1 **Early post-zygotic mutations contribute to congenital heart disease**
2

3 Alexander Hsieh[1,*], Sarah U. Morton[2,3,*], Jon A.L. Willcox[3,*], Joshua M. Gorham[3], Angela C. Tai[3], Hongjian
4 Qi[1], Steven DePalma[3], David McKean[3], Emily Griffin[1], Kathryn B. Manheimer[4], Daniel Bernstein[5], Richard
5 W. Kim[6], Jane W. Newburger[2], George A. Porter Jr.[7], Deepak Srivastava[8], Martin Tristani-Firouzi[9],
6 Martina Brueckner[10], Richard P. Lifton[14], Elizabeth Goldmuntz[13], Bruce D. Gelb[4], Wendy K. Chung[1,#],
7 Christine E. Seidman[3,11,12#], J. G. Seidman[3,#], Yufeng Shen[1,#^]
8
9 [1] Columbia University Medical Center, New York, NY, USA
10 [2] Boston Children's Hospital, Boston, Massachusetts, USA.
11 [3] Harvard Medical School, Boston, Massachusetts, USA.
12 [4] Icahn School of Medicine at Mount Sinai, New York, New York, USA
13 [5] Stanford University, Palo Alto, California, USA.
14 [6] Children's Hospital Los Angeles, Los Angeles, California, USA.
15 [7] University of Rochester Medical Center, Rochester, New York, USA.
16 [8] Gladstone Institutes and University of California San Francisco, San Francisco, California, USA.
17 [9] University of Utah School of Medicine, Salt Lake City, Utah, USA.
18 [10] Yale University School of Medicine, New Haven, Connecticut, USA.
19 [11] Brigham and Women's Hospital, Boston, Massachusetts, USA
20 [12] Howard Hughes Medical Institute, Harvard University, Boston, Massachusetts, USA
21 [13] Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA
22 [14] Rockefeller University, New York, NY, USA
23
24 *,# Equal contribution.
25 ^Corresponding author. Email: ys2411@cumc.columbia.edu (Y.S.).
26
27
28 alh2194@cumc.columbia.edu (A.H.)
29 Sarah.Morton@childrens.harvard.edu (S.U.M.)
30 Jon_Willcox@hms.harvard.edu (J.A.L.W.)
31 jgorham@genetics.med.harvard.edu (J.M.G.)
32 angela_tai@hms.harvard.edu (A.C.T.)
33 hq2130@columbia.edu (H.Q.)
34 depalma@genetics.med.harvard.edu (S.D.)
35 dmckean@genetics.med.harvard.edu (D.M.)
36 eg2871@cumc.columbia.edu (E.G.)
37 kathryn.manheimer@icahn.mssm.edu (K.B.M.)
38 danb@stanford.edu (D.B.)
39 rikim@chla.usc.edu (R.W.K.)
40 jane.newburger@cardio.chboston.org (J.W.N.)
41 george_porter@urmc.rochester.edu (G.A.P.)
42 deepak.srivastava@gladstone.ucsf.edu (D.S.)
43 martin.tristani@utah.edu (M.T.F.)
44 martina.brueckner@yale.edu (M.B.)
45 rickl@mail.rockefeller.edu (R.P.L.)
46 Goldmuntz@email.chop.edu (E.G.)
47 bruce.gelb@mssm.edu (B.D.G.)
48 wkc15@cumc.columbia.edu (W.K.C.)
49 cseidman@genetics.med.harvard.edu (C.E.S.)
50 seidman@genetics.med.harvard.edu (J.G.S.)
51 ys2411@cumc.columbia.edu (Y.S.)
52

**Abstract**

*Background*

The contribution of somatic mosaicism, or genetic mutations arising after oocyte fertilization, to congenital heart disease (CHD) is not well understood. Further, the relationship between mosaicism in blood and cardiovascular tissue has not been determined.

*Results*

We developed a computational method, Expectation-Maximization-based detection of Mosaicism (EM-mosaic), to analyze mosaicism in exome sequences of 2530 CHD proband-parent trios. EM-mosaic detected 326 mosaic mutations in blood and/or cardiac tissue DNA. Of the 309 detected in blood DNA, 85/97 (88%) tested were independently confirmed, while 7/17 (41%) candidates of 17 detected in cardiac tissue were confirmed. MosaicHunter detected an additional 64 mosaics, of which 23/46 (50%) among 58 candidates from blood and 4/6 (67%) of 6 candidates from cardiac tissue confirmed. Twenty-five mosaic variants altered CHD-risk genes, affecting 1% of our cohort. Of these 25, 22/22 candidates tested were confirmed. Variants predicted as damaging had higher variant allele fraction than benign variants, suggesting a role in CHD. The frequency of mosaic variants above 10% mosaicism was 0.13/person in blood and 0.14/person in cardiac tissue. Analysis of 66 individuals with matched cardiac tissue available revealed both tissue-specific and shared mosaicism, with shared mosaics generally having higher allele fraction.

*Conclusions*

We estimate that ~1% of CHD probands have a mosaic variant detectable in blood that could contribute to cardiac malformations, particularly those damaging variants expressed at higher allele fraction compared to benign variants. Although blood is a readily-available DNA source, cardiac tissues analyzed contributed ~5% of somatic mosaic variants identified, indicating the value of tissue mosaicism analyses.

*Keywords*

Mosaic, Somatic, Congenital Heart Disease, Exome Sequencing

1

81  **Background**

82      Mosaicism results from somatic mutations that arise post-zygotically in an early embryonic cell,

83  resulting in two or more cell populations with distinct genotypes in the developing embryo {Biesecker

84  2013}. The developmental status of the early embryonic cell at the time of mutagenesis determines the

85  proportion of variant-carrying cells and the tissue distribution of these cells in the post-natal child {Acuna-

86  Hidalgo 2015}. While germline variants have a variant allele frequency (VAF) of 0.5, somatic mosaic

87  variants have a significantly lower VAF.

88      Post-zygotic mosaic mutations have been implicated in several diseases including non-malignant

89  developmental disorders such as overgrowth syndromes {Poduri 2013; Lindhurst 2012; Kurek 2016},

90  structural brain malformations {Poduri 2012; Jamuar 2014; Riviere 2012; Lee 2012}, epilepsy {Stosser

91  2018}, and autism spectrum disorder {Lim 2017; Krupp 2017; Freed 2016; Dou 2017}. Recent analyses

92  also identified mosaic variants in a cohort of patients with congenital heart disease (CHD) {Manheimer

93  2018}, but the prevalence of these was far less than germline variants (CHD) {Zaidi 2013; Homsy 2015;

94  Jin 2017; Zaidi 2017}.

95      Assessment of the frequency of mosaicism in human disease is confounded by technical issues,

96  including differences in sequencing depth, DNA sources, and variant assessment pipelines. Low levels of

97  mosaicism can escape the detection threshold of traditional sequencing methods with standard read

98  depths, while post-zygotic mutations with a higher percentage of affected cells are difficult to discriminate

99  from germline *de novo* mutations {Acuna-Hidalgo 2015}. All of these issues can lead to substantially

100  different conclusions. For example, analyses of mosaicism in autism spectrum disorder was recently

101  assessed from whole exome sequence (WES) data from whole blood DNA from 2506 families (proband,

102  parents and unaffected sibling; trios and quads) in the Simons Simplex Collection (SSC) {Fischbach

103  2010}.  The primary sequence data were analyzed by three groups; one that identified a protein-coding

104  somatic mosaic variant rate of 0.074 per individual {Freed 2016}, another that found a mosaic rate of

105  0.059 per individual {Lim 2017}, and a third group that reported a mosaic rate of 0.125 per individual

106  {Krupp 2017}. This disparity suggests the need for more systematic mosaic mutation detection methods

107  that account for dataset-specific confounding factors.

108    By contrast, analyses of affected tissues can improve the sensitivity and specificity of detection of

109    somatic mosaicism. In cancer, methods to detect these events, such as MuTect {Cibulskis 2013},

110    compare tumor and benign tissues from the same patient. Mosaicism has also been demonstrated from

111    the analyses of unpaired samples with cancer and other pathologies {Sun 2018; Huang 2017; Smith

112    2015} by the demonstration of variants in affected tissues that are absent from blood-derived DNA

113    {Symoens 2017; McDonald 2018}. With access to cardiac tissues from patients with CHD obtained during

114    surgical repair, we hypothesized that analyses of mosaicism in cardiac tissue might improve insights into

115    the causes of this common congenital anomaly. As many cardiomyocyte lineages share a mesodermic

116    origin with blood cells but exit the cell cycle during embryogenesis, we also sought to determine if

117    mosaicism in the heart exhibited distinct patterns of mosaicism with regard to variant frequency and allele

118    fractions.

119    In this study, we developed a computational method, EM-mosaic (Expectation-Maximization-

120    based detection of Mosaicism), to detect mosaic single nucleotide variants (SNVs) using WES of proband

121    and parent DNA. To optimize this method, we measured mosaic detection power as a function of

122    sequencing depth. We compared EM-mosaic and MosaicHunter {Huang 2014} to investigate mosaicism

123    in 2530 CHD proband-parent trios from the Pediatric Cardiac Genomics Consortium (PCGC) {Jin et al

124    2017}, using exome sequences derived from blood-derived DNA. We detected predicted deleterious

125    mosaic mutations in genes involved in known biological processes relevant to CHD or developmental

126    disorders in 1% of probands. The accuracy of these mosaic variant detection algorithms was assessed

127    using an independent re-sequencing method. We found that among high-confidence mosaic mutations in

128    CHD-relevant genes, likely-damaging variants tended to have higher VAF than likely-benign variants.

129    In parallel we assessed mosaicism by EM-mosaic and MosaicHunter in 70 discarded tissues from

130    several heart regions obtained from 66 probands who underwent cardiac surgical repairs. While VAF

131    varied significantly (>3 fold) between blood and cardiovascular tissue at about 60% of sites, in general

132    mosaic variants with high (>15%) VAF were more likely shared between blood and cardiac tissue than

133    variants with lower VAF.

134

135    **Results**

136   *High-accuracy detection of mosaic mutations in WES data using EM-mosaic*

137   We analyzed whole exome sequence (WES) data from 2530 CHD proband-parent trios {Homsy

138   2015; Jin 2017} **(Table S1)**. Among this cohort, 1205 probands had CHD with neurodevelopmental

139   disorders (NDD) and/or extracardiac manifestations (EM), 788 had isolated CHD at the time of

140   enrollment, 539 had undetermined NDD status due to young neonatal age at the time of enrollment, and

141   9 subjects had incomplete data (**Table S2**).

142   Previous WES analyses {Jin et al 2017} identified 1742 germline *de novo* SNVs among 838

143   cases with NDD and/or EM, 516 isolated cases, 644 cases of unknown NDD status, and 7 with

144   incomplete data. These *de novo* variants were identified using the Genome Analysis Toolkit (GATK)

145   pipeline {McKenna 2010; DePristo 2011} assuming a germline diploid model in which the expected VAF

146   is 0.5. This model has limited sensitivity to detect mosaic mutations for which the fraction of alternative

147   allele reads is significantly below 0.5, especially because *de novo* variants with VAF<0.2 were excluded

148   to reduce false discovery.

149   To efficiently capture mosaic variants with VAF<0.4, we developed a new method (EM-mosaic) to

150   detect mosaic variants in WES sequence of a proband and parents (trios). Potential mosaic variants were

151   identified in WES sequence data using SAMtools *mpileup* {Li 2009} with settings designed to capture

152   sites with VAF between 0.1-0.4 and merged with the variants found by the GATK pipeline {Jin et al 2017}

153   (**Fig S1**) to create a union variant set. To reduce the elevated false positive rate inherent in low-VAF calls,

154   we applied a set of empirical filters to remove likely technical artifacts due to sequencing errors

155   associated with repetitive and/or low complexity sequences. We then manually inspected *de novo* SNVs

156   with VAF<0.3 (n=582) using IGV and filtered out an additional 188 likely false positives. After

157   preprocessing and outlier removal, the remaining 2971 *de novo* SNVs were used as input to our mosaic

158   detection model.

159   Among the 2971 *de novo* SNVs, this pipeline identified 309 sites as candidate mosaics based on

160   posterior odds score **(Fig 1A-B, Table S3)**, including 50 sites that were previously reported as germline

161   *de novo* variants {Jin et al 2017}. An additional 86 sites were identified as having posterior odds below

162   our threshold of 10 but greater than 1 **(Fig S2A, S2B)**, including a *ZEB2* variant with posterior odds 4.7

163   that was previously confirmed via ddPCR {Manheimer 2018}. Among these 86 variants, 53 are likely

4

164    mosaic and 33 are likely germline **(Fig S2B)**. We chose not to include these sites since there was

165    insufficient evidence to confidently resolve them individually as mosaic or germline.

166

167    *Mosaic mutations found in blood derived DNA with MosaicHunter*

168         We also employed MosaicHunter, which uses a Bayesian genotyping algorithm with a series of

169    stringent filters (see Supplemental Methods) for discovering mosaic variants using WGS genotype

170    information from trios. {Huang 2017} Among the 2530 CHD trios, MosaicHunter identified an initial set of

171    58976 sites showing evidence of mosaicism, including 214 high-confidence variants located in coding

172    regions. (**Fig S3).** After applying a minimum likelihood ratio (LR) cutoff of 80 for distinguishing mosaic

173    from germline mutation, and additional heuristic filters **(Supplemental Methods)**, MosaicHunter identified

174    116 coding sites **(Table S4)** or 0.05 mosaics /individual.

