies a beta-binomial model to the genome-wide distribution of read depth  
121 data, aiming to compare a test sample to a similar reference set selected by the tool. For the  
122 implementation of the ExomeDepth workflow, the generation of read count data was separated  
123 from that of identifying candidate CNVs. Thus for each experiment, read depth was initially  
124 calculated for all targets of the respective capture kit, and stored as a Bioconductor iRanges  
125 object<sup>10</sup>. In the second step, all iRanges objects from experiments generated using the same  
126 enrichment kit were analysed as a batch to generate raw CNV call sets. In this second step  
127 ExomeDepth automatically identifies an independent background reference set for each test  
128 sample by selecting the most closely correlated samples in terms of coverage from within the  
129 batch. Copy-Number prediction is provided by the ratio of observed/expected reads over a set  
130 of targets. We interpreted these ratios in diploid chromosomes as follows:

131

- 132 • O/E ratio <0.10 - Likely Homozygous Deletion i.e. Copy Number (CN) = 0
- 133 • 0.10 < O/E ratio <0.75 - Likely Heterozygous Deletion; CN=1
- 134 • 0.75 < O/E ratio <1.25 - Likely Copy Number neutral; CN=2 *i.e.* No CNV to report
- 135 • 1.25 < O/E ratio <1.75 - Likely Heterozygous Duplication; CN=3
- 136 • 1.75 < O/E ratio <2.25 - CN=4
- 137 • O/E ratio >2.25 - CN OTHER

138

139 ExomeDepth provides two indicators of quality. The first is a sample-level indicator of the  
140 correlation between the test sample and the background reference, which should be >0.97 for  
141 the results to be regarded as reliable. Secondly, regarding call quality, ExomeDepth provides  
142 a Bayes Factor (BF) based on the ratio of observed/expected reads over a set of apparently  
143 copy-number variant targets. Experiments with a correlation <0.97 were considered as failing  
144 QC, and any calls with a BF <0.15 were discarded as being unreliable.

## 145 CNV Classes

146 To aid downstream interpretation, each CNV call was categorised into one of six classes.

147

148 1) Putative CNVs longer than 500kb in length were initially identified regardless of the  
149 presence of absence of genes of interest in the ERN gene lists. The recent release of  
150 large CNVs catalogues such as DECIPHER , as well as the presence of a large number  
151 of case reports with chromosomal changes of this size and larger, allowed us to  
152 hypothesise that such variants could be interpreted successfully, even if the reported  
153 phenotypes of the patients exhibiting such variants may differ from the phenotypes  
154 expected for affected genes.

155 2) Homozygous deletions are generally rare, and the presence of a homozygous deletion  
156 needs to be interpreted very cautiously due to potentially incorrect enrichment kit  
157 reporting, or poor-quality library preparation. An important indicator that a putative  
158 homozygous deletion call is likely to be *bona fide* is the consanguinity status of the  
159 patient.

160 3) Heterozygous CNVs occurring in genes with a described autosomal-dominant mode  
161 of inheritance reported in OMIM.

162 4) Duplications with apparent copy number >3. These may represent cases where alleles  
163 on both chromosomes are duplicated, or cases where only the allele on one  
164 chromosome has been duplicated but multiple times.

165 5) Gonosomal CNVs. As gonosomal CNVs require a mixed workflow depending on the  
166 sex of the participant, a separate set of calls was generated for CNV calls on  
167 chromosomes X and Y. In the case of the Y-Chromosome, only “Long” CNVs that  
168 would fall into category 1 above were reported for interpretation, since there were no  
169 genes of interest on the Y-Chromosome on any of the ERN gene lists.

170 6) Potential compound heterozygote SNV/CNV “double-hits”. For a short list of  
171 experiments in which a single candidate SNV had been identified by the Solve-RD

172 SNV working group, which was either listed in ClinVar as Pathogenic/Likely Pathogenic  
173 or predicted to have a high impact in a gene of interest, affecting an individual where  
174 the mode of inheritance was suspected to be recessive, (see Laurie et al, 2023, under  
175 review) we investigated whether a potentially pathogenic CNV affecting the second  
176 allele of the same gene could explain the case as a compound heterozygote.

