

1 Assessing the contribution of rare variants to congenital heart disease through a large-  
2 scale case-control exome study

3

4 Enrique Audain<sup>1,2#</sup>, Anna Wilsdon<sup>3#</sup>, Gregor Dombrowsky<sup>1,2</sup>, Alejandro Sifrim<sup>4</sup>, Jeroen Breckpot<sup>4</sup>,  
5 Yasset Perez-Riverol<sup>5</sup>, Siobhan Loughna<sup>3</sup>, Allan Daly<sup>6</sup>, Pavlos Antoniou<sup>6</sup>, Philipp Hofmann<sup>1</sup>, Amilcar  
6 Perez-Riverol<sup>2</sup>, Anne-Karin Kahlert<sup>1</sup>, Ulrike Bauer<sup>7</sup>, Thomas Pickardt<sup>7</sup>, Sabine Klaassen<sup>8,9</sup>, Felix  
7 Berger<sup>9</sup>, Ingo Daehnert<sup>10</sup>, Sven Dittrich<sup>11</sup>, Brigitte Stiller<sup>12</sup>, Hashim Abdul-Khaliq<sup>13</sup>, Frances  
8 Bu'lock<sup>14</sup>, Anselm Uebing<sup>1,15</sup>, Hans-Heiner Kramer<sup>1</sup>, Vivek Iyer<sup>6</sup>, Lars Allan Larsen<sup>16</sup>, J David  
9 Brook<sup>3\*</sup>, Marc-Phillip Hitz<sup>1,2,6,15\*</sup>

10 # These authors contributed equally to this work.

11 \* Joint corresponding authors.

12

- 1 Department of Congenital Heart Disease and Pediatric Cardiology, University Hospital of Schleswig-Holstein, Kiel, Germany,
- 2 Institute for Medical Genetics, Klinikum Oldenburg, Oldenburg, Germany.
- 3 School of Life Sciences, University of Nottingham, University Park, Nottingham, United Kingdom
- 4 Department of Human Genetics, University of Leuven, KU Leuven, Leuven, Belgium
- 5 European Bioinformatics Institute (EBI), Wellcome Trust Genome Campus, Hinxton, Cambridge, United Kingdom
- 6 Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, United Kingdom
- 7 Competence Network for Congenital Heart Defects, Berlin, Germany
- 8 Experimental and Clinical Research Center (ECRC), Charité - Universitätsmedizin Berlin and Max Delbrück Center, Berlin, Germany
- 9 Deutsches Herzzentrum der Charité, Dept. of Congenital Heart Disease-Pediatric Cardiology, Berlin, Germany
- 10 Department of Pediatric Cardiology and Congenital Heart Disease, Heart Center, University of Leipzig, Leipzig, Germany
- 11 Department of Pediatric Cardiology, University Hospital Erlangen, Friedrich-Alexander-University Erlangen-Nürnberg (FAU), Erlangen, Germany
- 12 Department of Congenital Heart Disease and Pediatric Cardiology, University Heart Center Freiburg—Bad Krozingen, Freiburg, Germany
- 13 Department of pediatric Cardiology, Saarland University Hospital, Homburg, Germany
- 14 Congenital and Paediatric Cardiology, East Midlands Congenital Heart Centre and University of Leicester, Glenfield Hospital, United Kingdom
- 15 German Centre for Cardiovascular Research (DZHK), Partner Site Kiel, Germany.
- 16 Department of Cellular and Molecular Medicine, University of Copenhagen, Copenhagen, Denmark

13  
14

15 **ABSTRACT**

16

17 Several studies have demonstrated the value of large-scale human exome and  
18 genome data analysis, to maximise gene discovery in rare diseases. Using this  
19 approach, we have analysed the exomes of 4,747 cases and 52,881 controls, to  
20 identify single genes and digenic interactions which confer a substantial risk of  
21 congenital heart disease (CHD). We identified both rare loss-of-function and missense  
22 coding variants in ten genes which reached genome-wide significance (Bonferroni  
23 adjusted  $P < 0.05$ ) and an additional four genes with a significant association at a false  
24 discovery rate (*FDR*) threshold of 5%. We highlight distinct genetic contributions to  
25 syndromic and non-syndromic CHD at both single gene and digenic level, by  
26 independently analysing probands from these two groups. In addition, by integrative  
27 analysis of exome data with single-cell transcriptomics data from human embryonic  
28 hearts, we identified cardiac-specific cells as well as putative biological processes  
29 underlying the pathogenesis of CHD. In summary, our findings strengthen the  
30 association of known CHD genes, and have identified additional novel disease genes  
31 and digenic interactions contributing to the aetiology of CHD.

32 **Keywords:** Congenital Heart Disease; UK Biobank; genetic variants; digenic  
33 interactions.

34

35

## 36 INTRODUCTION

37

38 Congenital Heart Disease (CHD) is a global health challenge, affecting ~1-2% of live  
39 births worldwide<sup>1</sup>. However, despite advances in our understanding of the underlying  
40 disease aetiology in recent years, a significant proportion of CHD cases remains  
41 unexplained, suggesting that genetic mechanisms and other risk factors remain poorly  
42 understood<sup>2,3</sup>. Recent advances in exome and genome sequencing technologies have  
43 opened up new avenues of study and have resulted in novel insights into the genetic  
44 and epigenetic mechanisms underlying rare diseases, such as CHD<sup>4,5</sup>.

45 Previous studies have defined the association of inherited and *de novo* variations as  
46 a cause of CHD<sup>6,7</sup>. In addition, these studies have highlighted the differences between  
47 the genetic architecture of syndromic (with extracardiac malformations and/or  
48 neurodevelopmental delay) and non-syndromic (isolated) CHD<sup>6,7</sup>. Continuing  
49 collaboration between the scientific community and healthcare teams has driven  
50 efforts to integrate and analyse larger cohorts of patients, and has demonstrated the  
51 potential of this approach to uncover novel variants and genes associated with CHD<sup>6-</sup>  
52 <sup>8</sup>.

53 Here, we present a whole exome sequencing analysis of 4,747 CHD cases and 52,881  
54 controls. This is one of the largest cohorts of non-syndromic CHD cases (n=2,929)  
55 studied so far, meaning that we are in an advantageous position to refine our  
56 understanding of the genetic mechanisms underlying non-syndromic CHD specifically.  
57 This is especially important given that the vast majority of individuals with CHD, have  
58 non-syndromic CHD.

59 We used the case-control cohort to investigate both single genes and digenic  
60 interactions contributing to CHD. We integrated data obtained in the case-control

61 study with single-cell transcriptome data obtained from human embryonic hearts<sup>9</sup>. This  
62 complementary analysis identified biological processes enriched for genes  
63 differentially expressed in cardiac-specific cells, found also significant in our case-  
64 control analysis. Importantly, the data suggest a difference in cardiac developmental  
65 mechanisms between syndromic and non-syndromic CHD.

66 Taken together, we have identified ten genome-wide significant (*Bonferroni adjusted*  
67  $P < 0.05$ ) genes, and an additional four genes at *FDR* 5%, which are associated with  
68 CHD, as well as a larger contribution of digenic interactions to non-syndromic  
69 compared to syndromic CHD.

70

## 71 **RESULTS**

72

### 73 **Cohort description and analysis workflow**

74 We combined and analysed the exomes of 4,747 CHD cases (aCHD, refers to all CHD  
75 cases) and 52,881 controls. CHD cases were further classified into syndromic CHD  
76 (sCHD, individuals with extracardiac malformations or neurodevelopmental disability,  
77  $n=1,818$ ) and non-syndromic CHD (nsCHD, individuals with isolated CHD,  $n=2,929$ ).

78 All samples and genetic variants were subjected to a sequence of quality control steps  
79 to obtain a final cohort of unrelated and matched-ancestry individuals, as well as a set  
80 of high-confidence variants for downstream analysis (see **Methods, Supplemental**  
81 **Information 1**).

82 We evaluated the distribution of high-confidence loss-of-function (hcLOF) and  
83 missense constrained variants (missC) across a spectrum of LOF and missense  
84 constrained genes (**Methods**). In addition, we performed gene-based burden testing  
85 to identify genes conferring a high risk of CHD, as well as the expression pattern at

86 single-cell resolution. Lastly, we evaluate the contribution of digenic interactors to  
87 syndromic and non-syndromic CHD. **Figure 1** simplifies the workflow followed in this  
88 study to discover novel associations with CHD.

89

## 90 **Distinct pattern of loss-of-function constrained genes identified between sCHD** 91 **and nsCHD**

92 Previous studies have suggested a greater contribution of loss-of-function (LOF)  
93 variants to sCHD, compared to non-syndromic forms<sup>6,7</sup>. To determinate if this holds  
94 true in this present cohort, we evaluated the burden of rare variants in the sCHD and  
95 nsCHD cohort, compared with controls across the per gene spectrum of loss-of-  
96 function intolerance. Following the approach proposed by the gnomAD consortium<sup>10</sup>,  
97 we divided 19,923 protein-coding genes into ten bins (~1,900 genes per bin) based  
98 on its observed/expected LOF ratio upper fraction (termed LOEUF) and applied a  
99 logistic regression model (see **Methods**) to each bin (i.e., gene-set). This allowed us  
100 to assess enrichment across three different functional categories of variants (hcLOF,  
101 missC and synonymous), stratified by CHD probands (aCHD, sCHD and nsCHD).  
102 The highest enrichment was observed in the most LOF constrained genes (bin 1) for  
103 hcLOF variants (**Figure 2**). These variants provided a major contribution to sCHD  
104 cases ( $OR = 2.27$ ,  $P < 2 \times 10^{-16}$ ), and much less so for nsCHD ( $OR = 1.52$ ,  $P = 1.2 \times$   
105  $10^{-13}$ ). A moderate enrichment was observed for missC variations, suggesting that this  
106 class of variants could have a similar (although smaller) functional impact compared  
107 to hcLOF variants. Although reduced in magnitude, this same pattern was also  
108 observed in the set of genes in the second LOEUF constraint bin, whereas no  
109 enrichment was observed towards less LOEUF constrained bins (**Figure 2**). No

110 enrichment of synonymous variants was observed across the bins, providing a  
111 negative control set (**Figure 2**).

112 When the same analysis was performed across the missense constraint spectrum,  
113 assessed by the observed/expected missense ratio upper fraction gene-based metric  
114 (termed MOEUF), a similar pattern as described above (higher enrichment in the most  
115 missense-constrained genes) was observed (**Supplemental Figure 1**).

116 These results demonstrate a larger effect of hcLOF compared to missC variants  
117 across the LOEUF and MOEUF spectrum, with the major contribution observed in  
118 sCHD, compared with nsCHD. Nevertheless, the results suggest that both hcLOF and  
119 missC variants are important genetic components contributing to CHD development.

120

## 121 **Gene-based enrichment analysis**

122 To identify genes that confer a significant risk of CHD, we performed a case-control  
123 burden analysis by combining rare variants ( $MAF < 0.001$ ) at the gene level. It has  
124 been demonstrated that following a method of collapsing variants within specific  
125 genomic regions (e.g., genes), increases the power to discover new associations at  
126 low allele frequencies<sup>11</sup>. Following this principle, we conducted a Fisher's exact test to  
127 identify genes with a significant burden of non-synonymous variants in CHD cases  
128 compared to controls, and evaluated them independently for sCHD and nsCHD.

129 As with earlier comparable case-control exome studies<sup>12-14</sup>, the burden test was  
130 performed separately for hcLOF ( $P_{lof}$ ) and missC ( $P_{miss}$ ), and the minimal *p-value*  
131 observed per gene between these two variant categories was selected as the study-  
132 wide *p-value* ( $P$ ). hcLOF variants were defined using the LOFTEE tool<sup>10</sup>, whereas  
133 missense variants were defined based on different missense deleteriousness  
134 prediction scores (see **Methods, Supplemental Figure 2**). Ten genes were identified

135 with significant  $P$ , after correcting for multiple testing using the Bonferroni method  
136 (**Table 1, Supplemental Table 1**). Eight genes were associated with sCHD (*KMT2A*,  
137 *SMAD4*, *PTPN11*, *TAB2*, *NSD1*, *BCOR*, *KAT6A*, *PBX1*) and two were identified  
138 through the nsCHD (*FLT4* and *NOTCH1*) analysis. In addition, four genes showed  
139 significant associations with CHD at  $FDR$  5% (*CTCF*, *KAT6B*, *SHOX2*, *HCAR1*). The  
140 evaluation of the set of synonymous variants showed a similar distribution of expected  
141 vs observed  $p$ -values, suggesting no genomic inflation of the test statistic  
142 (**Supplemental Figure 3**).

143 Of the genes identified as significant in sCHD, *KMT2A* (AD Wiedemann-Steiner  
144 syndrome OMIM 159555) showed the highest enrichment (**Figure 3a**). *NOTCH1* (AD  
145 Adams-Oliver syndrome 5, Aortic valve disease 1 OMIM 190198) showed the highest  
146 number of variations in the nsCHD cohort (**Figure 3b**) and warranted further  
147 investigation (companion manuscript).

148 Other genes reaching a significant level of association included *NSD1* (AD Sotos  
149 Syndrome OMIM 606681), *TAB2* (AD Non-Syndromic CHD 2, OMIM 605101), *KAT6A*  
150 (AD Arboleda-Tham Syndrome OMIM 601408), *PTPN11* (AD Noonan Syndrome  
151 OMIM 176876), *SMAD4* (AD Mhyre Syndrome OMIM 600993), *FLT4* (AD Congenital  
152 heart defects, multiple types, 7 OMIM 136352), and the X-linked gene *BCOR* (XLD  
153 Syndromic Microphthalmia OMIM 300485). They have all been previously described in  
154 the context of CHD, and our results corroborate these findings.

155 The association of *PBX1* (AD Congenital anomalies of kidney and urinary tract  
156 syndrome with or without hearing loss, abnormal ears, or developmental delay OMIM  
157 176310), *CTCF* (AD Intellectual developmental disorder, autosomal dominant 21  
158 OMIM 604167) and *KAT6B* (AD Genitopatellar syndrome and SBBYSS syndrome  
159 OMIM 605880) with CHD (**Table 1**) have been previously reported in isolated cases

160 or small patient cohorts, and our results add further evidence for an association with  
161 CHD.

162 *HCAR1* (OMIM 606923) and *SHOX2* (OMIM 602504) have not previously been  
163 associated with CHD at a genome-wide level. However, both genes were significantly  
164 associated with nsCHD at *FDR* 5% (**Figure 3b**).

165

### 166 **Differentially expressed genes in cardiac-specific cells show a distinct** 167 **enrichment pattern in syndromic and non-syndromic CHD**

168 Previous studies have revealed significant levels of expression in the heart of genes  
169 associated with CHD<sup>7,8</sup>. By using publicly accessible bulk RNAseq data<sup>15</sup> (**Methods**),  
170 we consistently showed that genes with significant level of association in our case-  
171 control analysis also showed high expression in cardiac tissues (**Supplemental**  
172 **Figure 4**). Moreover, syndromic CHD genes showed a systematic elevated  
173 expression in other tissues (e.g., brain and kidney), compared to non-syndromic CHD  
174 (**Supplemental Figure 5**). The difference in expression patterns between these two  
175 groups was negligible in the heart, though ( $P > 0.05$ , Wilcoxon test; **Supplemental**  
176 **Figure 5**). Despite its relevance, bulk RNAseq data analysis does not stretch as far  
177 as the delineation of expression patterns at the cellular level.

