



Validation of copy number variation analysis for next-generation sequencing diagnostics

DOI:

[10.1038/ejhg.2017.42](https://doi.org/10.1038/ejhg.2017.42)

Document Version

Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Ellingford, J., Campbell, C., Barton, S., Bhaskar, S., Gupta, S., Taylor, R., Sergouniotis, P. I., Horn, B., Lamb, J., Michaelides, M., Webster, A. R., Newman, W., Panda, B., Ramsden, S., & Black, G. (2017). Validation of copy number variation analysis for next-generation sequencing diagnostics. *European Journal of Human Genetics*, 25, 719-724. <https://doi.org/10.1038/ejhg.2017.42>

Published in:

European Journal of Human Genetics

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



Research Article

Validation of copy number variation analysis for next-generation sequencing diagnostics

Jamie M Ellingford^{1,2},
Christopher Campbell¹,
Stephanie Barton¹,
Sanjeev Bhaskar¹,
Saurabh Gupta³,
Rachel L Taylor^{1,2},
Panagiotis I Sergouniotis¹,
Bradley Horn¹,
Janine A Lamb⁴,
Michel Michaelides^{5,6},
Andrew R Webster^{5,6},
William G Newman^{1,2},
Binay Panda³,
Simon C Ramsden¹,
Graeme CM Black^{1,2}

¹Manchester Centre for Genomic Medicine, Central Manchester University Hospitals NHS Foundation Trust, Manchester Academic Health Sciences Centre, St Mary's Hospital, Manchester M13 9WL, UK.

²Institute of Human Development, University of Manchester, Oxford Road, Manchester, M13 9WL, UK.

³Ganit Labs, Bio-IT Centre, Institute of Bioinformatics and Applied Biotechnology, Bangalore, 560100, India.

⁴Institute of Population Health, University of Manchester, Oxford Road, Manchester, M13 9PT, UK.

⁵Moorfields Eye Hospital NHS Foundation Trust, London, EC1V 2PD, UK

⁶UCL Institute of Ophthalmology, Department of Genetics, London, EC1V 9EL, UK.

Corresponding Author: Jamie M Ellingford,
Jamie.ellingford@postgrad.manchester.ac.uk, 0161 276 8703.

Conflict of Interest: The authors declare no conflict of interest.

Running Title: CNV detection in targeted NGS diagnostics

1 **Abstract**

2 Although a common cause of disease, copy number variants (CNVs) have not
3 routinely been identified from next-generation sequencing (NGS) data in a
4 clinical context. This study aimed to examine the sensitivity, and specificity of a
5 widely used software package, ExomeDepth, to identify CNVs from targeted NGS
6 datasets. We benchmarked the accuracy of CNV detection using ExomeDepth
7 v1.1.16 applied to targeted NGS datasets, through comparison to CNV events
8 detected through whole genome sequencing (WGS) for 25 individuals, and
9 determined the sensitivity and specificity of ExomeDepth applied to these
10 targeted NGS datasets to be 100% and 99.8%, respectively. To define quality
11 assurance metrics for CNV surveillance through ExomeDepth, we undertook
12 simulation of single exon ($n=1000$) and multiple-exon heterozygous deletion
13 events ($n=1749$), determining a sensitivity of 97% ($n=2749$). We identified that
14 the extent of sequencing coverage, the inter- and intra-sample variability in the
15 depth of sequencing coverage, and the composition of analysis regions are all
16 important determinants of successful CNV surveillance through ExomeDepth.
17 We then applied these quality assurance metrics during CNV surveillance for 140
18 individuals across 12 distinct clinical areas, encompassing over 500 potential
19 rare disease diagnoses. All 140 individuals lacked molecular diagnoses after
20 routine clinical NGS testing, and through application of ExomeDepth we
21 identified 17 CNVs contributing to the cause of a Mendelian disorder. Our
22 findings support the integration of CNV detection using ExomeDepth v1.1.16
23 with routine targeted NGS diagnostic services for Mendelian disorders.
24 Implementation of this strategy increases diagnostic yields and enhances clinical
25 care.

- 26 **Key words:** Copy Number Variation, Targeted Next-Generation Sequencing,
27 Mendelian disorders, Medical Genetics

28 **Introduction**

29 Molecular diagnostic services available for patients with genetically
30 heterogeneous Mendelian disease have been transformed by the adoption of
31 next-generation DNA sequencing (*NGS*) within the clinical setting.^{1,2} At present,
32 diagnostic services facilitated by *NGS* are frequently limited to targeted capture
33 techniques, including custom gene panels^{3, 4} and whole exome sequencing
34 (*WES*).^{5, 6} These techniques have demonstrated tremendous power to identify
35 rare and private single nucleotide variation and small insertions/deletions
36 underpinning disease onset.

