

1 **Systems Analysis of de novo Mutations in Congenital Heart Diseases Identified a**  
2 **Molecular Network in Hypoplastic Left Heart Syndrome**

3 Yuejun Jessie Wang<sup>1</sup>, Xicheng Zhang<sup>2</sup>, Chi Keung Lam<sup>3,4</sup>, Hongchao Guo<sup>3,4</sup>, Cheng  
4 Wang<sup>1</sup>, Sai Zhang<sup>2</sup>, Joseph C. Wu<sup>3,4,5</sup>, Michael Snyder<sup>2,3,\*</sup>, and Jingjing Li<sup>1,\*</sup>

5  
6 <sup>1</sup>the Eli and Edythe Broad Center of Regeneration Medicine and Stem Cell Research,  
7 the Bakar Computational Health Sciences Institute, the Parker Institute for Cancer  
8 Immunotherapy, and the Department of Neurology, School of Medicine, University of  
9 California, San Francisco, 35 Medical Center Way, San Francisco, CA 94143

10 <sup>2</sup>Department of Genetics and the Center for Genomics and Personalized Medicine,  
11 School of Medicine, Stanford University, 291 Campus Dr., Stanford, CA 94305

12 <sup>3</sup>Stanford Cardiovascular Institute, School of Medicine, Stanford University, 265  
13 Campus Dr., Stanford, CA 94305

14 <sup>4</sup>Department of Medicine, Division of Cardiology, School of Medicine, Stanford  
15 University, 265 Campus Dr., Stanford, CA 94305

16 <sup>5</sup>Department of Radiology, Stanford University School of Medicine, Stanford University,  
17 265 Campus Dr., Stanford, CA 94305

18

19 \* correspondence should be addressed to:

20 MS: [mpsnyder@stanford.edu](mailto:mpsnyder@stanford.edu)

21 JL: [Jingjing.Li@ucsf.edu](mailto:Jingjing.Li@ucsf.edu)

22

23

24 **Abstract**

25 Congenital heart diseases (CHD) are a class of birth defects affecting ~1% of live births.  
26 These conditions are hallmarked by extreme genetic heterogeneity, and therefore,  
27 despite a strong genetic component, only a very handful of at-risk loci in CHD have  
28 been identified. We herein introduced systems analyses to uncover the hidden  
29 organization on biological networks of genomic mutations in CHD, and leveraged  
30 network analysis techniques to integrate the human interactome, large-scale patient  
31 exomes, the fetal heart spatial transcriptomes, and single-cell transcriptomes of clinical  
32 samples. We identified a highly connected network in CHD where most of the member  
33 proteins had previously uncharacterized functions in regulating fetal heart development.  
34 While genes on the network displayed strong enrichment for heart-specific functions, a  
35 sub-group, active specifically at early developmental stages, also regulates fetal brain  
36 development, thereby providing mechanistic insights into the clinical comorbidities  
37 between CHD and neurodevelopmental conditions. At a small scale, we experimentally  
38 verified previously uncharacterized cardiac functions of several novel proteins  
39 employing cellular assays and gene editing techniques. At a global scale, our study  
40 revealed developmental dynamics of the identified CHD network and observed the  
41 strongest enrichment for pathogenic mutations in the network specific to hypoplastic left  
42 heart syndrome (HLHS). Our single-cell transcriptome analysis further identified  
43 pervasive dysregulation of the network in cardiac endothelial cells and the conduction  
44 system in the HLHS left ventricle. Taken together, our systems analyses identified novel  
45 factors in CHD, revealed key molecular mechanisms in HLHS, and provides a  
46 generalizable framework readily applicable to studying many other complex diseases.

## 47 **Introduction**

48 Congenital heart diseases (CHD), broadly defined by the structural and functional  
49 abnormalities in fetal heart, are the most common forms of birth defects and affects ~1%  
50 of live births<sup>1,2</sup>. Although CHD has a strong genetic component<sup>3,4</sup>, the underlying genetic  
51 basis has largely remained elusive. Like many other pediatric diseases, large-scale  
52 copy number variations (together with aneuploidies) potentially explain ~20% of CHD  
53 cases<sup>4,5</sup>, and cases with monogenic causes could be solved by familial analyses<sup>6,7</sup>.  
54 However, the genetic basis of sporadic cases, accounting for the majority of CHD  
55 probands, has largely remained unclear<sup>3,8</sup>. The Pediatric Cardiac Genomics Consortium  
56 (PCGC) aims to fill in the knowledge gap by performing whole exome/genome  
57 sequencing on large-scale patient samples representing major sporadic CHD sub-  
58 types<sup>9</sup>. The latest PCGC study analyzed de novo and rare variants in the whole-exome  
59 data from 2,871 CHD probands and identified seven genes achieving genome-wide  
60 significance, together with a handful of genes showing suggestive associations with  
61 CHD<sup>10</sup>. These analyses collectively explained ~10% of CHD cases in the cohort<sup>10</sup>.  
62 Targeting common variants, the latest genome-wide association analysis only identified  
63 one SNP reaching genome-wide significance<sup>11</sup>. Given the strong genetic basis in CHD,  
64 its complete genetic architecture has been yet to be discovered.

65

66 The existing analytical frameworks have been largely based on mutational recurrence  
67 analysis, where genes recurrently affected in probands than controls were identified for  
68 disease associations. However, in real clinical situations, different patients usually carry  
69 different sets of clinical mutations, and genes are more often individually than

70 recurrently affected in patient populations. Importantly, these seemingly heterogeneous  
71 mutations usually functionally conserved onto common molecular pathways, giving rise  
72 to similar clinical phenotypes<sup>12-15</sup>. This is particularly the case in CHD, where its risk  
73 factors were more likely to affect different components in a shared molecular network,  
74 as opposed to recurrently affecting the same genes among patient populations<sup>16,17</sup>. As  
75 such, it is not individual genes or mutations, but their structural organization on  
76 biological networks that defines the complete genetic architecture of the disease.

77

78 Our recent work has proposed a series of theoretic models to dissect convergent  
79 pathways in complex diseases from biological networks<sup>12</sup>. We herein leverage this  
80 system thinking to study CHD genomes. We developed a new framework to integrate  
81 network biology, genome analysis, spatial transcriptomics, and single-cell analysis for a  
82 direct revelation of the genetic basis in CHD. Combining computational and  
83 experimental approaches, we identified a highly connected cluster from the global  
84 human protein interaction network that was strongly implicated in fetal heart  
85 development and displayed significant enrichment for pathogenic mutations in CHD  
86 probands. Analyzing different CHD sub-types, the identified network was strongly  
87 associated with the hypoplastic left heart syndrome (HLHS) and displayed substantial  
88 dysregulation in cardiac endothelial cells and the conduction system in the under-  
89 developed left ventricle of the HLHS heart. We particularly note that the majority of the  
90 newly identified genes in this study had previously uncharacterized cardiac functions,  
91 nor their implications in CHD. Overall, our work provides a new systems framework for