## 177 CNV Visualisation

178 To provide support for interpretation of the technical validity of CNV calls, screenshots for  
179 regions containing CNV calls were generated automatically using the Integrative Genomics  
180 Viewer<sup>11</sup> (IGV), incorporating a variety of custom-built tracks (see **Supplementary Figure 3**).  
181 These included call tracks for each of the three callers in SEG format, normalised coverage  
182 tracks for ClinCNV and Conifer, beta-allele frequency, BAM DoC, Institute of Medical Genetics  
183 and Applied Genomics (Tübingen) in-house polymorphic CNV regions, and gene tracks from  
184 RefSeq genes, ERN candidate genes, and DECIPHER microdeletion and duplication  
185 syndromes<sup>12</sup>.

186  
187 For each CNV returned for interpretation, we generated IGV screenshots of both the whole  
188 sample (chr1-22 and chrX/Y) to allow evaluation of overall sample quality, and the region  
189 around the individual CNV (+/-10kb). Specifically in the case of long CNVs, the observation of  
190 clear deviations from the expected ratio of 50/50 in beta-allele frequencies provided strong  
191 additional support of variant validity. For rare cases in which a signal of unusual read pairing  
192 was observed, suggesting that a breakpoint may have been captured, a screenshot was  
193 generated including the suspected breakpoint.

## 194 Diagnostic interpretation

195 To facilitate diagnostic interpretation, we used AnnotSV<sup>13</sup> (version 3.0.7) to add a range of  
196 useful annotations to the reports returned to Clinical experts for interpretation as listed in



197 **Supplementary Table 2.** Diagnostic interpretation was undertaken by expert clinicians and  
198 clinical researchers from the respective ERNs. Each ERN prioritised the calls for further  
199 investigation according to their own strategy, based on their expert knowledge of underlying  
200 disease mechanisms in their respective patients. Some annotations, such as that of the  
201 ENCODE blacklist for high-signal regions were used to quickly discard overlapping CNVs by  
202 all ERNs, whereas other information, such as evidence of consanguinity, provided further  
203 support that homozygous deletions were likely to be relevant in affected cases. For the  
204 interpretation of heterozygous deletions, pLI scores from GnomAD<sup>14</sup>, and haploinsufficiency  
205 gene lists from the DDD project<sup>15</sup>, aided interpretation. The full workflow is illustrated in  
206 **Supplementary Figure 4.**

207

## 208 ERN EURO-NMD

209 The filtering strategy undertaken by EURO-NMD was determined per analysis (see section ‘CNV  
210 filtering’). In general, a balance had to be upheld whereby submitting clinical researchers would interpret  
211 as many CNVs as possible while maintaining a feasible interpretation load. Thus the following analyses  
212 were shared directly given the relative number of CNVs to be analysed: homozygous deletions, high  
213 copy number duplications, gonosomal CNVs, and potential compound heterozygote second hits,  
214 whereas heterozygous CNVs were split between CNVs of copy number one (CN1, i.e. deletions) and  
215 those of copy number three (CN3 i.e. duplications).

216

217 For CN1, CNVs for genes with DDD Haploinsufficiency scores > 90 or a GnomAD pLi < 0.1 were  
218 discarded, as these indicate that the gene is likely tolerant of heterozygous deletions. For both CN1 and  
219 CN3, CNVs identified through ClinCNV with a loglikelihood < 30 were discarded, as these are likely  
220 false positives. CNVs identified in genes only known to have recessive inheritance patterns were  
221 discarded, as were CNVs reported in Conrad *et al*<sup>16</sup>. For long CNVs, CNVs found in the Encode blacklist  
222 were discarded. Following these filtering steps, experts from the submitting groups applied a phenotype-  
223 first approach. If the phenotype could potentially match with the gene affected by the CNV call, IGV  
224 tracks were checked to evaluate the likelihood of the called CNV being a true CNV.

