1   **Robust identification of deletions in exome and genome sequence data based on**

2   **clustering of Mendelian errors**

3

4   **Authors**: Kathryn B. Manheimer[1], Nihir Patel[1], Felix Richter[1], Joshua Gorham[4], Angela

5   C. Tai[4], Jason Homsy[4,12], Marko T. Boskovski[13], Michael Parfenov[4], Elizabeth

6   Goldmuntz[5,6], Wendy K. Chung[7,8], Martina Brueckner[15,16], Martin Tristani-Firouzi[9],

7   Deepak Srivastava[10,11], Jonathan G. Seidman[4], Christine E. Seidman[4,14], Bruce D.

8   Gelb[1,2,3], Andrew J. Sharp[1,2*]

9

10

11  **Affiliations:**

12

13  [1]Mindich Child Health and Development Institute and Departments of [2]Genetics and

14  Genomic Sciences and [3]Pediatrics, Icahn School of Medicine at Mount Sinai, New York,

15  NY, USA, [4]Department of Genetics, Harvard Medical School, Boston MA, USA,

16  [5]Department of Pediatrics, The Perelman School of Medicine, University of

17  Pennsylvania, Philadelphia, PA, USA, [6]Division of Cardiology, The Children's Hospital

18  of Philadelphia, The University of Pennsylvania Perelman School of Medicine,

19  Philadelphia, PA, USA, Departments of [7]Pediatrics and [8]Medicine, Columbia University

20  Medical Center, New York, NY, USA, [9]Department of Pediatric Cardiology, University of

21  Utah, Salt Lake City, UT, USA, [10]Department of Pediatrics, UCSF, San Francisco, CA,

22  USA, [11]Gladstone Institutes, San Francisco, CA, USA, [12]Cardiovascular Research

23  Center, Massachusetts General Hospital, Boston, MA, USA, [13]Division of Cardiac

24  Surgery, The Brigham and Women's Hospital, Harvard Medical School, Boston, MA,

25  USA, [14]Department of Medicine (Cardiology), Brigham and Women's Hospital, Boston,

26  MA and the Howard Hughes Medical Institute, Chevy Chase, MD, USA, [15]Genetics and

27  [16]Pediatrics, Yale University School of Medicine, New Haven, CT, USA

28

29  *Address for correspondence: Andrew J. Sharp, Department of Genetics and Genomic

30  Sciences, Mount Sinai School of Medicine, Hess Center for Science and Medicine,

31  1470 Madison Avenue, Room 8-116, Box 1498, New York, NY 10029 USA. Telephone:

32  +1-212-824-8942, Fax: +1-646-537-8527, Email: andrew.sharp@mssm.edu

33

34

35    **Abstract**

36    Multiple tools have been developed to identify copy number variants (CNVs) from whole

37    exome (WES) and whole genome sequencing (WGS) data. Current tools such as

38    XHMM for WES and CNVnator for WGS identify CNVs based on changes in read depth.

39    For WGS, other methods to identify CNVs include utilizing discordant read pairs and

40    split reads and genome-wide local assembly with tools such as Lumpy and SvABA,

41    respectively. Here, we introduce a new method to identify deletion CNVs from WES and

42    WGS trio data based on the clustering of Mendelian errors (MEs). Using our Mendelian

43    Error Method (MEM), we identified 127 deletions (inherited and *de novo*) in 2,601 WES

44    trios from the Pediatric Cardiac Genomics Consortium, with a validation rate of 88% by

45    digital droplet PCR. MEM identified additional *de novo* deletions compared to XHMM,

46    and also identified sample switches, DNA contamination, a significant enrichment of

47    15q11.2 deletions compared to controls and eight cases of uniparental disomy. We

48    applied MEM to WGS data from the Genome In A Bottle Ashkenazi trio and identified

49    deletions with 97% specificity. MEM provides a robust, computationally inexpensive

50    method for identifying deletions, and an orthogonal approach for verifying deletions

51    called by other tools.

52

53    **Keywords:** copy number variant identification, whole exome sequencing, whole

54    genome sequencing, UPD

55

56

57

58

**Introduction**

60

61  Structural variation (SV), particularly *de novo* deletions, has been implicated in many

62  human diseases including autism spectrum disorders, developmental delay,

63  schizophrenia and congenital heart disease (Weischenfeldt et al., 2013; Gilissen et al.,

64  2014; Glessner et al., 2014; Szatkiewicz et al., 2014; Brandler et al., 2015). Previously

65  identified using microarrays, many tools have been developed in the past ten years to

66  identify SV from next generation sequencing (NGS) data (Tattini et al., 2015). These

67  tools utilize three main lines of evidence to detect SV: changes in read depth,

68  discordant read pairs and split reads. Assembly methods including genome-wide local

69  assembly and *de novo* assembly are also available (Weisenfeld et al., 2014; Wala et al.,

70  2017).

71

72  With respect to whole exome sequencing (WES) data, one tool to identify copy number

73  variants (CNVs) is XHMM, which identifies changes in normalized read depth within a

74  cohort (Fromer and Purcell, 2014). Although widely used for identifying CNVs from WES

75  data, XHMM has several limitations, including a minimum cohort size and the

76  requirement that CNVs must include at least three exons. Typically, ~20% of putative

77  CNVs identified by XHMM fail to be confirmed, and its sensitivity is limited (Glessner et

78  al., 2014). For example, one study that used both XHMM and SNP arrays to identify *de*

79  *novo* CNVs found that XHMM failed to detect 63% of CNVs identified by the SNP array

4

80  (Glessner et al., 2014). The limited sensitivity of XHMM stems from the limitations of

81  WES, some of which can be overcome with whole genome sequencing (WGS).

82

83  Multiple tools have been developed to identify SV from WGS data including CNVnator

84  and Lumpy (Abyzov et al., 2011; Layer et al., 2014). While CNVnator identifies CNVs

85  based on changes in normalized read depth (Abyzov et al., 2011), Lumpy utilizes

86  discordant read pairs and split reads to identify deletions, duplications and other types

87  of SVs (Layer et al., 2014). Lumpy is often used in combination with CNVnator to take

88  into account changes in read depth. In order to estimate the sensitivity and false

89  discovery rate (FDR), SVs identified by CNVnator and Lumpy were both compared to

90  SVs identified in the 1000 Genomes Project by other SV callers (*e.g.*, Delly, Pindel).

91  Although both tools are reported to have a low FDR (0.4 – 3%) and high sensitivity (60 –

92  90%) (Abyzov et al., 2011; Layer et al., 2014), the accuracy of these tools diminishes

93  when used for identifying *de novo* SV (Kloosterman et al., 2015). This problem results

94  from a lack of sensitivity when identifying SVs: false negatives in parental samples lead

95  to a high false positive rate for calling *de novo* SV, creating a significant challenge when

96  attempting to identify *de novo* events that are potentially pathogenic.

97

98  Here, we describe a novel approach called the Mendelian Error Method (MEM) to

99  identify and/or validate deletion SV in trios with WES and WGS data. MEM is based on

100  the principle described in McCarroll *et al.* 2006 (McCarroll et al., 2006), where the

101  presence of a heterozygous deletion reduces the underlying genotype to a hemizyous

102  state. As genotype callers such as GATK assign diploid genotypes to autosomal loci,

103    regions of heterozygous deletion are erroneously assigned homozygous genotypes. In

104    the context of a trio design, variants within heterozygous deletions frequently display

105    Mendelian errors as a result of this genotype mis-assignment (illustrated in Figure 1).

106    We, therefore, hypothesized that clusters of Mendelian errors could be used as a robust

107    signal for the presence of underlying deletions in sequencing data from trios. We

108    applied MEM to both WES and WGS trio data from the Pediatric Cardiac Genomic

109    Consortium (PCGC) and compared results to deletions identified by XHMM, CNVnator

110    and Lumpy. Overall, our results show that MEM identifies both inherited and *de novo*

111    deletions with a positive predictive value (PPV) exceeding 90%, and identifies additional

112    *de novo* deletions that are missed by other SV callers.