178 To accomplish this, we assessed the mutational burden of rare non-synonymous  
179 variants (hcLOF and missC) within differentially expressed genes (DEGs) in cardiac-  
180 specific cells. We meta-analysed the exome data with a publicly available human heart  
181 transcriptomic dataset generated from early developmental stages of the human heart  
182 (6.5 and 7 weeks post-conception)<sup>9</sup>. Using the logistic regression framework  
183 mentioned above, we performed gene-set enrichment analysis on DEGs defined on  
184 15 distinct cardiac cell clusters (C0-C14) reported by Asp *et al*<sup>9</sup>. Both hcLOF and



185 missC mutations were evaluated independently and the analysis was stratified further  
186 by proband CHD status versus controls (aCHD, sCHD and nsCHD, **Figure 4**).

187 Five cardiac-specific cell clusters were found significantly enriched (Bonferroni  
188 adjusted  $P < 0.05$ ) for hcLOF variations when analysing aCHD probands vs controls  
189 (**Figure 4**): Smooth muscle cells (C5), Cardiac neural crest cells (C14), Epicardium-  
190 derived cells (C3), Capillary endothelium (C0) and Atrial cardiomyocytes (C7).  
191 Enrichment of hcLOF variants for DEGs in Smooth muscle cells (C5) showed a  
192 significant contribution to both sCHD and nsCHD. Cardiac neural crest cells (C14) and  
193 Atrial cardiomyocytes (C7) contributed to sCHD in the main, whereas the cluster of  
194 Capillary endothelium cells was significantly enriched in nsCHD versus controls  
195 (**Figure 4**).

196 A similar enrichment pattern was observed when analysing the set of missC variants  
197 (**Supplemental Figure 6**). In addition to the Capillary endothelium (C0), Smooth  
198 muscle cells (C5), Atrial cardiomyocytes (C7) and Cardiac neural crest cells (C14)  
199 clusters; which were also found significantly enriched for hcLOF variants; two other  
200 cardiac-specific cell clusters showed a significant burden of missC variants in CHD  
201 cases (aCHD) compared to controls: Endothelium/pericytes cells (C10) and Fibroblast  
202 cells (C2).

203 The synonymous variants set was used as a negative control and did not identify  
204 enrichment in any clusters evaluated (Bonferroni adjusted  $P > 0.05$ , **Supplemental**  
205 **Figure 7**).

206 Together, these results provide valuable evidence regarding the possible mechanisms  
207 involved in the pathogenesis of CHD.

208

209

## 210 **Gene Ontology (GO) enrichment analysis**

211 To provide additional supporting evidence for our previous findings, we performed  
212 Gene Ontology (GO) enrichment analysis to identify relationships between the  
213 enriched DEGs in cardiac-specific cell clusters to biological processes. We analysed  
214 the set of DEGs with an unadjusted  $P < 0.01$  (Fisher Exact test) identified in the case-  
215 control burden analysis within the cell clusters showing enrichment in either the aCHD,  
216 sCHD or nsCHD analysis (**Supplemental Figure 8**).

217 Among the DEGs in cardiac-specific cells evaluated with the Enrichr tool<sup>16</sup> (see  
218 **Methods**), four clusters showed at least one GO term with  $FDR < 1\%$ . The data  
219 suggested that cell cluster C7 (Atrial cardiomyocytes, **Supplemental Figure 8a**) was  
220 mainly associated with biological processes involved in developing cardiac muscle  
221 tissue, and the observed signal was driven by *NKX2-5*, *MYH6*, *MYOCD*, *PKP2*, *BMP7*,  
222 *ANKRD1* and *ACTC1*. DEGs in C0 (Capillary endothelium, **Supplemental Figure 8b**)  
223 showed enrichment for vasculogenesis, with contribution from *KDR*, *NOTCH1* and  
224 *RASIP1*. C5 (Smooth muscle cells, **Supplemental Figure 8c**) was associated with  
225 extracellular matrix organisation processes, with a noteworthy contribution of genes  
226 that contain a collagen-like domain (e.g., *COL14A1* and *COL1A2*), as well as *ELN* and  
227 *FBN2*. DEGs in C10 (Endothelium and pericyte cells, **Supplemental Figure 8d**),  
228 demonstrated the higher enrichment of missC variants (**Supplemental Figure 6**), and  
229 enrichment of biological process involved in the cellular response to vascular  
230 endothelial growth factor stimulus and the regulation of cell migration as part of  
231 sprouting angiogenesis. *DLL4*, *FLT4*, *KDR*, *MEOX2* and *NOTCH1* all contributed to  
232 this cluster.

233

234

## 235 **Contribution of digenic interactions to syndromic and non-syndromic CHD**

236 Next, we studied the contribution of digenic interactions to CHD using the RareComb<sup>17</sup>  
237 framework (**see Methods**). Our analysis revealed a total of 2,083 digenic pairs  
238 significantly enriched for hcLOF and/or missC variants in CHD cases (aCHD)  
239 compared to controls at *FDR* 1% (**Figure 5a, Supplemental Table 2**). The data  
240 suggested that a significantly higher proportion of digenic interactors contributed to  
241 non-syndromic (n=810) forms of CHD ( $P = 6.7 \times 10^{-3}$ , proportion Z-test, **Figure 5a**)  
242 compared to syndromic forms (n=433).

243 The rate of a gene being observed in a digenic pair showed no correlation ( $r_2 < 0.2$ )  
244 with its coding sequence (CDS) length (**Supplemental figure 9a**). Thus, the digenic  
245 interactions implicated by our data does not appear to be biased by an increased  
246 mutation rate of larger genes.

247 A significantly lower LOEUF was observed in the genes forming digenic pairs  
248 contributing to syndromic CHD compared to non-syndromic CHD (**Figure 5b**). A  
249 similar pattern was observed by comparing the distribution of MOEUF (**Supplemental**  
250 **figure 9b**), whereas no significant differences were observed when analysing the  
251 observed/expected synonymous upper fraction ratio metric (SOEUF, **Supplemental**  
252 **figure 9c**). Thus, our data suggest that digenic pairs show similar correlation between  
253 loss-of-function intolerance and sCHD as we observed for single genes (**Figure 2b**).

254 We hypothesized that if genes in the digenic gene-sets are causative of CHD in our  
255 patient cohort, we would expect that they are enriched for known CHD genes. To test  
256 this hypothesis, we calculated the overlap between genes in the digenic lists aCHD,  
257 sCHD and nsCHD and curated lists of genes known to cause CHD in patients<sup>8</sup>. We  
258 observed significant enrichment of known CHD disease genes in all three digenic lists  
259 (**Figure 5c**).

260 To investigate if digenic pairs interact at a systems level, we generated protein-protein  
261 interaction (PPI) networks using data from STRING<sup>18</sup>. Networks of 1,929 nodes/3,587  
262 edges, 1,042 nodes/882 edges and 610 nodes/307 edges were generated for aCHD,  
263 nsCHD and sCHD digenic lists, respectively (**Figure 5d, Supplemental Figure 10**).  
264 Analysis of the networks showed that nsCHD and sCHD digenic gene-lists interact in  
265 physical networks with PPI enrichment values of 1.78 ( $P < 1.0 \times 10^{-16}$ ) and 1.51 ( $P =$   
266  $8.3 \times 10^{-12}$ ), respectively.

267 To test if the genes in digenic pairs share the same function, we performed a simple  
268 overlap test to see how many digenic pairs show direct (1<sup>st</sup> degree) interaction in the  
269 network. The results show that only very few pairs interact directly at the protein level  
270 (**Figure 5e**). Finally, to determine the degree of separation between digenic pairs in  
271 each network, we calculated the length of the shortest path between genes in each  
272 pair. This analysis showed that the median of the shortest path between digenic pairs  
273 is 5 and 6 for the sCHD and nsCHD networks, respectively (**Figure 5f, Supplemental**  
274 **Figure 11**).

275 In summary, our analyses suggest that the digenic pairs interact at systems level in  
276 complex PPI networks, with a high degree of separation between the genes in each  
277 digenic pair.

278

## 279 **DISCUSSION**

280

281 In this study, we amassed 57,628 human exomes and conducted both a gene- and  
282 gene-set centred case-control burden analysis to increase our understanding of the  
283 genetic causes of CHD. After quality control at the sample and variant level, we  
284 provide a comprehensive CHD case-control cohort with unrelated and ancestry-

285 matched individuals. Specifically, the availability of detailed phenotype data allowed  
286 us to explore the differences between syndromic and non-syndromic forms of CHD.  
287 By utilising gene-level constraint information<sup>10</sup>, we investigated the contribution and  
288 properties of loss-of-function and missense constraint variants independently for all  
289 CHD cases, as well as syndromic and non-syndromic CHD independently. Like earlier  
290 comparable studies<sup>6,7</sup>, our results revealed a higher contribution of LOF variants to  
291 CHD compared to missense variants, confirming that this type of variation represents  
292 the largest driver. Subsequently, the analysis of syndromic cases revealed a higher  
293 burden of LOF mutations when compared with the non-syndromic cohort. This effect  
294 was mainly a result of the contribution of genes with a higher intolerance to loss-of-  
295 function variations. This same pattern was also observed when analysing the genes  
296 based on missense constraint.

297

298 We next assessed the contribution to CHD at the gene level, by performing a gene-  
299 based case-control burden analysis. Our analysis revealed ten genes that reached  
300 genome-wide significant levels of association with CHD (*NSD1*<sup>19</sup>, *TAB2*<sup>20</sup>, *KAT6A*<sup>21</sup>,  
301 *PTPN11*<sup>22</sup>, *CTCF*<sup>23</sup>, *SMAD4*<sup>24</sup>, *FLT4*<sup>25</sup>, *NOTCH1*<sup>26,27</sup>, *BCOR*<sup>28</sup> and *KMT2A*<sup>29</sup>).  
302 Previous studies have associated these genes as a cause of CHD, and our results  
303 confirm this association (**Table 1**). Furthermore, four candidate genes (*PBX1*, *SHOX2*,  
304 *KAT6B*, *HCAR1*) were found contributing to both syndromic and non-syndromic CHD  
305 at *FDR* 5%. To our knowledge, these genes have been either not previously  
306 associated with, or have only been briefly described in the context of CHD.

307

308 *PBX1* has been primarily associated with congenital abnormalities of the kidney and  
309 urinary tract (CAKUT)<sup>30</sup>; however, previous studies have reported isolated cases

310 carrying *de novo* missense variations leading to syndromic CHD<sup>30,31</sup>. In line with these  
311 early reports, our analysis revealed a significant burden of missense constrained  
312 variants in *PBX1* in syndromic CHD patients (**Table 1**). It has also been demonstrated  
313 that deficiency of *Pbx1* impacts branchial arch artery patterning and results in the  
314 failure of cardiac outflow tract septation<sup>32</sup>. Interestingly, this gene was also found to  
315 be differentially expressed in Epicardium and Smooth muscle cells (**Supplemental**  
316 **Figure 12**). Together our findings suggest that *PBX1* contributes significantly to  
317 syndromic forms of CHD.

318  
319 *SHOX2* was significantly enriched (at *FDR* 5%) in the nsCHD cohort, for missC  
320 variants (**Table 1**). Recent studies in animal models have demonstrated that the *Shox2*  
321 null mice are embryonic-lethal<sup>33</sup>. Cardiovascular defects identified in these mice  
322 included an abnormally low heartbeat rate, a severely hypoplastic Sinoatrial Node  
323 (SAN), hypoplastic or absent sinus valves<sup>33</sup>, and other atrial abnormalities (e.g.,  
324 enlarged atrial chamber and thinner atrial wall). Subsequently, *SHOX2* has been  
325 described as playing a key role in developing the Sinoatrial Node<sup>33,34</sup>. In addition,  
326 *SHOX2* was identified as a significant DEG in atrial cardiomyocytes (**Supplemental**  
327 **Figure 12**), providing further supporting evidence of its role in heart development,  
328 most likely by regulating the activity of *NXK2-5*<sup>33,35</sup> and *TBX5*<sup>36</sup>. These results imply  
329 that *SHOX2* is a plausible novel non-syndromic CHD gene.

330  
331 Truncating variants in the Lysine Acetyltransferase 6B gene (*KAT6B*) have been  
332 associated with Say–Barber–Biesecker–Young–Simpson Syndrome (SBBYSS,  
333 OMIM 603736) and Genitopatellar Syndrome (GTPTS, OMIM 606170). Heart defects  
334 have been reported as part of the phenotypic spectrum of SBBYSS<sup>37</sup>. In a recent study

335 of 32 individuals with *KAT6B* disorder, 47% showed cardiovascular anomalies, mainly  
336 atrial septal defects, ventricular septal defects, and patent ductus arteriosus<sup>38</sup>. Our  
337 results have identified that *KAT6B* was differentially expressed in the cluster of atrial  
338 cardiomyocytes cells (**Supplemental Figure 12**), which suggests a possible role in  
339 the early cardiac development program. Our analysis extends previous findings  
340 associating loss-of-function variations in *KAT6B* to sCHD.

341

342 The Hydroxycarboxylic Acid Receptor 1 (*HCAR1*) does not appear to have been  
343 associated with CHD thus far, however our findings suggest this gene may be a novel  
344 candidate CHD gene. It was not differentially expressed in any of the cardiac-specific  
345 cell clusters analysed.

346

347 By meta-analysing the genomic data with heart single-cell transcriptomic data, we  
348 investigated the pattern of expression of DEGs for aCHD, sCHD and nsCHD in a range  
349 of cardiac-specific cells. Using Gene Ontology enrichment as a complementary  
350 analysis, we identified key gene markers and biological processes associated with  
351 CHD. Unlike previous studies<sup>7,39</sup>, which focused on whole heart bulk-RNA sequencing  
352 data, the use of transcriptomic data at a single-cell resolution allowed the analysis of  
353 candidate gene expression patterns in specific cardiac cell clusters important for early  
354 cardiac development. Our analysis highlighted distinct cardiac cell clusters  
355 contributing to sCHD and nsCHD. In addition, we demonstrated that missense  
356 constrained variants could have a similar functional impact compared to loss-of-  
357 function variants, although to a lesser degree. For instance, the significant enrichment  
358 of sCHD in cardiac neural crest cells (cNCCs) suggests a broader contribution of  
359 patients affected by syndromic occurrences, not limited to heart development only.

360 Perturbations in the cNCCs migration process can lead to a wide spectrum of human  
361 cardio-craniofacial syndromes, including DiGeorge Syndrome (22q11.2 Deletion  
362 Syndrome, OMIM 188400) and CHARGE (OMIM 214800). The enrichment observed  
363 in capillary endothelium and pericyte cells in nsCHD, associated with the  
364 vasculogenesis process, suggests that the phenotypic occurrence in these patients is  
365 limited to the cardiovascular system rather than affecting a broader spectrum of cells.  
366 Whilst the results are promising, they are limited because the currently available  
367 human heart single-cell map<sup>9</sup> is incomplete (e.g., only a few early developmental time  
368 points). Therefore, future studies integrating mouse and human single-cell heart and  
369 whole-embryo data are warranted.