225

## 226 ERN GENTURIS

227 Due to the small size of the ERN GENTURIS cohort, and the short gene list, only limited further  
228 filtering of calls was necessary. No additional filters were applied to call sets from Conifer. In  
229 the case of heterozygote deletions and duplications, specific filtering criteria were applied  
230 separately for ClinCNV and ExomeDepth. For ClinCNV, we first interpreted all events  
231 identified by more than one tool, independent of the ClinCNV loglikelihood value. After this,  
232 we proceeded to analyse all events called only by ClinCNV with a loglikelihood of at least 20.  
233 For ExomeDepth, we first interpreted all events called by more than one tool, independently  
234 of the Bayes Factor (BF), and subsequently considered events called only by ExomeDepth  
235 with a BF of at least 15. For long CNVs, we first discarded all those events found in the encode  
236 blacklist and analysed the rest. For all datasets, following IGV visualization, only CNVs  
237 observed to be rare in control populations were considered for further interpretation.

238

## 239 ERN ITHACA

240 For ERN ITHACA, as a first step, we discarded variants that were annotated to have low QC,  
241 had been previously annotated as benign, or occurred in regions on the Encode Blacklist, as  
242 provided by the AnnotSV annotation. Additionally, to reduce the proportion of false positives,  
243 we discarded deletions shorter than 10kb and duplications shorter than 20kb in length, with  
244 the exception of homozygous deletion calls and variants in parent-offspring trios identified as  
245 being *de novo* by ClinCNV. Following this, visual inspection of each of the remaining CNV  
246 calls in IGV images was undertaken to assess technical validity, using reads and coverage  
247 supporting the call and B-allele frequency. Based on this visual assessment, apparently real  
248 biological CNVs were defined. For detailed clinical interpretation, prioritisation was  
249 subsequently guided by genes present on the ERN ITHACA gene list with a disease-  
250 association validity score  $\geq 3$ , see Laurie et al, 2023 (under review), consistent with the

251 expected mode of inheritance. Of note, CNVs  $\geq 200$  kb were also investigated regardless of  
252 the presence or absence of a gene on the ERN ITHACA gene list, given the prior knowledge  
253 of large CNVs being involved in ITHACA-associated phenotypes. All CNVs passing the above  
254 criteria were returned to the submitting groups from DITF-ITHACA, for diagnostic interpretation  
255 based on the clinical relevance to the phenotype observed in the affected individual.

## 256 ERN RND

257 The filtering strategy of ERN RND was predominantly based on tool-specific metrics. In  
258 general, the goal was to exclude calls with a high likelihood of being false positives. For  
259 ClinCNV we discarded all calls with a loglikelihood  $< 30$  and first prioritised calls with a  
260 loglikelihood  $> 200$ . As Conifer provides no metrics for filtering, all Conifer calls were analysed.  
261 For ExomeDepth, we discarded all calls affecting less than three targets and those with a  
262 Bayes factor  $< 30$ , unless there was an overlapping CNV identified by one of the other tools.  
263 Following these filtering steps, the clinical researchers who submitted the case applied a  
264 phenotype-first approach. If the phenotype could potentially match that of the called CNV, IGV  
265 tracks were checked visually to evaluate the likelihood that the called CNV was *bona fide*.