113

114    **Methods**

115

116    *WES and WGS in cases with CHD*

117    Probands were recruited from 10 centers in the United States and United Kingdom as

118    part of the Congenital Heart Disease Genetic Network study of the PCGC as described

119    previously (Homsy et al., 2015). Cases (n=2,601) were subject to WES at the Yale

120    Center for Genome Analysis as described previously (Homsy et al., 2015), with a mean

121    depth of 107x.  All genomic coordinates quoted are based on human genome

122    hg19/build 37. Variants were called following the n+1 protocol from GATK.

123

124    Three hundred and fifty probands and their parents from the PCGC were selected for

125    WGS; of note 332 also have WES data. Cases were sequenced at the Broad Institute

126   (n=25), New York Genome Center (n=25) and Baylor College of Medicine Human

127   Genome Sequencing Center (n=300). Samples were sequenced with PCR-free library

128   preparation (n=325) or with SK2-IES (n=25) to a mean depth of 30x on Illumina HiSeq X

129   Ten sequencers. Variants were called by GATK HaplotypeCaller (version 3.3.2)

130   following GATK best practices for n+1 joint calling

131   (https://software.broadinstitute.org/gatk/best-practices/).

132

133   *WES and WGS of healthy population cohort*

134   Trios representing a typical population cohort (n=1,683) were provided by the Simons

135   Foundation Autism Research Initiative Simplex Collection. Simplex families (two

136   unaffected parents, one child with autism spectrum disorder, and one unaffected sibling)

137   underwent WES using DNA extracted from peripheral blood cells, with a mean depth of

138   117x (O'Roak et al., 2011; Sanders et al., 2012; Iossifov et al., 2014). Trios of

139   unaffected siblings and parents served as a typical population cohort for comparison.

140

141   Five hundred and nineteen quartet families selected from the Simons Simplex

142   Collection (SSC) underwent WGS at the New York Genome Center. Samples were

143   sequenced with either a PCR-based library preparation on an Illumina Hi-Seq 2000

144   (n=39) or PCR-free library preparation on an Illumina HiSeq X Ten (n=480). Sequencing

145   was performed with 150-bp paired reads with median coverage of 37.8x per individual.

146   Detailed information regarding this cohort can be found in Werling *et al.* (Werling et al.,

147   2017)

148

149    Variants were called using GATK HaplotypeCaller (version 3.1-1-g07a4bf8, n=19,

150    version 3.2-2-gec30ce, n=21, version 3.4-0-g7e26428, n=479). GATK best practices

151    (https://software.broadinstitute.org/gatk/best-practices/) were followed. Trios comprising an

152    unaffected sibling and their parents were used as a typical population cohort for

153    comparison in this study with permission from the SSC.

154

155    *Genome in a Bottle (GIAB) WGS with Illumina*

156    The GIAB Ashkenazi Jewish (AJ) trio was subject to WGS using both short and long

157    read methodologies. 148-bp paired-end reads were generated with an Illumina Hiseq

158    instrument. Reads were aligned with BWA-mem (details in Zook *et al.*, 2016) (Zook et

159    al., 2016). Variants were called by GATK HaplotypeCaller (version 3.3.2) following

160    GATK best practices using n+1 joint calling.

161

162    *GIAB deletions for AJ trio*

163    GIAB provided draft benchmark structural variants (SVs) for the AJ trio (v0.3.0a). SVs

164    from 119 different tools were compared and merged using the tool SURVIVOR (Jeffares

165    et al., 2017), which required the breakpoints to be within 1000 bp. Deletions identified

166    by a minimum of two tools were compared to deletions identified by MEM using

167    bedtools and required a 20% reciprocal overlap.

168

169    *Mendelian Error Method (MEM) Pipeline – Figure 2*

170

171    *1. Extract Mendelian errors (MEs) from WES and WGS VCFs*

8

172    MEs were extracted based on genotypes reported in the joint VCF produced by GATK

173    best practices, using in-house perl scripts or vcftools. Table S1 includes the eight

174    scenarios considered as MEs that could represent a deletion.

175

176    *2. Filtering*

177    Variants in PCGC, GIAB and SSC trios were filtered using the following criteria: read

178    depth ≥10, genotype quality >60 for WES and >30 for WGS (Table S2). B allele

179    frequency (BAF, defined as the alternate allele depth/total depth) was calculated for

180    heterozygous SNVs, and those with a BAF <0.25 or >0.75 were excluded. Regions

181    overlapping segmental duplications obtained from the UCSC Genome Browser track

182    were excluded. CNVs with a minor allele frequency ≥0.05 in European, African or East

183    Asian ancestry as identified in Conrad *et al.* were excluded (Conrad et al., 2012). For

184    WGS, SNVs with a mappability score <1 were excluded, based on the UCSC Genome

185    Browser track "Alignability of 100mers by GEM from ENCODE/CRG(Guigo)". Regions

186    with tandem repeats, taken from the UCSC Genome Browser track "Simple Repeats"

187    and expanded ±5 bp, were excluded. The Hardy Weinberg equilibrium (HWE) statistic

188    was calculated using vcftools for SNVs with a minimum allele frequency of 0.01 in

189    parents. Any SNVs with a HWE p-value equal to zero were removed.

190

191    *3. Sliding window analysis*

192    We generated 2-Mb windows with 95% overlap for WES analysis and 100-kb windows

193    with 90% overlap for WGS analysis using Bedtools (version 2.26.0) makewindows. In

194    house bash scripts utilizing Bedtools intersect were used to calculate the number of

195    MEs for each window. This was applied to each sample in the PCGC and SSC cohorts

196    separately.

197

198    For each unique window, the number of probands with MEs, the minimum number of

199    MEs, the maximum number of MEs and the average number of MEs per proband were

200    calculated for PCGC and SSC probands. We filtered for windows where the average

201    number of MEs per proband was >2 MEs.

202

203    *4. Comparison to population cohort*

204    Windows with MEs in PCGC cases were compared to corresponding windows in the

205    SSC population cohort. Windows with a ME cluster in three or more SSC probands

206    were excluded, except if the maximum number of MEs in cases was >5.

207

208    *5. Merge windows*

209    For each sample overlapping windows with MEs were merged to identify putative

210    deletion regions. The minimum, maximum and average number of MEs per window was

211    calculated for each region. The number of MEs in each putative deletion region was

212    calculated in SSC probands and regions with ME clusters as described in Step 4 were

213    removed from further analysis.

214

215    *6. Filter for ME clusters*

216    Finally, we filtered for regions with an average number of MEs per window >2 in cases.

217    We identified the first and last ME within each region and used these as the coordinates

218    for the putative deletions.

219

220    *Visualization*

221    1. XHMM

222    For putative deletions identified with MEM from the PCGC WES cohort, we extracted z-

223    scores of the PCA-normalized read depth for each exon from XHMM (Fromer and

224    Purcell, 2014). Putative deletions were inspected visually (Figure S1) and exons with z-

225    scores <-2 were considered candidates for deletions.

226

227    2. IGV

228    Integrated Genomics Viewer (IGV, version 2.3.34) pileup visualization was used as one

229    method for deletion validation. Variants were visualized in the proband and parents.

230    Deletions were excluded if any of the following aspects were detected: multiple reads

231    with quality scores of zero in child or parents, no clear drop of coverage in the proband,

232    or the presence of heterozygous SNVs in the proband.

233

234    *CNVnator*

235    CNVnator identifies CNVs in WGS data based on changes in normalized read depth

236    (Abyzov et al., 2011). Deletions were called for each case proband and the GIAB

237    proband with CNVnator (version 0.3.2) and genotyped for putative copy number within

238    the CNV regions on a scale from 0 – 3. We considered scores between 0.7 – 1.4 as

239    indicating a heterozygous deletion. *De novo* deletions were identified by filtering for a

240    score <1.4 in the child and >1.4 in the parents. We overlapped putative deletions in

241    WGS cases identified using MEM with *de novo* deletions identified by CNVnator using

242    Bedtools intersect, requiring a 25% reciprocal overlap. In the AJ trio, we overlapped

243    putative deletions identified with MEM with both inherited (proband genotype <1.4) and

244    *de novo* deletions called by CNVnator, and considered all intersections with at least 1

245    bp of overlap.