37 The identification of large structural variants and copy number variants (*CNVs*)
38 encapsulating the regions targeted by *WES* and custom gene panel assays have
39 proved challenging in a clinical context. While whole genome sequencing (*WGS*)
40 techniques have the potential to address this gap in diagnostic *NGS* services,^{7, 8}
41 the cost and data burdens remain substantial. Consequently, the application of
42 *CNV* detection algorithms in targeted *NGS* diagnostic services can facilitate
43 immediate improvement in clinical care for individuals with heterogeneous
44 Mendelian disorders. However, such techniques require formal assessment to
45 demonstrate accuracy, reliability and repeatability.

46 Here, we assess a framework for the implementation of *CNV* detection with
47 targeted *NGS* diagnostic services applied across a range of highly heterogeneous
48 Mendelian disorders.

49 **Methods**

50 **Study Design**

51 High coverage targeted NGS data was generated in a United Kingdom Accredited
52 Clinical Laboratory. We applied a CNV detection algorithm to validate the
53 sensitivity for (i) known CNV events, and (ii) simulated CNV events (Figure 1).

54 We assessed a number of factors to determine whether they influenced
55 successful CNV surveillance. We selected two key factors identified from
56 assessments of simulated and known CNVs (inter-sample variability and
57 insufficient coverage) as quality assurance metrics during prospective CNV
58 detection for individuals without molecular diagnoses through clinical NGS
59 testing (Figure 1).

60 Our analyses included individuals referred for diagnostic testing for four highly
61 heterogeneous disorders where targeted gene panel NGS is a routine diagnostic
62 service, specifically: inherited retinal dystrophies (IRD), congenital cataracts,
63 cardiac disorders and metabolic disorders.

64 **Sequencing & Variant Analysis**

65 **Whole Genome and Targeted Next-Generation Sequencing**

66 WGS data was generated for 25 individuals by Complete Genomics (Mountain
67 View, CA, USA) using a mate-paired sequencing technique, as described
68 previously.⁹ Read alignment and variant calling was performed using version 2.5
69 of the Complete Genomics pipeline.¹⁰

70 For targeted NGS, enrichments were performed on DNA extracted from
71 peripheral blood using Agilent SureSelect Custom Design target-enrichment kits

72 (Agilent, Santa Clara, CA, USA). Enrichment kits were designed to capture known
73 pathogenic intronic variants and the protein-coding regions +/-50 nucleotides of
74 selected NCBI RefSeq transcripts; conditions tested included IRD (105 genes or
75 180 genes), congenital cataracts (114 genes), cardiac disorders (72 genes
76 comprised of 10 sub-panels) and metabolic disorders (226 genes comprised of 6
77 sub-panels). The genes and transcripts included in the targeted capture regions
78 for each disease referral are available online (Supp Tables S1-S4) and through
79 the UK Genetic Testing Network ([http://ukgtn.nhs.uk/find-a-test/search-by-](http://ukgtn.nhs.uk/find-a-test/search-by-laboratory/laboratory/manchester-rgc-36/)
80 [laboratory/laboratory/manchester-rgc-36/](http://ukgtn.nhs.uk/find-a-test/search-by-laboratory/laboratory/manchester-rgc-36/)). Samples were pooled and paired-
81 end NGS was performed using the manufacturer protocols for the Illumina HiSeq
82 2000/2500 platform (Illumina, Inc., San Diego, CA, USA). Sequencing reads were
83 demultiplexed with CASAVA v.1.8.2. and aligned to the *hg19* reference genome
84 using Burrows-Wheeler Aligner short read (BWA-short v0.6.2) software¹¹ before
85 duplicate reads were removed using samtools v0.1.18. 10.3 million unique NGS
86 reads were generated, on average, per sample ($n=170$, $\text{min}=1,241,785$,
87 $\text{max}=23,240,481$, $\text{median}=10,812,279$), with an average coverage of 880 unique
88 reads per nucleotide ($n=70,514,012$, $\text{min}=0$, $\text{max}=7956$, $\text{median}=783$, $\text{sd}=515.4$)
89 and 2155 unique reads per exon ($n=388,974$, $\text{min}=0$, $\text{max}=317678$,
90 $\text{median}=1561$, $\text{sd}=3309.8$) within the complete region enriched for analysis. The
91 detection and clinical analysis of single nucleotide variants and small
92 insertions/deletions was performed as described previously.^{4, 12}