92 CHD genome analysis and has significantly expanded our knowledge about the genetic  
93 architecture of CHD.

## 94 **Methods and Materials**

### 95 **An overview of the genomic resources**

96 We analyzed 2,990 de novo mutations in cases and 1,830 de novo mutations in controls.  
97 These mutations were identified from the previous whole-exome sequencing study  
98 encompassing 2,871<sup>10</sup> probands and 1,789 control individuals<sup>18</sup>, where the control  
99 sibling subjects were unrelated individuals to the probands. Among all the de novo  
100 mutations, we considered 323 and 129 loss-of-function (LoF) mutations (i.e. stop gain  
101 and loss mutations, and frameshift indels) in cases and controls, respectively given their  
102 clear functional consequences compared with missense mutations. At the gene level,  
103 we considered those intolerant to LoF mutations with a pLI score<sup>19</sup> greater than 0.8, and  
104 therefore the presence of de novo LoF mutations in these genes is a strong indicator of  
105 mutational pathogenicity. With this, we identified 120 and 35 genes in the proband and  
106 control cohorts, respectively (**Table S1**), for subsequent analyses. Gene functional  
107 enrichment test confirmed a significant overrepresentation of genes regulating heart  
108 development in the 120 proband genes but not in the control genes, and we therefore  
109 considered them as CHD candidate genes. Gene function enrichment analyses  
110 throughout this work were performed using Enrichr<sup>20</sup> (Database: as of Oct, 2021).

111

### 112 **Analysis of the protein interaction network**

113 We seeded the 120 proband proteins in the human protein interaction network. The  
114 network was analyzed in our previous publication<sup>13</sup> with 16,085 unique proteins (or

115 genes) and 217,605 interactions. We implemented the random walk algorithm  
116 (personalized page rank<sup>21</sup>) on the network by setting the restart probability at 0.1 and  
117 the maximum number of iterations at 500. For each node on the network, we derived its  
118 probabilities of visiting all other nodes on the network and a greater probability indicates  
119 greater reachability between two nodes, thereby increased topological distance. To  
120 identify those topologically closer to the 120 CHD candidate proteins seeded on the  
121 network, we performed Wilcoxon rank-sum test to determine whether a given node is  
122 more reachable to the 120 CHD candidate proteins relative to its reachability to all other  
123 proteins on the network. We performed the test on each of the nodes on the network  
124 (excluding the 120 CHD proteins), and the derived P values were corrected by  
125 Benjamini-Hochberg (false discovery rates less than 0.05). Therefore, on the proteome  
126 scale, we agnostically identified another set of 120 new proteins forming a highly  
127 connected network the 120 seeded PCGC proteins. To determine the identification of  
128 the new proteins was not expected by chance, we performed degree-preserving  
129 shuffling<sup>22</sup> to permute the protein interaction network, and recorded the number of  
130 nodes in the largest connected component in each permutation simulation. The  
131 observation from the real dataset cannot be observed from the permutation test, thereby  
132 statistical significance suggesting biological functions.

133

### 134 **Gene expression analysis**

135 We analyzed the time-course gene expression data in the mouse cardiogenesis  
136 process<sup>23</sup>, where human-mouse orthology mapping was based on the Ensembl Biomart  
137 annotation, and we only considered those with one-to-one orthology mapping.

138 Expression of gene symbols mapped onto multiple probesets identifiers were averaged.  
139 Gene expression was then normalized across time points followed by a hierarchical  
140 clustering, revealing two expression clusters on the network, Group-I (G-I) and Group-II  
141 (G-II), in **Figure 2B**. We also analyzed gene expression in the developing fetal heart  
142 from postconceptional day 96 (gestational week 13.7) to 147 (gestational week 21)  
143 using microarray data from the NIH Roadmap Epigenomics Mapping Consortium<sup>24</sup>  
144 (GSE18927 in GEO). All probesets intensities were normalized onto a logarithm scale,  
145 and signals from probesets mapped onto the same gene symbol were averaged. At  
146 each time point, we compared expression of genes of interest against the transcriptome  
147 background to determine their molecular activities and only protein coding genes were  
148 used in this comparison. To confirm our observation, we also performed analysis using  
149 RNA-seq data in fetal heart samples in postconceptional weeks 19 and 28<sup>25</sup>  
150 (ENCSR000AEZ from the ENCODE consortium: <https://www.encodeproject.org>).  
151 Statistical comparisons were determined by the Wilcoxon rank-sum test.

152

### 153 **Spatial transcriptome analysis**

154 We analyzed the spatial transcriptome data in the fetal heart from a recent publication<sup>26</sup>,  
155 and analyzed spatial RNA-seq data in 3,115 tissue spots (across sections) at 4.5-5, 6.5  
156 and 9 postconceptional weeks (PCW). We standardized gene expression on spots from  
157 all three postconceptional weeks by NormalizeData (settings: normalization.method =  
158 "LogNormalize", scale.factor = 10000) from Seurat package<sup>27</sup> (version 4.0.4). We then  
159 performed quantile normalization across all tissue spots across all sections. In each  
160 tissue section, the original study clustered tissue spots into groups with shared

161 transcriptome profiles, which corresponded to different anatomical compartments in the  
162 developing heart. To quantify region-specific gene expression in the heart, we averaged  
163 expression of each gene across spots within the same anatomical compartments of a  
164 fetal heart across all sections. To determine statistical significance, we compared genes  
165 of interest against the transcriptome backgrounds (protein-coding genes) in each  
166 annotated anatomical compartments using Wilcoxon rank-sum test. For visualization,  
167 we used one representative tissue section with the most spots at each PCW.

168

### 169 **Whole-exome analysis of PCGC proband cohort**

170 To determine the enrichment of pathogenic mutations on this network genes among  
171 CHD proteins, we examined the PCGC whole-exome-sequencing data in dbGAP (as of  
172 Feb, 2021, dbGAP-24034, gap\_accession: phs000571, gap\_parent\_phs: phs001194,  
173 SRP025159). We used the same control subjects as described in the original study<sup>10,18</sup>,  
174 and downloaded the whole exome sequencing data from SFARI BASE  
175 (<https://www.sfari.org/resource/sfari-base/>). The analyzed individual identifiers were  
176 separately listed in **Table S6**. We downloaded the raw FastQ data files from dbGAP,  
177 and performed variant calls following the Best Practice procedure recommended by  
178 GATK. We utilized the Sentieon toolkit which substantially increases the computational  
179 efficiency while keeping genotyping accuracy<sup>28</sup>. We performed independent quality  
180 control analyses to ensure high quality of the called variants. We utilized VCFtools  
181 (<http://vcftools.sourceforge.net>) to compute the distribution of Ti/Tv ratios across all  
182 analyzed individuals (shown in **Figure S5**).