246

247    *Lumpy*

248    Lumpy identifies SVs based on discordant read pairs and split-reads (Layer et al.,

249    2014). Deletions were called for each case proband and the GIAB proband with Lumpy

250    (version 0.2.13) and genotyped using SVtyper (version 0.0.4). *De novo* deletions were

251    identified based on proband and parent genotypes. We overlapped PCGC WGS MEM

252    deletions with Lumpy *de novo* deletions in the same manner as CNVnator. In the AJ

253    trio, we overlapped putative deletions identified with MEM with both inherited and *de*

254    *novo* deletions by Lumpy, and considered all intersections with at least 1 bp of overlap.

255

256    SvABA

257    Deletions were called with SvABA from 350 WGS trios based on genome-wide local

258    assembly (Wala et al., 2017). Default parameters were employed to identify putative

259    copy number variants, which were further validated by IGV visualization prior to digital

260    droplet PCR analyses.

261

262    *Deletion validation*

263    Digital droplet PCR (ddPCR) was used to validate MEM WES deletions and WGS *de*

264    *novo* deletions identified by CNVnator and Lumpy, as previously reported (Mazaika and

265    Homsy, 2014) with the following modification.  PCR primers that amplified a portion of

266    the putative CNV were designed to avoid homopolymer runs or probes that begin with

267    G.  PCR-positive droplets were identified by EvaGreen dye (DNA-bound emission at

268    500/533 nm).  CNV product positive droplets were EvaGreen dye positive, VIC

269    negative. A VIC probe targeting the RPP30 gene was used as reference. Reaction

270    mixtures of 20μL volume comprising ddPCR Master Mix (Bio-Rad), relevant forward and

271    reverse primers and probe(s) and 50ng of DNA were prepared, ensuring that<40% of

272    the 5000-10000 droplets ultimately produced were positive for Evagreen dye and/or VIC

273    signal. For *de novo* CNV confirmations, DNA from the subject with CHD and parents

274    was used. After thermal cycling, plates were transferred to a droplet reader (Bio-Rad)

275    that flows droplets single-file past a 2-color fluorescence detector. Differentiation

276    between droplets that contain target and those that did not was achieved by applying a

277    global fluorescence amplitude threshold in QuantaSoft (Bio-Rad). The threshold was set

278    manually based on visual inspection at approximately the mid-point between the

279    average fluorescence amplitude of positives and negative droplet clusters on each of

280    the EvaGreen dye and VIC channels. Confirmed CNV duplications had ≈ 50% increase

281    in the ratio of positive to negative droplets, as did the reference channel. Conversely,

282    confirmed CNV deletions had approximately half the ratio of positive to negative

283    droplets, as did the reference channel. CNVs that were called, but were unable to be

13

284    confirmed or rejected due to ddPCR technical failure or DNA unavailability were

285    excluded from analysis.

286

287    **Results**

288

289    *MEM identifies inherited and de novo deletions from WES trios*

290    The MEM pipeline was used to analyze WES data from 2,601 PCGC trios and 1,683

291    healthy trios from the SSC. Windows with ME clusters in SSC probands were removed

292    as described in Methods in order to limit our findings to those of likely relevance to the

293    pathogenesis of congenital heart disease (CHD). MEM identified a final set of 171

294    merged and filtered regions containing putative deletions in the PCGC probands (Table

295    S3). We used the location of the first and the last ME in each region with a ME cluster to

296    define the minimal coordinates for the deletion. We utilized XHMM read depth data to

297    perform an initial assessment of the accuracy of our MEM deletion calls. The proband's

298    normalized XHMM z-scores for each exon within the deletion identified by MEM were

299    compared to the rest of the cohort (Figure S1). The presence of outlier negative z-

300    scores in the proband suggested a deletion. The parents' z-scores were also compared

301    to the rest of the cohort to determine if the deletion was inherited or *de novo.* In this

302    manner, 58 deletions were determined to be *de novo*, and 79 were noted to be

303    inherited. Of note, the exons in 13 ME clusters did not have negative normalized z-

304    scores, and seven ME clusters showed inconsistent scores, with some exons showing

305    reduced XHMM z-scores, while other exons were within the normal range (z-score >-2),

306    suggesting that these 20 calls could be false positives.

14

307

308 We directly compared the performance of MEM for the detection of *de novo* deletions

309 with that of XHMM. Fifty deletions were called by both tools, 46 by XHMM alone, and 25

310 by MEM alone (Figure 3A). Of note, the 25 MEM-exclusive deletions included 13 that

311 showed no reduction in z-scores with XHMM for proband or parents and, thus, could

312 represent either *de novo* deletions or false positives. We considered the size of the

313 deletions that MEM did and did not identify. For deletions ≥200 kb, MEM identified

314 100% of deletions, however for deletions <200 kb MEM identified 24% of deletions

315 (Figure 3B). The 46 XHMM-exclusive deletions had a mean size of 35 kb and, therefore

316 due to an insufficient number of SNPs within them, could not be identified by MEM with

317 high recall.

318

319 From the 171 MEM deletions, 36 overlapped with deletions previously confirmed by

320 digital droplet PCR (ddPCR). For the remaining 135 deletions, we performed ddPCR,

321 which was successful for 109 deletions. Ninety-six out of 109 were confirmed as true

322 deletions, achieving a positive predictive value (PPV) of 88.1%. Surprisingly, the results

323 from ddPCR indicated that five of the regions with the ME cluster were inherited

324 duplications. Thus, overall 137/145 (94.5%) of ME clusters identified by MEM were

325 confirmed as true CNVs. Deletions identified as inherited by inspection of XHMM z-

326 score plots confirmed with a PPV of 86% (49/57 inherited, 3/57 *de novo*). From the

327 possible false positives, two out of eight deletion regions without negative normalized z-

328 scores in XHMM were confirmed, and four of six regions with inconsistent loss of exons

15

329     confirmed. Finally, 26 *de novo* deletions were confirmed, four exclusively identified by

330     MEM.

331

332     *Enrichment of deletions on chromosome 15q11.2*

333     With MEM, we identified 15 deletions (13 inherited, 2 *de novo*) ranging from 11 kb to 1

334     MB in the chromosome region 15q11.2 in PCGC probands. These deletions fall in a

335     known microdeletion region between breakpoints (BP) 1 and 2, with a population

336     frequency of 0.25% (Cafferkey et al., 2014). Deletions in this region occurred at a

337     frequency of 0.58% (15/2,601) in the PCGC cohort, and are therefore enriched

338     compared to the reported population frequency (binomial, p=0.004) and to SSC

339     probands, which had a deletion frequency of 0.24% (4/1,683) deletions in this region

340     (binomial, p=0.002).

341

342     *Identification of uniparental disomy (UPD) in WES trios by MEM*

343     Following ME extraction and applying quality filters (Table S1), the majority of trios had

344     between 0.6 – 2% of loci that were scored as MEs (Figure 4A).  We identified eight

345     probands with an elevated rate of MEs distributed across an entire chromosome,

346     suggestive of possible uniparental disomy (UPD). Prior microarray experiments noted

347     UPD of chromosome 15 for one proband, and an extended region of homozygosity on

348     chromosome 16 for a second proband. However, there was no prior indication of UPD

349     in the other six cases.

350

351  All eight instances of UPD were classified as maternal heterodisomy, based on the

352  presence of heterozygous maternal SNPs. The heterodisomic inheritance was for

353  chromosomes 4 (x2), 8, 9, 14, 15 and 16 (x2). UPD was not found in any SSC

354  probands, and was therefore enriched in cases (binomial, p=0.026).

355

356  *MEs identify irregularities in WES trios*

357  We identified two other distinct ME patterns that were informative. Twenty trios had a

358  dramatically higher rate of MEs (~50% of all SNVs), which were distributed across every

359  chromosome (Figure 4D). Nearly all of the MEs were attributable to lack of inheritance

360  from one parent, suggesting either a sample switch or incorrect paternity.

361

362  Similarly, we observed an elevated, but lower, rate (20-30%) of MEs distributed across

363  the entire genome in six other probands (Figure 4C). We hypothesized that this pattern

364  might be due to DNA contamination, which was confirmed with the program

365  VerifyBamID (Jun et al., 2012).

366

367  All samples with likely sample mix-ups, DNA contamination or UPD were excluded from

368  further analysis.