93 **Copy Number Variant Detection**

94 For the 25 samples with WGS data, CNVs were identified using version 2.5 of the
95 Complete Genomics pipeline.¹⁰ Briefly, this strategy incorporates an assessment

96 of (i) sequencing read depth, and (ii) discordant mate-pairs. For each tested
97 individual, sequencing read depth was normalized for GC content and genomic
98 positional effects, and CNV status was calculated for non-overlapping 2Kb
99 genomic intervals through comparison to a baseline sample set – comprised of
100 52 unrelated individuals. To identify the location of breakpoints and insertion
101 points of CNV events, genomic regions where mate-pairs aligned to the reference
102 genome displayed abnormal genomic intervals between the two reads were
103 flagged. Within these flagged regions, local *de-novo* assembly was then
104 performed for sequencing reads where only one of the two reads within the
105 mate-pair aligned to the reference genome. Where possible, the genomic location
106 of breakpoints and insertion points was identified and reported.

107 For targeted NGS samples, CNV detection was performed using ExomeDepth
108 v1.1.6.¹³ For each tested individual, the ExomeDepth algorithm builds the most
109 suitable reference set from the BAM files of a presented group of potential
110 reference samples. We presented ExomeDepth with BAM files for >20
111 individuals that had been generated by identical laboratory and computational
112 procedures. All potential reference samples were individuals referred for
113 genomic diagnostic testing who were not knowingly related to the tested
114 individual, and had been obtained from the same sequencing run on the Illumina
115 HiSeq platform, where possible. The reference sample sets selected by
116 ExomeDepth are referred to as ‘reference samples’ herein.

117 **Accuracy of ExomeDepth in comparison to WGS and MLPA**

118 For 25 individuals with IRD we generated gene panel NGS and WGS datasets
119 (Figure 1). We used the variant detection techniques applied to the WGS datasets

120 as a reference standard for CNV detection and then assessed, at the exon level
121 ($n=1590$ exons per sample), the sensitivity and specificity of ExomeDepth
122 applied to gene panel NGS datasets. We defined sensitivity as the capability of
123 ExomeDepth to identify exons with abnormal CNV, and specificity as the
124 capability to identify exons with a normal CNV status. For a further five
125 individuals with cardiac disorders we generated gene panel NGS and MLPA
126 datasets (Figure 1) and then assessed the sensitivity of ExomeDepth applied to
127 gene panel NGS datasets for these individuals.

128

129 **Assessment of ExomeDepth to identify simulated CNV events**

130 Simulated CNV events were introduced into targeted NGS data for the 25 IRD
131 patients with complementary WGS data. The enrichment region for targeted NGS
132 for the 25 samples encapsulates 1590 protein-coding exons for 105 genes
133 associated with IRD. Importantly, we had previously defined and reported the
134 copy number status for each exon included within the targeted enrichment
135 through the analysis of WGS data.⁷

136 Simulation was performed using a random sample and exon selector, bedtools
137 v2.16.2 intersect, and software within the PicardTools v1.75 java package:
138 DownsampleSam and MergeSamFiles (Supp Figure S1). Exons were excluded
139 from analysis if they overlapped with known heterozygous deletion events in the
140 selected sample. We simulated deletion events for 1000 single exons and 1749
141 multiple exons (2, 3 and 4 exon events). In all cases, we assume that the intronic
142 breakpoints of the deletion event are not captured through NGS. Deletion events
143 are not expected to be detected above a test:reference sample read ratio of 0.7

144 (see supplemental results and methods). We created three discrete groups for
145 simulated deletion events, with the extent of sequencing reads randomly
146 removed indicated in parentheses: (i) control events (0%), (ii) deletion with
147 amplification bias (40%) and (iii) deletion without amplification bias (50%).
148 Further details on the simulation methodology are provided in the supplemental
149 results and methods.