183

184 Using all the called exonic variants, we performed principal component analysis (PCA)  
185 by aggregating individuals from case and control cohorts. The PCA analysis ensured  
186 almost identical population structure between cases and controls. For rare variants, we  
187 only considered those absent from the 1000 Genome Database  
188 (<https://www.internationalgenome.org>). We elected to use 1000 Genome as our  
189 reference, instead of the Genome Aggregation Database (gnomAD<sup>29</sup>) or TOPMed<sup>30</sup>,  
190 because a significant portion of samples in gnomAD and TOPMed were individuals with  
191 cardiovascular diseases, such as the those in Atrial Fibrillation Genetics Consortium  
192 (AFGen) and Myocardial Infarction Genetics Consortium. We annotated the called  
193 variants using the ANNOVAR package<sup>31</sup>, which provided CADD (Combined Annotation  
194 Dependent Deletion) phased scores (Database: CADD v1.6 as of Feb, 2021) for each  
195 of the called variants<sup>32,33</sup>. For each personal exome in this study (in case and control  
196 cohorts), we computed the mean CADD scores for non-synonymous (LoF and  
197 missense) variants affecting the network genes, and then compared the mean CADD  
198 score distribution among all probands in each CHD subtype against individuals in the  
199 control cohort. The same comparison was also performed on synonymous variants as a  
200 set of negative controls. As another set of negative control, we obtained 62 lung-  
201 specific protein-coding genes from a previous publication<sup>34</sup>. To determine whether the  
202 observed mutational pathogenicity was specific to Group-II genes, we performed a  
203 permutation analysis, where, in each permutation, we randomly sampled rare non-  
204 synonymous variants from the exome background in the HLHS cohort, matching the  
205 number of rare non-synonymous variants in Group-II genes in the same HLHS cohort.  
206 We performed the same sampling on the control cohort, and then compared CADD

207 scores between the two randomly sampled variant list. We performed the permutation  
208 100 times, and confirmed that CADD scores were not statistically significant between  
209 the randomly sampled variant sets from case and control cohorts, respectively ( $p = 0.98$ ,  
210 permutation test).

211

### 212 **Analysis of single-cell RNA-seq data from an HLHS heart**

213 We re-analyzed the published single-cell RNA-seq data in an HLHS fetal heart<sup>35</sup>, and  
214 specifically performed our comparison on the under-developed left ventricle against that  
215 of the matched control sample. We followed the cell type clustering described in the  
216 original publication<sup>35</sup>. For each gene, we averaged its expression across all cells in a  
217 given cell type, and for each cell type, we asked whether the network genes (Group-I  
218 and II genes) tended to have increased expression relative to the transcriptome  
219 background. Wilcoxon rank-sum test was used to determine statistical significance.

220

### 221 **iPSC conversion to cardiomyocytes**

222 All induced pluripotent stem cell (iPSC) lines were reprogrammed by the Sendai virus  
223 expressing 4 Yamanaka factors (CytoTune®-iPS Sendai Reprogramming Kit, Invitrogen)  
224 as previously described<sup>36</sup>. Protocols in the study were approved by Human Subjects  
225 Research Institutional Review Board (IRB) at Stanford University. Human iPSCs were  
226 maintained in Essential 8™ Medium (Gibco®, Life Technology). For passaging, iPSC  
227 culture was dissociated with Accutase (Innovative Cell Technologies) at 37°C for 15-20  
228 min, and suspended iPSCs were reseeded on Matrigel-coated plates (BD Biosciences,  
229 San Jose, CA) at a density of 500 K cells per well in 6-well plates.

230

231 Beating induced pluripotent stem cell-derived cardiomyocytes (iPSC-CMs) were  
232 generated using the RPMI + B27 method as described<sup>37,38</sup>. Briefly, human iPSCs were  
233 kept in 6-well plates until passage 20. Once they reached > 80% confluence, the  
234 medium was switched to RPMI 1640 with Insulin minus B27 supplement (Gibco®, Life  
235 Technology) and 6  $\mu$ M CHIR99021 (Selleckchem) for 2 days. They were allowed for 1-  
236 day recovery with RPMI + B27 (minus insulin) medium. Cells were then treated with 5  
237  $\mu$ M IWR-1 (Sigma) for 2 days, and then fresh RPMI + B27 (minus insulin) medium for  
238 another 2 days. Cells were then switched to RPMI + B27 medium for 3 days. Beating  
239 cardiomyocytes were purified for 2-3 rounds of 3-day glucose-free RPMI + B27 medium  
240 treatment.

241

242 To downregulate the expression of our target genes, we performed siRNA transfection  
243 in iPSC-CMs as described previously<sup>39</sup>. 80  $\mu$ l of siRNA (Dharmacon) was added into a  
244 master mix of 3.2  $\mu$ l DharmaFECT1 (Dharmacon) transfection reagent and 236.8  $\mu$ l  
245 OptiMEM (Thermo Fisher Scientific) and incubated for 20 min before addition to a 6-well  
246 plate of iPSC-CMs. The cell media was then changed after 24 hr. Cells were then  
247 subjected to contractility assays or harvested 48 hr after medium change.

248

### 249 **Cellular contractility assays**

250 To assess iPSC-CM contractility, iPSC-CMs were re-seeded on Matrigel-coated 24-well  
251 plates and cultured for 7 days to recover their spontaneous beating, as previously  
252 described<sup>40</sup>. Contraction of monolayer cardiomyocytes was recorded with high

253 resolution motion capture tracking using the SI8000 Live Cell Motion Imaging System  
254 (Sony Corporation). During data collection, cells were maintained under controlled  
255 humidified conditions at 37°C with 5% CO<sub>2</sub> and 95% air in a stage-top microscope  
256 incubator (Tokai Hit). Functional parameters were assessed from the averaged  
257 contraction-relaxation waveforms from 10-sec recordings.

258

### 259 **Western blotting**

260 Human iPSCs and iPSC-CMs grown in 6-wells plates were harvest and lysated in RIPA  
261 buffer with Complete Mini, EDTA-free Protease inhibitor cocktail tablets (Roche). The  
262 lysates were placed on ice for 30 minutes, followed by centrifuging at 14000 rpm for 20  
263 min, the supernatants were then collected as proteins. BCA Protein Assay kit (Thermo  
264 Fisher Scientific) was used to measure the protein concentration. Western blot was  
265 performed according to the standard protocol. Briefly, Equal amounts of protein was  
266 treated by SDS-PAGE electrophoresis and transferred to a nitrocellulose membrane.  
267 After nonfat milk blocking, the membrane was incubated with the following primary  
268 antibodies at 4°C overnight, respectively: TLK1 (Cell Signaling Technology, 4125S;  
269 1:1000 dilution), TEAD2 (Proteintech, 21159-1-AP; 1:500 dilution), RBBP5 (Bethyl  
270 Laboratories, A300-109A-M; 1:1000 dilution), ASH2L (Bethyl Laboratories, A300-107A-  
271 M; 1:1000 dilution) and GAPDH (Abcam, ab8245; 1:1000 dilution). Subsequently, the  
272 membrane was incubated in protein-specific HRP conjugated secondary antibody for 1  
273 hr at room temperature. Restore Western Blot Stripping Buffer (Thermo Fisher Scientific)  
274 was used to clean RBBP5 antibody for detecting ASH2L protein. The blots were  
275 visualized using chemiluminescence (Thermo Fisher Scientific).