370

371 Contrasting with the study of monogenic causes of CHD, oligogenic factors underlying  
372 the disease have been explored to a lesser extent. We took advantage of a newly  
373 developed method to study the contribution of digenic interactions (the simplest form  
374 of oligogenic) to CHD, in a case-control setting<sup>17</sup>. Interestingly, we observed a higher  
375 proportion of digenic interactions contributing to non-syndromic compared to  
376 syndromic CHD. These results contrast with those observed at the gene level, where  
377 16 out of 21 genes (~76%) found significant at *FDR* 10% (**Table 1**) were associated  
378 with syndromic CHD. These findings imply that sCHD is more likely to have a  
379 monogenic aetiology, and oligogenic interactions may be a more important component  
380 in the development of non-syndromic forms of CHD.

381 Functional annotation of subclusters in the networks, generated from nsCHD and  
382 sCHD digenic pairs, indicate that the identified digenic pairs encode proteins involved  
383 in transcriptional regulation, signalling pathways (e.g., BMP/TGF beta signalling and



384 NOTCH signalling) and tissue structures (e.g. sarcomere and extracellular matrix)  
385 which are important in heart development<sup>40–42</sup>.

386 Our network analyses suggest that rather than interacting within the same subcluster,  
387 the digenic pairs interact at systems level, with a high degree of separation between  
388 the genes in each pair, thus supporting previous results which suggest that CHD risk  
389 factors converge in higher-order developmental networks<sup>43</sup>.

390 In summary, we analysed ~57,000 exomes, and complemented this with  
391 transcriptomic data at single-cell resolution. The findings have strengthened the  
392 association of previously described genes with CHD, identified novel candidate genes,  
393 and provide a deeper understanding of the pathophysiological mechanisms underlying  
394 CHD at gene and digenic level and the potential different aetiologies between  
395 syndromic and non-syndromic CHD.

396

## 397 **METHODS**

398

### 399 **Cohort description**

400 To create a comprehensive CHD case-control cohort, exome sequencing data from  
401 multiple individuals was combined in a unique reference dataset. CHD cases were  
402 mainly sequenced as part of an initiative from the German Competence Network for  
403 Congenital Heart Defects, the Deciphering Developmental Disorder (DDD) project and  
404 the University of Nottingham (UK); controls were sequenced as part of the UK Biobank  
405 (UKBB). Samples from the UKBB dataset with phenotype description labelled as  
406 Schizophrenia (SCZ), bipolar disorder (BP) or developmental delay (DD) were  
407 excluded from the analysis. Accordingly, a small fraction of samples in the UKBB  
408 cohort (127 samples), labelled as CHD cases, were included in the analysis. In total,

409 we assembled an exome dataset consisting of 57,628 samples (4,747 CHD cases and  
410 52,881 controls).

411

### 412 **Alignment, quality control and variant annotation**

413 The assembled dataset was processed and harmonized using the same alignment  
414 (BWA v0.3), calling (GATK v4.0), annotation (VEP v95) and quality control (Hail v0.2)  
415 pipelines. **Supplemental Information 1** describes extensively the implementation and  
416 results of these methods.

417

### 418 **Defining a set of loss-of-function and missense constraint variants**

419 We enriched the dataset for high-confidence loss-of-function (hcLOF) variants and  
420 missense constrained (missC) variants. hcLOF variants were annotated as indicated  
421 by the LOFTEE tool (<https://github.com/konradjk/loftee>) with its default parameters  
422 and included stop-gained, essential splice and frameshift variants. To define a set of  
423 missC variants, we evaluated four state-of-art pathogenicity prediction scores:  
424 CADD<sup>44</sup>, MPC<sup>45</sup>, REVEL<sup>46</sup> and MVP<sup>47</sup>. Specifically, the performance of these scores  
425 was assessed by classifying benign and pathogenic missense variants (accessed  
426 through the ClinVar database, <https://www.ncbi.nlm.nih.gov/clinvar>) in the context of  
427 known CHD genes. In brief, receiver operating characteristic (ROC) analysis was  
428 conducted for benign and pathogenic variants within known CHD genes. The analysis  
429 was further stratified by splitting the gene set into LOF constraint (LOEUF < 0.35) and  
430 LOF non-constraint (LOEUF >= 0.35) genes. A score was defined as a 'good predictor'  
431 if achieved an area-under-curve (AUC) > 90% in both evaluated scenarios. Three of  
432 these scores (CADD, REVEL and MVP) met this criterion. A missense variant was  
433 defined as missC if it was predicted as likely deleterious by at least two of these scores

434 based on the optimal threshold suggested by the ROC analysis (**Supplemental**  
435 **Figure 2**).

436

### 437 **Defining rare variants**

438 Variants were filtered based on the cohort-specific allelic frequency ('internal' AF) as  
439 well as using external datasets. A variant was defined as rare if AF was lower than  
440 0.001 (MAF < 0.001) in the gnomAD database<sup>10</sup> (both exomes v2.1.1 and genomes  
441 v3.0.0), the RUMC cohort<sup>48</sup>, as well as AFs from an *in-house* German exome  
442 sequencing cohort.

443

### 444 **Gene-set enrichment analysis**

445 *Generation of gene sets.* Gene set-level association analysis was performed to assess  
446 whether an excess of the possible pathogenic variants was enriched for a particular  
447 category of genes (as described below). This procedure was executed for the following  
448 gene sets:

449 a) LOEUF gene bins: Constraint loss-of-function (LOF) metrics per protein-coding  
450 genes were accessed through gnomAD resource<sup>10</sup>. Genes were ranked by their  
451 observed/expected LOF mutation ratio upper fraction (termed LOUEF), and ten bins  
452 with an equal number of genes (~1,900 genes per bin) were defined. Lower values of  
453 LOEUF (e.g., bins 1 and 2) denote most LOF-constrained genes.

454 b) MOEUF gene bins: Similar as described above for LOEUF genes, but genes were  
455 binned based on their observed/expected missense mutation ratio upper fraction  
456 (termed MOEUF).

457 c) Differentially expressed genes (DEGs) in cardiac-specific cells: DEGs identified in  
458 15 distinct cardiac cell clusters reported by Asp *et al*<sup>9</sup>. In brief, genes were determined

459 as significantly differentially expressed in a particular cardiac cell cluster if the  
460 averaged log-fold change ( $\log_{2}FC$ ) > 0 (upregulated) at FDR 1%.

461

462 *Gene set-based association analysis.* For each sample within the filtered dataset, we  
463 generate a Minimal Allele Count (MAC) metric by aggregating high confidence  
464 Genotypes (DP  $\geq$  10, GQ  $\geq$  20 and allelic balance heterozygous > 0.2) across the  
465 genes within the gene set. Then, a burden logistic regression test was performed using  
466 CHD case/control status as response and the five first ancestry principal component  
467 and sex as covariates using the Hail function *hl.logistic\_regression\_rows*. The analysis  
468 was stratified at the sample and variant level. At the sample level, the data was divided  
469 based on the syndromic status; three categories were tested: aCHD (all CHD cases  
470 vs control), nsCHD (non-syndromic CHD cases vs control) and sCHD (syndromic CHD  
471 cases vs control). At variant level, three different groups were evaluated based on the  
472 predicted severity of the variants: hcLOF (most severe), missC and synonymous. The  
473 synonymous variant set was used as a negative control set at the variant level to  
474 evaluate for potential artefacts. The odds ratio ( $\exp(\beta)$ ), 95% confidence  
475 interval and *p-value* metrics were used to evaluate significant enrichment.

476

### 477 **Gene-based burden testing**

478 We performed case-control gene-centred burden test analysis to assess genes with  
479 significant association with CHD. Fisher Exact test was performed independently for  
480 rare (MAF < 0.001) hcLOF and missC variants. To define the significant study-wide *p-*  
481 *value*, the minimal *p-value* (*P*) per gene between these two categories was chosen.  
482 The analysis was further stratified by syndromic status to assess the distinct  
483 contribution of these categories to CHD. A gene was defined as genome-wide

484 significant if it reached a Bonferroni corrected  $P < 0.05$  and suggested significant if  
485  $FDR < 5\%$ . In addition, the set of synonymous variants was used as a negative control  
486 set since no difference between cases/control is expected on this set of variations  
487 (quantile-quantile plots, **Supplemental Figure 3**).

488

### 489 **Expression analysis using bulk RNAseq data**

490 A publicly available human transcriptomic dataset previously described by Cardoso-  
491 Moreira *et al*<sup>15</sup> was used to complement this study. To assess the gene expression  
492 levels in the heart, kidney, brain, and liver; the RPKM matrix hosted in ArrayExpress  
493 (E-MTAB-6814) was used. Gene expression levels were averaged among samples in  
494 the early developmental stages (4-8 weeks-post-conception). Percentile rank per gene  
495 was computed based on the mean expression.

496

### 497 **Gene Ontology enrichment analysis**

498 The R-package *Enrichr* (with the *Biological\_Process\_2018* database) was used to  
499 perform Gene Ontology (GO) enrichment analysis. The analysis was conducted on  
500 the differentially expressed genes (DEGs) in cardiac-specific cell clusters, which also  
501 showed unadjusted  $P < 0.01$  (Fisher Exact test) from the case-control burden analysis.  
502 The evaluated DEGs were previously reported by Asp *et al*<sup>9</sup> with no additional  
503 processing. GO terms with only one overlapping gene were filtered out. A biological  
504 process term was considered significant if  $FDR < 1\%$  as reported by the *Enrichr* tool<sup>15</sup>.

505

### 506 **Digenic analysis**

507 The digenic analysis was performed using the R-package *RareComb*<sup>17</sup>. *RareComb*  
508 combines inferences statistics with an a priori algorithm to elucidate digenic/oligogenic

509 combinations that are enriched for rare genetic variants. Specifically, we implemented  
510 the test '*enrichment\_depletion*', to assess digenic pairs significantly enriched with rare  
511 (MAF < 0.001) hcLOF and/or missC variations in CHD cases compared to controls  
512 (depleted). The analysis was further stratified by syndromic status (aCHD, sCHD or  
513 nsCHD vs. controls). Digenic pairs were defined as significant if  $FDR < 1\%$ .

514 Enrichment of known CHD genes was determined by calculating overlap between  
515 gene lists. Significant overlap was calculated using hypergeometric statistics. A  
516 representation factor was calculated as the number of overlapping genes, divided by  
517 the expected number of overlapping genes drawn from two independent groups:  
518  $RF = x / ((n * D) / N)$ , where  $x$  = number of overlapping genes,  $n$  = genes in group 1,  $D$  = genes  
519 in group 2,  $N$  = protein-coding genes in genome (20,000).

520 Protein-protein interactions (PPIs) were obtained from STRING v.11.5 using genes in  
521 the digenic pairs to query the database. The following parameters were used; network  
522 type: physical subnetwork, active interaction sources: experiments, databases (text  
523 mining data excluded), minimum required interaction score: 0.400 (medium  
524 confidence). PPIs were visualized as a network using CytoScape v3.9.1; nodes  
525 represent proteins and edges represent interactions between these proteins. PPI  
526 enrichment was analysed using the Analysis tool available in the online version of  
527 STRING (<https://string-db.org>). PPI enrichment was calculated as (observed number  
528 of edges) / (expected number of edges) and PPI enrichment P-values were obtained  
529 directly from STRING. Length of the shortest path between genes in each digenic  
530 pairs, within the nsCHD and sCHD networks was calculated using PesCa<sup>49</sup> v3.0.8.

531

532

533 **Table 1.** Top 21 genes in the case-control burden analysis using the Fisher Exact test stratified by  
 534 syndromic status (sCHD and nsCHD). A total of 16,351 genes were tested per variant type (hcLOF and  
 535 missC). Analysis: sCHD or nsCHD vs controls. Consequence: denotes the consequence group with the  
 536 minimal p-value ( $P$ ). sCHD: number of syndromic cases (heterozygous). nsCHD: number of non-  
 537 syndromic cases (heterozygous). Controls: number of controls (heterozygous).  $P$ : the minimal p-value  
 538 per gene between  $P_{lof}$  and  $P_{miss}$ .  $P$  adj (FDR): Adjusted minimal p-value ( $P$ ) using the B-H method with  
 539  $n = 2 \times 16,351$ .  $P$  adj (Bonferroni): Adjusted minimal p-value ( $P$ ) using the Bonferroni method with  $n =$   
 540  $2 \times 16,351$ . In bold are highlighted the ten genes with *Bonferroni adjusted*  $P < 0.05$ . **Supplemental Table**  
 541 **1** contains the results for all protein-coding genes tested.

Genes	Analysis	Consequence	sCHD	nsCHD	Controls	$P$	$P$ adj (FDR)	$P$ adj (Bonferroni)
<b>KMT2A</b>	sCHD	hcLOF	8	0	0	9.76E-13	3.19E-08	3.19E-08
<b>SMAD4</b>	sCHD	missC	11	3	16	2.47E-10	4.04E-06	8.09E-06
<b>NOTCH1</b>	nsCHD	hcLOF	2	7	0	8.48E-10	2.77E-05	2.77E-05
<b>PTPN11</b>	sCHD	missC	11	5	25	8.78E-09	9.57E-05	2.87E-04
<b>TAB2</b>	sCHD	hcLOF	5	1	0	3.13E-08	2.56E-04	1.02E-03
<b>NSD1</b>	sCHD	hcLOF	5	1	1	1.83E-07	1.20E-03	5.98E-03
<b>BCOR</b>	sCHD	hcLOF	4	0	0	9.93E-07	4.06E-03	3.25E-02
<b>KAT6A</b>	sCHD	hcLOF	4	1	0	9.93E-07	4.06E-03	3.25E-02
<b>PBX1</b>	sCHD	missC	6	3	6	7.73E-07	4.06E-03	2.53E-02
<b>FLT4</b>	nsCHD	hcLOF	0	5	0	3.32E-07	5.43E-03	1.09E-02
<b>CTCF</b>	sCHD	missC	4	1	1	4.84E-06	1.58E-02	1.58E-01
<b>KAT6B</b>	sCHD	hcLOF	4	1	1	4.84E-06	1.58E-02	1.58E-01
<b>SHOX2</b>	nsCHD	missC	1	10	21	1.81E-06	1.98E-02	5.93E-02
<b>HCAR1</b>	nsCHD	missC	2	9	18	4.40E-06	3.60E-02	1.44E-01
<b>ADNP</b>	sCHD	hcLOF	3	0	0	3.15E-05	6.44E-02	1.00E+00
<b>CHD7</b>	sCHD	hcLOF	3	0	0	3.15E-05	6.44E-02	1.00E+00
<b>EP300</b>	sCHD	hcLOF	3	1	0	3.15E-05	6.44E-02	1.00E+00
<b>KMT2D</b>	sCHD	hcLOF	3	0	0	3.15E-05	6.44E-02	1.00E+00
<b>KRT25</b>	sCHD	missC	8	3	31	2.51E-05	6.44E-02	8.19E-01
<b>QRICH1</b>	sCHD	hcLOF	3	0	0	3.15E-05	6.44E-02	1.00E+00
<b>SLC38A9</b>	nsCHD	missC	0	6	6	1.19E-05	7.78E-02	3.89E-01