175         Of the mosaic candidates detected by MosaicHunter, 58/116 (50%) were also identified by EM-

176    mosaic while 58/116 (50%) candidates were unique to MosaicHunter **(Table 1; Fig S4)**. Of the 58

177    candidates unique to MosaicHunter, 35 were filtered out by EM-mosaic on the basis of insufficient

178    alternate allele read support, 16 had a non-zero allelic depth in the parents, and 7 failed quality filters.

179    The 251 candidates unique to EM-mosaic were discarded by the MosaicHunter pipeline during BAM

180    reprocessing (n=13), quality filtering (n=146), application of LR cutoff (22), or were not called due to

181    inadequate read depth (n=70) **(Fig S3)**.

182

**Table 1: Mosaic variant detection by EM-Mosaic, MosaicHunter and validated by PCR product sequencing**

| | | Union | Shared | Unique | |
| --- | --- | --- | --- | --- | --- |
| | | | | EM-Mosaic | MosaicHunter |
| **Mosaic Variants (total)*** | | **315 (332)** | **56 (57)** | **218 (240)** | **29 (35)** |
| **Mosaic Candidates** | | 367 | 58 | 251 | 58 |
| **Mosaic Candidate VAF mean (SD)** | | 0.13 (0.06) | 0.12 (0.05) | 0.13 (0.06) | 0.10 (0.05) |
| **MiSeq Confirmation** | Total Tested | 143 | 22 | 75 | 46 |
| | Mosaic | 108 | 21 | 64 | 23 |
| | Germline | 3 | 0 | 3 | 0 |
| | No Variant | 32 | 1 | 8 | 23 |
| **Validation Rate** | | 76% | 95% | 85% | 50% |

183    *Estimated number of mosaic variants found among 2530 CHD probands (total number of mosaic variants detected
184    by EM-mosaic and MosaicHunter).

5

185

186

187    *Sequence confirmation of candidate mosaic variants and estimation of mosaicism in CHD*

188    From the 367 high-confidence EM-mosaic and/or MosaicHunter SNVs, we selected 143

189    candidates (97 identified by EM-mosaic; 68 identified by MosaicHunter) for experimental confirmation

190    using MiSeq amplicon resequencing (**Table S5; Table S11 and S12; Methods**). DNA fragments

191    encompassing the putative mosaic variant were PCR-amplified from proband and each parent DNA,

192    sequenced on an Illumina MiSeq next generation sequencer and VAF was calculated for each individual.

193    These candidate mosaics included SNVs on the extremes of the VAF spectrum, as well as mosaics that

194    were flagged by MosaicHunter quality filters. Candidates mosaic variants were considered confirmed by

195    MiSeq analyses if they demonstrated an amplicon VAF exceeding 0.01 but less than 0.45, so as to

196    indicate a variant of post-zygotic origin. MiSeq VAF values closely correlated with those originally

197    determined by exome sequencing ($P$= $2.2 \times 10^{-16}$; **Fig S5)**.

198    We confirmed 85/97 (88%) EM-mosaic candidate mosaic variants**.** Three candidate variants were

199    likely germline *de novo* SNVs (VAF>0.45). Nine candidate variants were 'false positives' that were neither

200    germline *de novo* SNVs or mosaic SNVs since either no variant reads were detected by MiSeq

201    sequencing of the proband amplicon, or the same small fraction of variants were detected in proband

202    amplicon and one parent's amplicon.

203    Parallel analyses with MosaicHunter confirmed 44/68 (65%) candidate mosaic variants. There

204    were 23 sites for which no variant reads were detected by MiSeq amplicon sequencing (MiSeq

205    VAF<0.001) or in which the same small fraction of variant reads was detected in the proband amplicon as

206    in one parent's amplicon.

207    We considered whether estimates of mosaic variant frequency were sensitive to whole exome

208    sequencing depth by calibrating estimates of mosaic detection power using properties of the sequence

209    data (average read depth, prior mosaic fraction, and the value of our overdispersion parameter *θ*) (**Fig**

210    **S6; Supplemental Methods)**. Our projected mosaic detection power curves demonstrated more than a

211    doubling of power to detect mosaic variants with VAF 0.2 as sequencing depth increases from 40x to 80x

212 **(Fig 1C)**. Projected mosaic detection power curves for less stringent mosaic cutoffs showed similar

213 increases of power with increasing sequencing depth **(Fig S8).**

214 　　　To estimate the 'true' frequency of mosaicism per blood DNA exome, independent of average

215 coverage detection power constraints, we estimated the 'true' mosaic count in a VAF range by multiplying

216 the number of mosaics by the inverse of the detection power for each VAF bin. Applying this method to

217 the 184 of 309 high-confidence EM-mosaic variants with VAF>0.1, we estimated the adjusted number of

218 mosaics with VAF>0.1 to be 361 **(Fig S8A).** Thus, the true frequency of coding mosaics in the blood

219 (0.4>VAF >0.1) is 0.14 variants per individual, representing a non-negligible class of mutations with

220 potential contribution to genetic risk for congenital heart disease. The estimated true mosaic frequency

221 does not change significantly when using less stringent mosaic definitions **(Figure S8)**. In sum, we

222 identified 315 blood mosaic variants in 2530 CHD probands or 0.13 mosaic variants per subject with a

223 mean VAF of 0.13±0.06. We do not anticipate that doubling the sequencing depth would change

224 significantly this estimate.

225

226 *Mosaic variants occurred most frequently at CpG sequences.*

227 　　　Previous studies demonstrated a strong preference for *de novo* C>T mutations at CpG

228 dinucleotides compared to other dinucleotides due to the spontaneous deamination of 5-methylcytosine

229 {Fryxell 2005; Francioli 2015}. We asked whether the germline *de novo* variants observed in CHD

230 probands and the 332 mosaic sites demonstrated a similar sequence preference **(Fig 1, Table 1, S3, S4)**.

231 Of the 2662 germline *de novo* mutations identified in 2530 CHD probands, 979 variants (37% of all

232 variants) involved mutation of the cytosine of a CpG dinucleotide (**Fig 2A**). By contrast, 99 (29% of all

233 mosaic SNVs) of 332 mosaic SNVs altered the cytosine of a CpG dinucleotide; significantly more than

234 expected by chance (2.2x above expectation; p=$2x10^{-15}$). These observations suggest that somatic *de*

235 *novo* mutations were 1.4-fold less likely to involve a CpG dinucleotide than germline *de novo* mutations in

236 CHD probands (*P*=0.01; **Fig 2B)**. Even ignoring the high CpG mutation frequency, cytosines and

237 guanines were ~2-fold more likely to be mutated than adenines or thymidines both for germline mutations

238 and for mosaic variants. Surprisingly, somatic mutations of A>C/T>G transversions in ApC dinucleotides

239 were ~2-fold greater than the corresponding germline mutations (*P*=$5x10^{-8}$; **Fig 2B**).

240

241   *Detection of mosaic mutations in CHD tissues*

242   Using EM-mosaic and MosaicHunter we analyzed exome sequences from 70 cardiac tissues

243   derived from 66 subjects with CHD (**Table S6)** and paired blood samples**.** Among 57 *de novo* variants

244   (allele depth approximately 0.5) that were previously identified in blood-derived DNA, 54 were also found

245   in CHD tissues. Of the 3 *de novo* variants not present in cardiac tissue, 1 was outside of the tissue WES

246   capture region and 2 occurred in a single proband (Table 2). In addition, 23 distinct candidate mosaic

247   variants were detected by EM-mosaic (n=13), MosaicHunter (n=6), or by both algorithms (n= 4). All 23

248   candidates were tested via MiSeq amplicon sequencing of blood and cardiac tissue DNAs; 15 of 23

249   unique candidate mosaics were confirmed  (**Table 2, S7),** including a *CCNC* variant that was identified in

250   two different CHD tissues from proband 1-01684. Ten (86%) confirmed mosaic variants were detected in

251   blood and cardiac tissues (MAF>0.01), four were found only in cardiac tissue, and one was found only in

252   blood. Of the 7 mosaics detected by blood WES analysis, 4 were confirmed in the corresponding cardiac

253   tissue sample.  Remarkably, five confirmed cardiac tissue mosaic variants occurred in one proband (1-

254   07004), one of which was also present in blood DNA.

255   These analyses indicate a frequency of coding mosaics (0.4>VAF >0.1) in the cardiac tissues of

256   0.14 per individual (9 of 66 probands), which approximated our estimate of 0.14 blood mosaics per

257   individual **(Fig S8A)**. Despite these similar frequencies, multiple distinct mosaic variants were identified in

258   these tissues. Mosaics with highest VAF were more likely to be found in both tissues (Mann-Whitney U

259   Test *P*=0.019), presumably indicating that the mutation occurred earlier in lineage development (**Fig S9)**.

260

261
262   **Table 2. Mosaics detected in individuals with matched cardiovascular tissue and blood**
263   **< Insert Table 2; see end of document >**
264   Characteristics of mosaic variants predicted for individuals with blood and cardiovascular tissue WES data available.
265   Among 15 mosaics, 5 were detected via analysis of blood WES, 8 were detected from cardiovascular tissue WES,
266   and 2 were detected by both approaches. Six of 7 (86%) mosaics detected from analysis of blood were present in
267   both DNA sources with MiSeq VAF≥0.01. Two additional variants previously identified as *de novo* germline variants in
268   blood WES were absent from CHD tissue WES. Minimum 1023 MiSeq reads used to determine VAF. Abbreviations:
269   AD, allelic depth (reference, alternate); AO, aorta; AtrSpt, atrial septum; Bmis, benign missense; Dmis, deleterious
270   missense; LOF, Loss of function variant; LV, left ventricle; RV, right ventricle; VAF, variant allele fraction.
271
272

273   *Blood and Cardiac Tissue Mosaics Likely to Contribute to CHD*

274     Our prior genetic studies of CHD studies showed that damaging *de novo* variants typically

275     occurred in genes highly expressed in the top quartile of the developing E9.5 mouse heart (HHE).{Zaidi

276     2013; Homsy 2015} or contributed to CHD in mouse models {Jin 2017}. Among the 342 mosaic variants

277     identified from blood or cardiac tissue analyses that were not false by MiSeq, 65 altered these HHE

278     and/or mouse CHD genes (n=4558, **Table S8**). RefSeq functional annotation predicted 52 variants as

279     likely-damaging variants (LOF, Dmis), and 46 as likely benign, missense **(Table S8, S9)**. In total, we

280     observed potentially CHD-causing mosaic mutations in 25 participants, representing 1% of the 2530 total

281     participants in our CHD cohort. Among these 25 mosaics, we confirmed 22/22 (100%) candidates tested

282     via MiSeq.  Notably, multiple likely-damaging mosaic variants altered genes (*ISL1*, *SETD2*, *NOVA2*,

283     *SMAD9*, *LZTR1*, *KCTD10*, *KCTD20*, *FZD5,* and *QKI* ) involved in key developmental pathways, which

284     may account for the extra-cardiac phenotypes observed in these patients (**Table 3, S10**). There was no

285     difference in the proportion of individuals with extracardiac features among those with damaging mosaic

286     variants compared to the overall cohort (11/25 vs 909/ 2521, *P*=0.68), and there was a wide range of

287     CHD subtypes.  Five subjects carried additional *de novo* LoF or Dmis variants (1-06216, *TYRP1*; 1-

288     04046, *KRT13*; 1-06677, *TRIP4;* 1-05011, *KDM5B*; 1-00018, *SBF1*) and 4 genes harbored *de novo* LoF

289     or Dmis variants other than those listed in Table 3 (*FBN1*; *PKD1*; *LZTR1*; *PIK3C2G*). No CNVs were

290     detected in these subjects, with the exception of 1-00192 (duplication at chr15:22062306-23062355; non-

291     overlapping with the *GLYR1* mosaic).

292     If mosaic variants were unrelated to CHD, we would expect similar allelic fractions between

293     mosaics with variants predicted as likely damaging or likely benign. However, we found that the allele

294     fraction of likely damaging variants was significantly higher (Mann-Whitney U Test *P*=0.001, **Fig 4A**).