150 **Assessment of factors influencing successful identification of CNV events**

151 We assessed a number of criteria for known and simulated CNV events in order
152 to assess whether they are key determinants of successful CNV surveillance
153 through ExomeDepth, including: (i) the intra-sample variation in coverage, using
154 the normalized read count (reads-per-kilobase-per-million, *rpkm*) coefficient of
155 variation (CV) for surveyed genes in test samples (Supp Figure S2), (ii) the inter-
156 sample variation in coverage, using the *rpkmCV* for surveyed exons across
157 reference samples selected by ExomeDepth (Supp Figure S2), (iii) the percentage
158 of nucleotides and the number of exons containing nucleotides with appropriate
159 sequencing depth for in-house diagnostic surveillance (>50x unique sequencing
160 reads), (iv) the total and normalized read depth across surveyed exons, (v) the
161 GC content of the surveyed regions, (vi) the size of exons, and (vii) the distance
162 between neighbouring exons. All statistical analyses were performed in R v3.2.1
163 software.

164 **Integration of CNV detection during clinical NGS testing**

165 We integrated CNV detection using ExomeDepth into the NGS workflow for 140
166 individuals from 12 distinct referral groups (Supp Table S5). The reasons for
167 assessment of CNV events were (i) an assessment of whether a heterozygous

168 CNV event was *in-trans* to a clearly or likely pathogenic variant, or (ii) an
169 assessment of whether a heterozygous CNV event was present in a gene highly
170 specific to an individual's clinical presentation. In accordance with the
171 recommendations of the ExomeDepth developers, test samples with an overall
172 correlation to selected reference samples <0.97 were repeated with an
173 alternative set of reference samples or excluded from analysis. Clinical
174 interpretation of CNVs was restricted to genes relevant to their referral on a
175 case-by-case basis. We performed additional assays to confirm the presence of
176 all identified CNVs before they were clinically reported. Where kits designed and
177 created by MRC-Holland (Amsterdam, Netherlands) were available, we carried
178 out multiplex ligation-dependent probe amplification (MLPA) assays. In the
179 absence of a suitable MLPA kit, we validated CNVs using droplet digital PCR or a
180 bespoke multiplex quantitative fluorescence methodology (see Supplemental
181 Methods). Validated CNV events were submitted to the ClinVar database.¹⁴

182 Results

183 Accuracy of ExomeDepth in comparison to WGS and MLPA

184 To establish the accuracy and reliability of ExomeDepth when applied to
185 targeted NGS data, we analysed targeted NGS datasets for 30 individuals in
186 whom CNV detection had been performed using either WGS ($n=25$) or MLPA
187 ($n=5$). This allowed calculation of the sensitivity and specificity for identified
188 deletions and duplications. Overall, we found a sensitivity of 92.9% and
189 identified that variable and insufficient coverage within surveyed genes reduces
190 the capability of ExomeDepth to identify single exon deletions.

191 In comparison to WGS, we determined that ExomeDepth applied to targeted NGS
192 datasets (encompassing 1590 exons from 105 genes) has a sensitivity of 100%
193 and a specificity of 99.8% (Supp Table S6) at the exon level. True positive events
194 included a single exon deletion in *GPR98*, a 2 exon deletion in *USH2A*, and a 6
195 exon deletion in *PCDH15* (Supp Table S7). In comparison to MLPA, we identified
196 3 out of 4 single exon deletions and one single exon duplication (Supp Table S7).
197 We assessed a number of key factors, and observed that the sequencing data for
198 the individual in whom a single exon deletion was erroneously not identified,
199 showed the highest intra-sample variation (62%) and the highest level of
200 insufficient coverage (9.5% of exons and 0.86% of nucleotides; sample
201 *14011718*, Supp Table S8).

202 We assessed metrics calculated by ExomeDepth for the 8 previously identified
203 deletions and duplication events, observing that the average confidence (Bayes
204 factor, *BF*) determined by ExomeDepth for true positive CNV events was 45.04
205 (Supp. Table S7, min=6.4, max=76.8) and the average ratio of sequencing reads

206 between test and reference samples for deletions was 0.61 (Supp. Table S7,
207 min=0.539, max=0.745) and 1.4 for the sole duplication event.

208 **Capability of ExomeDepth to identify simulated CNV events**

209 In order to assess factors that influence the successful identification of CNV
210 events in targeted NGS data using ExomeDepth, we introduced simulated events,
211 *in-silico*, into the targeted NGS datasets created in a clinical setting for the 25
212 individuals for whom we held complementary WGS data. We found a 97%
213 sensitivity for simulated events when 50% of the NGS reads were removed from
214 selected exons ($n=2749$), and identified that inter-sample variation – a measure
215 of consistency of NGS read coverage across reference samples (Supp Figure S2) –
216 and insufficient coverage were key determinants of whether simulated events
217 were missed or identified by ExomeDepth (Tables 1 & S9).