542  
543

544 **References**

- 545
- 546 1. van der Linde, D. *et al.* Birth Prevalence of Congenital Heart Disease Worldwide.  
547 *J. Am. Coll. Cardiol.* **58**, 2241–2247 (2011).
  - 548 2. Pierpont, M. E. *et al.* Genetic Basis for Congenital Heart Disease: Revisited: A  
549 Scientific Statement from the American Heart Association. *Circulation* **138**, e653–  
550 e711 (2018).
  - 551 3. Morton, S. U., Quiat, D., Seidman, J. G. & Seidman, C. E. Genomic frontiers in  
552 congenital heart disease. *Nat. Rev. Cardiol.* 1–17 (2021) doi:10.1038/s41569-021-  
553 00587-4.
  - 554 4. Homsy, J. *et al.* De novo mutations in congenital heart disease with  
555 neurodevelopmental and other congenital anomalies. *Science* **350**, 1262–6  
556 (2015).
  - 557 5. Izarzugaza, J. M. G. *et al.* Systems genetics analysis identifies calcium-signaling  
558 defects as novel cause of congenital heart disease. *Genome Med.* **12**, 76 (2020).
  - 559 6. Sifrim, A. *et al.* Distinct genetic architectures for syndromic and nonsyndromic  
560 congenital heart defects identified by exome sequencing. *Nat. Genet.* **48**, 1060–5  
561 (2016).
  - 562 7. Jin, S. C. *et al.* Contribution of rare inherited and de novo variants in 2,871  
563 congenital heart disease probands. *Nat. Genet.* (2017) doi:10.1038/ng.3970.
  - 564 8. Audain, E. *et al.* Integrative analysis of genomic variants reveals new associations  
565 of candidate haploinsufficient genes with congenital heart disease. *PLOS Genet.*  
566 **17**, e1009679 (2021).
  - 567 9. Asp, M. *et al.* A Spatiotemporal Organ-Wide Gene Expression and Cell Atlas of  
568 the Developing Human Heart. *Cell* **179**, 1647-1660.e19 (2019).
  - 569 10. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from  
570 variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
  - 571 11. Povysil, G. *et al.* Rare-variant collapsing analyses for complex traits: guidelines  
572 and applications. *Nat. Rev. Genet.* **20**, 747–759 (2019).
  - 573 12. Guo, M. H., Plummer, L., Chan, Y.-M., Hirschhorn, J. N. & Lippincott, M. F. Burden  
574 Testing of Rare Variants Identified through Exome Sequencing via Publicly  
575 Available Control Data. *Am. J. Hum. Genet.* **103**, 522–534 (2018).
  - 576 13. Flannick, J. *et al.* Exome sequencing of 20,791 cases of type 2 diabetes and  
577 24,440 controls. *Nature* **570**, 71–76 (2019).
  - 578 14. Singh, T., Neale, B. M. & Daly, M. J. Exome sequencing identifies rare coding  
579 variants in 10 genes which confer substantial risk for schizophrenia on behalf of  
580 the Schizophrenia Exome Meta-Analysis (SCHEMA) Consortium\*. *medRxiv*  
581 2020.09.18.20192815 (2020) doi:10.1101/2020.09.18.20192815.
  - 582 15. Cardoso-Moreira, M. *et al.* Gene expression across mammalian organ  
583 development. *Nature* **571**, 505–509 (2019).
  - 584 16. Kuleshov, M. V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web  
585 server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).



- 586 17. Pounraja, V. K. & Girirajan, S. A general framework for identifying oligogenic  
587 combinations of rare variants in complex disorders. *Genome Res.* gr.276348.121  
588 (2022) doi:10.1101/gr.276348.121.
- 589 18. Szklarczyk, D. *et al.* The STRING database in 2021: customizable protein-protein  
590 networks, and functional characterization of user-uploaded gene/measurement  
591 sets. *Nucleic Acids Res.* **49**, D605–D612 (2021).
- 592 19. Tatton-Brown, K. *et al.* Genotype-Phenotype Associations in Sotos Syndrome: An  
593 Analysis of 266 Individuals with NSD1 Aberrations. *Am. J. Hum. Genet.* **77**, 193–  
594 204 (2005).
- 595 20. Thienpont, B. *et al.* Haploinsufficiency of TAB2 Causes Congenital Heart Defects  
596 in Humans. *Am. J. Hum. Genet.* **86**, 839–849 (2010).
- 597 21. Urreiziti, R. *et al.* Five new cases of syndromic intellectual disability due to KAT6A  
598 mutations: widening the molecular and clinical spectrum. *Orphanet J. Rare Dis.*  
599 **15**, 44 (2020).
- 600 22. Sarkozy, A. *et al.* Correlation between PTPN11 gene mutations and congenital  
601 heart defects in Noonan and LEOPARD syndromes. *J. Med. Genet.* **40**, 704–8  
602 (2003).
- 603 23. Gregor, A. *et al.* De novo mutations in the genome organizer CTCF cause  
604 intellectual disability. *Am. J. Hum. Genet.* **93**, 124–131 (2013).
- 605 24. Lin, A. E. *et al.* Gain-of-function mutations in SMAD4 cause a distinctive repertoire  
606 of cardiovascular phenotypes in patients with Myhre syndrome. *Am. J. Med.*  
607 *Genet. A.* **170**, 2617–31 (2016).
- 608 25. Reuter, M. S. *et al.* Haploinsufficiency of vascular endothelial growth factor related  
609 signaling genes is associated with tetralogy of Fallot. *Genet. Med.* **21**, 1001–1007  
610 (2019).
- 611 26. Kerstjens-Frederikse, W. S. *et al.* Cardiovascular malformations caused by  
612 NOTCH1 mutations do not keep left: data on 428 probands with left-sided CHD  
613 and their families. *Genet. Med.* **18**, 914–923 (2016).
- 614 27. Garg, V. *et al.* Mutations in NOTCH1 cause aortic valve disease. *Nature* **437**, 270–  
615 274 (2005).
- 616 28. Fan, Z. *et al.* BCOR regulates mesenchymal stem cell function by epigenetic  
617 mechanisms. *Nat. Cell Biol.* **11**, 1002–1009 (2009).
- 618 29. Baer, S. *et al.* Wiedemann-Steiner syndrome as a major cause of syndromic  
619 intellectual disability: A study of 33 French cases. *Clin. Genet.* **94**, 141–152 (2018).
- 620 30. Arts, P. *et al.* Paternal mosaicism for a novel PBX1 mutation associated with  
621 recurrent perinatal death: Phenotypic expansion of the PBX1-related syndrome.  
622 *Am. J. Med. Genet. A.* **182**, 1273–1277 (2020).
- 623 31. Alankarage, D. *et al.* Functional characterization of a novel PBX1 de novo  
624 missense variant identified in a patient with syndromic congenital heart disease.  
625 *Hum. Mol. Genet.* **29**, 1068–1082 (2020).
- 626 32. CP, C. *et al.* Pbx1 functions in distinct regulatory networks to pattern the great  
627 arteries and cardiac outflow tract. *Dev. Camb. Engl.* **135**, 3577–3586 (2008).
- 628 33. Espinoza-Lewis, R. A. *et al.* Shox2 is essential for the differentiation of cardiac  
629 pacemaker cells by repressing Nkx2-5. *Dev. Biol.* **327**, 376–385 (2009).

- 630 34. Munshi, N. V. Gene regulatory networks in cardiac conduction system  
631 development. *Circ. Res.* **110**, 1525–1537 (2012).
- 632 35. Yang, T., Huang, Z., Li, H., Wang, L. & Chen, Y. P. Conjugated activation of  
633 myocardial-specific transcription of *Gja5* by a pair of *Nkx2-5-Shox2* co-responsive  
634 elements. *Dev. Biol.* **465**, 79–87 (2020).
- 635 36. Puskaric, S. *et al.* *Shox2* mediates *Tbx5* activity by regulating *Bmp4* in the  
636 pacemaker region of the developing heart. *Hum. Mol. Genet.* **19**, 4625–4633  
637 (2010).
- 638 37. Gannon, T. *et al.* Further delineation of the KAT6B molecular and phenotypic  
639 spectrum. *Eur. J. Hum. Genet.* **23**, 1165–1170 (2015).
- 640 38. Zhang, L. X. *et al.* Further delineation of the clinical spectrum of KAT6B disorders  
641 and allelic series of pathogenic variants. *Genet. Med. Off. J. Am. Coll. Med. Genet.*  
642 **22**, 1338–1347 (2020).
- 643 39. Sevim Bayrak, C., Zhang, P., Tristani-Firouzi, M., Gelb, B. D. & Itan, Y. De novo  
644 variants in exomes of congenital heart disease patients identify risk genes and  
645 pathways. *Genome Med.* **12**, 9 (2020).
- 646 40. Lage, K. *et al.* Dissecting spatio-temporal protein networks driving human heart  
647 development and related disorders. *Mol. Syst. Biol.* **6**, 381 (2010).
- 648 41. Del Monte-Nieto, G., Fischer, J. W., Gorski, D. J., Harvey, R. P. & Kovacic, J. C.  
649 Basic Biology of Extracellular Matrix in the Cardiovascular System. *J. Am. Coll.*  
650 *Cardiol.* **75**, 2169–2188 (2020).
- 651 42. Kathiriya, I. S., Nora, E. P. & Bruneau, B. G. Investigating the transcriptional control  
652 of cardiovascular development. *Circ. Res.* **116**, 700–714 (2015).
- 653 43. Lage, K. *et al.* Genetic and environmental risk factors in congenital heart disease  
654 functionally converge in protein networks driving heart development. *Proc. Natl.*  
655 *Acad. Sci. U. S. A.* **109**, 14035–14040 (2012).
- 656 44. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD:  
657 predicting the deleteriousness of variants throughout the human genome. *Nucleic*  
658 *Acids Res.* **47**, D886–D894 (2019).
- 659 45. Samocha, K. E. *et al.* Regional missense constraint improves variant  
660 deleteriousness prediction. *bioRxiv* 148353 (2017) doi:10.1101/148353.
- 661 46. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the  
662 Pathogenicity of Rare Missense Variants. *Am. J. Hum. Genet.* **99**, 877–885 (2016).
- 663 47. Qi, H. *et al.* MVP predicts the pathogenicity of missense variants by deep learning.  
664 *Nat. Commun.* **12**, 510 (2021).
- 665 48. Kaplanis, J. *et al.* Evidence for 28 genetic disorders discovered by combining  
666 healthcare and research data. *Nature* **586**, 757–762 (2020).
- 667 49. Scardoni, G., Tosadori, G., Pratap, S., Spoto, F. & Laudanna, C. Finding the  
668 shortest path with PesCa: a tool for network reconstruction. *F1000Research* **4**, 484  
669 (2015).
- 670  
671  
672

673 **Data availability**

674 The CRAM-level data from CHD patients used in this study can be accessed under  
675 the following accession codes (European Genome-phenome Archive):  
676 EGAD00001002200, EGAD00001000796, EGAD00001000797, EGAD00001000800,  
677 EGAS00001000544, EGAS00001000775, EGAS00001000762. UK Biobank 50K  
678 WES dataset freeze was accessed under the application number 44165.

679

680 **Code availability**

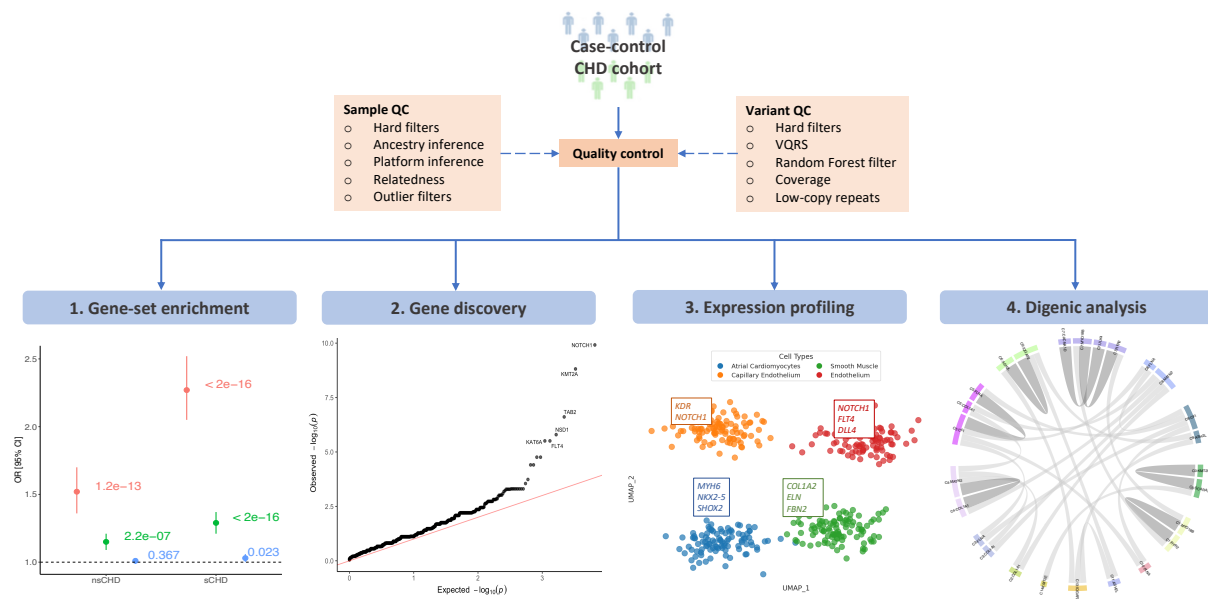
681 Pipelines for sample/variant quality control (QC), annotation, and burden testing are  
682 available on GitHub: [https://github.com/enriquea/wes\\_chd\\_ukbb](https://github.com/enriquea/wes_chd_ukbb).