295     Moreover, among mosaic variants in genes that are not included among HHE or mouse CHD genes, we

296     found no significant difference of allele fraction (*P*=0.985, **Fig 4B**). We repeated these analyses using

297     less stringent posterior odds cutoffs of 2 and 5 and found the same result **(Fig S10)**. Together these data

298     support our conclusion that at least some likely-damaging mosaic variants identified here contribute to

299     CHD.  These results were determined independently of MiSeq validation results.

300
301
302     **Table 3. Damaging Mosaics in CHD-relevant genes**

9

303 **< Insert Table 3; see end of document>**
304 There were 25 potentially-pathogenic mosaic mutations based on known gene function and patient phenotype. Some
305 of these probands have previously described rare LoF/Dmis variants, though none are likely pathogenic for CHD {Jin
306 2017}. Additionally, some genes were previously found to have LoF/Dmis variants among other individuals in this
307 CHD cohort. Abbreviations: ASD, atrial septal defect; BAV, bicuspid aortic valve; Dmis, deleterious missense;
308 episcore, haploinsufficiency score (percentile rank) {Han 2018}; Heart Exp, heart expression percentile rank; LoF,
309 loss-of-function; pLI, probability of loss-of-function intolerance {gnomAD}; PCGC, Pediatric Cardiac Genomics
310 Consortium; VAF, variant allele fraction; VSD, ventricular septal defect. *VAF refers to CHD tissue WES.
311
312
313

314 **Discussion**

315 Distinguishing mosaic mutations from constitutional mutations has both clinical management and

316 reproductive implications for proband and parents. Individuals with mosaic mutations are generally

317 clinically less severely affected for conditions that affect multiple parts of the body {Happle 1986; Wallis

318 1990; Cohn 1990; Etheridge 2011; Donkervoort 2015; Weinstein 2016}. Mutations that occur post-

319 zygotically should have no recurrence risk for the parents and could have a recurrence risk of less than

320 50% for the proband depending on gonadal involvement. This study is among the first investigations of

321 the role of post-zygotic mosaic mutations in CHD. We developed a new computational method to robustly

322 detect mosaic single nucleotide variants from blood WES data at standard read depth. Applying this

323 method to a cohort of 2530 CHD patients, we detected 309 high-confidence mosaics (with a confirmation

324 frequency of 88% in a subset of variants assessed) or 0.12 variants per proband. Sequencing of cardiac

325 tissue to greater depth identified an additional 8 mosaic variants that had not been detected in blood

326 WES, 6 of which are present in cardiac tissue but not blood. We found significantly more variants per

327 proband in cardiac tissue DNA (0.23 variants per proband) than in blood DNA (0.12 variants per proband;

328 p=0.02). While the increased numbers of mosaic variants in cardiac tissue DNA vs blood DNA may reflect

329 technical differences such as sequencing read depth of cardiac tissue DNA vs blood DNA, it is possible

330 that somatic variation occurs more frequently in cardiac tissue of CHD probands than in their blood.

331 Whether or not there are more cardiac tissue mosaic variants in CHD probands than blood DNA variants,

332 we found 10 mosaic variants among 66 CHD proband cardiac tissues with a higher VAF in tissue than in

333 blood **(Table 2)** and 5 variants among these individuals with a higher VAF in blood than in tissue.

334 In total, we observed potentially CHD-causing mosaic mutations in 25 participants, representing

335 1% of the 2530 total participants in our CHD cohort. Among these 25 mosaics, we confirmed 22/22

336 (100%) candidates tested. We found that in CHD-related genes, likely-damaging mosaic mutations have

337  significantly greater alternative allele fraction than likely-benign mosaics, suggesting that some of these

338  variants contribute to CHD. Comparison of blood and cardiovascular tissues demonstrated tissue-specific

339  mosaic variants, though those variants with a higher VAF were more likely to be shared between tissues.

340  Due to limitations of conventional clinical interpretation for both mosaic and constitutional CHD variants

341  **(Supplemental Methods)**, we cannot know with complete certainty which among these 25 variants is

342  pathogenic and instead propose that, among our detected mosaics, the 23 detected from blood WES

343  data provide an estimate of the disease-causing mosaics detectable in blood with standard exome-

344  sequencing read depth. Nine of these variants affect genes known to have a role in cardiac development:

345  *ISL1*, *SETD2*, *NOVA2*, *QKI*, *SMAD9*, *LZTR1*, *KCTD10*, *KCTD20*, and *FZD5*.

346  The mosaic LOF mutation in *ISL1* is likely to be the cause of CHD in participant 1-05095. *ISL1* is

347  a transcription factor essential to normal cardiac development that regulates expression of *NKX*, *GATA*,

348  and *TBX* family genes {Golzio 2012; Colombo 2018} and controls secondary heart field differentiation and

349  atrial septation {Colombo 2018; Briggs 2012}. *ISL1* deficiency has been shown to lead to severe CHD in

350  mice {Cai 2003; Golzio 2012}. Participant 1-05095 has an isolated atrial septal defect consistent with a

351  secondary heart field defect phenotype {Stevens 2010} and has no other previously reported damaging

352  germline variants in CHD-related genes.

353  Damaging germline *de novo* variants in CHD subjects are enriched in genes related to chromatin

354  modification and RNA processing {Homsy 2015; Jin 2017}. Three genes with damaging mosaic variants

355  discovered here have related functions. *SETD2* is a histone methyltransferase required for embryonic

356  vascular remodeling {Hu 2010}; it is both sensitive to haploinsufficiency and highly expressed in the heart

357  during development. *NOVA2* is a key alternative-splicing regulator involved in angiogenesis that has been

358  shown to disrupt vascular lumen formation when depleted {Giampietro 2015}. *QKI* encodes an RNA-

359  binding protein that regulates splicing, RNA export from the nucleus, protein translation, and RNA stability

360  {Lauriat 2008}. *QKI* is also highly expressed in the heart during development and has been shown to

361  cause CHD and other blood vessel defects in mice when dysregulated {Noveroske 2002}.

362  Other damaging mosaic variants affect processes known to be relevant to CHD. *SMAD9* is

363  involved in the TGF-beta signaling pathway. TGF-beta signaling plays a critical role in cardiac

364  development and cardiovascular physiology, leading to pulmonary arterial hypertension and cardiac

365    abnormalities in mice when dysregulated {Drake 2015; Soubrier 2013}. *LZTR1* encodes a member of the

366    BTB-Kelch superfamily that is highly expressed in the heart during development and has been associated

367    with Noonan {Yamamoto 2015; Ghedira 2017} and DiGeorge Syndromes {Kurahashi 1995}, both of which

368    are characterized by CHD. *KCTD10* binds to and represses the transcriptional activity of *TBX5* (T-box

369    transcription factor), which plays a dose-dependent role in the formation of cardiac chambers {Tong

370    2014}. *KCTD10* is highly expressed in the heart during development and has been shown to produce

371    CHD in mice when dysregulated {Ren 2014}. *KCTD20* is a positive regulator of Akt {Nawa 2013} also

372    highly expressed in the heart during development. *FZD5* is haploinsufficient and encodes a

373    transmembrane receptor involved in Wnt, mTOR, and Hippo signaling pathways and has been shown to

374    play a role in cardiac development {Dawson 2013}.

375        Finally, two mosaic variants found in cardiac tissue, genes encoding *RFX3* and *PIK3C2G,* may be

376    disease-relevant. *PIK3C2G* is a signaling kinase involved in cell proliferation, survival, and migration, as

377    well as oncogenic transformation and protein trafficking {OMIM: 609001; RefSeq}. The effects of

378    *PIK3C2G* haploinsufficiency during cardiac development has not been characterized. *RFX3* is a highly-

379    constrained ciliogenic transcription factor that leads to pronounced laterality defects {Rasmdell 2005 } and

380    disruption of *RFX3* leads to congenital heart malformations in mice {Lo 2011 MGI: 5560494}. Notably the

381    RFX3 LoF variant has a 4-fold higher VAF in cardiac tissue than in blood.

382        Several investigators, who studied cancer and diseases with cutaneous manifestations, proposed

383    that the VAF correlates with time of mutation acquisition and disease burden {Belickova 2016; Sallman

384    2016; Happle 1986}. In this study, we used VAF as a proxy for cellular percentage and mutational timing,

385    with increasing VAF corresponding to events occurring earlier in development. Thus, we assume that

386    CHD-causal mosaic events identified in blood-derived DNA occurred during or shortly after the

387    gastrulation process (3[rd] week of development) {Moorman 2003} in the mesodermal progenitor cells that

388    differentiate into both heart precursor cells (cardiogenic mesoderm) and blood precursor cells

389    (hemangioblasts). We found that in CHD-relevant genes, mosaic sites predicted to be damaging tended

390    to have higher VAF than sites predicted to be likely benign, consistent with the hypothesis that these

391    mutations arose early in fetal development and play significant roles in CHD. However, additional

392    functional studies are necessary to fully assess causality. .

393     Finally, we recognize that while our method is able to detect a large fraction of mosaic variants in

394     blood, our calibrated estimates for the true number of mosaics suggest there are a non-negligible number

395     of additional mutations that were not identified by our method. At our current average sequencing depth

396     of 60x, we have limited sensitivity in the low VAF (<0.05) range. To reliably identify these low allelic

397     fraction sites, ultra-deep sequencing will be critical to distinguishing true variants from noise. At 500x, we

398     estimate detection sensitivity for mosaic events at VAF 0.05 to be above 80%. We also recognize age-

399     related clonal hematopoiesis {Jaiswal 2014; Genovese 2014} as a potential confounding factor in somatic

400     mutation detection; however, our study cohort includes mostly pediatric cases and we did not observe

401     mosaic mutations in genes related to clonal expansion (e.g. *ASXL1, DNMT3A, TET2, JAK2*) nor did we

402     observe a relationship between proband age and mosaic rate **(Fig S11, S12)**, suggesting minimal impact

403     from this process.

404

405     **Conclusions**

406     This study is among the first investigations of the role of post-zygotic mosaic mutations in CHD.

407     Despite limitations in sequencing depth and sample type, EM-mosaic was able to detect 309 high-

408     confidence mosaics with resequencing confirmation in 88% of cases assessed. Using MosaicHunter, an

409     additional 64  candidate mosaic sites were identified, of which 23/46 (50%) candidates from blood DNA

410     and 4/6 (67%) from CHD tissue DNA validated. In total, we observed potentially CHD-causing mosaic

411     mutations in 25 participants, representing 1% of our CHD cohort, and propose that these 25 cases

412     provide an estimate of the disease-causing mosaics detectable in blood with standard exome-sequencing

413     read depth. Additionally, we found that in CHD-related genes, likely-damaging mosaics have significantly

414     greater alternative allele fraction than likely benign mosaics, suggesting that many of these variants

415     cause CHD and occurred early in development. In the subset of our cohort for which cardiovascular

416     tissue samples were available, we show that mosaics detected in blood can also be found in the disease-

417     relevant tissue and that, while the VAF for mosaic variants often differed between blood and

418     cardiovascular tissue DNA, variants with higher VAF were more likely to be shared between tissues.

419     Given current limitations in sequencing depth and on the availability of relevant tissues, particularly for

420     conditions impacting internal organs like the heart, the full extent of the role of mosaicism in many

421     diseases remains to be explored. However, as datasets containing larger numbers of blood and other

422     tissue samples sequenced at higher depths become increasingly available, we will be able to more fully

423     characterize the biological processes underlying post-zygotic mutation and, by extension, the contribution

424     of mosaicism to disease using the methods presented here.

425

426     **Methods**

427     *Samples and Sequencing Data*

428         We analyzed WES data from 2530 Congenital Heart Disease (CHD) proband-parents trio families

429     who were recruited as part of the Pediatric Cardiac Genomics Consortium (PCGC) study {Homsy 2015;

430     Jin 2017}. Genomic DNA from venous blood or saliva was captured using Nimblegen v.2 exome capture

431     reagent (Roche) or Nimblegen SeqCap EZ MedExome Target Enrichment Kit (Roche) followed by

432     Illumina DNA sequencing (paired-end, 2x75bp) {Jin 2017, Zaidi 2013}. Genomic DNA from 70 surgically-

433     discarded cardiovascular tissue samples (2-10mg) was isolated using DNeasy Blood & Tissue Kit

434     (QIAgen), then captured using xGen Exome Research Panel v1.0 reagent (IDT) followed by Illumina DNA

435     sequencing (paired-end, 2x75bp). Sequence reads were mapped to the hg19 human reference genome

436     with BWA-MEM and BAM files were further processed following GATK Best Practices, which included

437     duplication marking, indel realignment, and base quality recalibration steps. Blood and saliva samples

438     had sample average depth 60x and cardiovascular tissue samples had sample average depth 160x.

439

440     *De Novo Variant Calling and Annotation*

441         We processed our sample BAMs and called variants on a per-trio basis using SAMtools (v1.3.1-

442     42) and BCFtools (v1.3.1-174). Pileups were generated using samtools '*mpileup'* command with mapQ

443     20 and baseQ 13 to minimize the effect of poorly mapped reads on variant allele fraction, followed by

444     bcftools '*call'* using a cutoff of 1.1 for the posterior probability of the homozygous reference genotype

445     parameter (-p) to capture additional sites with variant allele fraction suggestive of post-zygotic origin that

446     would otherwise be excluded under the default threshold of 0.01. To identify *de novo* mutations from trio

447     VCF files, we selected sites with (i) a minimum of 6 reads supporting the alternate allele in the proband

448     and (ii) for both parents, a minimum depth of 10 reads and 0 alternate allele read support. Variants were

449    then annotated using ANNOVAR (v2017-07-17) to include information from refGene, gnomAD (March

450    2017), 1000 Genomes (August 2015), ExAC, genomicSuperDups, CADD (v1.3) COSMIC (v70), and

451    dbSNP (v147) databases, as well as pathogenicity predictions from a variety of established methods

452    included as part of the dbNSFP (v3.0a) database or generated in-house (MCAP, REVEL, MVP, MPC).