276

## 277 **Cell sorting by FACS**

278 Fluorescent activated cell sorting (FACS) was performed to determine the  
279 cardiomyocyte differentiation efficiency. iPSC-CMs were dissociated and stained with  
280 cardiac troponin T antibody (ab45932, Abcam), and goat anti-Rabbit secondary  
281 antibody (A32731, Thermo Scientific) using Fixation/Permeabilization Kit (554714, BD  
282 Biosciences). The stained cells were filtered through a 35 µm cell strainer snap cap and  
283 collected in a 5 ml FACS tube (Corning). The cells were analyzed on a BD Biosciences  
284 FACS Aria II instrument using FACSDiva software. The cells were gated on the basis of  
285 forward scatter and side scatter. Flow cytometric gates were set using parental cells.

286

## 287 **siRNA experiments and RNA-seq**

288 To generate ASH2L knockout iPSCs, CRISPR/Cas9 gene editing was performed using  
289 two single-guide RNA (sgRNAs) flanking the exon 4-5 region of ASH2L. The guide DNA  
290 oligos were designed using a web-based tool ([crispr.mit.edu/](http://crispr.mit.edu/)) and chosen based on a  
291 high score for on-target binding and the lowest off-target score. The gRNAs were cloned  
292 into the pSpCas9(BB)-2A-GFP vector (PX458; a gift from Feng Zhang; Addgene  
293 plasmid #48138) using annealed reverse complementary guide DNA oligos. The  
294 sequences of the sgRNAs were gRNA\_3': GGCAGAGACGGATGCAACAG and  
295 gRNA\_3': GTGGTTGTATAACAGAATAT. Two CRISPR/Cas9 vectors (1 µg each) were  
296 transfected in hiPSCs (SCVI480 using the Lipofectamine Stem Transfection Reagent  
297 (Thermo Fisher Scientific). The cells were dissociated using TrypLE express 1x  
298 (Thermo Fisher Scientific) and GFP+ cells were sorted by flow cytometry 24 hr post-

299 transfection. GFP+ cells were seeded at a density of 1,000 cells per well in a 6-well  
300 plate to generate clonal isolates. Ten to fourteen days after seeding, 48 individual  
301 clones were picked for genotypic screening by PCR. (FW:  
302 AGCTAGTTTCAGAGTCCAAGATAAA, RV: GATGGAGAAAGAAGCTATAGAGGAG)  
303 Knock-out clones were confirmed by DNA agarose gel. Heterozygous clones F5 and  
304 H12 are selected for subsequent experiments.

305

306 For RNA-seq, NEBNext Ultra II RNA kit with PolyA Selection was used to construct  
307 RNA library. Two replicates were completed for each condition. Then standard RNA-seq  
308 protocols were followed and we generated > 28 million 76 bp single end reads per  
309 sample. The original reads were mapped to human genome (GRCh37) with STAR  
310 2.7.3a<sup>41</sup> with default settings, then reads with MAPQ > 20 were collected for further  
311 analysis. FeatureCounts<sup>42</sup> was applied to count the number of reads mapping to  
312 genomic feature. Differentially expression was detected by DESeq2<sup>43</sup>. GO terms were  
313 enriched with clusterProfiler<sup>44</sup>.

314

## 315 **Results**

### 316 ***NetWalker identified a sub-network in CHD***

317 We analyzed de novo mutations identified from the PCGC whole-exome sequencing  
318 data from 2,871 CHD probands<sup>10</sup>. We agnostically identified 120 genes from the PCGC  
319 CHD probands that were affected by de novo loss-of-function (LoF) mutations (stop  
320 gain and loss mutations and frameshift indels), and LoF mutations in these genes were  
321 highly depleted in natural human populations (see **Methods and Materials** and **Table**

322 **S1A**), suggesting their sensitivity to gene copy loss. The 120 genes displayed an overall  
323 functional enrichment for abnormal heart development (FDR =  $4.7e-4$ ), decreased fetal  
324 cardiomyocyte proliferation (FDR =  $2.2e-3$ ) and abnormal cardiovascular system  
325 morphology (FDR =  $6.5e-3$ ) (**Table S1B**), well representing CHD candidate genes. With  
326 exactly the same procedure, we also identified 35 genes from the accompanied control  
327 cohort from the original publication<sup>10</sup>, which, as expected, did not show enrichment for  
328 cardiac functions (**Table S1C and D**). We examined the topological occupancies of  
329 these identified PCGC candidate genes on the high-quality experimentally derived  
330 human protein interaction network that was compiled and quality-checked in our  
331 previous publication<sup>13</sup>, encompassing 16,085 unique proteins (or genes) and 217,605  
332 interactions (see **Methods and Materials, Table S2**). We observed that these PCGC  
333 candidate proteins tended to have significantly increased connectivity compared with  
334 those identified from unaffected sibling controls ( $p = 0.02$ , Wilcoxon rank-sum test,  
335 **Figure 1A**), thereby occupying the central topological positions on the interaction  
336 network. We confirmed this observation using genes affected by de novo synonymous  
337 mutations in the PCGC probands and unaffected siblings, respectively: as expected, no  
338 significant difference was observed ( $p = 0.24$ , **Figure 1A**). This comparison suggested a  
339 global impact of the identified 120 PCGC proteins given their central topological  
340 positions on the network. We next asked whether these identified proteins were  
341 topological neighbors on the interaction network, which would help reveal convergent  
342 molecular pathways in CHD. We calculated the fraction of the identified 120 proteins  
343 that were direct interacting partners on the network, which indeed displayed a significant  
344 enrichment compared to the genes identified from the unaffected siblings using the

345 same procedure ( $p = 6.8e-3$ , Fisher's exact test, **Figure 1B**), and the pattern was  
346 absent from control genes affected by de novo synonymous mutations ( $p = 0.19$ ,  
347 Fisher's exact test, **Figure 1B**). Taken together, although the 120 genes were  
348 individually and agnostically identified from the PCGC CHD exomes, their topological  
349 positions on the global interaction network were not random: (1) they tended to occupy  
350 central positions on the network to exert their global impact on biological processes; (2)  
351 they were more likely to interact with each other on the biological network, suggesting  
352 their formation of a functional convergent sub-network underlying heart development.

353

354 We developed the NetWalker algorithm to dissect the complete structure of the  
355 underlying convergent molecular network in CHD seeded by the 120 candidate proteins  
356 (termed the PCGC proteins/genes). NetWalker is essentially a random walk scheme on  
357 a complex network which calculates the probability of visiting one node from any other  
358 nodes assuming stochastic flow on the network (the "reachability" between any two  
359 nodes on the network, **Figure 1C**)<sup>45,46</sup>. Specifically, for each protein on the protein  
360 interaction network, we computed its reachability to each of the PCGC proteins, as well  
361 as to the entire collection of the human proteins on the network. At a false discovery  
362 rate of 0.05 and fold change of 2, we agnostically identified another set of 120 new  
363 proteins (the NetWalker proteins/genes) topologically more reachable to the 120 PCGC  
364 proteins than to the global proteome background (**Figure 1D**). Most of these 120 newly  
365 identified NetWalker proteins had not been reported by previous CHD mutational  
366 analyses<sup>10,47</sup>, and we next sought to functionally characterize their potential functions in  
367 heart development.