683

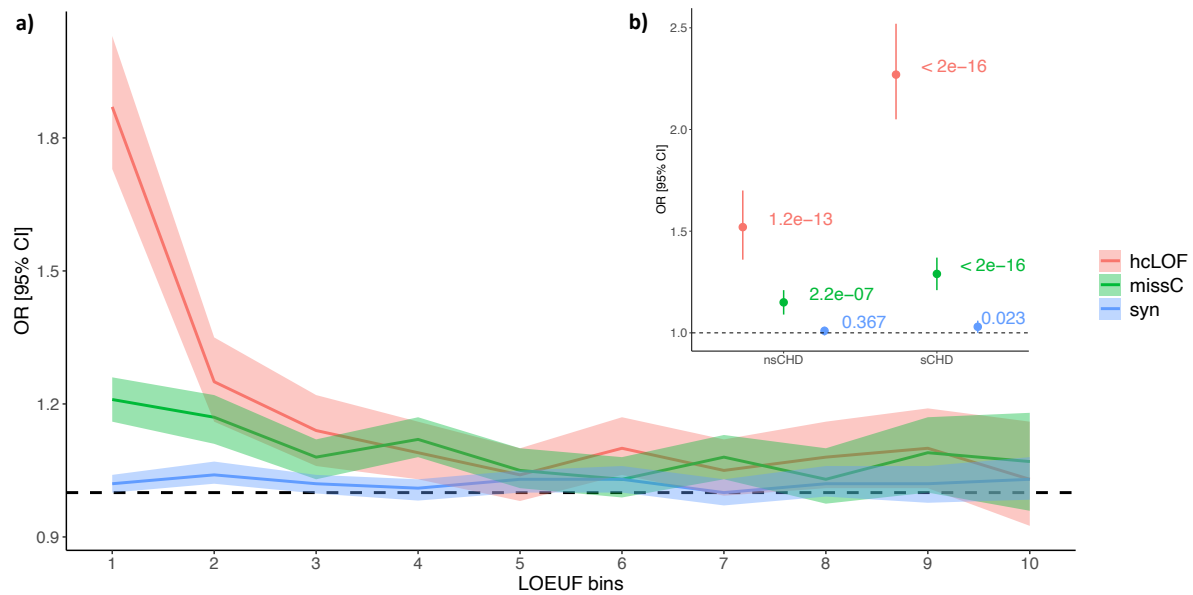
## 684 **Acknowledgements**

685 This research was conducted using the UKBB Resource under application number  
686 44165. We used data from the Deciphering Developmental Disorders (DDD) study.  
687 The DDD study presents independent research commissioned by the Health  
688 Innovation Challenge Fund, a parallel funding partnership between the Wellcome  
689 Trust and the UK Department of Health, and the Wellcome Trust Sanger Institute. The  
690 views expressed in this publication are those of the author(s) and not necessarily those  
691 of the Wellcome Trust or the UK Department of Health. The authors wish to thank Prof.  
692 Matthew Hurles (Sanger Institute, UK) for his significant contribution to this study. We  
693 thank the KinderHerzen e. V. for providing research funding for this study. We thank  
694 Prof. Dr. Christian Gilissen (RadboudUMC), Prof. Dr. Peter Krawitz (University Boon),  
695 and collaborators from the Universitaetsklinikum Tuebingen (Prof. Dr. Stephan  
696 Ossowski, Prof. Dr. Olaf Horst Rieß, Prof. Dr. Tobias Haack), for providing us with  
697 Central European allele frequencies. This work was partly funded by PROCEED  
698 project ERA PerMED joint Translational Call Initiative (DLR Funding reference  
699 number: 01KU1919).

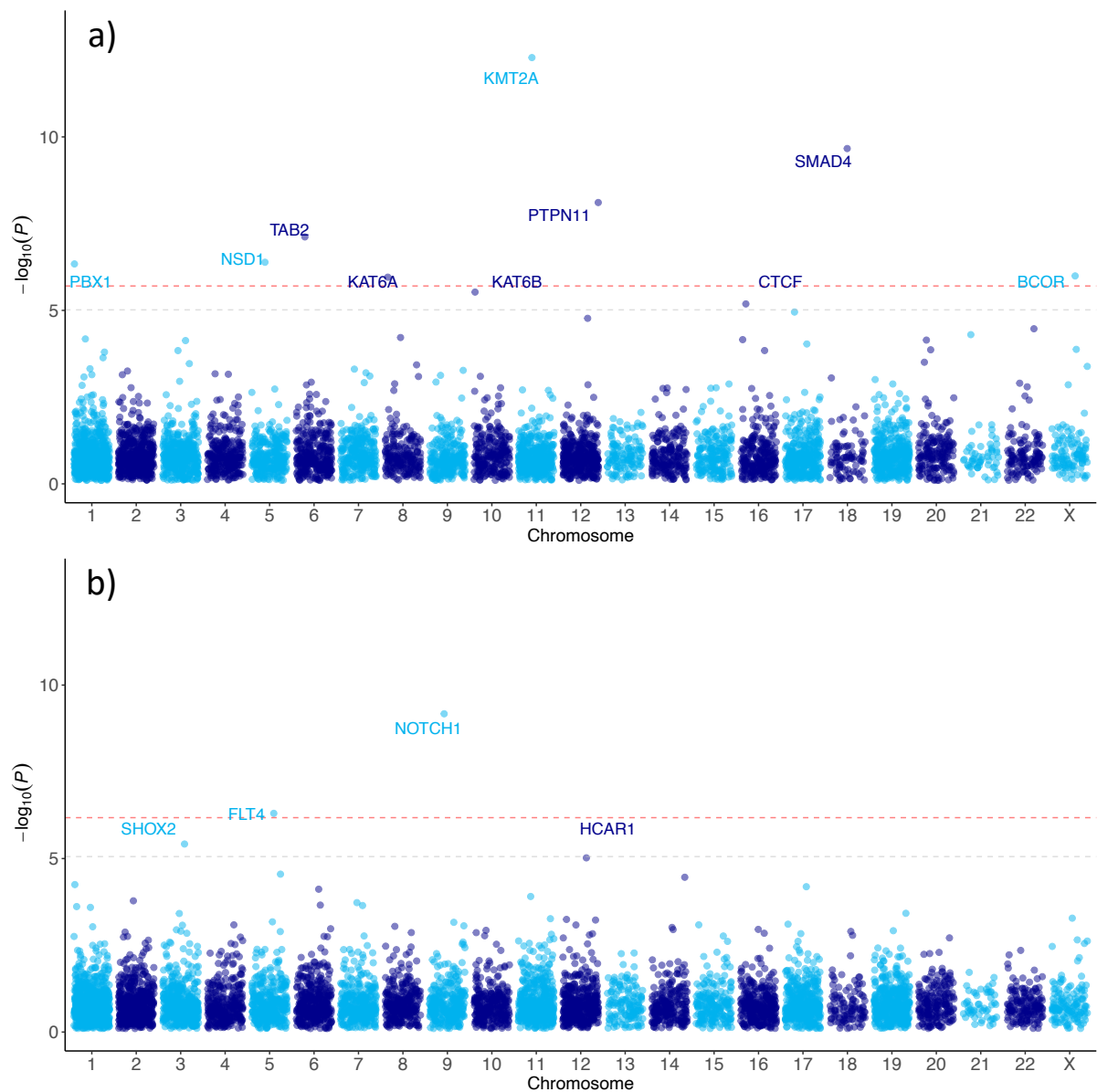
## Main figures



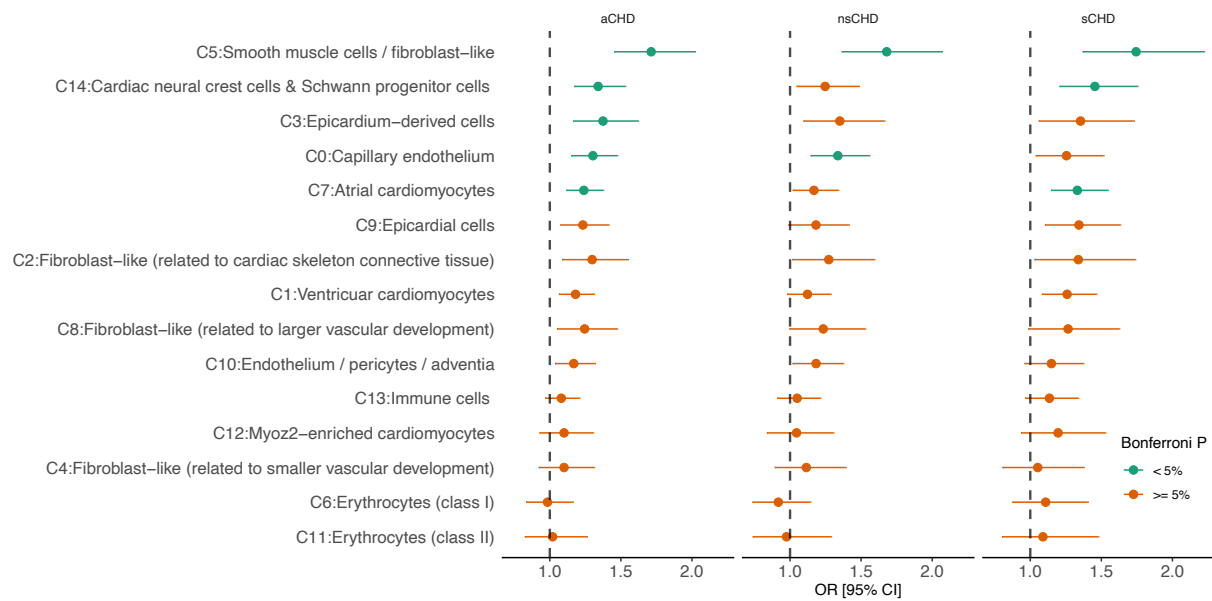
**Figure 1.** Analysis workflow for disease gene discovery. Quality control processes were conducted at the sample and variant levels. **(1)** Gene-set enrichment analysis was performed on the gene intolerance to missense and loss-of-function spectrum. **(2)** Gene-based case-control burden testing (Fisher's Exact test) was performed for high-confidence loss-of-function (hcLOF) and missense constrained variants (missC) independently. The per gene minimal  $p$ -value ( $P$ ) from both analyses was set as the study-wide  $p$ -value, corrected for multiple testing using the Bonferroni and B-H methods. **(3)** Expression profiling of significant CHD genes differentially expressed on cardiac specific cell clusters **(4)** Digenic analysis was conducted by comparing the rate of mutations observed on cases compared to controls. All analysis were stratified by syndromic status (aCHD, sCHD and nsCHD) vs control.



**Figure 2.** Enrichment analysis across the LOF constraint gene spectrum. Protein-coding genes were binned based on the LOEUF metric as proposed by gnomAD. Every bin contains ~1,900 genes. Top bins (1, 2) contain genes with the highest intolerance to loss-of-function. **a)** Enrichment analysis comparing aCHD vs controls. **b)** Enrichment analysis stratified by syndromic status (sCHD and nsCHD) vs controls in the top constraint LOF bin (1). The x-axis indicates the constraint bins; the y-axis shows the Odd Ratios (OR) and the 95% confidence interval.

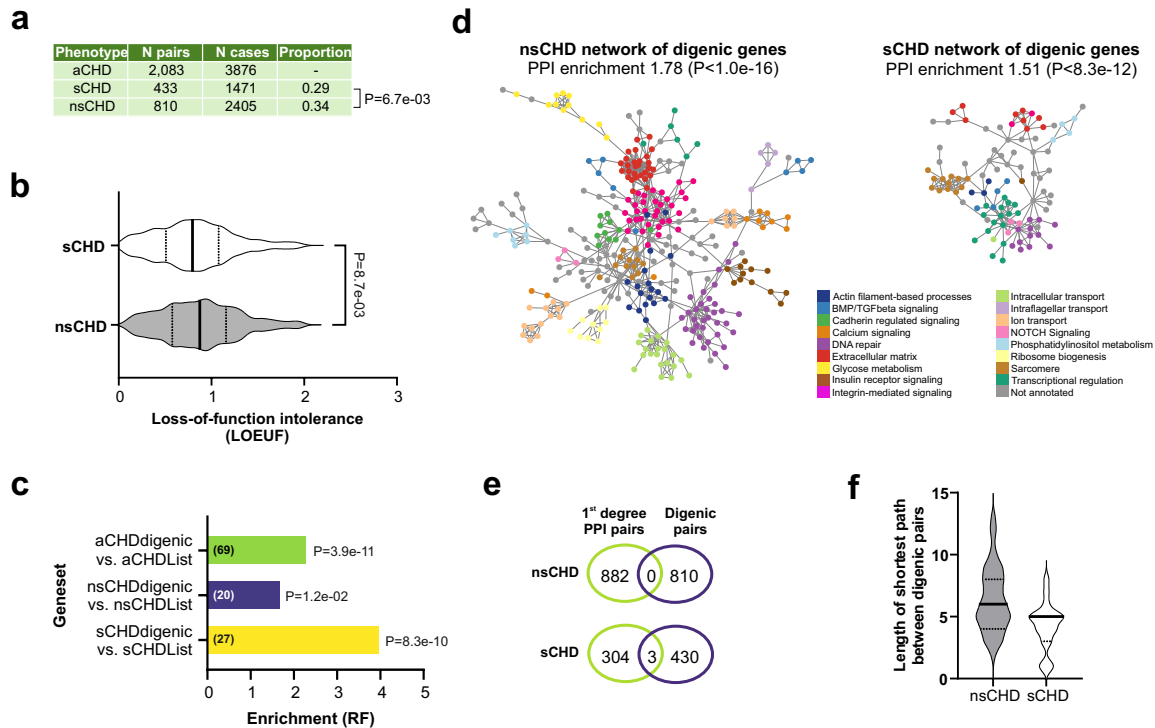


**Figure 3.** Log-transformed minimal p-value ( $P$ ) per gene (y-axis) against its chromosomal location (x-axis). Red dashed line denotes the threshold for genes reaching exome-wide significance (Bonferroni adjusted  $P < 0.05$ ); grey dashed line marks the threshold for genes reaching suggestive exome-wide significance (FDR 5%). a) Burden analysis of sCHD vs controls; b) burden analysis of nsCHD vs controls.



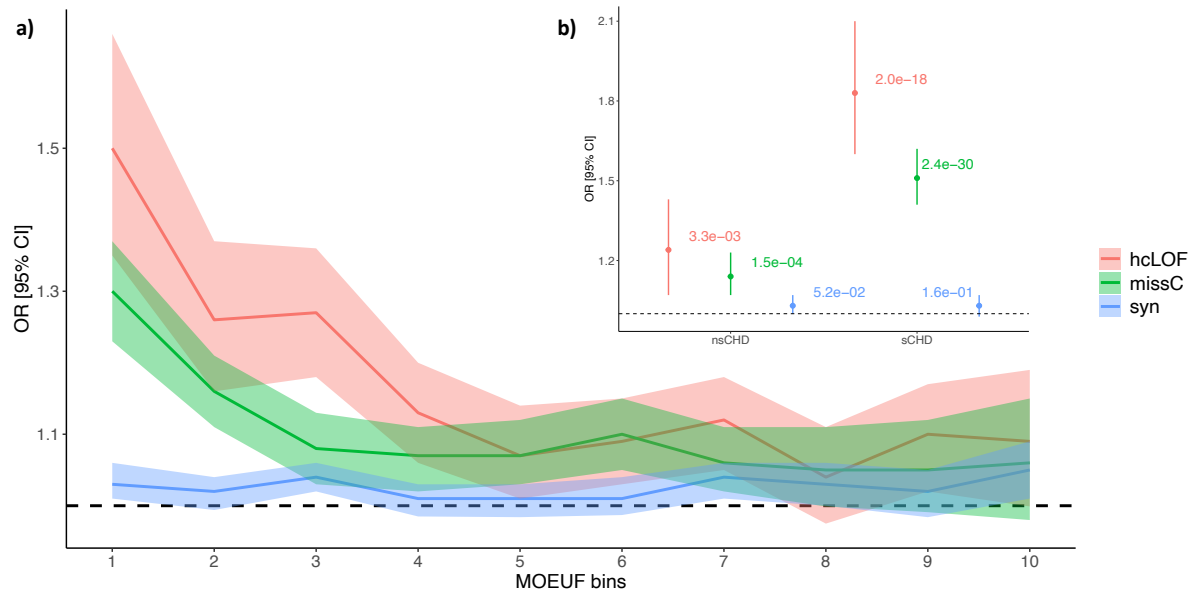
**Figure 4.** Logistic regression-based enrichment analysis of differentially expressed genes (DEGs) in cardiac-specific cell clusters for high-confidence loss-of-function variants (hcLOF). The analysis was stratified by syndromic status (aCHD, sCHD and nsCHD). The x-axis denotes the Odds Ratio (OR) and the 95% confidence interval. *P-values* were adjusted using the Bonferroni method (0.05 / 45 tests) to assess significant enrichment. Cardiac cell clusters C0, C3, C5, C7 and C14, show significant enrichment when analysing aCHD vs controls. The enrichment observed in clusters C7 and C14 showed a major contribution of sCHD. In comparison, cluster C0 provided the major contribution to nsCHD.