266

267

268

## 269 References

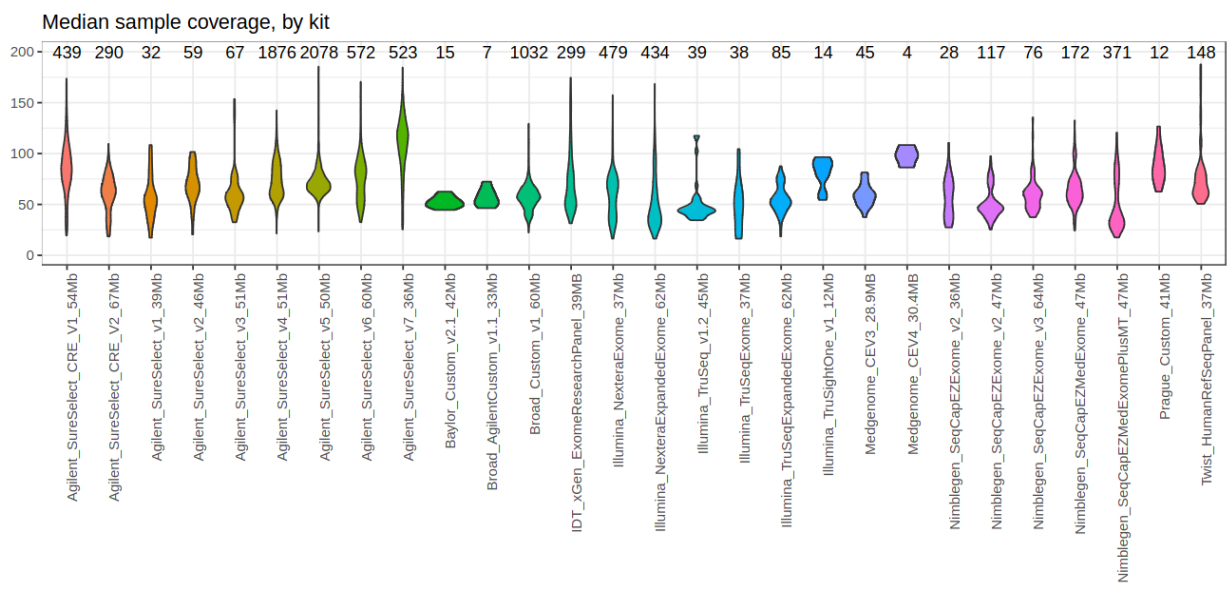
- 270 1. Auton, A. *et al.* A global reference for human genetic variation. *Nature* vol. 526 68–74  
271 Preprint at <https://doi.org/10.1038/nature15393> (2015).
- 272 2. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-  
273 MEM. *arXiv preprint arXiv* **00**, 3 (2013).
- 274 3. Parrish, A. *et al.* An enhanced method for targeted next generation sequencing copy  
275 number variant detection using ExomeDepth. *Wellcome Open Research* **2017 2:49 2**,  
276 49 (2017).
- 277 4. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and  
278 Projection for Dimension Reduction. (2018).
- 279 5. Hahsler, M., Piekenbrock, M. & Doran, D. dbscan: Fast Density-Based Clustering with  
280 R. *J Stat Softw* **91**, 1–30 (2019).
- 281 6. Olshen, A. B., Venkatraman, E. S., Lucito, R. & Wigler, M. Circular binary segmentation  
282 for the analysis of array-based DNA copy number data. *Biostatistics* **5**, 557–572 (2004).
- 283 7. Bentley, J. L. & Chan, P. *Programming Pearls*. (Addison-Wesley, 1989).
- 284 8. Krumm, N. *et al.* Copy number variation detection and genotyping from exome  
285 sequence data. *Genome Res* **22**, 1525–1532 (2012).
- 286 9. Plagnol, V. *et al.* A robust model for read count data in exome sequencing experiments  
287 and implications for copy number variant calling. *Bioinformatics* **28**, 2747–2754 (2012).
- 288 10. Lawrence, M. *et al.* Software for Computing and Annotating Genomic Ranges. *PLoS*  
289 *Comput Biol* **9**, e1003118 (2013).
- 290 11. Robinson, J. T., Thorvaldsdóttir, H., Wenger, A. M., Zehir, A. & Mesirov, J. P. Variant  
291 review with the integrative genomics viewer. *Cancer Res* **77**, e31–e34 (2017).
- 292 12. Firth, H. V. *et al.* DECIPHER: Database of Chromosomal Imbalance and Phenotype in  
293 Humans Using Ensembl Resources. *Am J Hum Genet* **84**, 524–533 (2009).
- 294 13. Geoffroy, V. *et al.* AnnotSV: an integrated tool for structural variations annotation.  
295 *Bioinformatics* **34**, 3572–3574 (2018).
- 296 14. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in  
297 141,456 humans. *Nature* **581**, (2020).
- 298 15. Huang, N., Lee, I., Marcotte, E. M. & Hurles, M. E. Characterising and predicting  
299 haploinsufficiency in the human genome. *PLoS Genet* **6**, e1001154–e1001154 (2010).
- 300 16. Conrad, D. F. *et al.* Origins and functional impact of copy number variation in the human  
301 genome. *Nature* **464**, 704–12 (2010).