368

369 These newly identified 120 proteins have extensive interactions with the candidate  
370 PCGC proteins, where 179 (88/120 from the PCGC proteins and 91/120 from the  
371 NetWalker proteins) out of the 240 proteins (120 PCGC+120 NetWalker) formed a  
372 highly connected network (**Figure 2A**). To demonstrate statistical significance of the  
373 identified network, we performed degree-preserving shuffling<sup>22</sup> on the network, where  
374 we randomly permuted the interacting partners for each protein on the interaction  
375 network while maintaining their respective connectivity. With the same set of seed  
376 proteins, implementing NetWalker 100 times on these permuted networks only identified  
377 highly fragmented sub-networks (**Figure S1**), suggesting that the identified highly  
378 connected network (**Figure 2A**) was not expected by chance. Excluding the seed  
379 PCGC candidate proteins, throughout our analyses we sought to characterize the newly  
380 identified NetWalker proteins (the orange nodes in **Figure 2A**) for their biological  
381 significance in regulating heart development.

382

### 383 ***Functional characterization of the CHD network***

384 Although the newly identified NetWalker proteins had no significant mutations from  
385 previous mutational analyses in the PCGC cohort (the orange nodes in **Figure 2A**),  
386 several known CHD risk factors (most from previous clinical studies) could be  
387 immediately recognized, including TBX5<sup>48</sup>, NKX2-5<sup>49</sup>, CITED2<sup>50</sup>, IFT80<sup>51</sup>, ZFPM2<sup>52,53</sup>  
388 and ACVR2B<sup>54</sup>. Additionally, our algorithm also identified MSX1, whose association with  
389 CHD was recently suggested by a CHD GWAS study<sup>11</sup>. Despite these known proteins,  
390 to the best of our knowledge, the majority of the proteins identified by NetWalker

391 **(Figure 2A)** had uncharacterized function in heart development nor in CHD. We  
392 investigated their gene expression in the developing heart using mouse models and the  
393 human fetal heart. We first considered the time-course transcriptomic data during  
394 cardiogenesis in mouse, where transcriptomes in the developing heart were densely  
395 sampled and profiled across key heart developmental stages, from embryonic stem  
396 cells to fetal, postnatal and adult stages<sup>23</sup>. Utilizing one-to-one unambiguous mouse-  
397 human orthology mapping, we observed that the NetWalker genes formed two  
398 expression clusters across the heart developmental stages: Group-I (G-I) genes  
399 displayed preferential expression from embryonic stem cells (ESCs) to E7.3, whereas  
400 Group-II (G-II) genes exhibited substantial expression enrichment from E7.3 to  
401 postnatal and adult stages (**Figure 2B**). Note that the heart tube forms at ~E7.3  
402 (corresponding to ~2.5 postconceptional weeks, PCW, in humans). Close examination  
403 of Group-II genes further revealed two subcluster structure, where Group-II-A (G-II-A)  
404 was preferentially expressed across fetal developmental stages and Group-II-B (G-II-B)  
405 was more specific in postnatal stages, particularly strong in the adult heart (**Figure 2B**).  
406 It is also important to note that the observed expression propensities were relative  
407 comparisons across developmental stages, and therefore the increased expression in  
408 particular time points do not necessarily preclude its molecular activities in other time  
409 points.

410

411 We further leveraged the recently published spatial transcriptomic data in the fetal  
412 heart<sup>26</sup> to investigate the molecular activities of the identified genes in specific heart  
413 compartments. The original experiments evenly sampled tissue spots (each containing

414 ~30 cells) from the fetal heart of serial sections in postconceptional weeks (PCW) 4.5-5,  
415 6.5 and 9, and determined the transcriptome in each tissue spot so as to gain insights  
416 into the spatial effects on modulating gene expression during heart development<sup>26</sup>.  
417 These tissue spots were then clustered by their transcriptome similarity to reveal cell  
418 architecture defining spatial compartments of the developing fetal heart. We analyzed  
419 our identified genes in these spatial compartments, and did not observe significant  
420 expression enrichment for Group-I genes across all heart compartments in PCW 4.5-5,  
421 6.5 and 9. The observation was expected given their early embryonic expression in the  
422 mouse data (**Figure 2B** before ~E7.3). The lack of expression enrichment was also  
423 observed for Group-II-B genes (**Figure 2B**), consistent with their postnatal and adult  
424 expression in the mouse data. However, the Group-II-A genes displayed strong  
425 expression enrichment across all fetal heart compartments in both PCW 6.5 and 9  
426 (**Figure 3A and B**), but not in PCW 4.5-5. This again demonstrated the overall  
427 concordance with their fetal expression in the developing mouse heart (**Figure 2B**), but  
428 more precisely suggested their developing timing after PCW 4.5-5. Overall, these  
429 spatial transcriptome data supported our observation from the mouse data.

430

431 Leveraging resources in the Epigenome Roadmap Project<sup>24</sup>, we further analyzed the  
432 fetal heart transcriptome data in 7 discrete time points from postconceptional day 96  
433 (gestational week 13.7) to day 147 (gestational week 21), and again we observed  
434 strong molecular activities of the Group-II-A genes across all time points relative to the  
435 transcriptome background (**Figure 3C**). The significance was absent for Group-I and  
436 Group-II-B genes. We further examined RNA-seq data from the ENCODE project<sup>25</sup> for

437 the fetal heart in gestational weeks 19 and 28, again confirmed strong activities of the  
438 Group-II-A genes in both gestational weeks (**Figure 3D**). Taken together, all the human  
439 data confirmed molecular dynamics of the identified genes as we observed from the  
440 mouse data: Group-I genes were specific to early embryonic stages; Group-II-A genes  
441 were specific for fetal development, whereas Group-II-B genes were more specific for  
442 the postnatal and adult heart.

443  
444 We performed gene ontology analysis on the 120 NetWalker proteins (the orange nodes  
445 in **Figure 2A**) to determine their molecular functions in modulating heart development.  
446 Consistent with the above observation on gene expression, Group-II genes displayed  
447 significant functional enrichment for outflow tract septum morphogenesis, ventricular  
448 septum development and regulation of cardiac muscle cell proliferation (**Table 1**). The  
449 enrichment for heart functions was also observed when splitting Group-II genes into two  
450 sub-groups (Group-II-A and Group-II-B). For Group-I genes, their functional enrichment  
451 was significant for heart development, particularly in right atrial isomerism and abnormal  
452 cardiovascular development; however, unexpectedly, their gene ontology also displayed  
453 unexpected enrichment for neural functions, especially for open neural tube defects and  
454 exencephaly (**Table 1**). Given the significant comorbidities between CHD and  
455 neurodevelopmental diseases<sup>55,56</sup>, particularly with neural tube malformations (e.g. 40.6%  
456 of probands with spina bifida aperta develop CHD<sup>57</sup>), this observation likely suggested  
457 the underlying etiological causes (which will be experimentally verified below). Because  
458 Group-I genes were specific to early developmental stages and the enrichment for brain  
459 functions was completely absent from Group-II genes (**Table 1**), the shared molecular



460 etiologies between CHD and neurodevelopmental disorders should occur only at the  
461 early developmental stages.