369

370  *ME clusters are non-random in the genome*

371  Before applying MEM to WGS data, we first needed to determine if the increased SNV

372  density in WGS data relative to WES data could lead to ME clusters by chance alone.

373  To test this, we generated a null model of SNV clusters across the genome. We only

374    considered heterozygous SNVs, and also applied additional filters for genotypes

375    generated from WGS as shown in Table 1. After applying these quality filters, the

376    median number of MEs per proband among the 350 PCGC WGS trios was 317. We

377    then ran 1000 permutations of selecting 317 informative SNV positions from one trio,

378    assuming those were MEs, and implemented MEM with a 100-kb window and 10-kb

379    slide. We calculated the number of windows with SNV clusters divided by the number of

380    windows with at least 1 SNV.  The null model had a mean of 0.3% of windows with a

381    SNV cluster (Figure S2). In contrast, 21.4% of windows with at least 1 ME among the

382    PCGC WGS probands had a ME cluster and they were infrequent across the genome

383    (Figure S2). From these results, we inferred that ME clusters in WGS were likely non-

384    random and were likely identifying underlying deletions.

385

386    *Mendelian error clusters identify deletions from GIAB Ashkenazi trio*

387    To test the robustness of MEM for calling deletions from WGS, we identified putative

388    deletions using MEM based on genotypes generated using Illumina short read WGS

389    data for an Ashkenazi Jewish (AJ) trio sequenced by the GIAB consortium (Zook et al.,

390    2016). We processed filtered SNV genotypes from the Illumina WGS data in this trio

391    using the parameters listed in Table 1 and searched for ME clusters. Using the MEM

392    pipeline we identified 32 putative deletions (Table S3) that contained an average of 9.4

393    MEs, with a mean size of 31.5 kb.

394

395    To determine the accuracy of the MEM deletion calls, we intersected them with draft

396    benchmark deletions provided by GIAB. Requiring a 20% reciprocal overlap between

18

397  deletions, 27/32 MEM deletions overlapped with those from GIAB. After removing the

398  20% overlap requirement 31/32 MEM deletions overlapped. The five deletions that did

399  not overlap by 20% were visualized in IGV, where we found evidence for a deletion in

400  4/5. Therefore, MEM identified deletions with 97% precision from WGS for the GIAB AJ

401  proband. Of note, one 215-kb MEM deletion overlapped two GIAB deletions.

402  Visualization in IGV confirmed the presence of two separate deletion events at this

403  locus, which the distribution of MEs also supports (Figure S3).

404

405  Next, we looked at the deletions identified by GIAB that MEM did not identify

406  (n=24,090). These do not include deletions in segmental duplication regions but do

407  include 14,690 deletions at tandem repeat loci. Due to the challenges of sequencing

408  tandem repeats with short read sequencing we would not expect MEM to accurately

409  identify deletions with tandem repeats, as variant calling is unreliable in these regions.

410  The MEM false negatives (FNs) had a median size of 39 bp and a mean size of 306 bp

411  and were attributable to inadequate number of MEs in those deletions as 93.5% did not

412  include any MEs before filtering. Only 1% of the MEM FNs were related, at least in part,

413  to the filtering of MEs, having >2 MEs prior to filtering.

414

415  We also compared the MEM calls for the AJ trio to calls from CNVnator and Lumpy. Of

416  the 32 MEM deletion calls, 27 (84%) and 23 (72%) overlapped with calls from CNVnator

417  and Lumpy, respectively. There were many calls from CNVnator and Lumpy that were

418  not made by MEM, however most of them contained no MEs. ME filtering accounted for

419  21% of FNs from CNVnator calls and 6% of FNs from Lumpy calls.

19

420

421    *MEM identifies deletions from WGS trios*

422    Based on the promising results from GIAB, we proceeded to apply the MEM pipeline to

423    identify deletions from 350 WGS case trios from the PCGC, and 517 healthy trios from

424    the SSC. From the PCGC trios, MEM identified 6,645 regions with ME clusters

425    (mean=19.1/proband) that ranged in size from 3 bp to 9 Mb, with a median size of 2.9

426    kb and a mean size of 20 kb (Table S3). Eleven percent of regions included exons. We

427    used the first and last MEs as coordinates for the putative deletions. For 332 PCGC

428    trios that have both WES and WGS data we compared the deletions identified by MEM

429    from both data sets. MEM identified 11 deletions from WES, all of which were detected

430    by MEM with WGS. All of the deletions were the same size or larger when detected by

431    WGS except for one. This is expected as the increased SNP density of WGS provides

432    more informative sites for MEM, thus facilitating a better estimate of the deletion size.

433

434    To determine if the ME clusters in WGS data identified true deletions, we integrated

435    normalized read depth data from CNVnator. Each region was labeled with a CNVnator

436    score where 0 corresponds to a homozygous deletion, 0.7-1.5 to a heterozygous

437    deletion, 1.5-2.4 to being normally diploid and >2.4 to a duplication. The vast majority

438    (97%) of MEM deletions had a CNVnator score between 0.7 – 1.5 suggesting MEM was

439    identifying true heterozygous deletions (Figure S4). We visualized MEM deletion calls

440    with a CNVnator score >1.5 in IGV. Based on this manual curation, we concluded that

441    the majority (66%) were false positives, but 34% were heterozygous deletions: 10%

442    covering the entire region and 24% being either a deletion of a portion of the region or

20

443    two smaller deletions located close together. In addition, we visualized in IGV a test set

444    of MEM deletions with a range of CNVnator scores. The vast majority of false positives

445    (93.5%) had a score of 1.5 or greater, while 100% of the true or possible deletions had

446    a score between 0.7 and 1.5 (Figure S5). Overall, our comparison with read depth data

447    supports a PPV of 92% (Supplementary Formula 1) for identifying heterozygous

448    deletions from WGS with MEM.

449

450    Next, we identified which MEM deletions were *de novo* based on the proband and

451    parents' CNVnator scores. We used two sets of filters (Table S4) and identified 37

452    putative *de novo* deletion calls (mean = 0.12 *de novo* deletions/proband) After

453    visualization in IGV, we determined that 20/37 represented likely true *de novo* deletions,

454    while 17 were inherited. We compared these to *de novo* deletions identified by

455    CNVnator, Lumpy and a third WGS tool called SvABA that uses genome-wide local

456    assembly to identify SV (Wala et al., 2017). The deletions called by the other SV tools

457    were confirmed by ddPCR. Of the 20 *de novo* deletions found by MEM, five were also

458    identified by CNVnator, Lumpy, and SvABA, three were identified by CNVnator and

459    SvABA but not Lumpy, and 12 were not found by the three other tools. Thirteen

460    additional *de novo* deletions were identified with a combination of CNVnator, Lumpy and

461    SvABA: all three tools but not MEM (n=7), CNVnator and SvABA (n=2), CNVnator and

462    Lumpy (n=1), CNVnator only (n=2), and SvABA only (n=1). None of these deletions,

463    which had a median size of 6.5 kb, included any MEs, suggesting MEM is less sensitive

464    for deletions smaller than ~10 kb in WGS.

465

466   *MEM is computationally efficient*

467   We compared the computational resources required for MEM and the other CNV

468   detection tools used in this study for deletion identification in one trio (Table 1). Runtime

469   and memory for all tools were based on the use of an Intel Haswell 2.4 GHz processor

470   with 64 GB memory and Cray nodes. We did not utilize parallelization for any of the

471   tools. Runtime and memory for MEM was calculated for Step 1 of the MEM pipeline (ME

472   extraction). All other steps in the MEM pipeline can be performed on the command line

473   and do not require significant time or memory. Of note, resources required for the

474   preliminary steps for all tools (DepthOfCoverage for XHMM, Samblaster for Lumpy, and

475   variant calling for MEM) were not included.