218 Single exon deletions ($n=1000$) were introduced into 101 of 105 genes enriched
219 during NGS and we observed that the sensitivity of ExomeDepth for simulated
220 events was 93.5%, with 930 deletions precisely detected at the exon level and 5
221 included in deletion events erroneously identified as spanning to adjacent exons.
222 This sensitivity is reduced to 79.5% when accounting for amplification bias in
223 simulated events (Supp Tables S10 & S11), with an additional 140 false negative
224 events identified when only 40% of the original NGS reads were removed from
225 the selected exon. Interestingly, 51% (36/70) of the false negative simulated
226 events without amplification bias (50% of NGS reads removed) were exons
227 flanked by neighbouring exons within 250 nucleotides of the canonical donor or
228 acceptor sites. Further, all of these 36 events could be identified if the
229 neighbouring exon boundaries were merged into a single analysis region for

230 simulations, increasing the overall sensitivity of ExomeDepth for simulated
231 events to 97.1% (Supp Table S11).

232 Multiple exon deletions ($n=1749$) – where 50% of the NGS reads were randomly
233 removed from adjacent exons – were introduced into all of the 105 genes
234 enriched during targeted NGS for all 25 individuals. We observed sensitivity
235 rates of 96.6% ($n=620$), 95.9% ($n=586$) and 97.1% ($n=543$) for 2 exon, 3 exon
236 and 4 exon deletions, respectively.

237 To ensure that the process of introducing simulated events into targeted NGS
238 data did not influence the performance of ExomeDepth, we performed the same
239 computational processes of the simulation technique for each event, without
240 removing any NGS reads. No single exon or multiple exon simulated deletion
241 events were identified by ExomeDepth in any of these control simulation
242 experiments.

243 **Integration of CNV detection during clinical NGS testing**

244 Following assessment of the accuracy and the reliability of ExomeDepth applied
245 to targeted NGS datasets, we then integrated CNV detection using ExomeDepth
246 into the NGS workflow for 140 individuals from 12 distinct referral groups to
247 assess specific clinical evaluations. These included either (i) an assessment of
248 whether a heterozygous CNV event was *in-trans* to a clearly or likely pathogenic
249 variant, or (ii) an assessment of whether a heterozygous CNV event was present
250 in a gene highly specific to an individual's clinical presentation. This analysis
251 strategy led to the surveillance of a single gene for 128 individuals, two genes for
252 10 individuals and three genes for 2 individuals.

253 **Confirmation of molecular diagnoses for 17 individuals**

254 Analysis on a gene-by-patient basis identified 17 heterozygous CNV events (15
255 deletions, 1 duplication and 1 complex event; Supp Table S12; Supp Figure S3).
256 All events were verified through an alternative technique, were concluded to
257 contribute to the molecular diagnosis for referred individuals and have been
258 submitted to the ClinVar database (Submission number: SUB2171211). The
259 heterozygous CNV events identified by ExomeDepth ranged from a 20 exon
260 deletion in *PCDH15* (NG_009191.2, NM_001142770.1; >600Kb) to single exon
261 deletions in *RPGRIP1* (NG_008933.1, NM_020366.3), *BEST1* (NG_009033.1,
262 NM_004183.3) and *NMNAT1* (NG_032954.1, NM_022787.3). For a single
263 individual referred with a provisional clinical diagnosis of Marfan syndrome, we
264 identified a complex event in *FBN1* (NG_008805.2, NM_000138.4): a 3-exon
265 deletion (chr15:48737523-48741140, c.(5545+1_5546-1)_(5917+1_5918-1)del)
266 and a 2-exon duplication (chr15:48720493-48723049, c.(6739+1_6740-
267 1)_(6997+1_6998-1)dup), consistent with a clinical diagnosis of Marfan
268 syndrome (Figure 2).

269 We assessed metrics calculated by ExomeDepth for identified deletion and
270 duplication events, observing that the average confidence score (BF) attributed
271 to identified CNV events by the ExomeDepth algorithm was 87 (Supp Table S12,
272 min=22, max=321) and the average read count ratio between test and selected
273 reference samples was 0.56 (min=0.518, max=0.637) and 1.35 (min=1.31,
274 max=1.38), respectively.