462

463 ***Experimental validation confirmed the role of the novel proteins in heart***  
464 ***development***

465 Given the clear and strong implication of Group-II genes in heart development (**Figure 3**  
466 and **Table 1**), we next closely examined how Group-I genes modulate heart  
467 development, particularly those also implicated in neurodevelopmental disorders. We  
468 followed our previous practice<sup>15</sup> and partitioned the identified network (**Figure 2A**) into  
469 33 non-overlapping topological clusters (**Table S3**), where in each cluster, proteins were  
470 densely connected with their interacting partners but sparsely interacted with proteins  
471 outside of their respective clusters. This approach has enabled us to understand each  
472 protein's function in the context of its own interaction module. Cluster #4 is an excellent  
473 example to demonstrate the convergence of CHD-associated mutations (**Figure 4A**),  
474 where 12 out of the 15 member proteins were affected by de novo LoF mutations  
475 identified from CHD proband (the PCGC candidate genes). These genes were heralded  
476 by FOXM1, whose mouse mutants displayed ventricular hypoplasia, cardiomyocyte  
477 deficiencies<sup>58</sup> and many other cardiac anomalies<sup>59</sup>. Although these de novo mutations  
478 were individually identified from different CHD probands, their topological clustering in  
479 the same module demonstrated the mutational convergence in CHD onto a common  
480 component. As such, for the remaining three novel proteins in the same cluster  
481 (SLC6A2, DRP2 and MLC1, **Figure 4A**) identified by our NetWalker algorithm, it is  
482 reasonable to expect their function in modulating heart development. Indeed, SLC6A2

483 mouse mutants display smaller heart sizes with increased heart rate (Mouse Genome  
484 Informatics<sup>60</sup>), and DRP2 is specific to heart-derived endothelial cells among many  
485 other endothelial cell types<sup>61</sup>.

486

487 The identified network also encompassed a cluster structure implicating NOTCH1-  
488 mediated signaling (cluster #2 in **Table S3, Figure 4B**), a master regulator of numerous  
489 development processes, including both heart and brain development<sup>62</sup>. NOTCH1 is a  
490 known CHD factor and was affected by a de novo LoF mutation in this PCGC cohort.  
491 We prioritized to experimentally validate two Group-I genes in this cluster (given clear  
492 cardiac functions for Group-II genes, **Table 1**): TEAD2 (TEA Domain Transcription  
493 Factor 2) and TLK1 (Tousled Like Kinase 1). TEAD2 is a transcription factor implicated  
494 in Hippo signaling, while TLK1 is a kinase regulating chromatin assembly. Notably,  
495 TEAD2, as a member of the Group-I gene (**Figure 2B and Table 1**), has been  
496 characterized as a regulator of neural development<sup>63</sup>, where *Tead2*<sup>-/-</sup> mouse mutants  
497 displayed exencephaly and open neural tube defects<sup>64</sup>. However, both TEAD2 and  
498 TLK1 have uncharacterized function in regulating heart development. We examined  
499 their cardiac functions using induced pluripotent stem cell-derived cardiomyocytes  
500 (iPSC-CMs). We converted human iPSCs into cardiomyocytes (**Methods and**  
501 **Materials**), and subsequently determined high protein abundance of TEAD2 and TLK1  
502 in the iPSC-CMs (**Figure S2A**). We performed siRNA knockdown to suppress TLK1 and  
503 TEAD2 expression, respectively, followed by RNA-seq and cellular assays to determine  
504 their regulatory and phenotypic effects in cardiomyocytes. Confirming the knockdown  
505 efficiencies by their respective siRNAs (**Figure S3**), we observed that massive genes

506 associated with cardiac muscle contraction exhibited differential expression upon  
507 TEAD2 or TLK1 knockdown in cardiomyocytes (**Figure 4D-F** and **Table S4**).  
508 Specifically, for TEAD2, despite its known function in regulating neural development, its  
509 knockdown in cardiomyocytes has clearly perturbed numerous genes specific for  
510 regulating cardiac muscle contractility, including the cardiac muscle myosin factors  
511 (MYL3, MYH6/7), and the troponin complex subunits (TNNI3 and TNNT1), the muscle  
512 intermediate filament protein (DES) and the connexin gap junction protein (GJA5,  
513 **Figure 4E**). Performing cellular contractility assay in the iPSC-CMs, as expected, we  
514 observed that the siRNA knockdown against TEAD2 resulted in a marked reduction of  
515 the cardiomyocyte beating rate relative to the siRNA control ( $p = 1.1e-3$ , **Figure 4G**),  
516 accompanied with significantly increased contraction velocity, contraction deformation  
517 distance, relaxation velocity and relaxation deformation distance ( $P < 0.05$ , **Figure 4H-**  
518 **K**). Similar observation was also made from the TLK1 knockdown experiments, where  
519 numerous muscle contractility genes displayed differential expression in the iPSC-CMs  
520 (**Figure 4D** and **F**). In the meantime, knocking down TLK1 expression resulted in  
521 reduced beating rate of cardiomyocytes (**Figure 4G**), confirming the role of TLK1 in  
522 regulating heart muscle contraction.

523

524 The CHD network also encompassed a cluster structure mediated by key factors  
525 (GATA4, TBX5, and NKX2-5) regulating heart development (cluster #3 in **Table S3**,  
526 **Figure 4C**). However, our analysis now revealed their extensive interactions with  
527 numerous histone modification proteins (e.g., KDM6A, KDM4B, ASH2L, RBBP5 and  
528 JADE1). We particularly note that many member proteins in this cluster are also

529 implicated in brain development despite three key regulators of heart development  
530 GATA4, TBX5, and NKX2-5. For example, GATA4, TBX5, and NKX2-5 all interacted  
531 with KDM6A, the causal gene of Kabuki syndrome<sup>65</sup>, characterized by intellectual  
532 disability, distinctive facial features and growth delay, etc. Notably, coarctation of the  
533 aorta and ventricular/atrial septal defects are also common clinical manifestations of  
534 Kabuki syndrome<sup>66</sup>. Indeed, KDM6A mouse mutants also displayed open neural tube  
535 defects accompanied with multiple cardiac anomalies including failure of heart looping,  
536 thin ventricular wall and myocardial trabeculae hypoplasia (Mouse Genome  
537 Informatics<sup>60</sup>). Therefore, this local network structure suggested a potential mechanism  
538 of the heart defects in the Kabuki syndrome by perturbing GATA4/TBX5/NKX2-5-  
539 mediated cardiac functions through their interactions with KDM6A. More interestingly,  
540 this cluster structure also encompassed CHD7, the causal gene for CHARGE syndrome  
541 where heart defects and growth retardation are common among patients. This cluster  
542 therefore represent a convergent structure underlying several distinct but related  
543 monogenic disorders.