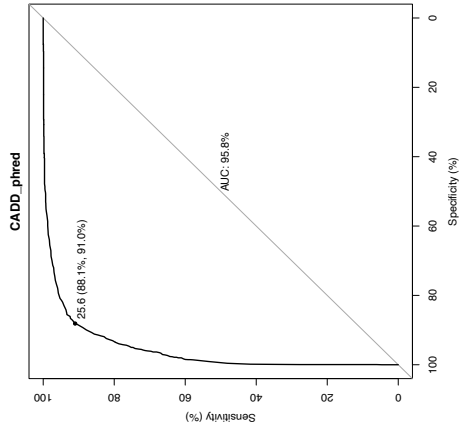
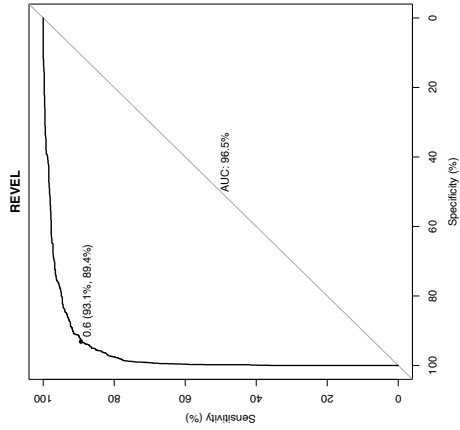
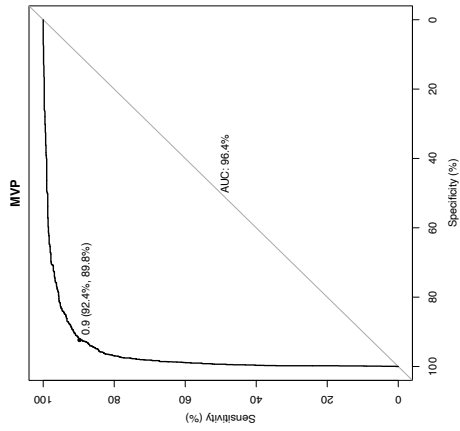
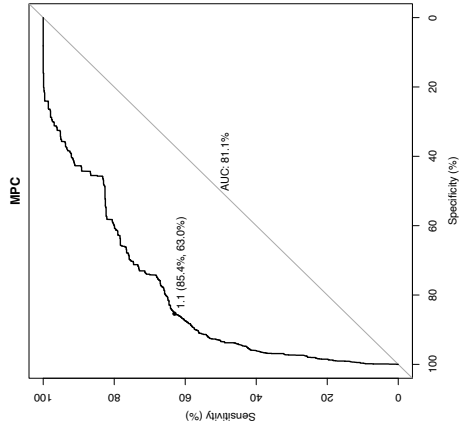
**Figure 5.** Case-control enrichment analysis at the digenic level using the *RareComb* framework. **(a)** Proportion of digenic pairs contributing to syndromic and non-syndromic CHD identified at FDR 1%. **(b)** Comparison of the distribution of LOEUF metric (at gene level) between syndromic and non-syndromic. The analysis was performed on genes observed in just one digenic pair. **(c)** Overlap of genes forming digenic interactions with known CHD genes. **(d)** Protein-protein interaction (PPI) network analysis of digenic genes for sCHD and nsCHD (PPI networks with annotated gene names are shown in **Supplemental Figure 10**). **(e)** Overlap between digenic pairs and first-degree interactors for sCHD and nsCHD in its respective PPI networks. **(f)** Length of the shortest path observed between genes forming a digenic pair in a PPI network.

## Supplemental figures

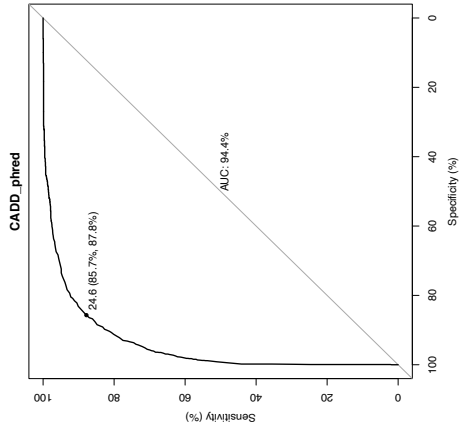
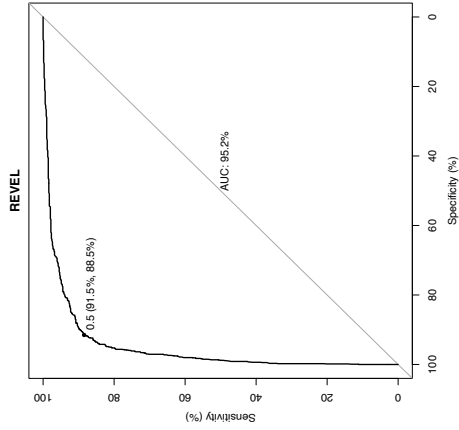
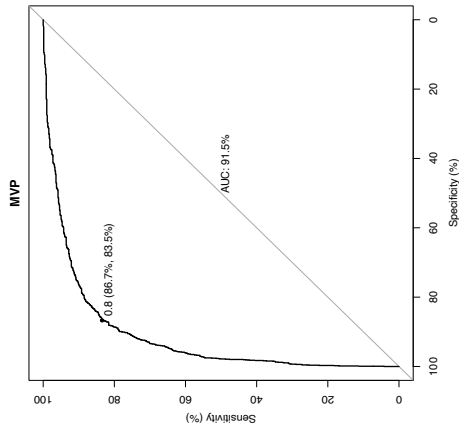
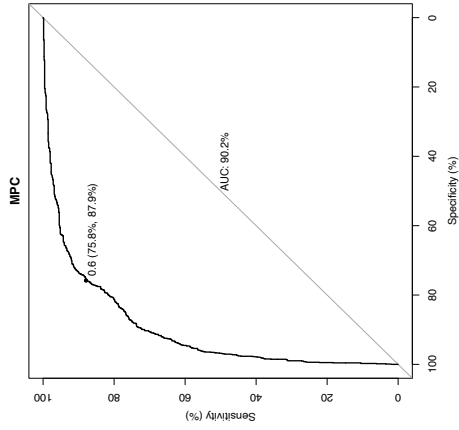


**Supplemental Figure 1.** Enrichment analysis across the missense constraint gene spectrum. Protein-coding genes were binned based on the MOEUF metric as proposed by gnomAD. Every bin contains ~1,900 genes. Top bins (1, 2) contain the genes with the highest intolerance to missense variation. **a)** Enrichment analysis per bin for aCHD vs controls are shown. **b)** Enrichment analysis stratified by syndromic status (sCHD and nsCHD) vs controls in the top constraint MOEUF bin (1). The x-axis indicates the constraint bins; the y-axis shows the Odd Ratios (OR) and the 95% confidence interval. hcLOF: high-confidence loss-of-function variants; missC: missense constrained variants; syn: synonymous variants.

LOF constraint genes

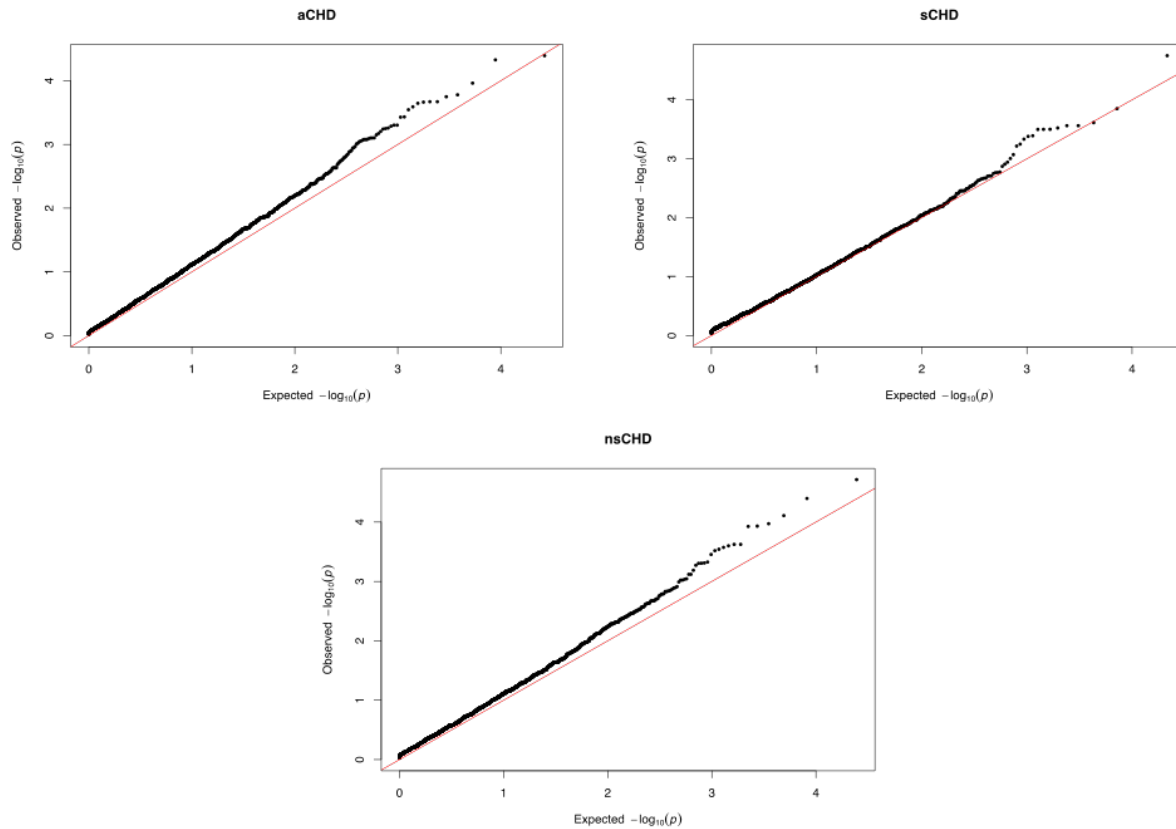


LOF non-constraint genes

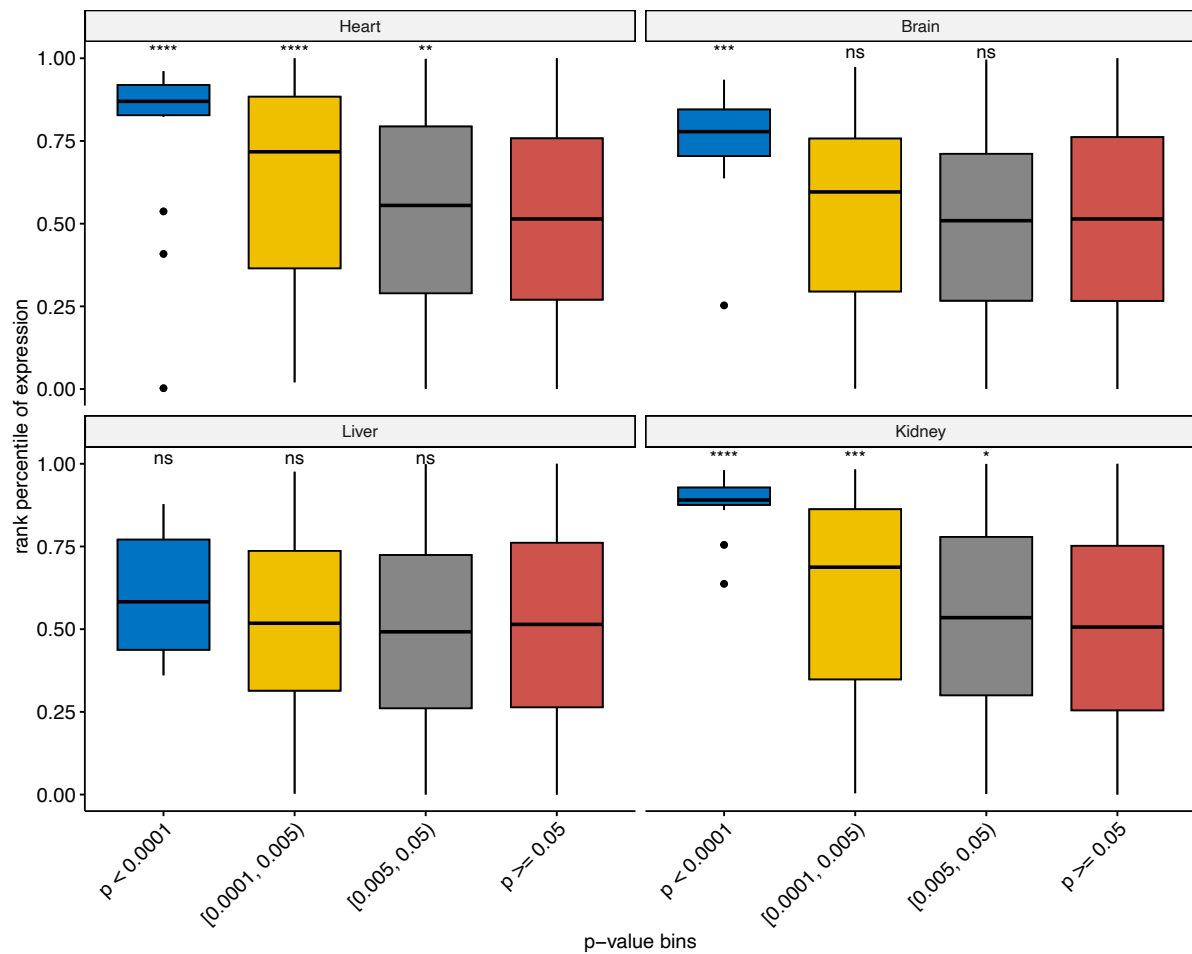


(Supplemental Figure 2. Legend on next page)