275 **Accuracy of ExomeDepth applied in a clinical context**

276 To estimate the accuracy of ExomeDepth applied to targeted NGS datasets for the
277 123 individuals determined to be absent of CNV events, we assessed (i) copy
278 number variant status through orthogonal techniques, and (ii) two key factors
279 identified through assessments of simulated and known CNV variants: inter-
280 sample variation and insufficient coverage (Table 1).

281 We calculated the sequencing coverage for each individual, and identified that
282 3% (135/4551) of the surveyed exons contained at least one nucleotide with less
283 than 50 unique NGS reads. Nine of these exons were found in individuals with a
284 confirmed CNV event in the gene, and 28 were in a gene confirmed to be absent
285 of a CNV event through orthogonal techniques (MLPA; Supp Figure S4). Of the
286 remaining 97 exons, 34 were unique patient-exon combinations and 63 were
287 accounted for by 12 exons with insufficient coverage across multiple samples. On
288 average, 4.6% of the nucleotides within these 97 poor coverage exons received
289 less than 50 unique NGS reads ($n=97$, $\text{min}=0.1\%$, $\text{max}=40.9\%$, $\text{median}=3.6\%$),
290 and all exons were within the range of insufficient coverage values observed for true
291 positive simulated deletion events (Table 1).

292 To estimate the accuracy of ExomeDepth in relation to reference samples, we
293 calculated the variability of sequencing coverage across the selected references
294 for each individual, and identified an average inter-sample variation for
295 surveyed exons of 5.1% ($n=4551$, $sd=3.4\%$), with average minimum and
296 maximum values observed per-individual of 2.4% ($sd=1.9\%$) and 9.9%
297 ($sd=5.5\%$), respectively. In comparison to simulated single exon deletions, these
298 data are consistent with an average sensitivity of 98.7% ($sd=1.5\%$, $\text{min}=88.7\%$,
299 $\text{max}=100\%$; Figure 3).

300 For 6 individuals, data from MLPA analyses provided additional support for the
301 absence of a CNV event (Supp Figure S4). For a single individual, we identified a
302 false negative event after subsequent MLPA analysis of the *DSP* gene. We found
303 that alteration of the analysis region, to survey 5 sub-exonic regions enriched by
304 non-overlapping probes though ExomeDepth identified a partial exon
305 duplication event within the *DSP* gene which complemented the result from
306 MLPA (Supp Figure S5).

307 **Discussion**

308 Copy number variants (CNVs) are an important and common form of genomic
309 variation in the general population,^{15, 16} and are implicated in many Mendelian
310 disorders.^{7, 8, 17} An ability to accurately survey for CNV events, in particular in
311 targeted NGS datasets, therefore has the power to increase diagnostic yields and
312 enhance clinical care. While it has already been shown that read count CNV detection
313 algorithms can be successfully applied to targeted NGS data in a research context,^{13,}
314 ¹⁸⁻²⁰ their integration within diagnostic services has been slower due to a lack of
315 validation parameters. In this study, we have identified key factors which can
316 facilitate the successful application of a widely used bioinformatics tool,
317 ExomeDepth,¹³ for CNV surveillance of targeted NGS datasets within the clinical
318 environment.

319

320 CNV detection tools used in a diagnostic context must be able to identify deletion and
321 duplication events that encapsulate single targets/exons included within the targeted
322 enrichments of custom gene panel and WES techniques, which is a known limitation
323 of some publically available algorithms. Since large datasets of known true positive
324 single exon CNV events do not exist, we have developed and applied a
325 computational simulation technique which permits extended assessment of single
326 exon CNV events. As a result, we have been able to perform an assessment of trends
327 in large and controlled datasets (Table 1), We have then used real-time comparison
328 between WGS and targeted NGS data to assess their applicability to real datasets.
329 Using this combined approach we have shown that amplification bias within NGS
330 assays and the distance between exons enriched during NGS influences the overall
331 sensitivity of ExomeDepth (Supp Table S11). After accounting for these dominating