453    We used REVEL {Ionnidis 2016} to evaluate missense variant functional consequence, using the

454    recommended threshold of 0.5 corresponding to sensitivity of 0.754 and specificity of 0.891. We used

455    spliceAI {Jaganathan 2019} to predict the variant functional impact on splicing using the delta score

456    thresholds of 0.2 for likely pathogenic (high recall), 0.5 for pathogenic (recommended), and 0.8 for

457    pathogenic (high precision).  We considered sites predicted to be Likely Gene-Disrupting (LOF) (stopgain,

458    stoploss, frameshift indels, splice-site), Deleterious Missense (Dmis; nonsynonymous SNV with

459    REVEL>0.5), or splice-damaging (Benign Missense or synonymous SNV with delta score > 0.5) to be

460    damaging and likely disease causing. We considered sites predicted to be Synonymous (delta score ≤

461    0.5) or Benign missense (Bmis; nonsynonymous SNV with REVEL ≤ 0.5 and delta score ≤ 0.5) to be non-

462    damaging.

463

464    *Pre-processing and QC*

465         To reduce the number of low VAF technical artifacts introduced by our variant calling approach,

466    we pre-processed our variants using a variety of filters. We first excluded indels from further analysis, as

467    their downstream model parameter estimates were less stable than those of SNVs. We then filtered our

468    variant call set for rare heterozygous coding mutations (Minor Allele Frequency (MAF) ≤ $10^{-4}$ across all

469    populations represented in gnomAD and ExAC databases). To account for regions in the reference

470    genome that are likely to affect read-depth estimates, we removed variant sites found in regions of non-

471    unique mappability (score<1; 300bp), likely segmental duplication (score>0.95), and known low-

472    complexity {Li 2014}. We then excluded sites located in MUC and HLA genes and imposed a maximum

473    variant read depth threshold of 500. We used SAMtools PV4 to exclude sites with evidence of technical

474    issues using a cutoff of 1e-3 for baseQ Bias and Tail Distance Bias and a cutoff of 1e-6 for mapQ Bias.

475    To account for potential strand bias, we used an in-house script to flag sites that have either (1) 0

476    alternate allele read support on either the forward or reverse strand or (2) p<1e-3 and (Odds Ratio

477 (OR)<0.33 or OR>3) when applying a two-sided Fisher's Exact Test to compare proportions of reference

478 and alternate allele read counts on the forward and reverse strands. We also excluded sites with cohort

479 frequency>1%, as well as sites belonging to outlier samples (with abnormally high *de novo* SNV (dnSNV)

480 counts, cutoff = 8) and variant clusters (defined as sites with neighboring SNVs within 10bp). Finally, we

481 applied an FDR-based minimum $N_{alt}$ filtering step **(Fig S7)** to control for false positives caused purely by

482 sequencing errors.

483

484 *IGV Visualization of Low Allele Fraction de novo SNVs*

485       To reduce the impact of technical artifacts on model parameter estimation, we manually

486 inspected *de novo* SNVs with VAF<0.3 (n=558) using Integrative Genomics Viewer (v2.3.97) to visualize

487 the local read pileup at each variant across all members of a given trio family. We focused on the allele

488 fraction range 0.0-0.3 since this range is enriched for technical artifacts that could potentially impact

489 downstream parameter estimation. Sites were filtered out if (1) there are inconsistent mismatches in the

490 reads supporting the mosaic allele, (2) the site overlaps or is adjacent to an indel, (3) the site has low

491 MAPQ or is not Primary alignment, (4) there is evidence of technical bias (strand, read position, tail

492 distance), or (5) the site is mainly supported by soft-clipped reads.

493

494 *Expectation-Maximization to Estimate Prior Mosaic Fraction and Control FDR*

495       Current estimates for the fraction of *de novo* events occurring post-zygotically are unstable due to

496 differences in study factors such as variant calling methods, average sequencing depth, and paternal

497 ages. In order to use this fraction as a prior probability in our posterior odds and false discovery

498 calculations, we reason that this value must be estimated from the data itself. We used an expectation-

499 maximization algorithm to jointly estimate the fraction of mosaics among apparent *de novo* mutations and

500 to calculate a per-site likelihood ratio score. This initial mosaic fraction estimate gives us a prior

501 probability of mosaicism, independent of sequencing depth or variant caller, and allows us to calculate for

502 each variant the posterior odds that a given site is mosaic rather than germline. To control for false

503 discovery among our predicted mosaic candidates, we chose a posterior odds threshold of 10 to restrict

504 FDR to 9.1%.

505

506    *Mosaic Mutation Detection Model*

507    To distinguish variant sites that show evidence of mosaicism from germline heterozygous sites,

508    we modeled the number of reads supporting the variant allele ($N_{alt}$) as a function of the total site depth

509    ($N$). In the typical case, $N_{alt}$ follows a binomial model with parameters $N$ = site depth and $p$ = mean VAF.

510    However, we observed notable overdispersion in the distribution of variant allele fraction compared to the

511    expectations under this binomial model **(Fig S6)**. To account for this overdispersion, we instead modeled

512    $N_{alt}$ using a beta-binomial distribution {Heinrich 2012; Ramu 2013}. We estimated an overdispersion

513    parameter $\theta$ for our model as follows: for site depth values $N$ in the range 1 to 500, we (1) bin variants by

514    identifying all sites with depth $N$, (2) calculate a maximum-likelihood estimate $\theta$ value using $N$ and all $N_{alt}$

515    values for variants in a given bin, and (3) estimate a global $\theta$ value by taking the average of $\theta$ values

516    across all bins, weighted by the number of variants in each bin. We then used $\theta$ in our expectation-

517    maximization approach to jointly estimate prior mosaic fraction and to calculate per-site likelihood ratios.

518    To calculate the posterior odds that a given variant arose post-zygotically, we first calculated a

519    likelihood ratio (LR) of two models: $M_0$: germline heterozygous variant, and $M_1$: mosaic variant. Under our

520    null model $M_0$, we calculated the probability of observing $N_{alt}$ from a beta-binomial distribution with site

521    depth $N$, observed mean germline VAF $p$, and overdispersion parameter $\theta$. Under our alternate model $M_1$,

522    we calculated the probability of observing $N_{alt}$ from a beta-binomial distribution with site depth $N$,

523    observed site VAF $p=N_{alt}/N$, and overdispersion parameter $\theta$. Finally, for each variant, we calculated LR

524    by using the ratio of probabilities under each model and posterior odds by multiplying LR by our E-M

525    estimated prior mosaic fraction estimate. We defined mosaic sites as those with posterior odds greater

526    than 10 (corresponding to 9.1% FDR). We used posterior odds in this context to be able to control for

527    false discovery, but we output similarly valid p-value and likelihood ratio scores for each *de novo* SNV.

528

529    *Mutation Confirmation by MiSeq Amplicon Sequencing*

530    Chromosome coordinates were expanded 500 bp upstream and downstream of the candidate

531    mosaic variants in the UCSC Genome Browser. Primer 3 Plus software was used to design forward and

532    reverse primers to generate 150-300 bp amplimers containing the candidate site. PCR reactions

17

533    consisting of genomic DNA, primers, and Phusion polymerase were amplified by thermal cycling and

534    purified with AMPure XP beads. The purified PCR product was quantified, and 0.5-1.0 ng of product was

535    used to construct Nextera XT libraries according to the protocol published by Illumina. Libraries were

536    purified using AMPure XP beads, and final libraries were quantified and pooled to undergo sequencing

537    through Illumina MiSeq.

538    We experimentally tested for the presence our predicted post-zygotic sites in the original blood

539    DNA and cardiovascular tissue DNA samples using Illumina MiSeq Amplicon sequencing. The Amplicon

540    Deep Sequencing workflow, optimized for the detection of somatic mutations in tumor samples, offers

541    ultra-high sequencing depth (>1000x) that gives us the resolution to confirm low VAF variants, accurately

542    estimate site VAF, and to distinguish true variant calls from technical artifacts. Mosaic candidates were

543    considered validated if the variant allele matched the MiSeq call and both the mosaic VAF and MiSeq

544    VAF indicated post-zygotic origin (VAF<0.45).

545    Mosaic candidates were selected for confirmation on the basis of VAF, plausible involvement in

546    CHD (based on predicted pathogenicity and HHE status), and detection method **(Table S11; Table S12)**.

547    We sampled mosaics from both ends of the VAF spectrum to evaluate our ability to distinguish high VAF

548    mosaics (VAF>0.2; n=29) from germline variants and to distinguish low VAF mosaics (VAF<=0.1; n=52)

549    from technical artifacts.  Confirmation rate across different VAF bins is shown in **Figure S13**.  We also

550    selected for confirmation mosaics detected uniquely by either EM-mosaic or MosaicHunter, for the sake

551    of method comparison **(Table 1)**.

552    To examine a potential source of bias in our candidate selection process, we compared the

553    posterior odds distribution of selected candidate mosaics (n=97) against those not chosen (n=212). We

554    found that our tested candidates had lower posterior odds than untested mosaics (mean$_{tested}$=5.382,

555    mean$_{untested}$=7.050, log$_{10}$-scale; Mann Whitney *U P*=0.002) **(Fig 14)**, suggesting that our validation rate is

556    not buoyed by testing variants with the strongest evidence of mosaicism.  For method development

557    purposes, we intentionally focused on mosaics with lower posterior odds as these fall in the VAF range

558    for which it is most difficult to distinguish germline from mosaic.

559

560    *Investigating the relationship between VAF and pathogenicity*

561    We hypothesized that mosaic contribution to disease is positively correlated with cellular

562    percentage and by extension mutational timing. Here, we used variant allele fraction as a proxy for

563    cellular percentage. We grouped mosaics into likely-damaging and likely-benign and compared the

564    distribution of allele fraction in CHD-related genes. We defined likely-damaging variants as: (a) likely

565    gene-disrupting (LOF) variants (including premature stop-gain, frameshifting, and variants located in

566    canonical splice sites), (b) missense variants predicted to be damaging by REVEL {Ioannidis 2016} (with

567    score ≥ 0.5) or (c) missense variants and synonymous predicted to be splice-damaging by spliceAI (with

568    score > 0.5). One of the main findings from previous CHD studies is that damaging *de novo* variants in

569    genes highly expressed in the developing heart ("HHE", ranked in the top 25% by cardiac expression data

570    in mouse at E14.5 {Zaidi 2013; Homsy 2015}) contribute to non-isolated CHD cases that have additional

571    congenital anomalies or neurodevelopmental disorders. Therefore, we considered the union of HHE

572    genes and known candidate CHD genes {Jin 2017} as CHD-related genes (n=4558). For mosaics in

573    CHD-related genes and for mosaics in other genes, we used a Mann-Whitney *U* Test to compare the VAF

574    distributions of likely-damaging and likely-benign groups.

575

576    *Estimated contribution of mosaicism to CHD*

577    We identified likely disease-causing mosaic mutations on the basis of predicted pathogenicity and

578    presence in genes involved in biological processes relevant to CHD or developmental disorders. Each

579    mosaic mutation was annotated with gene-specific information, including heart expression percentile,

580    probability of loss-of-function intolerance (pLI) score {Lek 2016}, whether dysregulation causes CHD in

581    mice {Smith 2018; Finger 2017}, and gene function {NCBI RefSeq}. We focused on HHE genes, genes

582    with high pLI (pLI>0.9), genes that cause CHD phenotypes in mice, and genes involved in key

583    developmental processes such as Wnt, mTOR, and TGF-beta signaling pathways. Then, for each patient,

584    we used the clinical phenotype to further prioritize mosaic mutations most likely contributing to that

585    individual's clinical features. Detailed mutation annotation and clinical phenotypes for the mosaic carriers

586    described above can be found in **Table S10.** We estimate the contribution of mosaicism to CHD as the

587    percentage of individuals carrying likely disease-causing mosaic mutations among all individuals in our

588    CHD cohort.

19

589

590

591

**Abbreviations**

593 ASD: Autism Spectrum Disorder

594 CHD: Congenital Heart Disease

595 dnSNV: *de novo* SNV

596 Dmis: Deleterious Missense mutation

597 $DP_{site}$: Total Read Depth at a variant site

598 $DP_{sample}$: Sample-wide Average Read Depth

599 ExAC: Exome Aggregation Consortium

600 FDR: False Discovery Rate

601 gnomAD: Genome Aggregation Database

602 HHE: High Heart Expression

603 LOF: Loss-of-Function

604 LR: Likelihood Ratio

605 MAF: Minor Allele Fraction

606 N: Total Read Depth

607 $N_{alt}$: Alternate Allele Read Depth

608 OR: Odds Ratio

609 PCGC: Pediatric Cardiac Genomics Consortium

610 pLI: Probability of Loss-of-Function Intolerance

611 PV4: P-value for strand bias, baseQ bias, mapQ bias and tail distance bias (SAMtools)

612 SNV: Single Nucleotide Variant

613 VAF: Variant Allele Fraction

614 WES: Whole Exome Sequencing

615

616

617

618                                        **References**

619 Acuna-Hidalgo, R., Bo, T., Kwint, M. P., van de Vorst, M., Pinelli, M., Veltman, J. A., … Gilissen, C.

620       (2015). Post-zygotic Point Mutations Are an Underrecognized Source of De Novo Genomic

621       Variation. *The American Journal of Human Genetics*, *97*(1), 67–74.

622       http://doi.org/10.1016/J.AJHG.2015.05.008

623 Belickova, M., Vesela, J., Jonasova, A., Pejsova, B., Votavova, H., Merkerova, M. D., … Cermak, J.

624       (2016). TP53 mutation variant allele frequency is a potential predictor for clinical outcome of

625       patients with lower-risk myelodysplastic syndromes. *Oncotarget*, *7*(24), 36266–36279.