302

303

304 **Supplementary Figures**

305 **Supplementary Figure 1**



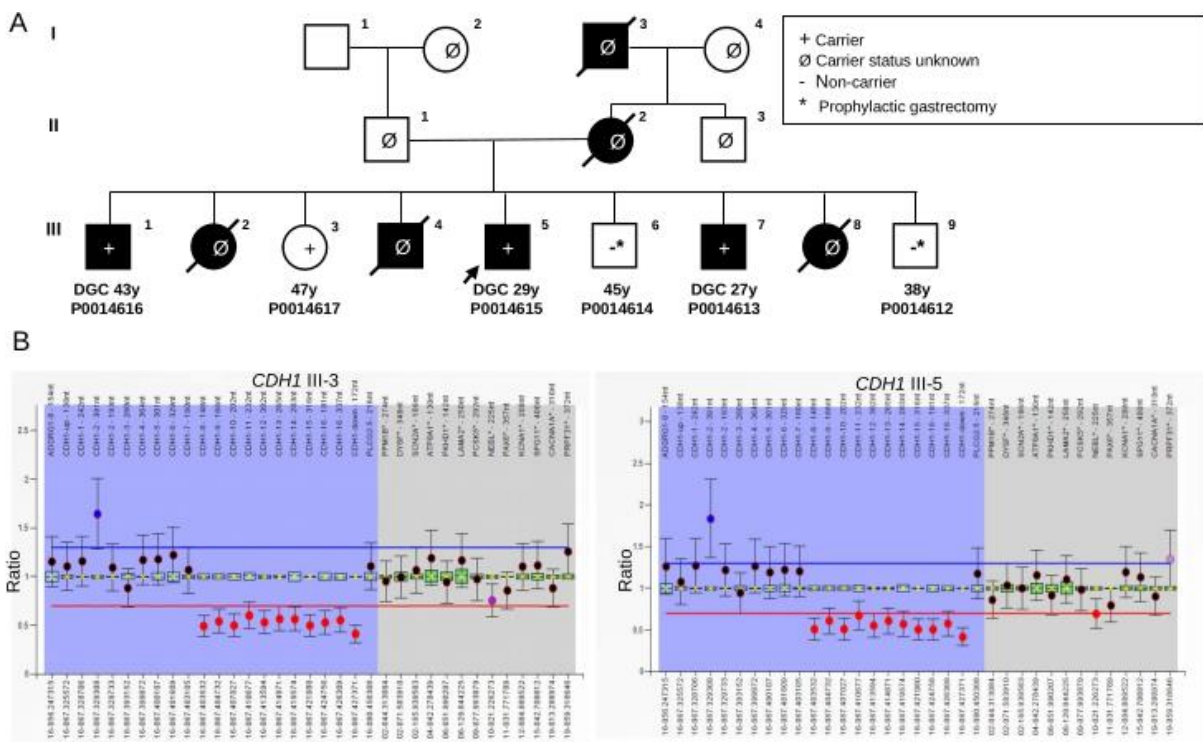
306

307 **Supplementary Figure 1.** Violin plot of median depth of coverage by kit for 9,351 ES  
 308 experiments pertaining to 28 different enrichment kits. The number of experiments pertaining  
 309 to each kit is shown above the plots. Coverage is shown on the Y-axis. Thickness of the plotted  
 310 shape indicates the proportion of experiments which have a particular coverage.