332 factors, we have demonstrated how variability of sequencing coverage between and
333 within samples, the extent of read depth, the size of surveyed exons and the level of
334 insufficient coverage are important determinants of successful identification of single
335 exon deletion events through ExomeDepth (Table 1, Table S9 and S10). Whilst all
336 these metrics are indicated as important quality assurance parameters for the accurate
337 detection of single exon CNVs, they are neither completely independent nor equally
338 applicable to real datasets on an individual basis. We therefore selected two key
339 metrics for routine incorporation into diagnostics: insufficient coverage (test sample
340 dependent) and inter-sample variability (reference sample dependent). This two-part
341 process firstly checks for the quantity of sequencing coverage over exons surveyed in
342 the tested sample, and second, assesses the consistency of NGS read coverage across
343 reference samples for each surveyed exon. We have assimilated this information to
344 successfully integrate surveillance of CNVs into the clinical bioinformatics pipeline
345 for 140 individuals in a clinical setting, achieving a definitive molecular diagnosis in
346 17 of 140 individuals. Importantly, we have shown that 97.2% of the exons surveyed
347 and determined to be absent of a CNV event have sufficient coverage, and none of the
348 insufficiently covered exons lie outside the range of true positives identified from
349 simulated experiments. Moreover, we have calculated the inter-sample variability for
350 surveyed exons on an individual basis, and through comparison to simulated single
351 exon events, estimated the accuracy of ExomeDepth to be 98.7% for the 123
352 individuals without an identified CNV (Figure 3). Both of these quality assurance
353 observations are supported by their integration with other CNV software tools²¹ and
354 the absence of CNV events in 6 individuals tested through MLPA.
355

356 Taken together, our data illustrate the utility of CNV assessments within a diagnostic
357 setting using the publically available ExomeDepth software, and support the
358 utilization of quality assurance parameters in complement to CNV detection
359 algorithms in targeted NGS diagnostic services. Whilst other types of software can be
360 routinely applied to WGS datasets to detect CNVs at single nucleotide resolution, we
361 expect that application of the approaches outlined in this study will improve the
362 utilization of read depth CNV tools in diagnostic environments across heterogeneous
363 targeted NGS gene panel approaches, including small and large gene panels, as
364 described here, and WES.

365 **Figure Legends**

366 **Figure 1. Study Design.** The approach taken in this study to assess the
367 'accuracy' and 'key factors' influencing the accuracy of ExomeDepth applied
368 to targeted next-generation sequencing datasets. The 'key factors' assessed
369 by application of ExomeDepth to datasets with known and simulated CNVs
370 are outlined in Table 1. *CNV*, copy number variation; *WGS*, whole genome
371 sequencing; *MLPA*, multiplex ligation-dependent probe amplification; *gene panel*,
372 next-generation sequencing data generated in a diagnostic environment after
373 enrichment for a set of genes known as a cause of specific Mendelian disorders.

374

375 **Figure 2. *FBN1* copy number variant.** A complex 3-exon deletion,
376 c.(5545+1_5546-1)_(5917+1_5918-1)del, and 2-exon duplication,
377 c.(6739+1_6740-1)_(6997+1_6998-1)dup event identified in *FBN1*
378 (NG_008805.2, NM_000138.4), confirming a clinical diagnosis of Marfan
379 syndrome for the referred individual. *Red crosshairs*, the ratio of reads between
380 test and reference samples; *grey bar*, the 95% confidence interval of expected
381 read ratios in comparison to reference samples.

382

383 **Figure 3. Inter-sample variation in sequencing coverage across surveyed**
384 **exons.** *Simulations*, the variability of sequencing coverage in selected reference
385 samples for 971 identified and 29 missed single exon simulated deletions.
386 *Diagnostic survey*, the variability of sequencing coverage in selected reference
387 samples for 4551 exons surveyed for copy number variants in a diagnostic
388 context.

389 **Declarations**

390

391 **Competing interests**

392 The authors of this manuscript have no competing interests to declare.

393

394 **Funding**

395 This work was supported by the Biotechnology and Biological Sciences Research
396 Council, the Manchester Biomedical Research Centre, The National Institute for
397 Health Research Biomedical Centre at Moorfields Eye Hospital and the UCL Institute
398 of Ophthalmology, the DST UK-India Education and Research Initiative, and an
399 independent research grant funded by the Manchester Academic Health Science
400 Centre.

401

402 **Acknowledgments**

403 We thank all patients, referring clinicians, clinical scientists and genetic counselors
404 involved in this study.

405

406 **Authors' contributions**

407 JME, GCMB, SCR and BP designed and coordinated the study. JME, CC, StB, SaB, SG,
408 RLG, PIS, BH, MM, ARW, WGN, SCR, GCMB contributed genetic and/or phenotypic
409 data. JME wrote the manuscript and all authors provided important revisions and
410 intellectual content.