626       http://doi.org/10.18632/oncotarget.9200

627 Biesecker, L. G., & Spinner, N. B. (2013). A genomic view of mosaicism and human disease. *Nature*

628       *Reviews Genetics*, *14*(5), 307–320. http://doi.org/10.1038/nrg3424

629 Briggs, L. E., Kakarla, J., & Wessels, A. (2012). The pathogenesis of atrial and atrioventricular

630       septal defects with special emphasis on the role of the dorsal mesenchymal protrusion.

631       *Differentiation*, *84*(1), 117–130. http://doi.org/10.1016/j.diff.2012.05.006

632 Cai, C.-L., Liang, X., Shi, Y., Chu, P.-H., Pfaff, S. L., Chen, J., & Evans, S. (2003). Isl1 identifies a

633       cardiac progenitor population that proliferates prior to differentiation and contributes a majority

634       of cells to the heart. *Developmental Cell*, *5*(6), 877–89. Retrieved from

635       http://www.ncbi.nlm.nih.gov/pubmed/14667410

636 Cibulskis, K., Lawrence, M. S., Carter, S. L., Sivachenko, A., Jaffe, D., Sougnez, C., … Getz, G.

637       (2013). Sensitive detection of somatic point mutations in impure and heterogeneous cancer

638       samples. *Nature Biotechnology*, *31*(3), 213–219. http://doi.org/10.1038/nbt.2514

639 Cohn, D. H., Starman, B. J., Blumberg, B., & Byers, P. H. (1990). Recurrence of lethal osteogenesis

640       imperfecta due to parental mosaicism for a dominant mutation in a human type I collagen gene

641       (COL1A1). *American Journal of Human Genetics*, *46*(3), 591–601. Retrieved from

642       http://www.ncbi.nlm.nih.gov/pubmed/2309707

643 Colombo, S., de Sena-Tomás, C., George, V., Werdich, A. A., Kapur, S., MacRae, C. A., & Targoff,

644       K. L. (2017). *nkx* genes establish SHF cardiomyocyte progenitors at the arterial pole and

645    pattern the venous pole through Isl1 repression. *Development*, dev.161497.

646    http://doi.org/10.1242/dev.161497

647    Dawson, K., Aflaki, M., & Nattel, S. (2013). Role of the Wnt-Frizzled system in cardiac

648    pathophysiology: A rapidly developing, poorly understood area with enormous potential.

649    *Journal of Physiology*, *591*(6), 1409–1432. http://doi.org/10.1113/jphysiol.2012.235382

650    DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V, Maguire, J. R., Hartl, C., … Daly, M. J.

651    (2011). A framework for variation discovery and genotyping using next-generation DNA

652    sequencing data. *Nature Genetics*, *43*(5), 491–8. http://doi.org/10.1038/ng.806

653    Donkervoort, S., Hu, Y., Stojkovic, T., Voermans, N. C., Foley, A. R., Leach, M. E., … Bönnemann,

654    C. G. (2015). Mosaicism for Dominant Collagen 6 Mutations as a Cause for Intrafamilial

655    Phenotypic Variability. *Human Mutation*, *36*(1), 48–56. http://doi.org/10.1002/humu.22691

656    Dou, Y., Yang, X., Li, Z., Wang, S., Zhang, Z., Ye, A. Y., … Wei, L. (2017). Postzygotic single-

657    nucleotide mosaicisms contribute to the etiology of autism spectrum disorder and autistic traits

658    and the origin of mutations. *Human Mutation*, *38*(8), 1002–1013.

659    http://doi.org/10.1002/humu.23255

660    Drake, K. M., Comhair, S. A., Erzurum, S. C., Tuder, R. M., & Aldred, M. A. (2015). Endothelial

661    Chromosome 13 Deletion in Congenital Heart Disease–associated Pulmonary Arterial

662    Hypertension Dysregulates SMAD9 Signaling. *American Journal of Respiratory and Critical*

663    *Care Medicine*, *191*(7), 850–854.

664    Etheridge, S. P., Bowles, N. E., Arrington, C. B., Pilcher, T., Rope, A., Wilde, A. A. M., … Tristani-

665    Firouzi, M. (2011). Somatic mosaicism contributes to phenotypic variation in Timothy

666    syndrome. *American Journal of Medical Genetics Part A*, *155*(10), 2578–2583.

667    http://doi.org/10.1002/ajmg.a.34223

668    Finger, J. H., Smith, C. M., Hayamizu, T. F., McCright, I. J., Xu, J., Law, M., … Ringwald, M. (2017).

669    The mouse Gene Expression Database (GXD): 2017 update. *Nucleic Acids Research*, *45*(D1),

670    D730–D736. http://doi.org/10.1093/nar/gkw1073

671     Fischbach, G. D., & Lord, C. (2010). The Simons Simplex Collection: A Resource for Identification of

672         Autism Genetic Risk Factors. *Neuron*, *68*(2), 192–195.

673         http://doi.org/10.1016/j.neuron.2010.10.006

674     Francioli, L. C., Polak, P. P., Koren, A., Menelaou, A., Chun, S., Renkens, I., … Sunyaev, S. R.

675         (2015). Genome-wide patterns and properties of de novo mutations in humans. *Nature*

676         *Genetics*, *47*(7), 822–826. https://doi.org/10.1038/ng.3292

677     Freed, D., & Pevsner, J. (2016). The Contribution of Mosaic Variants to Autism Spectrum Disorder.

678         *PLOS Genetics*, *12*(9), e1006245. http://doi.org/10.1371/journal.pgen.1006245

679     Fryxell, K. J., & Moon, W.-J. (2005). CpG Mutation Rates in the Human Genome Are Highly

680         Dependent on Local GC Content. *Molecular Biology and Evolution*, *22*(3), 650–658.

681         https://doi.org/10.1093/molbev/msi043

682     Genovese, G., Kähler, A. K., Handsaker, R. E., Lindberg, J., Rose, S. A., Bakhoum, S. F., …

683         McCarroll, S. A. (2014). Clonal Hematopoiesis and Blood-Cancer Risk Inferred from Blood

684         DNA Sequence. *New England Journal of Medicine*, *371*(26), 2477–2487.

685         http://doi.org/10.1056/NEJMoa1409405

686     Ghedira, N., Kraoua, L., Lagarde, A., Abdelaziz, R. Ben, Olschwang, S., Desvignes, J. P., … Mrad,

687         R. (2017). Further Evidence for the Implication of LZTR1, a Gene not Associated with the Ras-

688         Mapk Pathway, in the Pathogenesis of Noonan Syndrome. *Biology and Medicine*, *09*(06), 4–7.

689         http://doi.org/10.4172/0974-8369.1000414

690     Giampietro, C., Deflorian, G., Gallo, S., Di Matteo, A., Pradella, D., Bonomi, S., … Ghigna, C.

691         (2015). The alternative splicing factor Nova2 regulates vascular development and lumen

692         formation. *Nature Communications*, *6*, 1–15. http://doi.org/10.1038/ncomms9479

693     Golzio, C., Havis, E., Daubas, P., Nuel, G., Babarit, C., Munnich, A., … Etchevers, H. C. (2012).

694         ISL1 Directly Regulates FGF10 Transcription during Human Cardiac Outflow Formation. *PLoS*

695         *ONE*, *7*(1), e30677. http://doi.org/10.1371/journal.pone.0030677

696     Happle, R. (1986). The McCune-Albright syndrome: a lethal gene surviving by mosaicism. *Clinical*

697         *Genetics*, *29*(4), 321–4. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/3720010

698    Happle, R. (1993). Mosaicism in human skin. Understanding the patterns and mechanisms.

699         *Archives of Dermatology*, *129*(11), 1460–70. Retrieved from

700              http://www.ncbi.nlm.nih.gov/pubmed/8239703

701    Heinrich, V., Stange, J., Dickhaus, T., Imkeller, P., Krüger, U., Bauer, S., … Krawitz, P. M. (2012).

702         The allele distribution in next-generation sequencing data sets is accurately described as the

703         result of a stochastic branching process. *Nucleic Acids Research*, *40*(6), 2426–2431.

704              http://doi.org/10.1093/nar/gkr1073

705    Homsy, J., Zaidi, S., Shen, Y., Ware, J. S., Samocha, K. E., Karczewski, K. J., … Chung, W. K.

706         (2015). De novo mutations in congenital heart disease with neurodevelopmental and other

707         congenital anomalies. *Science*, *350*(6265), 1262–1266. http://doi.org/10.1126/science.aac9396

708    Hu, M., Sun, X.-J., Zhang, Y.-L., Kuang, Y., Hu, C.-Q., Wu, W.-L., … Chen, Z. (2010). Histone H3

709         lysine 36 methyltransferase Hypb/Setd2 is required for embryonic vascular remodeling.

710         *Proceedings of the National Academy of Sciences*, *107*(7), 2956–2961.

711              http://doi.org/10.1073/pnas.0915033107

712    Huang, A. Y., Zhang, Z., Ye, A. Y., Dou, Y., Yan, L., Yang, X., … Wei, L. (2017). MosaicHunter:

713         accurate detection of postzygotic single-nucleotide mosaicism through next-generation

714         sequencing of unpaired, trio, and paired samples. *Nucleic Acids Research*, *45*(10), e76–e76.

715              http://doi.org/10.1093/nar/gkx024

716    Ioannidis, N. M., Rothstein, J. H., Pejaver, V., Middha, S., McDonnell, S. K., Baheti, S., … Sieh, W.

717         (2016). REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare Missense

718         Variants. *The American Journal of Human Genetics*, *99*(4), 877–885.

719              http://doi.org/10.1016/j.ajhg.2016.08.016

720    Jaiswal, S., Fontanillas, P., Flannick, J., Manning, A., Grauman, P. V., Mar, B. G., … Ebert, B. L.

721         (2014). Age-Related Clonal Hematopoiesis Associated with Adverse Outcomes. *New England*

722         *Journal of Medicine*, *371*(26), 2488–2498. http://doi.org/10.1056/NEJMoa1408617

723    Jamuar, S. S., Lam, A.-T. N., Kircher, M., D'Gama, A. M., Wang, J., Barry, B. J., … Walsh, C. A.

724         (2014). Somatic Mutations in Cerebral Cortical Malformations. *New England Journal of*

725         *Medicine*, *371*(8), 733–743. http://doi.org/10.1056/NEJMoa1314432

726    Jin, S. C., Homsy, J., Zaidi, S., Lu, Q., Morton, S., DePalma, S. R., … Brueckner, M. (2017).

727        Contribution of rare inherited and de novo variants in 2,871 congenital heart disease probands.

728        *Nature Genetics*, *49*(11), 1593–1601. http://doi.org/10.1038/ng.3970

729    Krupp, D. R., Barnard, R. A., Duffourd, Y., Evans, S. A., Mulqueen, R. M., Bernier, R., … O'Roak, B.

730        J. (2017). Exonic Mosaic Mutations Contribute Risk for Autism Spectrum Disorder. *The*

731        *American Journal of Human Genetics*, *101*(3), 369–390.

732        http://doi.org/10.1016/j.ajhg.2017.07.016

733    Kurahashi, H., Akagi, K., Inazawa, J., Ohta, T., Niikawa, N., Kayatani, F., … Nishisho, I. (1995).

734        Isolation and characterization of a novel gene deleted in DiGeorge syndrome. *Human*

735        *Molecular Genetics*, *4*(4), 541–9. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/7633402

736    Kurek, K. C., Luks, V. L., Ayturk, U. M., Alomari, A. I., Fishman, S. J., Spencer, S. A., … Warman,

737        M. L. (2012). Somatic Mosaic Activating Mutations in PIK3CA Cause CLOVES Syndrome. *The*

738        *American Journal of Human Genetics*, *90*(6), 1108–1115.

739        http://doi.org/10.1016/j.ajhg.2012.05.006

740    Lauriat, T. L., Shiue, L., Haroutunian, V., Verbitsky, M., Ares, M., Ospina, L., & McInnes, L. A.

741        (2008). Developmental expression profile ofquaking, a candidate gene for schizophrenia, and

742        its target genes in human prefrontal cortex and hippocampus shows regional specificity.