476

477   For WES, MEM required 5.5 sec and an average of 12 MB of memory per trio. XHMM

478   required 453 sec and on average 81 MB of memory. For WGS, MEM required 407 sec

479   and an average of 7 MB of memory per trio. CNVnator required 77,629 sec and, on

480   average, 709 MB of memory. Lumpy/SVTyper required 4,238 sec and an average of

481   4,898 MB of memory. SVTyper produced genotypes for deletions only and not other

482   types of SV (duplications, translocations, inversions). For both WES and WGS, MEM

483   performed significantly faster and required significantly less memory compared to other

484   CNV detection tools. Of note, ME extraction execution time grows sub-linearly based on

485   the number of trios present in the VCF, however average memory required does not

486   increase significantly.

487

488   **Discussion**

489    A variety of tools have been developed to identify CNVs including XHMM and CoNIFER

490    for WES, and CNVnator, Lumpy and SvABA for WGS. Each of these tools has

491    limitations such as a requirement for 50 samples, the need for extensive computational

492    resources, or that up to 20% of CNVs will fail to confirm. In addition, false negative calls

493    in parents lead to a high false positive rate for *de novo* deletion CNV calls, making the

494    identification of true *de novo* CNVs difficult and time intensive. As documented in this

495    report, we developed a novel method, MEM: the Mendelian Error Method, to identify

496    deletion CNVs based on ME clustering. This orthogonal method identifies deletions with

497    a PPV >90% for both WES and WGS, and identifies additional *de novo* deletions

498    compared to other SV callers.

499

500    When used with WES, we demonstrate that MEM has several advantages compared to

501    XHMM. First, MEM can be used on a single trio, while XHMM requires a minimum of 50

502    samples to accurately normalize read depth and calculate z-scores. Second, MEM

503    requires substantially less memory and runtime compared to XHMM. Third, MEM can

504    be used as a method for quality control, as it can identify UPD, sample mix-ups and

505    DNA contamination. MEM is also a worthwhile complementary tool to XHMM as MEM

506    identified additional *de novo* deletions that XHMM missed due to spurious evidence of

507    inheritance or seemingly inconsistent loss of exons. In addition, MEM identified

508    deletions with less than 3 exons with high precision, albeit with low sensitivity. The

509    combination of evidence from both XHMM and MEM can increase our ability to identify

510    smaller deletions with high precision and increased sensitivity, as well as reducing the

511    need for PCR-based validation, which is expensive and time-consuming.

23

512

513    CNV identification from WGS data is still under development. We propose MEM as a

514    worthwhile addition to the WGS CNV identification toolbox as it can be efficiently

515    implemented in less than a day and identifies deletions with a >90% PPV. It can be

516    implemented on a large cohort without significantly increasing the computational

517    requirements, and identifies additional *de novo* deletions compared to CNVnator,

518    Lumpy and SvABA. While there are other SV tools for WGS data (*e.g.*, Delly, Pindel),

519    the methods utilized by CNVnator, Lumpy and SvABA, represent three primary ways to

520    identify CNVs: changes in read depth, discordant/split reads and local assembly, yet

521    MEM identified additional *de novo* deletions. Equally helpful is the orthogonal nature of

522    MEM, which may reduce the need for PCR validation for deletions identified by MEM

523    and a second tool.

524

525    MEM's primary limitation is the need for a complete trio, as many cohorts only recruit

526    singletons. The trio design is necessary in order to identify MEs and, therefore, cannot

527    be avoided. MEM is also limited regarding the size of the deletions it can detect with

528    high recall, which is a function of the SNV density in NGS data. WES deletions <200 kb

529    are identified with 24% recall, while deletions >200 kb are identified with 100% recall. Of

530    note, although the smaller deletions are not identified with high sensitivity, the PPV

531    remains high when they are called (78%). Based on deletions identified in GIAB AJ trio,

532    MEM identifies deletions from WGS with a range of sizes (100 – 660,000 bp); however,

533    we estimate that MEM has ~1% recall for deletions smaller than 3 kb and only 18%

534    recall for deletions 3-10 kb. Deletions >10 kb are identified with 45% recall. For this

24

535     reason MEM applied to WGS is particularly valuable as a secondary and orthogonal

536     method to confirm deletions identified by other tools, as the PPV is 92 - 97% with WGS

537     data.

538

539     MEM's sensitivity was also reduced by ME filtering, which accounted for ~5% of the

540     false negatives. Filtering is necessary in order to remove MEs caused by poor

541     genotyping or other errors and to achieve a high PPV. We suggest noting the number of

542     filtered MEs when verifying deletions identified with other tools, as even the presence of

543     1 or 2 MEs after filtering is evidence for a deletion in 88% of calls (data not shown).

544

545     Interestingly, 3.5% of regions with ME clusters identified with MEM were scored as

546     inherited duplications by ddPCR. Although ME genotypes are not indicative of a

547     duplication, it is has been noted that some CNVs are complex events with multiple

548     breakpoints comprising both deletions and duplications in close proximity (Quinlan et

549     al., 2010). We hypothesize that this phenomenon likely underlies our observations, and

550     that in these few cases the primer placement for ddPCR targeted a region of duplication

551     rather than the deletion found my MEM.

552

553     The pursuit of disease-causing CNVs in family trios often focuses on the identification of

554     *de novo* or rare CNVs. MEM identifies both inherited and *de novo* deletions, however

555     one is unable to distinguish between inherited and *de novo* deletions without the use of

556     a secondary tool that identifies deletions in parents. In order to identify rare CNVs from

557     a large cohort, one must eliminate regions with deletions in the general population. This

558   is included in the MEM pipeline in Steps 4 and 5. If population data are not available,

559   one could determine the number of samples with deletions in each region identified by

560   MEM as an alternative. Deletion regions found in multiple samples are less likely to be

561   disease-causing.

562

563   We applied MEM to trios from the PCGC to identify deletions that are causal for CHD

564   that had not been seen with previous studies (Glessner et al., 2014). With MEM, we

565   identified and quantified two genetic mechanisms associated with CHD; BP1-BP2

566   deletions in 15q11.2 and UPD. Deletions in the region 15q11.2 BP1-BP2 account for

567   ~0.3% of CHD cases in the PCGC cohort. Although 15q11.2 deletions are associated

568   with a wide range of phenotypic anomalies, CHD have been reported in ~9% of carriers

569   (Cox and Butler, 2015), which explains the presence of an inherited mutation present in

570   both a proband with CHD and their apparently unaffected parent.

571

572   Using MEM, we also identified whole-chromosome maternal heterodisomy in ~0.3% of

573   CHD cases in the PCGC cohort. The likely genetic mechanism for maternal

574   heterodisomic UPDs is non-disjunction and subsequent trisomy rescue. Thus, there is a

575   possibility that probands with UPD may be mosaic for trisomy of the UPD chromosome,

576   and this mosaic trisomy could be the underlying cause of the probands' CHD. UPD

577   could also lead to CHD due to changes in methylation of imprinted genes. One example

578   from the chromosomes affected in PCGC probands is chromosome 8, which harbors

579   the known CHD gene *CHD7* (MIM:608892) that is maternally methylated (Joshi et al.,

26

580 2016). Maternal heterodisomy would lead to hypermethylation and altered expression of

581 *CHD7.*

582

583 In conclusion, MEM is an orthogonal tool that identifies deletion CNVs with over 90%

584 PPV and is a valuable addition to CNV detection pipelines for both WES and WGS. As

585 NGS data becomes more accessible, the need to identify CNVs from WES and WGS

586 data will only increase. This is particularly true with relation to disease causing CNVs as

587 CNVs have been implicated in a number of different human diseases including

588 congenital heart disease, schizophrenia, developmental delay and autism spectrum

589 disorders. MEM helps overcome some of the challenges associated with identifying

590 pathogenic CNVs due to limited specificity of current SV tools.

591

592 **Acknowledgments**

615

616    **Conflict of Interest:** The authors have no conflicts of interest to declare.

617

618    **Ethical compliance:** All procedures performed in studies involving human participants

619    were in accordance with the ethical standards of the following Institutional Review

620    Boards: Boston Children's Hospital, Brigham and Women's Hospital, Great Ormond

621    Street Hospital, Children's Hospital of Los Angeles, Children's Hospital of Philadelphia,

622    Columbia University Medical Center, Icahn School of Medicine at Mount Sinai,

623    Rochester School of Medicine and Dentistry, Steven and Alexandra Cohen Children's

624    Medical Center of New York, and Yale School of Medicine. Informed consent was

625    obtained from all individual participants or their parent/guardian included in this study.