411 **Bibliography**

412

- 413 1. Baetens M, Van Laer L, De Leeneer K, et al. Applying Massive Parallel
414 Sequencing to Molecular Diagnosis of Marfan and Loeys-Dietz Syndromes.
415 *Human Mutation* 2011; **32**: 1053-62.
- 416 2. O'Sullivan J, Mullaney BG, Bhaskar SS, et al. A paradigm shift in the
417 delivery of services for diagnosis of inherited retinal disease. *J Med Genet* 2012;
418 **49**: 322-6.
- 419 3. Nishio SY, Hayashi Y, Watanabe M, Usami SI. Clinical Application of a
420 Custom AmpliSeq Library and Ion Torrent PGM Sequencing to Comprehensive
421 Mutation Screening for Deafness Genes. *Genet Test Mol Biomarkers* 2015.
- 422 4. Ellingford JM, Barton S, Bhaskar S, et al. Molecular findings from 537
423 individuals with inherited retinal disease. *J Med Genet* 2016.
- 424 5. Lee H, Deignan JL, Dorrani N, et al. Clinical Exome Sequencing for Genetic
425 Identification of Rare Mendelian Disorders. *JAMA* 2014.
- 426 6. Yang Y, Muzny DM, Xia F, et al. Molecular Findings Among Patients
427 Referred for Clinical Whole-Exome Sequencing. *JAMA* 2014.
- 428 7. Ellingford JM, Barton S, Bhaskar S, et al. Whole Genome Sequencing
429 Increases Molecular Diagnostic Yield Compared with Current Diagnostic Testing
430 for Inherited Retinal Disease. *Ophthalmology* 2016; **123**: 1143-50.
- 431 8. Gilissen C, Hehir-Kwa JY, Thung DT, et al. Genome sequencing identifies
432 major causes of severe intellectual disability. *Nature* 2014.
- 433 9. Drmanac R, Sparks AB, Callow MJ, et al. Human genome sequencing using
434 unchained base reads on self-assembling DNA nanoarrays. *Science* 2010; **327**:
435 78-81.
- 436 10. Carnevali P, Baccash J, Halpern AL, et al. Computational techniques for
437 human genome resequencing using mated gapped reads. *J Comput Biol* 2012; **19**:
438 279-92.
- 439 11. Li H, Durbin R. Fast and accurate short read alignment with Burrows-
440 Wheeler transform. *Bioinformatics* 2009; **25**: 1754-60.
- 441 12. Gillespie RL, O'Sullivan J, Ashworth J, et al. Personalized diagnosis and
442 management of congenital cataract by next-generation sequencing.
443 *Ophthalmology* 2014; **121**: 2124-37.
- 444 13. Plagnol V, Curtis J, Epstein M, et al. A robust model for read count data in
445 exome sequencing experiments and implications for copy number variant
446 calling. *Bioinformatics* 2012; **28**: 2747-54.
- 447 14. Landrum MJ, Lee JM, Benson M, et al. ClinVar: public archive of
448 interpretations of clinically relevant variants. *Nucleic Acids Res* 2016; **44**: D862-
449 8.
- 450 15. Iafrate AJ, Feuk L, Rivera MN, et al. Detection of large-scale variation in
451 the human genome. *Nat Genet* 2004; **36**: 949-51.
- 452 16. Sebat J, Lakshmi B, Troge J, et al. Large-scale copy number polymorphism
453 in the human genome. *Science* 2004; **305**: 525-8.
- 454 17. Glessner JT, Bick AG, Ito K, et al. Increased frequency of de novo copy
455 number variants in congenital heart disease by integrative analysis of single
456 nucleotide polymorphism array and exome sequence data. *Circ Res* 2014; **115**:
457 884-96.

- 458 18. Krumm N, Sudmant PH, Ko A, et al. Copy number variation detection and
459 genotyping from exome sequence data. *Genome Res* 2012; **22**: 1525-32.
- 460 19. Fromer M, Moran JL, Chambert K, et al. Discovery and statistical
461 genotyping of copy-number variation from whole-exome sequencing depth. *Am J*
462 *Hum Genet* 2012; **91**: 597-607.
- 463 20. Li J, Lupat R, Amarasinghe KC, et al. CONTRA: copy number analysis for
464 targeted resequencing. *Bioinformatics* 2012; **28**: 1307-13.
- 465 21. Johansson LF, van Dijk F, de Boer EN, et al. CoNVaDING: Single Exon
466 Variation Detection in Targeted NGS Data. *Hum Mutat* 2016; **37**: 457-64.
467
468