743        *Journal of Neuroscience Research*, *86*(4), 785–796. http://doi.org/10.1002/jnr.21534

744    Lee, J. H., Huynh, M., Silhavy, J. L., Kim, S., Dixon-Salazar, T., Heiberg, A., … Gleeson, J. G.

745        (2012). De novo somatic mutations in components of the PI3K-AKT3-mTOR pathway cause

746        hemimegalencephaly. *Nature Genetics*, *44*(8), 941–945. http://doi.org/10.1038/ng.2329

747    Lek, M., Karczewski, K. J., Minikel, E. V., Samocha, K. E., Banks, E., Fennell, T., … Consortium, E.

748        A. (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature*, *536*(7616),

749        285–291. http://doi.org/10.1038/nature19057

750    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … 1000 Genome Project Data

751        Processing Subgroup. (2009). The Sequence Alignment/Map format and SAMtools.

752        *Bioinformatics*, *25*(16), 2078–2079. http://doi.org/10.1093/bioinformatics/btp352

753      Lim, E. T., Uddin, M., De Rubeis, S., Chan, Y., Kamumbu, A. S., Zhang, X., … Walsh, C. A. (2017).

754      Rates, distribution and implications of postzygotic mosaic mutations in autism spectrum

755      disorder. *Nature Neuroscience*, *20*(9), 1217–1224. http://doi.org/10.1038/nn.4598

756      Manheimer, K. B., Richter, F., Edelmann, L. J., D'Souza, S. L., Shi, L., Shen, Y., … Gelb, B. D.

757      (2018). Robust identification of mosaic variants in congenital heart disease. *Human Genetics*,

758      *137*(2), 183–193. http://doi.org/10.1007/s00439-018-1871-6

759      McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., … DePristo, M. A.

760      (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation

761      DNA sequencing data. *Genome Research*, *20*(9), 1297–303.

762      http://doi.org/10.1101/gr.107524.110

763      Moorman, A., Webb, S., Brown, N. A., Lamers, W., & Anderson, R. H. (2003). Development of the

764      heart: (1) formation of the cardiac chambers and arterial trunks. *Heart (British Cardiac*

765      *Society)*, *89*(7), 806–14. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/12807866

766      Nawa, M., & Matsuoka, M. (2013). KCTD20, a relative of BTBD10, is a positive regulator of Akt.

767      *BMC Biochemistry*, *14*(1), 27. http://doi.org/10.1186/1471-2091-14-27

768      Noveroske, J. K., Lai, L., Gaussin, V., Northrop, J. L., Nakamura, H., Hirschi, K. K., & Justice, M. J.

769      (2002). Quaking is essential for blood vessel development. *Genesis (New York, N.Y. : 2000)*,

770      *32*(3), 218–30. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/11892011

771      Poduri, A., Evrony, G. D., Cai, X., Elhosary, P. C., Beroukhim, R., Lehtinen, M. K., … Walsh, C. A.

772      (2012). Somatic Activation of AKT3 Causes Hemispheric Developmental Brain Malformations.

773      *Neuron*, *74*(1), 41–48. http://doi.org/10.1016/j.neuron.2012.03.010

774      Ramsdell, A. F. (2005). Left–right asymmetry and congenital cardiac defects: Getting to the heart of

775      the matter in vertebrate left–right axis determination. *Developmental Biology*, *288*(1), 1–20.

776      http://doi.org/10.1016/J.YDBIO.2005.07.038

777      Ramu, A., Noordam, M. J., Schwartz, R. S., Wuster, A., Hurles, M. E., Cartwright, R. A., & Conrad,

778      D. F. (2013). DeNovoGear: de novo indel and point mutation discovery and phasing. *Nature*

779      *Methods*, *10*(10), 985–987. http://doi.org/10.1038/nmeth.2611

780    Ren, K., Yuan, J., Yang, M., Gao, X., Ding, X., Zhou, J., … Zhang, J. (2014). KCTD10 Is Involved in

781        the Cardiovascular System and Notch Signaling during Early Embryonic Development. *PLoS*

782        *ONE*, *9*(11), e112275. http://doi.org/10.1371/journal.pone.0112275

783    Rivière, J.-B., Mirzaa, G. M., O'Roak, B. J., Beddaoui, M., Alcantara, D., Conway, R. L., … Dobyns,

784        W. B. (2012). De novo germline and postzygotic mutations in AKT3, PIK3R2 and PIK3CA

785        cause a spectrum of related megalencephaly syndromes. *Nature Genetics*, *44*(8), 934–940.

786        http://doi.org/10.1038/ng.2331

787    Sallman, D. A., Komrokji, R., Vaupel, C., Cluzeau, T., Geyer, S. M., McGraw, K. L., … Padron, E.

788        (2016). Impact of TP53 mutation variant allele frequency on phenotype and outcomes in

789        myelodysplastic syndromes. *Leukemia*, *30*(3), 666–673. http://doi.org/10.1038/leu.2015.304

790    Sampson, J., Jacobs, K., Yeager, M., Chanock, S., & Chatterjee, N. (2011). Efficient study design

791        for next generation sequencing. *Genetic Epidemiology*, *35*(4), n/a-n/a.

792        http://doi.org/10.1002/gepi.20575

793    Smith, C. L., Blake, J. A., Kadin, J. A., Richardson, J. E., Bult, C. J., & Mouse Genome Database

794        Group. (2018). Mouse Genome Database (MGD)-2018: knowledgebase for the laboratory

795        mouse. *Nucleic Acids Research*, *46*(D1), D836–D842. http://doi.org/10.1093/nar/gkx1006

796    Smith, K. S., Yadav, V. K., Pei, S., Pollyea, D. A., Jordan, C. T., & De, S. (2016). SomVarIUS:

797        somatic variant identification from unpaired tissue samples. *Bioinformatics*, *32*(6), 808–813.

798        http://doi.org/10.1093/bioinformatics/btv685

799    Soubrier, F., Chung, W. K., Machado, R., Grünig, E., Aldred, M., Geraci, M., … Humbert, M. (2013).

800        Genetics and Genomics of Pulmonary Arterial Hypertension. *Journal of the American College*

801        *of Cardiology*, *62*(25), D13–D21. http://doi.org/10.1016/J.JACC.2013.10.035

802    Stevens, K. N., Hakonarson, H., Kim, C. E., Doevendans, P. A., Koeleman, B. P. C., Mital, S., …

803        Gruber, P. J. (2010). Common Variation in ISL1 Confers Genetic Susceptibility for Human

804        Congenital Heart Disease. *PLoS ONE*, *5*(5), e10855.

805        http://doi.org/10.1371/journal.pone.0010855

806    Stosser, M. B., Lindy, A. S., Butler, E., Retterer, K., Piccirillo-Stosser, C. M., Richard, G., &

807        McKnight, D. A. (2018). High frequency of mosaic pathogenic variants in genes causing

808     epilepsy-related neurodevelopmental disorders. *Genetics in Medicine*, *20*(4), 403–410.

809     http://doi.org/10.1038/gim.2017.114

810     Sun, J. X., He, Y., Sanford, E., Montesion, M., Frampton, G. M., Vignot, S., … Yelensky, R. (2018).

811     A computational approach to distinguish somatic vs. germline origin of genomic alterations

812     from deep sequencing of cancer specimens without a matched normal. *PLOS Computational*

813     *Biology*, *14*(2), e1005965. http://doi.org/10.1371/journal.pcbi.1005965

814     Tong, X., Zu, Y., Li, Z., Li, W., Ying, L., Yang, J., … Zhang, B. (2014). Kctd10 regulates heart

815     morphogenesis by repressing the transcriptional activity of Tbx5a in zebrafish. *Nature*

816     *Communications*, *5*, 1–10. http://doi.org/10.1038/ncomms4153

817     Wallis, G. A., Starman, B. J., Zinn, A. B., & Byers, P. H. (1990). Variable expression of osteogenesis

818     imperfecta in a nuclear family is explained by somatic mosaicism for a lethal point mutation in

819     the alpha 1(I) gene (COL1A1) of type I collagen in a parent. *American Journal of Human*

820     *Genetics*, *46*(6), 1034–40. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/2339700

821     Weinstein, M. M., Kang, T., Lachman, R. S., Bamshad, M., Nickerson, D. A., Krakow, D., & Cohn, D.

822     H. (2016). Somatic mosaicism for a lethal *TRPV4* mutation results in non-lethal metatropic

823     dysplasia. *American Journal of Medical Genetics Part A*, *170*(12), 3298–3302.

824     http://doi.org/10.1002/ajmg.a.37942

825     Yamamoto, G. L., Aguena, M., Gos, M., Hung, C., Pilch, J., Fahiminiya, S., … Bertola, D. R. (2015).

826     Rare variants in SOS2 and LZTR1 are associated with Noonan syndrome. *Journal of Medical*

827     *Genetics*, *52*(6), 413–421. http://doi.org/10.1136/jmedgenet-2015-103018

828     Zaidi, S., & Brueckner, M. (2017). Genetics and Genomics of Congenital Heart Disease. *Circulation*

829     *Research*, *120*(6), 923–940. http://doi.org/10.1161/CIRCRESAHA.116.309140

830     Zaidi, S., Choi, M., Wakimoto, H., Ma, L., Jiang, J., Overton, J. D., … Lifton, R. P. (2013). De novo

831     mutations in histone-modifying genes in congenital heart disease. *Nature*, *498*(7453), 220–

832     223. http://doi.org/10.1038/nature12141

833

834

835

836

837     **Fig 1. Mosaic detection by Expectation-Maximization. (A)** Expectation-Maximization (EM) Estimation

838     to decompose the variant allele fraction (VAF) distribution of our input variants into mosaic and germline

839     distributions. The EM-estimated prior mosaic fraction was 12.15% and the mean of the mosaic VAF

840     distribution was 0.15. **(B)** Read depth vs. VAF distrubution of individual variants. The blue line denotes

841     mean VAF (0.49) and the red lines denote the 95% confidence interval under our Beta-Binomial model.

842     Mosaic variants are defined as sites with posterior odds > 10, corresponding to a False Discovery Rate of

843     9.1%.Germline variants are represented in black and mosaic variants are represented in red. **(C)**

844     Estimated mosaic detection power as a function of average sample depth for values between 40x and

845     500x.

846

847     **Fig 2. Mutation spectrum of detected germline and mosaic variants**. Rates of specific mutations were

848     compared in germline **(A)**, blood mosaic **(B)**, and CHD tissue mosaic **(C)** variants. Transitions

849     predominated in both variant sets.

850

851     **Fig 3. Validated mosaics detected in probands with matched blood and cardiovascular tissue**

852     **samples available.** Validation VAF from blood compared to validation VAF from cardiovascular tissue

853     demonstrated tissue-specific mosaicism (red) as well as shared mosaicism (blue). Predicted effect of

854     mosaic variants corresponds to marker shape.

855

856     **Fig 4. Damaging mosaics in CHD-related genes have higher variant allele fraction than likely-**

857     **benign mosaics. (A)** Among the 76 mosaics in CHD-related genes, likely damaging variants have a

858     higher VAF than likely benign (Mann-Whitney *U* p=0.001). **(B)** Among the 233 mosaics in Other (non-

859     CHD-related) genes, there is no difference in VAF based on predicted effect (p=0.985).

860

861

862     **Supplement**

863     **Supplemental Methods**

864   *Union with Validated de novo SNVs from Jin et al. Nature Genetics 2017*

865        As part of the PCGC program, Jin et al. previously sequenced and processed a cohort of 2871

866   CHD probands – including 2530 parent-offspring trios used in this study – to investigate the contribution

867   of rare inherited and *de novo* variants to CHD. They called a total of 2992 proband *de novo* variants,

868   including 2872 SNVs and 118 indels, and Sanger confirmed a subset of the most likely-disease causing

869   variants. Since we processed the same proband-parent trios using different variant calling pipelines, we

870   combined the results of our two approaches to provide a more complete input *de novo* call set for mosaic

871   variant detection.

872        We first processed our SAMtools *de novo* calls using our upstream filters (*n*=2396 sites passing

873   all filters). We then applied the same upstream filters to the published dnSNVs from Jin et al. (*n*=2650

874   sites passing all filters) before finally taking the union of these two call sets (*n*=3192). There were 1814

875   sites in the intersection, with 836 sites unique to the Jin et al. calls and 542 sites unique to our SAMtools

876   calls. After preprocessing, outlier removal, and FDR-based minimum $N_{alt}$ filtering, the remaining 2971

877   dnSNVs were used as input to our mosaic detection model.

878

879   *Mutation Spectrum Analysis*

880        We compared the mutation spectrum – the frequencies of all possible base changes – of our

881   predicted mosaic candidates against the spectrum of our predicted germline heterozygous variants.

882   Under the assumption that that post-zygotic events occur randomly (i.e. due to errors in DNA replication

883   rather than a specific biological process), the mosaic mutation spectrum should not differ significantly

884   from the germline mutation spectrum. We used Pearson's Chi-square Test to test for a difference in

885   frequencies across all base changes between our predicted sets of variants. We interpreted large

886   qualitative differences in base change frequencies as evidence of technical artifacts and rejection of the

887   Chi-square null as evidence of systemic issues in our pipeline.

888

889   *Mosaic Detection Power Given Sample Average Coverage*

890        To model statistical power in the context of mosaic variant detection, we considered two

891   conditional probabilities: (i) the probability of detecting a mosaic event (i.e. the probability of a variant's

892    posterior odds exceeding a threshold) given site depth $DP_{site}$, VAF, and overdispersion parameter $\theta$ and

893    (ii) the probability of observing site depth $DP_{site}$, given sample-wide average coverage $DP_{sample}$.

894    (i) Pr(detect mosaic | $DP_{site}$, VAF, $\theta$) was calculated by first identifying the VAF range (and by extension,

895    the range of $N_{alt}$) over which posterior odds > cutoff, then by integrating the beta-binomial probability

896    mass function over this range, with considerations for the probability of strand bias (P(strand bias | $DP_{site}$)

897    ~ Binomial($N_{alt}$, $DP_{site}$, $p$=0.5)).