311

312 Supplementary Figure 2

313



314

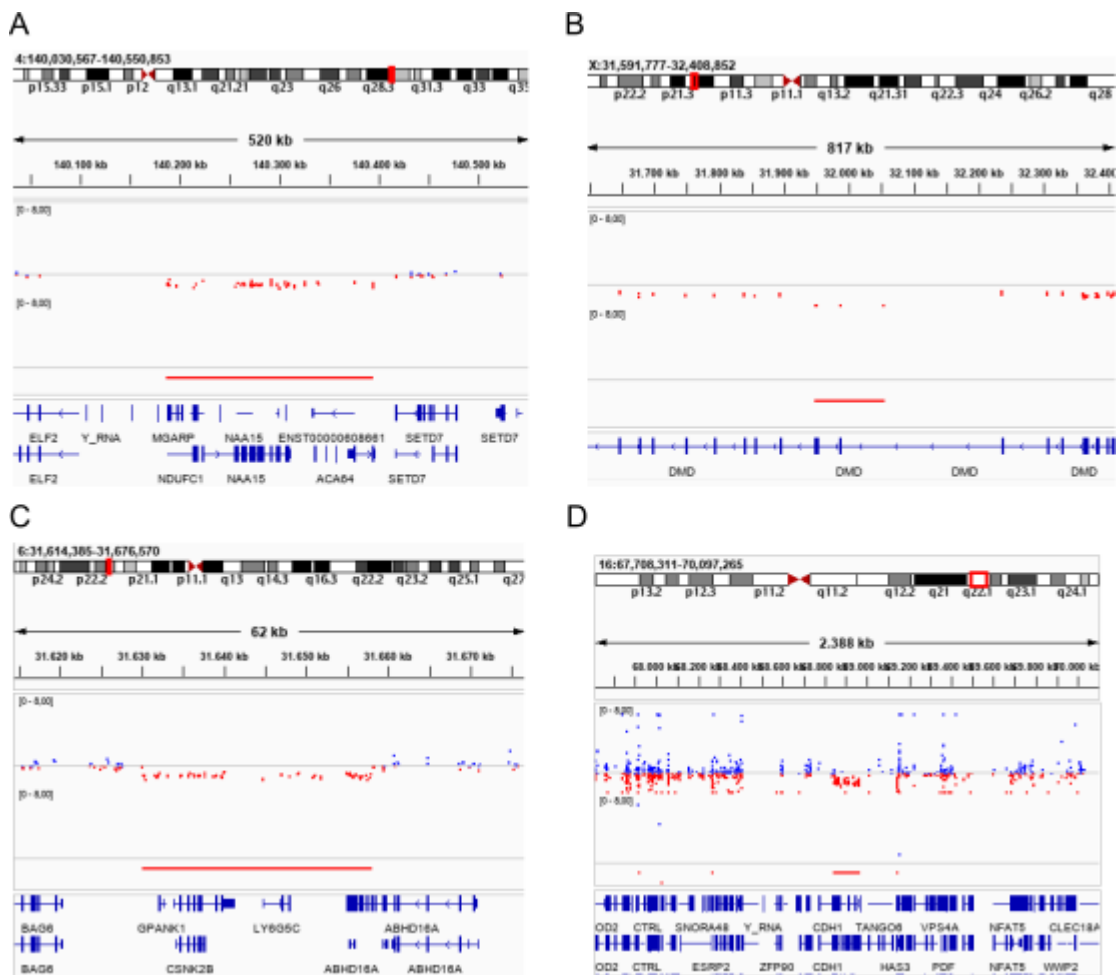
315 **Supplementary Figure 2.** Family pedigree and MLPA confirmation results for a Mexican  
 316 family extensively affected by Hereditary Gastric Cancer. **A)** Family tree of the family of  
 317 proband P0014615 (represented by an arrow). Exome Sequencing data from six individuals  
 318 of the family was submitted to Solve-RD for re-analysis, following prior analysis in 2015 for  
 319 both SNVs and CNVs which retrieved a negative result. Three of the sequenced family  
 320 members were affected by diffuse gastric cancer (DGC, black symbols: P0014616, P0014615,  
 321 P0014613) while the other three were unaffected (P0014617, P0014614, P0014612).  
 322 Individual III-3 (P0014617) is currently a healthy carrier, perhaps due to the incomplete  
 323 penetrance reported for *CHD1*. The age shown below affected individuals indicates age of  
 324 disease onset, while that below healthy individuals represents their current age. **B)** MLPA  
 325 validation results using SALSA MLPA-Probemix P083 *CDH1* (MRC Holland) in the healthy-  
 326 carrier III-3, and in the proband, III-5. A ratio above the blue line indicates elevated number of  
 327 copies, while a ratio below the red line indicates a decrease in copy number. The shaded blue  
 328 area represents position of probes for *CDH1* and two neighbouring genes while the grey area  
 329 represents reference probes.