626

**References**

628  Abyzov A, Urban AE, Snyder M, Gerstein M. 2011. CNVnator: An approach to discover,

629  genotype, and characterize typical and atypical CNVs from family and population

630  genome sequencing. Genome Res 21:974–984.

631  Brandler WM, Antaki D, Gujral M, Noor A, Rosanio G, Chapman TR, Barrera DJ, Lin

632  GN, Malhotra D, Watts AC, Wong LC, Estabillo JA, et al. 2015. Frequency and

633  complexity of de novo structural mutation in autism. bioRxiv 1–19.

634  Cafferkey M, Ahn JW, Flinter F, Ogilvie C. 2014. Phenotypic Features in Patients With

635  15q11.2(BP1-BP2) Deletion: Further Delineation of an Emerging Syndrome. Am J Med

636  Genet Part A 2:1916–1922.

637  Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD,

638  Barnes C, Campbell P, Hu M, Ihm CH, et al. 2012. Origins and functional impact of copy

639  number variation in the human genome. 464:704–712.

640  Cox DM, Butler MG. 2015. The 15q11.2 BP1-BP2 microdeletion syndrome: A review. Int

641  J Mol Sci 16:4068–4082.

642  Fromer M, Purcell SM. 2014. Using XHMM Software to Detect Copy Number Variation

643  in Whole-Exome Sequencing Data. Curr Protoc Hum Genet 81:7.23.1-7.23.21.

644  Gilissen C, Hehir-Kwa JY, Thung DT, Vorst M van de, Bon BWM van, Willemsen MH,

645  Kwint M, Janssen IM, Hoischen A, Schenck A, Leach R, Klein R, et al. 2014. Genome

646  sequencing identifies major causes of severe intellectual disability. Nature 511:344–

647  347.

648  Glessner J, Bick AG, Ito K, Homsy J, Rodriguez-Murillo L, Fromer M, Mazaika EJ,

649    Vardarajan B, Italia MJ, Leipzig J, DePalma S, Golhar R, et al. 2014. Increased

650    frequency of de novo copy number variations in congenital heart disease by integrative

651    analysis of SNP array and exome sequence data. Circ Res.

652    Homsy J, Zaidi S, Shen Y, Ware JS, Samocha KE, Karczewski KJ, Depalma SR,

653    Mckean D, Wakimoto H, Gorham J, Jin SC, Deanfield J, et al. 2015. De novo mutations

654    in congenital heart disease with neurodevelopmental and other congenital anomalies.

655    Science (80- ) 350:1262–1266.

656    Iossifov I, O'roak BJ, Sanders SJ, Ronemus M, Krumm N, Levy D, Stessman HA,

657    Witherspoon K, Vives L, Patterson KE, Smith JD, Paeper B, et al. 2014. The

658    contribution of de novo coding mutations to autism spectrum disorder. Nature 13:216–

659    221.

660    Jeffares DC, Jolly C, Hoti M, Speed D, Shaw L, Rallis C, Sedlazeck FJ. 2017. Transient

661    structural variations have strong effects on quantitative traits and reproduction isolation

662    in fission yeast. Nat Commun 1–11.

663    Joshi RS, Garg P, Zaitlen N, Lappalainen T, Watson CT, Azam N, Ho D, Li X,

664    Antonarakis SE, Brunner HG, Buiting K, Cheung SW, et al. 2016. DNA Methylation

665    Profiling of Uniparental Disomy Subjects Provides a Map of Parental Epigenetic Bias in

666    the Human Genome. Am J Hum Genet 99:555–566.

667    Jun G, Flickinger M, Hetrick KN, Romm JM, Doheny KF, Abecasis GR, Boehnke M,

668    Kang HM. 2012. Detecting and estimating contamination of human DNA samples in

669    sequencing and array-based genotype data. Am J Hum Genet 91:839–48.

670    Kloosterman WP, Francioli LC, Hormozdiari F, Marschall T, Hehir-kwa JY, Abdellaoui A,

671    Lameijer E, Moed MH, Koval V, Renkens I, Roosmalen MJ Van, Arp P, et al. 2015.

672    Characteristics of de novo structural changes in the human genome. Genome Res 792–

673    801.

674    Layer RM, Chiang C, Quinlan AR, Hall IM. 2014. LUMPY: a probabilistic framework for

675    structural variant discovery. Genome Biol 15:R84.

676    Mazaika E, Homsy J. 2014. Digital Droplet PCR: CNV Analysis and Other Applications.

677    Curr Protoc Hum Genet 82:7.24.1-7.24.13.

678    McCarroll S a, Hadnott TN, Perry GH, Sabeti PC, Zody MC, Barrett JC, Dallaire S,

679    Gabriel SB, Lee C, Daly MJ, Altshuler DM, The International HapMap Consortium.

680    2006. Common deletion polymorphisms in the human genome. Nat Genet 38:86–92.

681    Mccarroll SA, Hadnott TN, Perry GH, Sabeti PC, Zody MC, Barrett JC, Dallaire S,

682    Gabriel SB, Lee C, Daly MJ, Altshuler DM, Hapmap I. 2006. Common deletion

683    polymorphisms in the human genome. 38:86–92.

684    O'Roak BJ, Deriziotis P, Lee C, Vives L, Schwartz JJ, Girirajan S, Karakoc E,

685    Mackenzie AP, Ng SB, Baker C, Rieder MJ, Nickerson D a, et al. 2011. Exome

686    sequencing in sporadic autism spectrum disorders identifies severe de novo mutations.

687    Nat Genet 43:585–9.

688    Quinlan AR, Clark RA, Sokolova S, Leibowitz ML, Zhang Y, Hurles ME, Mell JC, Hall

689    IM. 2010. Genome-wide mapping and assembly of structural variant breakpoints in the

690    mouse genome. Genome Res 623–635.

691    Sanders SJ, Murtha MT, Gupta AR, Murdoch JD, Raubeson MJ, Willsey AJ, Ercan-

692    Sencicek AG, DiLullo NM, Parikshak NN, Stein JL, Walker MF, Ober GT, et al. 2012. De

693    novo mutations revealed by whole-exome sequencing are strongly associated with

694    autism. Nature 485:237–241.

695 Szatkiewicz JP, O'Dushlaine C, Chen G, Chambert K, Moran JL, Neale BM, Fromer M,

696 Ruderfer D, Akterin S, Bergen SE, Kähler a, Magnusson PKE, et al. 2014. Copy

697 number variation in schizophrenia in Sweden. Mol Psychiatry 19:762–73.

698 Tattini L, D'Aurizio R, Magi A. 2015. Detection of genomic structural variants from next-

699 generation sequencing data. Front Bioeng Biotechnol 3:.

700 Wala J, Bandopadhayay P, Greenwald N, Rourke RO, Stewart C, Schumacher S, Li Y,

701 Weischenfeldt J, Nusbaum C, Campbell P, Meyerson M, Zhang Z. 2017. SvABA:

702 Genome-wide detection of structural variants and indels by local assembly. bioRxiv 1–

703 40.

704 Weischenfeldt J, Symmons O, Spitz F, Korbel JO. 2013. Phenotypic impact of genomic

705 structural variation□: insights from and for human disease. Nat Rev Genet 14:125–138.

706 Weisenfeld NI, Yin S, Sharpe T, Lau B, Hegarty R, Holmes L, Sogoloff B, Tabbaa D,

707 Williams L, Russ C, Nusbaum C, Lander ES, et al. 2014. Comprehensive variation

708 discovery in single human genomes. Nat Publ Gr 46:1350–1355.

709 Werling DM, Brand H, An J-Y, Stone MR, Glessner JT, Zhu L, Collings RL, Dong S,

710 Layer RM, Markenscoff-Papadimitriou E, Farrell A, Schwartz GB, et al. 2017. Limited

711 contribution of rare, noncoding variation to autism spectrum disorder from sequencing of

712 2,076 genomes in quartet families. bioRxiv 1–45.