898    (ii) Pr($DP_{site}$ | $DP_{sample}$) follows an overdispersed poisson distribution that we approximated using a

899    negative binomial model with overdispersion parameter $\theta$ {Sampson 2011}. For each $DP_{sample}$ value, we

900    calculated a vector of weights corresponding to Pr($DP_{site}$ | $DP_{sample}$) for $DP_{site}$ values in the range (1,

901    1500).

902         Finally, we took the sum of the detection probabilities described in (i) multiplied by the weights

903    described in (ii) to determine the probability of detecting a mosaic variant given a sample average

904    coverage value – Pr(detect mosaic | $DP_{sample}$). Our estimated detection power curves for a range of

905    sample average coverage values typical of exome-sequencing studies are shown in **(Figure 4A).** Our

906    CHD cohort was sequenced to sample average depth of 60x, with prior mosaic fraction=0.121 and

907    estimated $\theta$=116.

908         To estimate the true rate of mosaicism per exome given sample average coverage, we first split

909    our set of predicted mosaics into VAF bins of size 0.05. For each bin above VAF 0.1, we multiplied the

910    number of mosaics by the inverse of the detection power for that given VAF bin to estimate the true count

911    of mosaic variants in that VAF range, assuming full detection power. Since EM-mosaic is underpowered

912    to detect mosaics with VAF < 0.1 in the blood and since this range is enriched for technical artifacts that

913    potentially affect our counts, we did not apply this scaling procedure to these bins to avoid over-inflating

914    our adjusted mosaic rate estimate **(Figure 4B)**.

915

916    *Filtering of MosaicHunter Candidate Variants **(Fig S3)***

917         MosaicHunter was used to identify candidate mosaic variants from blood exome-sequencing trio

918    data using default settings {Huang 2014}. Filtering of original MosaicHunter candidate variants excluded,

919    in order, any variant present in ExAC (46634), G to T mutations with fewer than $N_{alt}$<10 oxidative

31

920     indicating DNA damage {Costello 2013} (3995), non-uniquely called sites (4719), germline SNVs

921     previously called by GATK HaplotypeCaller (591), probands with >20 mosaic variants (1490 in 10

922     probands), mosaic log posterior likelihood ratio <10 (940), variants with >2 parental alternative allele

923     reads (244), variants with gnomAD population frequency > 1e-4 or located in MUC or HLA genes (40).

924

925     *Filtering of cardiovascular tissue Candidate Variants*

926     We used the MosaicHunter pipeline in trio mode to identify candidate variants in WES data from

927     70 cardiovascular tissue samples (belonging to 66 unique probands). From the list of variants initially

928     reported by the pipeline using default settings, we applied the same filtration steps listed for

929     MosaicHunter candidate variants in blood samples with the exception of the removal of G to T mutations

930     with fewer than 10 alternative allele reads and the mosaic log posterior likelihood ratio <10. Finally, we

931     removed variants that were identified in either parent or had a total read depth <10 in either parent.

932

933     *Clinical interpretation of mosaic variants – limitations*

934     We note that conventional clinical interpretation of mosaic mutations is challenging for several

935     reasons: (i) it is unclear in which tissues each mosaic mutation is expressed (ii) several study participants

936     were very young at time of clinical assessment and many classical disease features may not yet have

937     developed or been noted, and (iii) the absence of additional clinical features does not necessarily rule out

938     a mosaic mutation as being for the cause of the CHD. For the purposes of this study, we selected these

939     mosaic mutations on the basis of predicted pathogenicity and detection in genes involved in biological

940     processes relevant to CHD or developmental disorders

941

942     **Fig S1. EM-mosaic Flowchart.** We first processed our SAMtools *de novo* calls using our upstream filters

943     ($n$=2396 sites passing all filters). We then applied the same upstream filters to the published dnSNVs

944     from Jin et al. ($n$=2650 sites passing all filters) before finally taking the union of these two call sets

945     ($n$=3192). High-confidence mosaics (n=309) were defined as mosaics passing IGV inspection and having

946     posterior odds > 10. Grey text indicates which filters removed candidate mosaic variants called by

947     MosaicHunter but not by EM-mosaic.

948

949 **Fig S2. Blood variants with posterior odds between 1 and 10.(A)** Distribution of the 86 variants with

950 posterior odds between 1 and 10.**(B)** Histogram of counts by bin. To estimate the number of potential

951 mosaics mosaics missed by our threshold, counts of each bin were scaled by the estimated true positive

952 rate (TPR; posterior odds / 1+posterior odds). By our estimate, 54/86 variants were likely mosaic and

953 32/86 were likely germline.

954

955 **Fig S3. MosaicHunter workflow.** Quality Control filters excluded any sites that were (1) present in ExAC

956 (2) G>T with $N_{alt}$<10 (3) parent $N_{alt}$>2. Outliers were defined as probands carrying more than 20 mosaics,

957 or non-unique sites. We also removed sites called as germline by GATK Haplotype Caller. High-

958 confidence mosaics (n=116) were defined as having Likelihood Ratio > 80 and affecting coding regions

959 excluding MUC/HLA genes. Grey text indicates which filters removed variants called by EM-mosaic but

960 not by MosaicHunter.

961

962 **Fig S4. Comparison of variant allele fraction (VAF) and read depth of EM-mosaic and**

963 **MosaicHunter.** Candidate mosiac variants detected by the two pipelines had more overlapping variants

964 at low read depth and VAF values.

965

966 **Fig S5. Targeted sequencing to validate candidate blood mosaic variants.** (A) EM-mosaic and (B)

967 MosaicHunter variants were assayed using PCR followed by MiSeq for high-depth assessment of

968 mosaicism. Variants with x symbols were shared by both pipelines. Mosaic variants that validated are

969 black, while variants with VAF > 0.45 and therefore germline are red.

970 Validation VAF values demonstrated significant correlation with the original WES-derived VAF for EM-

971 mosaic (Pearson's correlation $P$=2.2x10$^{-16}$) and MosaicHunter ($P$=8.2x10$^{-11}$).

972

973 **Fig S6. Overdispersion.** Overdispersion is commonly seen in WES data {Heinrich 2012; Ramu 2013}

974 and is defined as observing variance (in terms of the VAF of variants with a given DP value) higher than

975　　expected across DP values, under a given statistical model. The blue line denotes the expectation under

976　　a Binomial model and the red line denotes the expectation under a Beta-Binomial model.

977

978　　**Fig S7. FDR-based minimum $N_{alt}$ threshold.** A FDR-based approach was used to determine a threshold

979　　for the minimum number of reads supporting the alternate allele for each site to avoid false positives

980　　caused purely by sequencing errors. Assuming that sequencing errors are independent and that errors

981　　occur with probability 0.005, with the probability of an allele-specific error being 0.005/3=0.00167, and

982　　given the total number of reads ($N$) supporting a variant site, we iterated over a range of possible $N_{alt}$

983　　values between 1 and 0.5*N and estimated the expected number of false-positives due to sequencing

984　　error, exome-wide ((1- $f_{poisson}$(x=$n$, λ=$N$*0.005/3))*3x10$^7$ ; where $f_{poisson}$ is the probability of x events in a

985　　Poisson process with mean λ). Assuming one coding *de novo* SNV per exome {Acuna-Hidalgo 2016} and

986　　that roughly 10% of *de novo* SNVs arise post-zygotically {Lim 2017; Krupp 2017; Freed 2016}, we used a

987　　conservative assumption of 0.1 mosaic mutation per exome. To constrain theoretical FDR to 10% we

988　　allowed a maximum of 0.01 false positives per exome and used the corresponding $N_{alt}$ value to define an

989　　FDR-based minimum $N_{alt}$ threshold for each variant. We then excluded variants with alternate allele read

990　　counts below this threshold.

991

992　　**Fig S8. Estimated mosaic detection power using less stringent mosaic definitions. (A)** Estimated

993　　true frequency of detectable coding mosaics (0.4>VAF>0.1) adjusted by detection power (n=341;

994　　0.135/exome) **(B)** Calibrated mosaic detection power and estimated true mosaic frequency of detectable

995　　coding mosaics, using posterior odds cutoff of 5 (n=361; 0.143/exome). **(C)** Calibrated mosaic detection

996　　power and estimated true mosaic frequency of detectable coding mosaics, using posterior odds cutoff of

997　　2 (0.4>VAF>0.1; n=424; 0.168/exome).

998

999　　**Fig. S9. Mosaic variants shared in blood and cardiovascular tissues have higher variant allele**

1000　　**fraction.** Validation VAF from **(A)** cardiovascular tissue and **(B)** blood had higher VAF for shared variants

1001　　compared to tissue-specific variants (p=0.101 and 0.015, respectively).

1002

1003 **Fig. S10. Damaging CHD-related mosaics have higher VAF under less stringent definitions of**

1004 **mosaicism. (A)** Using posterior odds cutoff of 5 (corresponding to 315 mosaics). Among 78 mosaics in

1005 CHD-related genes (left), there were 14 variants predicted as damaging, 63 variants predicted as likely-

1006 benign, and 1 variant of unknown functional consequence. Among 237 mosaics in non-CHD-related

1007 genes (right), there were 41 variants predicted as damaging, 184 variants predicted as likely-benign, and

1008 2 variants of unknown functional consequence. **(B)** Using posterior odds cutoff of 2 (corresponding to 352

1009 mosaics). Among 89 mosaics in CHD-related genes (left), there were 17 variants predicted as damaging,

1010 71 variants predicted as likely-benign, and 1 variant of unknown functional consequence. Among 263

1011 mosaics in non-CHD-related genes (right), there were 54 variants predicted as damaging, 206 variants

1012 predicted as likely-benign, and 3 variants of unknown functional consequence.

1013

1014 **Fig. S11. Mosaic rate by age**. (**A**) Age distribution for all 2530 probands in cohort. (**B**) Mosaic Rate

1015 across Age ranges. Rate = # mosaics/# probands in age bin. Note: 9/2530 probands missing Age

1016 information. 1/367 mosaic belong to a proband with missing Age.

1017

1018 **Fig. S12**. **Mosaic rate by parental age at birth.** Mosaic rate by age of father (blue) and mother (red) at

1019 birth. Rate = # mosaics/# probands in each parental age bin. Note: 9/2530 probands missing age

1020 information. 1/367 mosaic belong to a proband with missing age.

1021

1022 **Fig. S13**. **Confirmation rate across VAF bins.** The number of candidates for which we performed MiSeq

1023 resequencing among (**A**) the union set (n=143 tested) (**B**) all EM-mosaic calls (n=97) and (**C**) all

1024 MosaicHunter (n=68) calls vs. the number confirmed as mosaic for VAF ranges [0, 0.1), [0.1, 0.2), and

1025 [0.2, 0.3).

1026

1027 **Fig. S14**. **Posterior odds comparison for tested vs. untested mosaics.** Among 309 candidates with

1028 EM-mosaic posterior odds scores available, we compared the distribution of tested (n=97) vs. untested

1029 (n=212) mosaics. The $\log_{10}$-scaled posterior odds distribution for the tested group is shown in blue

1030 (mean=5.382). The $\log_{10}$-scaled mean posterior odds for the untested group is shown in red

35

1031    (mean=7.050).  The selected candidates had lower posterior odds than those not selected for

1032    confirmation (Mann Whitney *U* test *P*=0.002).

1033

1034    **Table S1**. **Proband Phenotypes**. Cardiac and neurodevelopmental phenotypes for CHD probands. NDD

1035    diagnosis is unknown for patients <1 year of age.

1036

1037    **Table S2**. **Cohort Summary**. Number of *de novo* and mosaic variants for probands with isolated CHD,

1038    extracardiac anomalies (ECA), neurodevelopmental delay (NDD) or unknown phenotypes.

1039

1040    **Table S3**. **EM-mosaic Mosaic Candidates**. Candidate mosaic variants identified by EM-mosaic and

1041    MiSeq validation results.

1042

1043    **Table S4**. **MosaicHunter Mosaic Candidates**. Candidate mosaic variants identified by MosaicHunter

1044    and MiSeq validation results.

1045

1046    **Table S5. Blood mosaic Candidate Validation by MiSeq.** 143 candidate mosaic sites were assessed

1047    using targeted PCR and deep sequencing. 85/97 (88%) of selected EM-mosaic sites and 44/68 (67%) of

1048    selected MosaicHunter sites were confirmed.

1049

1050    **Table S6**. **Cardiovascular Tissues with Whole Exome Sequencing**. 70 tissues from 66 CHD probands

1051    were assessed for mosaic variants.

1052

1053    **Table S7. CHD tissue mosaic Candidate Validation by MiSeq.** 24 candidate mosaic sites were

1054    assessed using targeted PCR and deep sequencing. 85/92 (92%) of selected EM-mosaic sites and 44/64

1055    (69%) of selected MosaicHunter sites were confirmed .

1056

1057    **Table S8**. **CHD related genes.** We considered the union of genes highly expressed in the developing

1058    heart (HHE) and known candidate CHD genes {Jin 2017} as CHD-related genes (n=4558).

1059

1060 **Table S9**. **All Protein-Altering Mosaics**. Detailed information for 398 mosaic variants predicted to affect

1061 protein sequence.

1062

1063 **Table S10**. **Damaging Mosaics in CHD-Relevant Genes**. Detailed information for 25 mosaic variants

1064 likely to contribute to CHD.