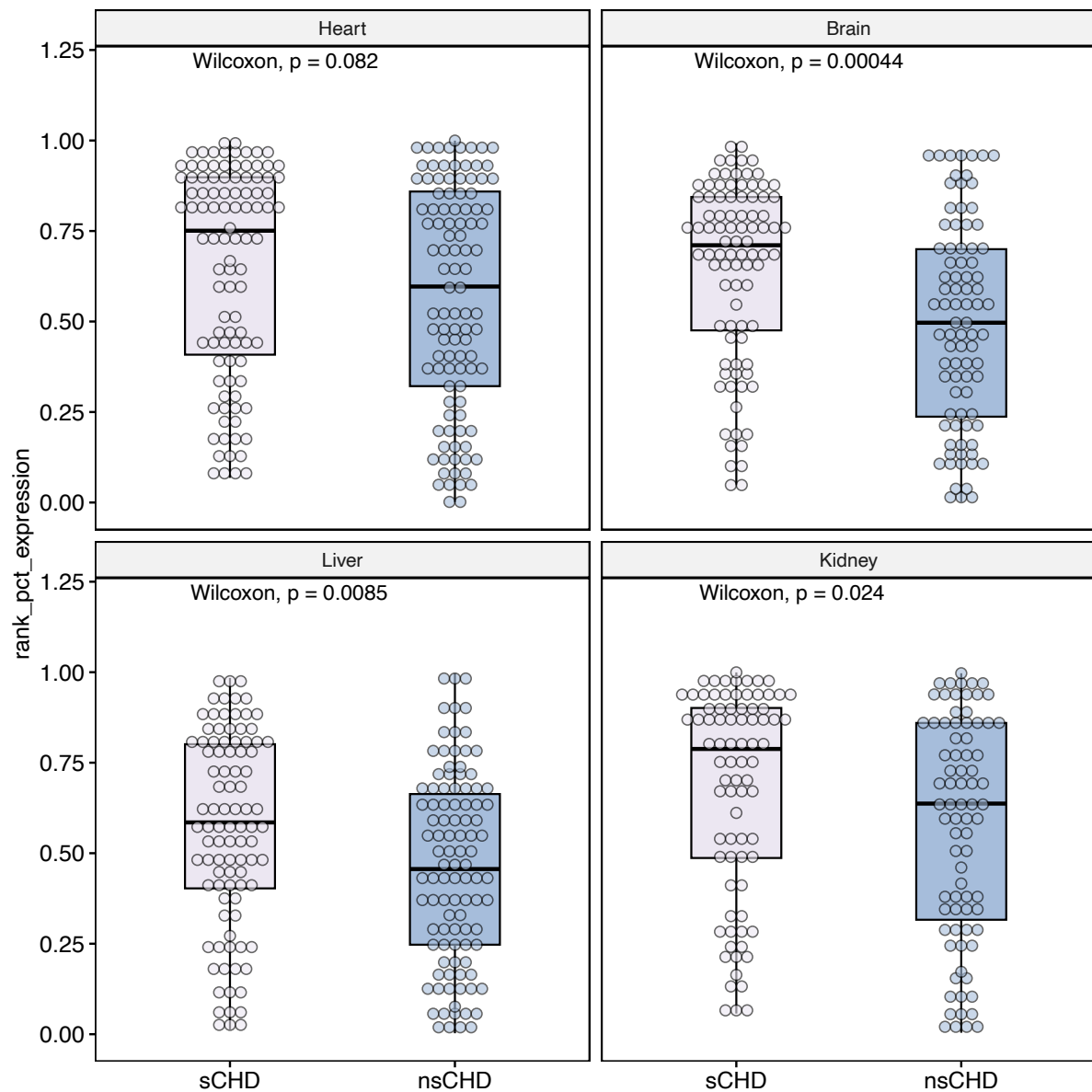
**Supplemental Figure 2.** ROC analysis of pathogenicity prediction scores (CADD, REVEL, MVP and MPC). The analysis was performed on a balanced set of benign (true negative) and likely pathogenic (true positive) variants from the ClinVar database within known CHD genes. The top panels show the results for LOF constraint genes (LOUEF < 0.35). The bottom panels show the results for LOF non-constraint genes (LOUEF >= 0.35).



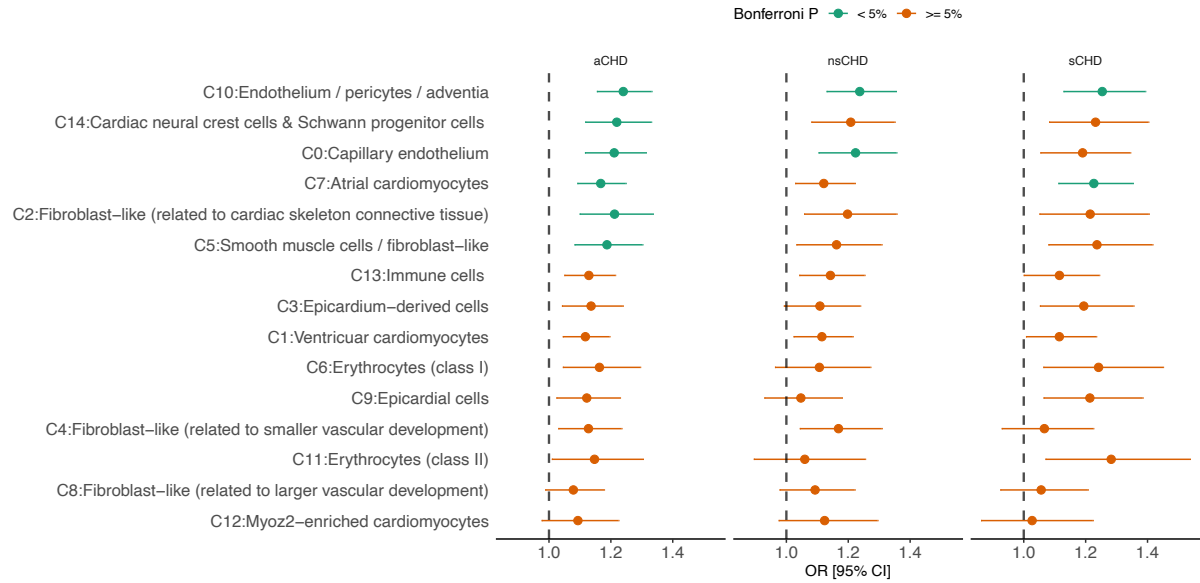
**Supplemental Figure 3.** Quantile-quantile plots. Expected vs observed p-values for synonymous variants stratified by syndromic status (MAF 0.1%). Q-Q plots for aCHD, sCHD and nsCHD vs controls are shown.



**Supplemental Figure 4.** Expression pattern of CHD genes in different tissues (Heart, Brain, Liver and Kidney). X-axis denotes gene p-value bins. P-value refers to the minimal p-value ( $P$ ) observed in the gene-based enrichment analysis for rare hcLOF and missC variants. The gene association analysis was performed by comparing all CHD probands (aCHD) vs controls. Y-axis denotes tissue-specific percentile rank of mean expression. Averaged expression was computed for samples between 4-8 weeks-post-conception (developmental stage). More significant genes (blue box) in the CHD case-control analysis showed the higher expression rank (e.g., Heart, Brain, and Kidney). Mean comparisons between bins were computed using the Wilcoxon test (alternative: greater; reference group (red box): genes with  $P > 0.05$  in the case-control analysis). ns:  $p > 0.05$ ; \*:  $p \leq 0.05$ ; \*\*:  $p \leq 0.01$ ; \*\*\*:  $p \leq 0.001$ ; \*\*\*\*:  $p \leq 0.0001$ .

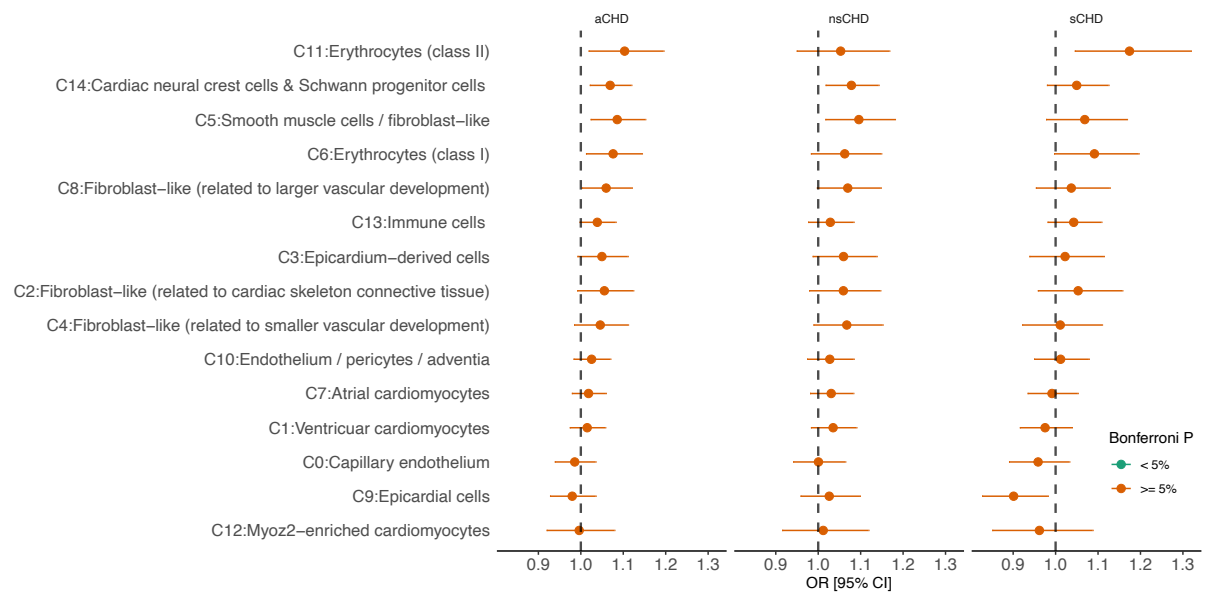


**Supplemental Figure 5.** Tissue-specific expression pattern of CHD genes identified in the case-control analysis, stratified by syndromic status (syndromic (sCHD) and non-syndromic (nsCHD) vs controls). Only genes with unadjusted  $P < 0.005$  in the case-control analysis are included. X-axis denotes the probands used in the case-control analysis (sCHD or nsCHD vs controls). Y-axis denotes tissue-specific percentile rank of mean expression. Averaged expression was computed among samples between 4-8 weeks-post-conception (developmental stage). Mean comparisons between groups were computed using the Wilcoxon test (two-sided). No significant difference was observed in the Heart for sCHD and nsCHD genes ( $P > 0.05$ ), compared to other tissues (e.g. brain, liver and kidney,  $P < 0.05$ ).

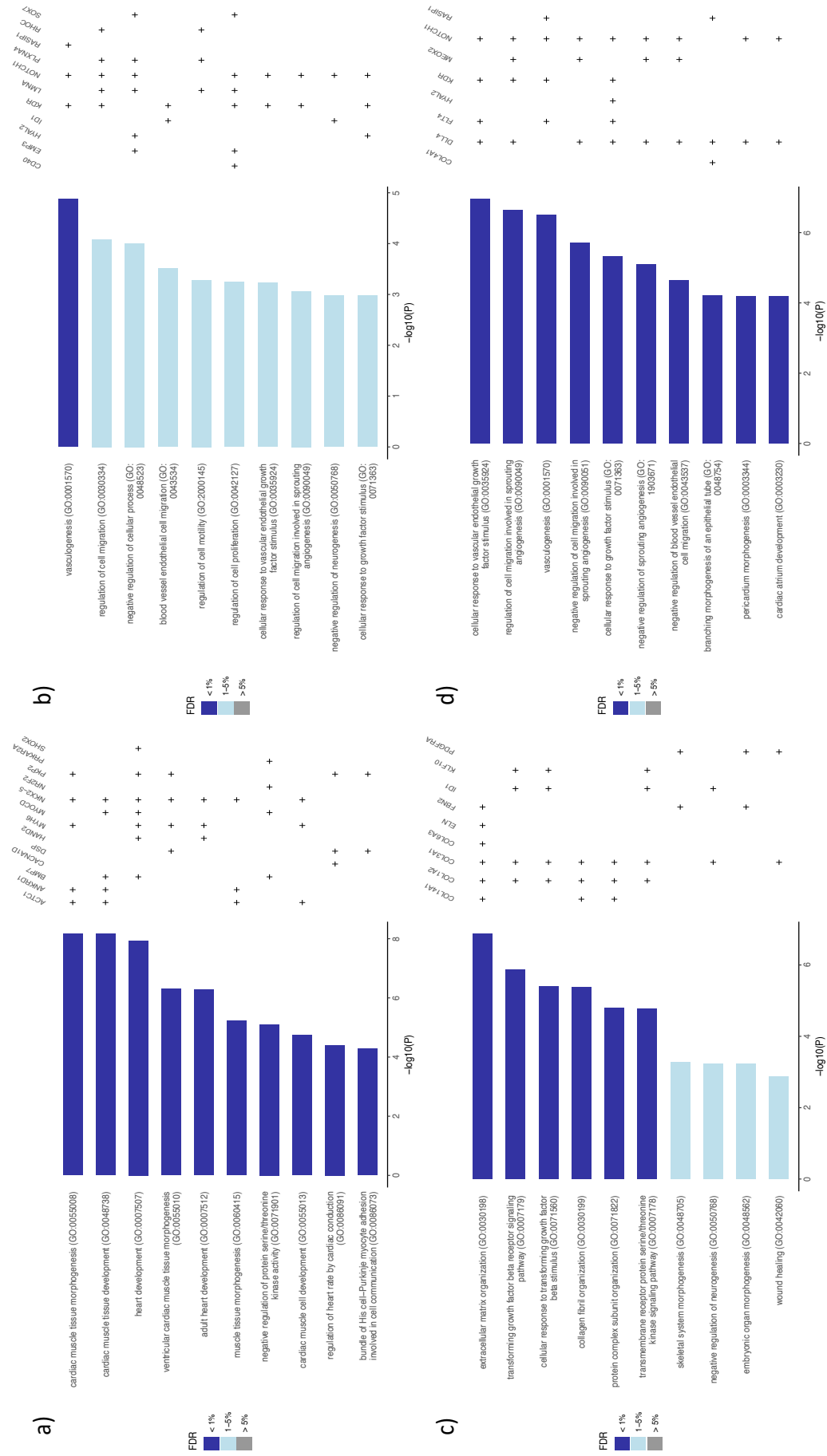


**Supplemental Figure 6.** Logistic regression-based enrichment analysis of differentially expressed genes (DEGs) in cardiac-specific cell clusters for missense constrained variants (missC). The analysis was stratified by syndromic status (aCHD, sCHD and nsCHD). The x-axis denotes the Odds Ratio (OR) and the 95% confidence interval. *P-values* were adjusted using the Bonferroni method (0.05 / 45 tests) to assess for significant enrichment.