330

331

332 Supplementary Figure 3

333

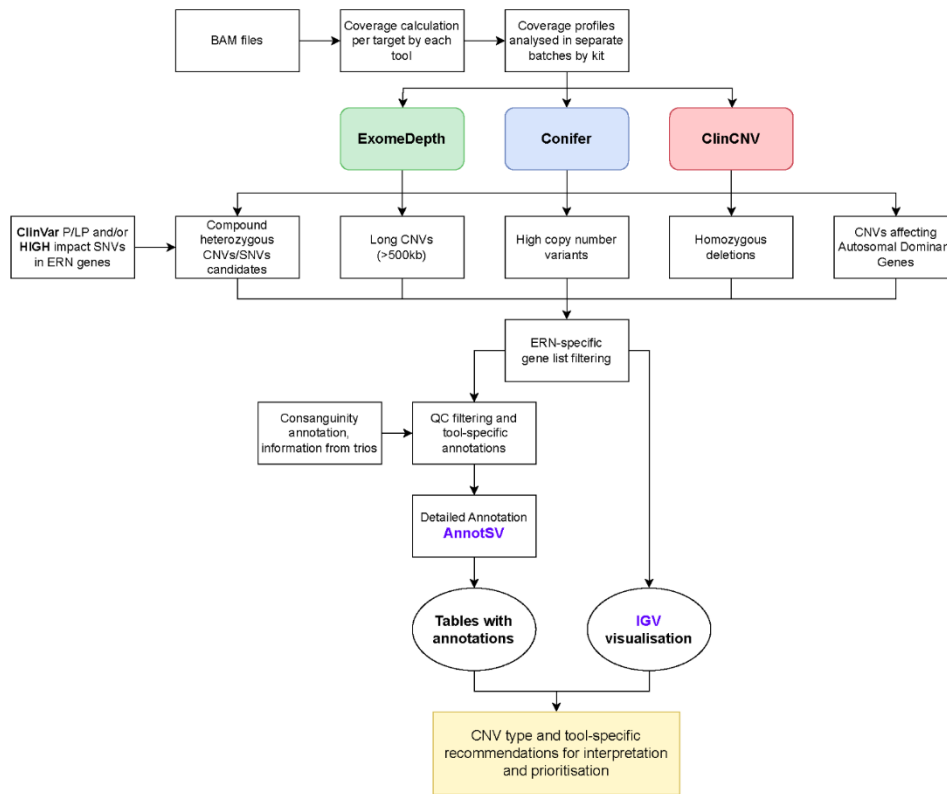


334

335 **Supplementary Figure 3.** IGV screenshots corresponding to the four illustrative newly  
 336 diagnosed individuals described in the main text, one from each ERN. **A)** RND: Heterozygous  
 337 deletion spanning *NAA15*, in an individual with intellectual disability, which was found to be  
 338 inherited from her paucisymptomatic mother. **B)** EURO-NMD: Hemizygous deletion of exons  
 339 45-47 of *DMD* resulting in Becker Muscular Dystrophy. **C)** ITHACA: Heterozygous *de novo*  
 340 deletion spanning *CSNK2B*, resulting in POBINDS **D)** GENTURIS: Inherited heterozygous  
 341 deletion affecting *CDH1* and *TANGO6*, resulting in autosomal dominant HDGC. Images show  
 342 customised coverage tracks and the position of the identified CNV (red bar). Blue dots above  
 343 the midline indicate elevated coverage, while red dots below the line indicate reduced  
 344 coverage. The position of genes is indicated at the bottom of the image, while the  
 345 chromosomal position is indicated at the top of the image.

346

347 Supplementary Figure 4



348

349 **Supplementary Figure 4.** Workflow used for CNV calling, filtering, and annotation, prior to  
 350 returning to calls to clinical experts for interpretation.

351

352 **Supplementary Tables (see Excel file)**

353 **Supplementary Table 1.** Curated ERN gene lists.

354 **Supplementary Table 2.** Annotations added to CNV calls to aid variant interpretation.

355 **Supplementary Table 3.** Number of families and affected individuals analysed per ERN, and  
 356 number of families and affected individuals with at least one CNV requiring interpretation.

357 **Supplementary Table 4.** Summary statistics regarding 7,849 CNVs initially returned for  
 358 interpretation.

359 **Supplementary Table 5.** Summary statistics regarding the length of 3,487 duplications and  
 360 4,362 deletions returned for interpretation.