1065

1066

| ID | Gene | Variant Class | Pipeline | CHD Tissue | | | | Blood WES VAF | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Location | WES AD | WES VAF | MiSeq VAF | WES AD | WES VAF | MiSeq VAF |
| 1-00543 | *CTCFL* | Bmis | EM-mosaic | AO | 138,36 | 0.21 | 0.32 | 29,8 | 0.22 | 0.19 |
| 1-00984 | *ZNF16* | syn | EM-mosaic | LV | 262,1 | 0.00 | 0.01 | 100,7 | 0.07 | 0.07 |
| 1-01282 | *GABRA6* | Dmis | MosaicHunter | RV | 104,1 | 0.01 | 0.01 | 55,12 | 0.18 | 0.18 |
| 1-01684 | *CCNC* | Bmis | Both | AoValve, RV | 36,7 | 0.16 | 0.17, 0.19 | 224,40 | 0.15 | 0.14 |
| 1-02672 | *TOR1A* | syn | Both | AtrSpt | 159,10 | 0.06 | 0.10 | 29,6 | 0.17 | 0.19 |
| 1-03512 | *RFX3* | LoF | MosaicHunter | RV | 156,15 | 0.09 | 0.08 | 39,0 | 0.00 | 0.03 |
| 1-04652 | *PCDH10* | syn | Both | AtrSpt | 154,19 | 0.11 | 0.14 | 15,1 | 0.06 | 0.10 |
| 1-07004 | *ANK2* | Bmis | MosaicHunter | SubAoMembr | 226,13 | 0.05 | 0.04 | 30,0 | 0.00 | 0.00 |
| 1-07004 | *MYH14* | Bmis | Both | SubAoMembr | 124,22 | 0.15 | 0.27 | 33,0 | 0.00 | 0.00 |
| 1-07004 | *NRG3* | Bmis | EM-mosaic | SubAoMembr | 152,30 | 0.16 | 0.24 | 43,0 | 0.00 | 0.00 |
| 1-07004 | *NUDT21* | Bmis | Both | SubAoMembr | 137,22 | 0.14 | 0.14 | 74,0 | 0.00 | 0.02 |
| 1-07004 | *TET3* | Dmis | MosaicHunter | SubAoMembr | 131,1 | 0.01 | 0.03 | 81,16 | 0.16 | 0.27 |
| 1-07299 | *RRS1* | syn | Both | RV, UNK | 160,25 | 0.14 | 0.25 | 22,2 | 0.08 | 0.14 |
| 1-09869 | *PIK3C2G* | LoF | MosaicHunter | LV | 126,9 | 0.07 | 0.10 | 31,0 | 0.00 | 0.00 |
| 1-11800 | *TMEM45A* | Bmis | MosaicHunter | RV | 213,0 | 0.00 | 0.00 | 32,7 | 0.18 | 0.06 |

1067 **Table 2. Mosaics detected in individuals with matched cardiovascular tissue and blood**
1068
1069 Characteristics of mosaic variants predicted for individuals with blood and cardiovascular tissue WES data available. Among 15 mosaics, 5 were detected via
1070 analysis of blood WES, 8 were detected from cardiovascular tissue WES, and 2 were detected by both approaches. Six of 7 (86%) mosaics detected from analysis
1071 of blood were present in both DNA sources with MiSeq VAF≥0.01. Two additional variants previously identified as *de novo* germline variants in blood WES were
1072 absent from CHD tissue WES. Minimum 1023 MiSeq reads used to determine VAF. Abbreviations: AD, allelic depth (reference, alternate); AO, aorta; AtrSpt, atrial
1073 septum; Bmis, benign missense; Dmis, deleterious missense; LOF, Loss of function variant; LV, left ventricle; RV, right ventricle; VAF, variant allele fraction

| ID | Gene | Variant Class | Blood VAF | pLI | Episcore | Heart Exp | Age (y) | Clinical Phenotype | | PCGC denovo LoF/Dmis Variants in Mosaic Gene |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Cardiac Abnormalities | Extracardiac Abnormalities | |
| 1-00761 | FBN1 | Dmis | 0.24 | 1.00 | 98 | 93 | 4.3 | Mitral stenosis | dysmorphic features, subglottic stenosis, hypoplasic left mainstem bronchus, short stature | 3 |
| 1-07004 | TET3 | Dmis | 0.16 | 1.00 | 7 | 87 | 10.3 | Subaortic stenosis | None | 0 |
| 1-05662 | SETD2 | LoF | 0.13 | 1.00 | 99 | 85 | 0.8 | Aortic coarctation, mitral valve hypoplasia | None | 0 |
| 1-00344 | UBR5 | splice | 0.27 | 1.00 | 95 | 90 | 16 | D-transposition of the great arteries, VSD, valvar and subvalvar pulmonary stenosis | None | 0 |
| 1-03512 | RFX3 | LoF | 0.09 | 1.00 | 100 | 46 | 0.4 | Tetralogy of Fallot with pulmonary stenosis | None | 0 |
| 1-06216 | ITSN1 | Dmis | 0.21 | 1.00 | 98 | 86 | 0.3 | ASD | plagiocephaly, rib anomaly, single kidney, dysmorphic facial features | 0 |
| 1-00363 | QSER1 | Dmis | 0.06 | 1.00 | 94 | 79 | 3.7 | Tetralogy of Fallot with pulmonary stenosis, VSD | inguinal hernia | 0 |
| 1-13185 | PKD1 | Dmis | 0.10 | 1.00 | 87 | 84 | 0.8 | VSD, partially anomalous pulmonary venous return | hemangioma | 1 |
| 1-00192 | GLYR1 | Dmis | 0.22 | 0.99 | 89 | 93 | 0.4 | ASD, VSD, interrupted aortic arch, hypoplastic tricuspid valve, BAV | None | 0 |
| 1-04046 | FZD5 | Dmis | 0.09 | 0.99 | 89 | 48 | 0.2 | Tetralogy of Fallot with pulmonary stenosis, VSD | None | 0 |
| 1-06649 | NOVA2 | Dmis | 0.15 | 0.95 | 75 | 56 | 0.6 | Tetralogy of Fallot with pulmonary stenosis | None | 0 |
| 1-05095 | ISL1 | LoF | 0.07 | 0.90 | 97 | 25 | 2.4 | ASD | None | 0 |
| 1-06677 | KCTD10 | Dmis | 0.16 | 0.84 | 75 | 91 | 10.1 | Aortic coarctation, pulmonary valve stenosis | dysmorphic facial features, hydrocephalus, pyloric stenosis, single kidney, imperforate/atretic anus | 0 |
| 1-05447 | HNRNPAB | Dmis | 0.09 | 0.76 | 72 | 99 | 7.8 | ASD, BAV, aortic coarctation | None | 0 |
| 1-00021 | QKI | LoF | 0.13 | 0.76 | 94 | 97 | 0.5 | Doublet outlet right ventricle, pulmonary stenosis, VSD | None | 0 |
| 1-11871 | FHOD3 | Dmis | 0.18 | 0.05 | 91 | 92 | 0.0 | Tetralogy of Fallot with pulmonary atresia | hypocalcemia, thrombocytopenia, lymphopenia | 0 |
| 1-01458 | HK2 | Dmis | 0.27 | 0.04 | 89 | 90 | 0.0 | Hypoplastic left heart with aortic and mitral atresia, aortic coarctation | None | 1 |
| 1-00669 | PRKD3 | splice | 0.19 | 0.02 | 77 | 82 | 0 | D-transposition of the great arteries, conal VSD, bilateral conus, interrupted aortic arch | None | 0 |
| 1-00524 | RNF20 | LoF | 0.10 | 0.00 | 55 | 83 | 23.7 | Left-dominant complete atrioventricular canal | Heterotaxy with situs inversus totalis, asplenia, duodenal atresia | 0 |
| 1-01851 | SUCLA2 | LoF | 0.11 | 0.00 | 72 | 89 | 15.5 | Balanced complete atrioventricular canal, aortic coarctation | None | 0 |
| 1-03885 | LZTR1 | Dmis | 0.20 | 0.00 | 31 | 84 | 14.7 | Abnormal pulmonary vein draining into the right atrium | left sided/midline liver, asplenia, maltrotation | 2 |
| 1-05011 | KCTD20 | Dmis | 0.26 | 0.00 | 76 | 77 | 24.5 | Transposition of the great arteries, Tricispid and Pulmonary valve atresia | left-sided/midline liver | 1 |
| 1-00018 | FIG4 | Dmis | 0.19 | 0.00 | 49 | 70 | 11.8 | BAV, mitral atresia, aortic coarctation, VSD, total anomalous pulmonary venous return | nephritis | 1 |
| 1-05661 | SMAD9 | Dmis | 0.06 | 0.00 | 84 | 39 | 9.3 | Common atrioventricular canal | None | 0 |
| 1-09869 | PIK3C2G | LoF | 0.07* | 0.00 | 73 | 28 | 6.3 | Common atrioventricular canal, aortic stenosis, aortic arch hypoplasia, VSDs | dysmorphic facial features, low-set ears, campomelic dysplasia | 1 |

**Table 3. Damaging Mosaics in CHD-relevant genes**
There were 25 potentially-pathogenic mosaic mutations based on known gene function and patient phenotype. Some of these probands have previously described rare LoF/Dmis variants, though none are likely pathogenic for CHD {Jin 2017}. Additionally, some genes were previously found to have LoF/Dmis variants among other individuals in this CHD cohort. Abbreviations: ASD, atrial septal defect; BAV, bicuspid aortic valve; Dmis, deleterious missense; episcore, haploinsufficiency score (percentile rank) {Han 2018}; Heart Exp, heart expression percentile rank; LoF, loss-of-function; pLI, probability of loss-of-function intolerance {gnomAD}; PCGC, Pediatric Cardiac Genomics Consortium; VAF, variant allele fraction; VSD, ventricular septal defect. *VAF refers to CHD tissue WES.

**Declarations**

*Ethics approval and consent to participate*

CHD subjects were recruited to the Congenital Heart Disease Network Study of the Pediatric Cardiac Genomics Consortium (CHD GENES: ClinicalTrials.gov identifier NCT01196182). The institutional review boards of Boston's Children's Hospital, Brigham and Women's Hospital, Great Ormond Street Hospital, Children's Hospital of Los Angeles, Children's Hospital of Philadelphia, Columbia University Medical Center, Icahn School of Medicine at Mount Sinai, Rochester School of Medicine and Dentistry, Steven and Alexandra Cohen Children's Medical Center of New York, and Yale School of Medicine approved the protocols. All subjects or their parents provided informed consent.

*Consent for publication*

See above.

*Availability of data and material*

EM-mosaic and code for analyzing data are available from https://github.com/ShenLab/mosaicism. The MosaicHunter software is available from http://mosaichunter.cbi.pku.edu.cn/. SAMtools is available from http://www.htslib.org/. ANNOVAR is available from http://annovar.openbioinformatics.org/en/latest/. Integrative Genomics Viewer (IGV) software is available from https://software.broadinstitute.org/software/igv/. Whole-exome sequencing data have been deposited in the database of Genotypes and Phenotypes (dbGaP) under accession numbers phs000571.v1.p1, phs000571.v2.p1 and phs000571.v3.p2. In-house pipelines are available from the corresponding authors on reasonable request.

*Competing interests*

The authors declare that they have no competing interests.

*Authors' contributions*

YS, JGS, CES, and WKC conceived and oversaw the study. AH, SUM, JALW, HQ, KBM, JGS, CES, YS, WKC analyzed the data. AH developed the EM-mosaic pipeline and wrote the statistical analysis code. SUM, JALW carried out MosaicHunter analyses of blood and tissue samples. JMG, AT, SD performed MiSeq experimental confirmation. AH, SUM, EG, CES, WKC interpreted the impact of mosaics on participant clinical phenotypes. DB, RWK, JWN, GAP, DS, MT-F, MB, RPL, EG, BDG, CES, JGS, WKC were involved in cohort ascertainment, phenotypic characterization, and recruitment. DM collected cardiovascular tissue samples. AH, SUM, JALW, YS, JGS, CES, WKC wrote the manuscript. All authors read and approved the manuscript.

# Fig 1

**Fig 2**

**Fig 3**

**Fig 4**

A

B

Likely Benign

Damaging

p-value=0.001

Likely Benign

Damaging

p-value=0.985

Variant Allele Fraction

Variant Allele Fraction

**Fig S1**

Parse validated DNMs from Jin et al.

n = 2650

Call *de novo* SNVs using SAMtools with settings optimized for mosaic sites

n = 17704

- 28 MH variants not called

Apply QC filters

n = 5717

- 11 MH variants removed due to parent alt reads >0
- 5 MH variants removed for other quality filters

Preprocessing and outlier sample removal

n = 2396

Create union set with Jin et al. validated DNMs

n = 3192

IGV visualization for variants VAF<0.3

n = 2971

Score variant sites and identify candidate mosaics

n = 309

high-confidence mosaic SNVs

**Fig S2**

**Fig S3**

**Fig S4**

**Fig S5**

**Fig S6**

**Fig S7**

**Fig S8**

Fig S9

# Fig S10

# Fig S11

**A**



**B**

**Fig S12**

Fig S13

# Fig S14