544

545 To demonstrate the overall implication of this cluster in regulating heart development,  
546 we prioritized two Group-I genes, RBBP5 and ASH2L, to experimentally validate their  
547 function in heart development given their known function in regulating brain  
548 development but uncharacterized cardiac functions. RBBP5 plays a key role in  
549 differentiating embryonic stem cells along the neural lineage (UniProtKB/Swiss-Prot  
550 Summary), whereas ASH2L regulates the corticogenesis process<sup>67</sup>, neuronal structure  
551 and behavior<sup>68</sup>. Both RBBP5 (RB Binding Protein 5) and ASH2L (ASH2 Like) encode

552 subunits of the histone lysine methyltransferase complex. We performed the same  
553 siRNA knockdown experiments and cellular contractility assays as described above. We  
554 first confirmed high protein abundance of RBBP5 and ASH2L in iPSC-CMs (**Figure**  
555 **S2B**). Our subsequent RNA-seq revealed differential expression of numerous cardiac  
556 muscle proteins in the cardiomyocytes upon RBBP5 knockdown (**Figure 4D and L**),  
557 resulting in substantially increased beating rate of the iPSC-CMs (**Figure 4N**). ASH2L  
558 knockdown led to fewer differentially expressed genes, but those affected genes were  
559 critical factors modulating heart contraction including the myosin light chain proteins  
560 MYL3/4, and the troponin subunit TNNC1 implicated in cardiomyopathy<sup>69</sup>. Knocking  
561 down ASH2L in iPSC-CMs did not significantly alter the beating rate, but substantially  
562 altered all other contractility parameters, including velocity and deformation distance for  
563 contraction and relaxation, respectively (**Figure 4O-R**). To gain deeper insights, we  
564 further asked whether ASH2L modulates cardiomyocytes differentiation from iPSCs. We  
565 generated heterozygous ASH2L knockout iPSC (ASH2L<sup>+/-</sup>) using the CRISPR-Cas9  
566 gene engineering technique. The editing effects were verified by Sanger sequencing  
567 and western blot in iPSCs for the two ASH2L<sup>+/-</sup> knockout clones (**Figure S4**). We then  
568 differentiated the edited iPSCs into cardiomyocytes. On day 11, we observed that in  
569 both ASH2L<sup>+/-</sup> knockout lines, the heterozygous ASH2L knockouts indeed significantly  
570 reduced the differentiation efficiencies into cardiomyocytes (TNNT2 positive cells) from  
571 iPSCs, thereby confirming the role of ASH2L in regulating the cardiogenesis process  
572 (**Figure 4S and T**).

573

574 Taken together, because Group-II genes were enriched for heart-specific functions  
575 (**Table 1**), our study prioritized four Group-I genes (TEAD2, TLK1, RBBP5 and ASH2L)  
576 for experimental validation, whose cardiac functions had not been previously  
577 characterized. We particularly note RBBP5 and ASH2L, where previous work in mouse  
578 suggested *Rbbp5* involvement in epigenetic regulation in murine cardiomyocytes<sup>70</sup> and  
579 *Ash2l* interaction with *Tbx1* in vitro<sup>71</sup>. Because RBBP5 is a subunit of histone lysine  
580 methyltransferase complex and is widely expressed across human tissues, its role of  
581 epigenetic regulation across many cell types (including cardiomyocytes) is anticipated.  
582 However, given its well characterized role in differentiating stem cells along the neural  
583 lineage, its dual function in regulating heart development was unexpected and is now  
584 revealed by our study. For *Ash2l* interaction with *Tbx1*<sup>71</sup>, because *Tbx1* is a key gene in  
585 22q11.2 deletion syndrome, interacting with *Tbx1* likely contribute to the role of *Ash2l* in  
586 regulating heart development. Nevertheless, our experimental data not only  
587 demonstrated the effectiveness of our network biology approach on identifying novel  
588 genes in regulating fetal heart development, but also provides direct mechanistic  
589 insights into the overlapping molecular basis between the brain and heart  
590 developmental programs.

591

### 592 ***The network is enriched for pathogenic mutations in HLHS probands***

593 Having established the implication of the identified network in regulating heart  
594 development, we sought to identify the etiological associations with CHD. Because  
595 these genes were not identified from the existing PCGC mutation analysis<sup>10</sup>, we  
596 examined the latest release of the PCGC sequencing data (as of Feb, 2021). We

597 performed analyses on patients' clinical records, and included patients with Tetralogy of  
598 Fallot (TOF, N = 328), ventricular septum defects (VSD, N = 776), atrial septum defects  
599 (ASD, N = 830), hypoplastic left heart syndrome (HLHS, N = 224), and transposition of  
600 the great arteries (TGA, N = 167). Note that all these PCGC proband samples were  
601 subjected to whole-exome sequencing on the same platform (Illumina 2000). We  
602 retrieved the deposited FASTQ data in PCGC from dbGAP (dbGAP-24034,  
603 gap\_accession: phs000571, gap\_parent\_phs: phs001194, SRP025159), and performed  
604 variant call and variant annotation (**Methods and Materials**). The variant calls were  
605 made by aggregating the whole-exome data for 1,817 unaffected siblings (control  
606 subjects, **Methods and Materials, Table S5**), which were also used as controls in the  
607 previous PCGC study<sup>10</sup>. The joint variant call minimized potential bias from different  
608 variant call platforms. We performed additional quality control procedures, which  
609 demonstrated high-quality of the identified variants from our variant call procedures  
610 (**Figure S5** and **Methods and Materials**).

611

612 Using all the called variants, we first confirmed similar population structure between  
613 cases and controls (**Figure S6**). Due to weak effect sizes of common variants, we  
614 focused our analyses on rare variants from the PCGC cohort. We considered rare  
615 variants that were not present in the 1000 Genome database. For each CHD sub-type,  
616 we compared deleteriousness of non-synonymous mutations (LoF and missense) in  
617 probands and controls. We used the well-known CADD phased scores<sup>32</sup> to quantify  
618 mutational deleteriousness in exonic regions, and computed the mean CADD scores for  
619 non-synonymous mutations affecting the network genes in each personal exome in

620 each CHD sub-type (TOF, VSD, ASD, HLHS and TGA) as well as in the control cohort  
621 (the unaffected siblings). We reasoned that if the network genes (the orange nodes in  
622 **Figure 2A**) were implicated in CHD, we would expect that CHD probands tend to carry  
623 more pathogenic mutations affecting the network genes relative to controls. Given  
624 differential expression enrichment of the network genes at specific heart development  
625 stages (**Figure 2B**), observation of excessive pathogenic mutations specifically affecting  
626 particular sub-groups (genes in Group-I or II, **Figure 2B**) among probands of a given  
627 CHD sub-type would indicate developmental timing. We separately analyzed Group-I  
628 and II genes on the network (**Figure 2B**), and compared the mean CADD scores  
629 affecting Group-I or II genes in each proband and control subject. While probands  
630 across different CHD subtypes did not display mutational enrichment affecting Group-I  
631 genes (**Figure S7**), we did observe significant elevation of mutational pathogenicity  
632 affecting Group-II genes, which was specific to HLHS probands ( $p = 9.7e-3$ , Wilcoxon  
633 rank-sum test, **Figure 5A**) but was insignificant for other CHD subtypes in comparison  
634 (**Figure 5A**). Because HLHS is typically comorbid with ASD, we indeed observed that  
635 56 HLHS probands (among total 224 HLHS cases) also received ASD diagnosis.  
636 Excluding these overlapping cases further boosted the statistical significance for the  
637 enrichment of pathogenic non-synonymous mutations in the identified network genes ( $p$   
638  $= 2.7e-3$ , Wilcoxon rank-sum test, **Figure 5C**). We performed a set of additional control  
639 experiments to confirm the observation: (1) we performed the same analysis on rare  
640 synonymous mutations affecting Group-II genes, but did not observe any statistical  
641 significance among all CHD sub-types, including HLHS (**Figure 5B**). This observation  
642 demonstrated that the observed statistical significance indeed implied functional