713 Zook JM, Catoe D, Mcdaniel J, Vang L, Spies N, Sidow A, Weng Z, Liu Y, Mason CE,

714 Alexander N, Henaff E, Mcintyre ABR, et al. 2016. Data Descriptor□: Extensive

715 sequencing of seven human genomes to characterize benchmark reference materials.

716 Nature 1–26.
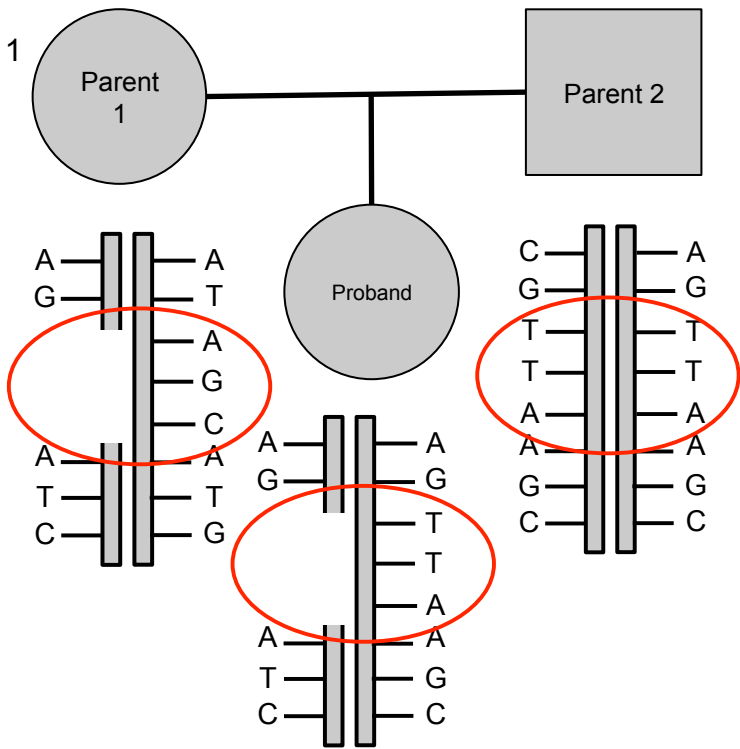
717

718    **Figure Legends**

719    **Figure 1: Schematic of MEM principle.** Diagram of trio where proband inherited a

720    deletion from parent 1. Tools report homozygous genotypes (red) that violate Mendelian

721    laws of segregation in the case of hemizygosity due to a heterozygous deletion.

722    Adapted from McCarroll *et al.* 2006 (McCarroll et al., 2006).

723

724    **Figure 2: MEM pipeline for WES and WGS data.**

725

726    **Figure 3:** A) Comparison of *de novo* deletions called by XHMM and MEM. B) Size

727    distribution of *de novo* deletions called by XHMM. Colors in stacked histogram indicate

728    which tools detected the deletion (red = MEM and XHMM detected, green = MEM

729    detected and not XHMM, blue = XHMM detected and not MEM).

730

731    **Figure 4: MEs plotted by chromosome** A) MEs in a trio after quality filtering. B)

732    Sample with UPD on chromosome 9. C) Trio with DNA contamination. D) Trio with a

733    sample mix.

734

735

736

737

738

739

Figure 1

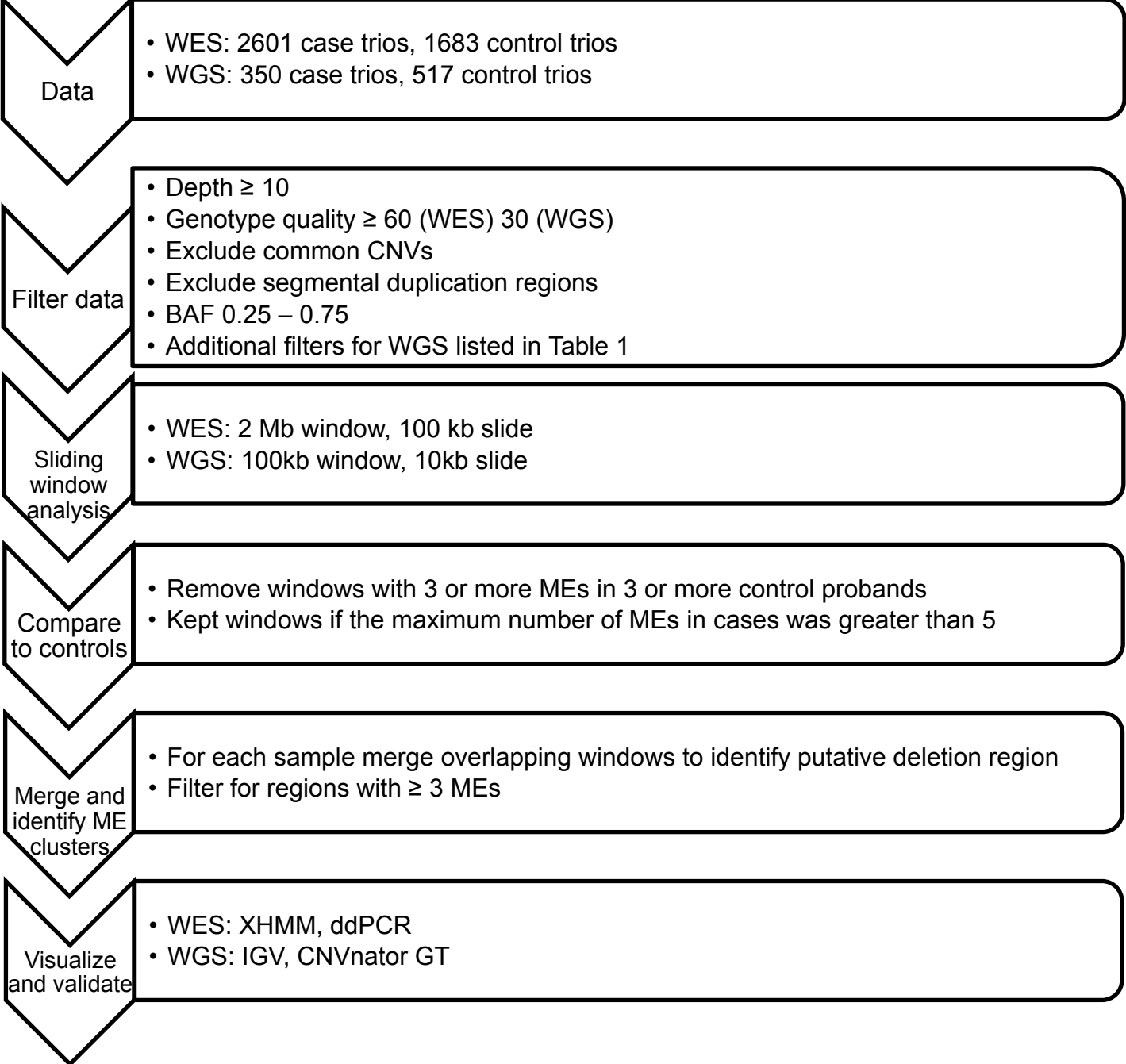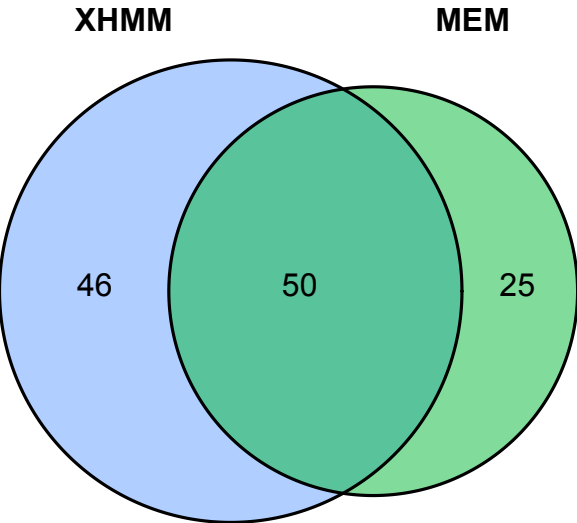| Parent 1 | Parent 2 | Expected in Child | Observed in Child |
|----------|----------|-------------------|-------------------|
| AA | TT | AT | TT |
| GG | TT | GT | TT |
| CC | AA | CA | AA |