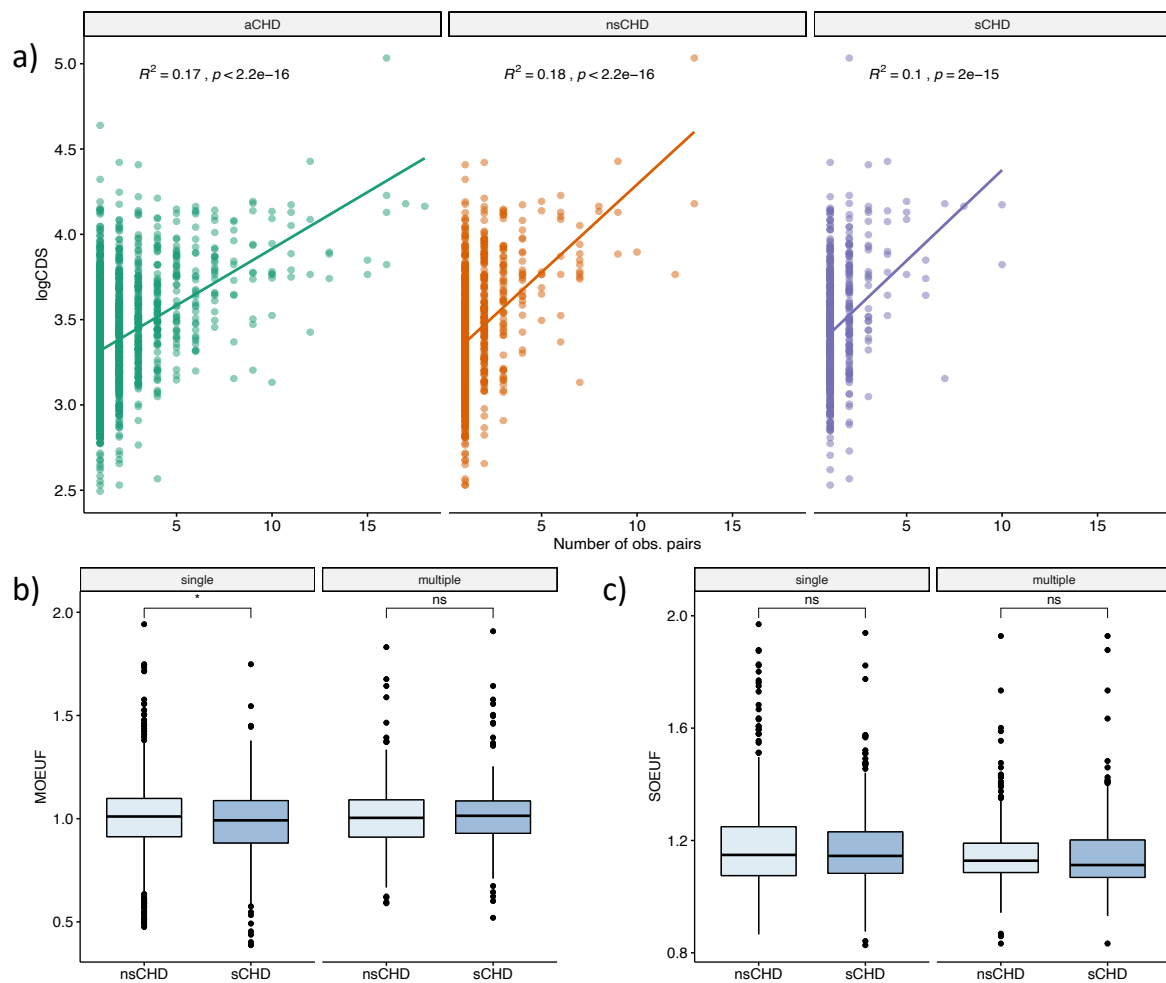


**Supplemental Figure 7.** Logistic regression-based enrichment analysis of differentially expressed genes (DEGs) in cardiac-specific cell clusters for synonymous variants. The analysis was stratified by syndromic status (aCHD, sCHD and nsCHD). The x-axis denotes the Odds Ratio (OR) and the 95% confidence interval. *P-values* were adjusted using the Bonferroni method (0.05 / 45 tests) to assess for significant enrichment.

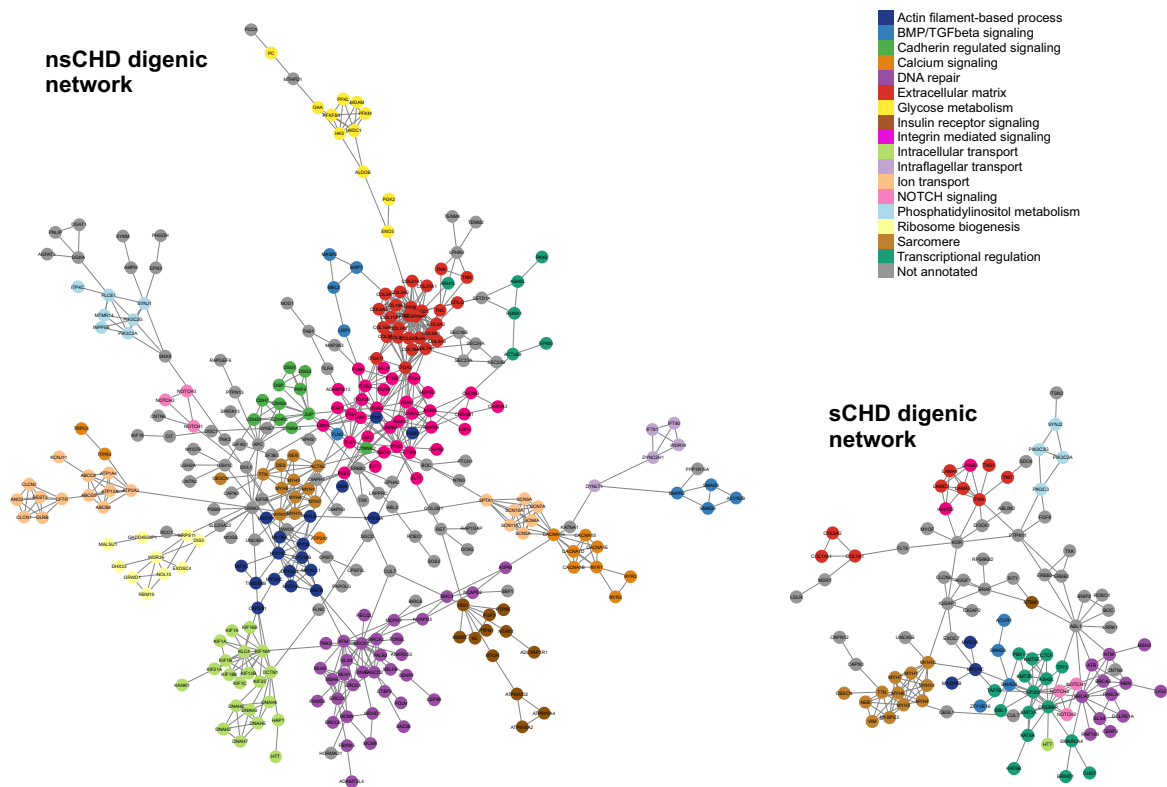


(Supplemental Figure 8. Legend on next page)

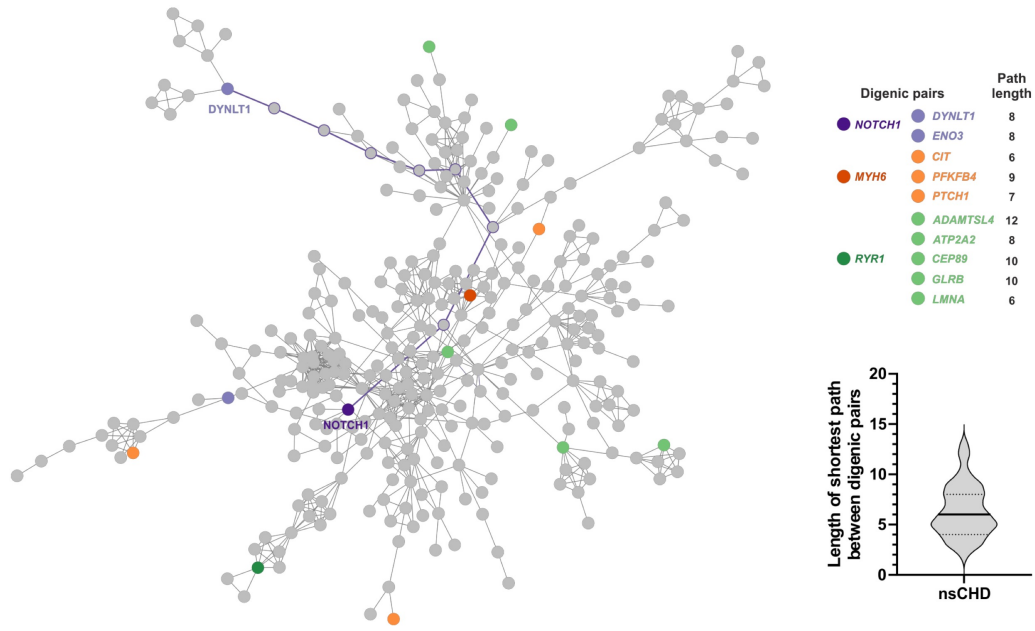
**Supplemental Figure 8.** Gene Ontology (GO) enrichment analysis of differentially expressed genes (DEGs) in cardiac-specific cells with unadjusted  $P < 0.01$  in the case-control burden analysis. a) C7: atrial cardiomyocytes cells, b) C0: capillary endothelium, c) C5: smooth muscle cells and d) C10: endothelium and pericytes cells. Only clusters with at least one GO term with  $FDR < 1\%$  are shown. For every GO term, the overlapping DE genes (+) are shown.



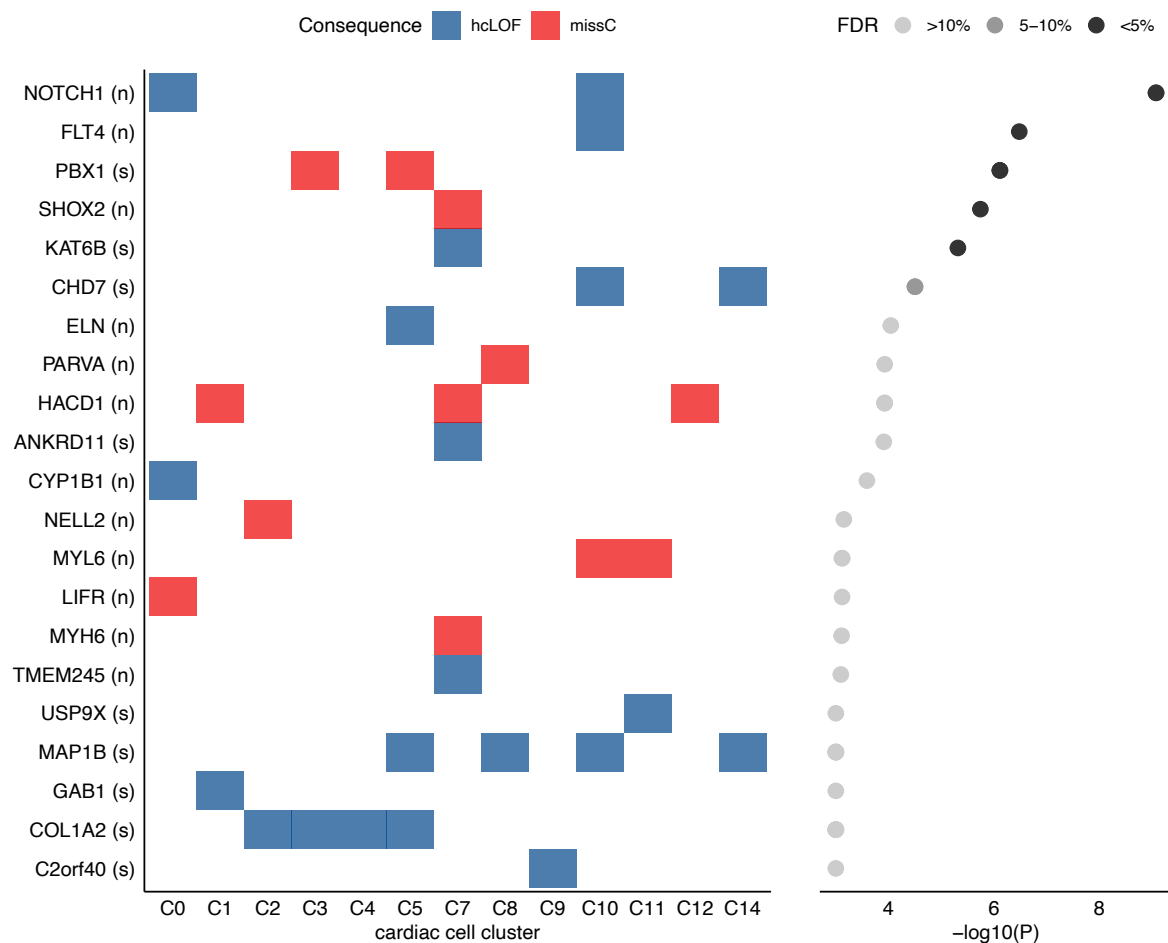
**Supplemental Figure 9.** a) Correlation between the frequency at which a gene is observed forming multiple digenic pairs (x-axis) vs. its log-transformed CDS length (y-axis). The correlation analysis was stratified by syndromic status (aCHD, sCHD and nsCHD). b) Comparison of the distribution of missense observed/expected ratio upper fraction (MOEUF) metric (at gene level) between syndromic and non-syndromic. The analysis is stratified further into single-genes (i.e., genes observed in just one digenic pairs) or multiple-genes (i.e., genes that appears two or more times forming digenic pairs). c) Same as (b), but for the synonymous observed/expected ratio upper fraction (SOEUF).



**Supplemental Figure 10.** Protein-protein interaction network for syndromic and non-syndromic CHD digenic genes. Nodes are labeled with the corresponding gene name and annotated with the specific biological process.



**Supplemental Figure 11.** Digenic pairs are scattered across protein-protein interaction (PPI) network. Examples of non-syndromic digenic pairs and its shortest paths are highlighted in the network. The median of the shortest path between non-syndromic digenic pairs was six (violin plot).



**Supplemental Figure 12.** Top enriched genes (unadjusted  $P < 0.001$ , case-control Fisher Exact test) found differentially expressed in at least one cardiac-specific cell cluster. The left plot shows the gene/cluster overlap and highlights the variant category with the highest enrichment (blue: hcLOF, red: missC). The x-axis denotes de cardiac clusters; the y-axis indicates the genes and the CHD category analysed (s: sCHD, n: nsCHD). The right plot shows the log-transformed  $P$  (x-axis) and the  $FDR$  significant level per gene. Six genes showed  $FDR < 10\%$ : *NOTCH1*, *FLT4*, *PBX1*, *SHOX2*, *KAT6B* and *CHD7*. C0: Capillary endothelium, C1: Ventricular cardiomyocytes, C2: Fibroblast-like (related to cardiac skeleton connective tissue), C3: Epicardium-derived cells, C4: Fibroblast-like (related to smaller vascular development), C5: Smooth muscle cells, C7: Atrial cardiomyocytes, C8: Fibroblast-like (related to larger vascular development), C9: Epicardial cells, C10: Endothelium/pericytes/adventia, C11: Erythrocytes (class II), C12: Myoz2-enriched cardiomyocytes, C14: Cardiac neural crest cells & Schwann progenitor cells.