Figure 2



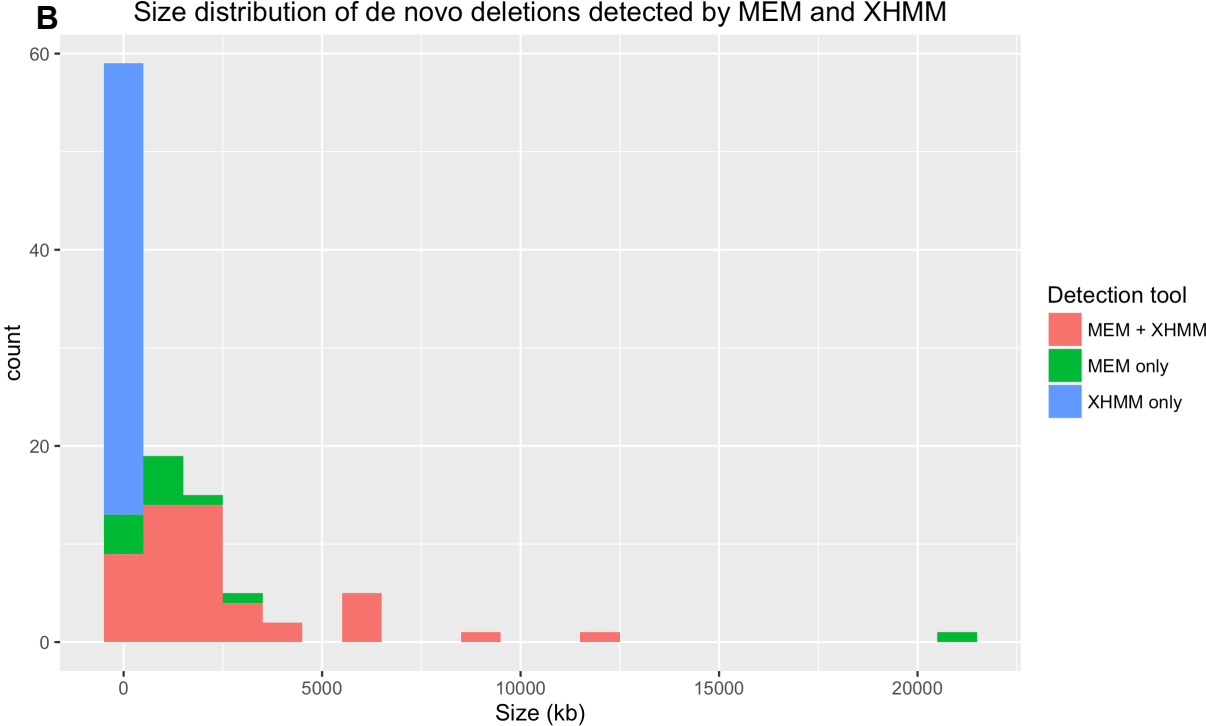| Data | • WES: 2601 case trios, 1683 control trios<br>• WGS: 350 case trios, 517 control trios |

| Filter data | • Depth ≥ 10<br>• Genotype quality ≥ 60 (WES) 30 (WGS)<br>• Exclude common CNVs<br>• Exclude segmental duplication regions<br>• BAF 0.25 – 0.75<br>• Additional filters for WGS listed in Table 1 |

| Sliding window analysis | • WES: 2 Mb window, 100 kb slide<br>• WGS: 100kb window, 10kb slide |

| Compare to controls | • Remove windows with 3 or more MEs in 3 or more control probands<br>• Kept windows if the maximum number of MEs in cases was greater than 5 |

| Merge and identify ME clusters | • For each sample merge overlapping windows to identify putative deletion region<br>• Filter for regions with ≥ 3 MEs |

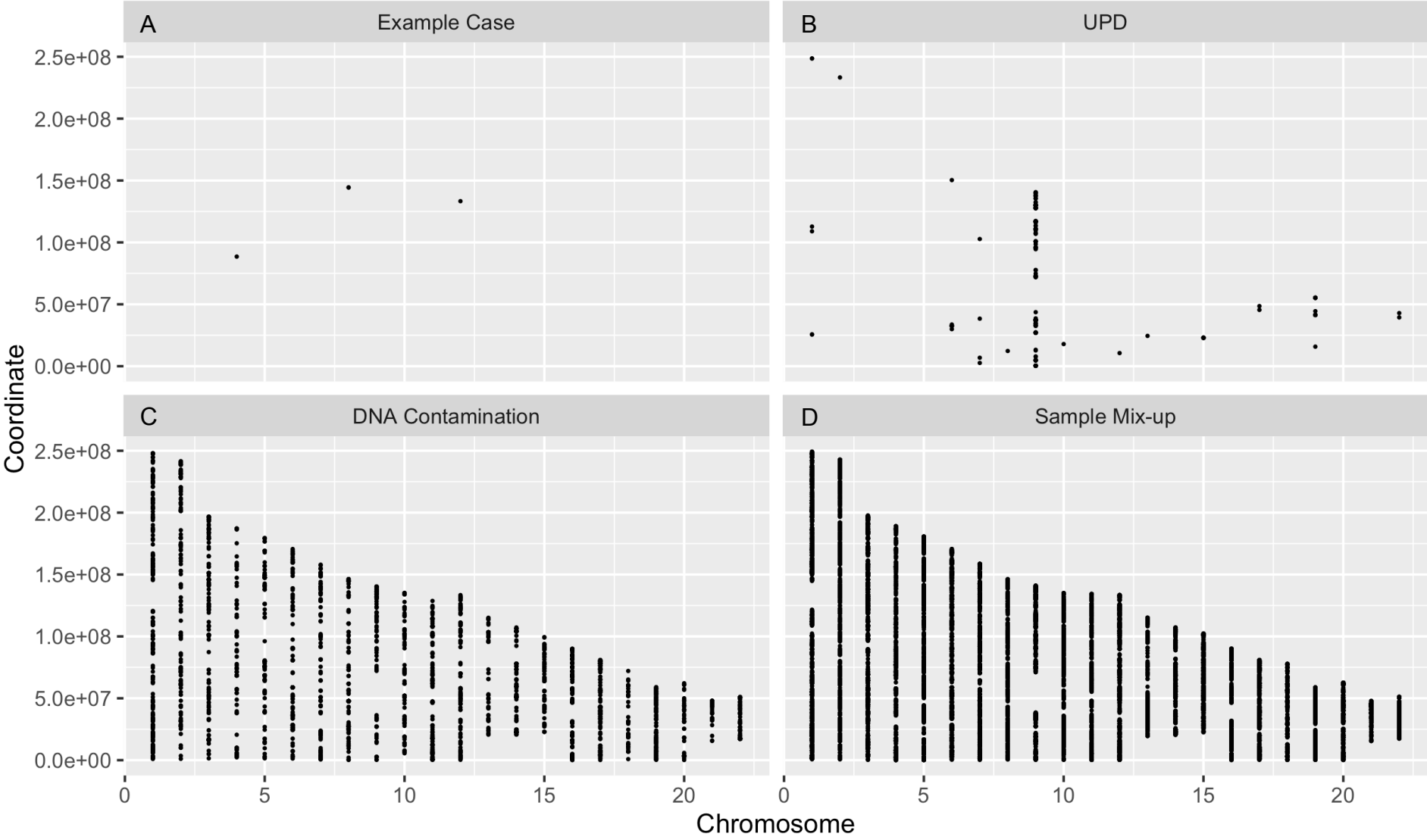| Visualize and validate | • WES: XHMM, ddPCR<br>• WGS: IGV, CNVnator GT |

Figure 3

Figure 4

**Tables**

**Table 1: Computational resources required for NGS CNV detection tools**

| Tool | Runtime (seconds) | Max Memory (MB) | Average Memory (MB) |
|---|---|---|---|
| MEM WES | 5.5 | 21 | 12 |
| XHMM | 453 | 278 | 81 |
| MEM WGS | 407 | 21 | 7 |
| CNVnator | 77,629 | 7,674 | 709 |
| Lumpy/SVTyper | 4,238 | 12,876 | 4,898 |