643 consequences. (2) We performed the same analysis on 62 lung-specific protein-coding  
644 genes<sup>34</sup>, and no significance was observed on both non-synonymous ( $p = 0.86$ ,  
645 Wilcoxon rank-sum test) and synonymous ( $p = 0.36$ , Wilcoxon rank-sum test) variants  
646 (**Figure 5D**). This comparison confirmed specificity of the identified genes in regulating  
647 heart development. (3) In each exome, we confirmed that the number of rare non-  
648 synonymous variants in each HLHS proband did not significantly differ from the  
649 unaffected siblings in the control cohort ( $p = 0.55$ , Wilcoxon rank-sum test, **Figure S8**),  
650 suggesting that the observed enrichment of pathogenic mutations cannot be explained  
651 by excessive mutations identified in the proband cohort than the control cohort, but by  
652 the increased mutational pathogenicity. (4) Lastly, we asked whether our observation  
653 could be merely explained by increased CADD scores of all rare non-synonymous  
654 variants across the exome background in cases relative to controls. We performed a  
655 permutation study, where, in each permutation, we randomly sampled rare non-  
656 synonymous variants from the exome backgrounds from cases and control cohorts,  
657 matching the number of rare non-synonymous variants in Group-II genes in cases and  
658 control cohorts, respectively. We performed the permutation 100 times, and confirmed  
659 that mutational pathogenicity scores were not significantly different between cases and  
660 controls when randomly sampling rare non-synonymous variants from the two cohorts  
661 (**Figure 5E**). Taken together, our comparisons collectively demonstrated the specificity  
662 of the identified network in the molecular etiologies of HLHS.

663

664 ***Single cell analysis of the CHD network in HLHS***

665 Since HLHS is a form of critical congenital heart defect (CCHD), we next sought to  
666 derive further mechanistic insights into the molecular etiologies of HLHS. HLHS affects  
667 blood flow through the heart due to the underdeveloped left ventricle accompanied with  
668 malformations of mitral and aortic valves<sup>72</sup>. To understand the association of the  
669 network with HLHS, we leveraged the recently published single-cell data from an HLHS  
670 fetal heart (day 84), and compared gene expression in the underdeveloped left ventricle  
671 of this HLHS heart against the left ventricle of a typically developing fetal heart on day  
672 83 (**Methods and Materials**)<sup>35</sup>. For Group-I and II genes in the network (**Figure 2A**), we  
673 compared their expression in the HLHS heart across all cell types. We did not observe  
674 significant expression differences in the HLHS heart for Group-I genes across all cell  
675 types, which was expected given their expression specificity in early developmental  
676 stages (compared with the HLHS heart from day 84). However, for Group-II genes, we  
677 observed their significant expression reduction only in endothelium cells of the HLHS  
678 left ventricle across all cell types ( $p = 5.8e-4$ , Wilcoxon rank-sum test, **Figure 6A**).  
679 Close examining the Group-II genes, we observed that the reduction was consistent  
680 between the two sub-groups (Group-II-A and Group-II-B, **Figure 2B**), but was  
681 particularly pronounced among Group-II-A genes (**Figure 6B**) whose expression were  
682 specific across fetal development stages (**Figure 3**), thereby highlighting the significant  
683 contribution to HLHS from fetal endothelium development. We also noted the  
684 conduction system in **Figure 6A**, where expression of Group-II genes displayed  
685 marginal statistical significance ( $p = 0.08$ , Wilcoxon rank-sum test). Close examination  
686 revealed a significant expression reduction only specific to Group-II-A genes in the  
687 conduction system in the HLHS heart (left ventricle) (**Figure 6C**), again suggesting a

688 dysregulation of the conduction system in HLHS. Taken together, our analysis thus  
689 strongly suggests the impaired endothelium and conduction system in HLHS.

690

## 691 **Discussion**

692 Disease-associated mutations are not randomly dispersed across the genome but affect  
693 common sets of molecular pathways on biological networks leading to common clinical  
694 manifestations<sup>12</sup>. The systems thinking has motivated us to develop a novel  
695 computational framework to integrate large-scale CHD genomes, the human  
696 interactome, the fetal heart spatial transcriptome as well as the single-cell transcriptome  
697 from clinical samples. This integrative strategy identified numerous novel proteins with  
698 previously uncharacterized roles in regulating fetal heart development, and our  
699 subsequent multi-omic analyses further demonstrated their function more specific to  
700 certain CHD subtypes with the strongest effect on HLHS. Overall, our integrative  
701 analysis significantly advanced our understanding of the genetic architecture in CHD,  
702 revealed the molecular etiologies in HLHS, and can be readily extended to study other  
703 complex diseases.

704

705 We started our analysis by seeding CHD candidate proteins on the human protein  
706 interaction network. These proteins were intolerant to copy losses and were affected by  
707 de novo LoF mutations in CHD probands. Their central positions and the clustering  
708 patterns on the network suggested their significant impacts and convergent functions in  
709 CHD, which has enabled us to implement the NetWalker approach to identify a  
710 connected network underlying fetal heart development. We observed two expression

711 clusters (Group-I and II genes, **Figure 2A**) in the identified CHD network. While Group-II  
712 genes displayed strong tissue specificity in the fetal heart, it was interesting that Group-I  
713 genes, showing the strongest expression at very early developmental stages (**Figure**  
714 **2A**), modulate both neurodevelopmental and heart developmental programs. This  
715 observation was concordant with previous work, where autism genes were also  
716 frequently identified as CHD candidate genes<sup>10</sup>; however, in this work, the shared  
717 molecular etiologies were only limited to Group-I genes, further demonstrating the  
718 functional convergence was specific to early developmental stages. Leveraging the  
719 global protein interaction network, we were able to further pinpoint the underlying  
720 mechanisms in context of the local interacting proteins. The FOXM1 cluster (**Figure 4A**)  
721 best demonstrated the notion of mutational convergence, where the majority of its  
722 member proteins were affected by de novo LoF mutations in CHD probands. In the  
723 same vein, our experimental validation in iPSC-CMs further confirmed four additional  
724 Group-I genes for their previously uncharacterized function in heart development,  
725 including the previously characterized factor in neural tube defects, TEAD2, and ASH2L  
726 regulating corticogenesis. Our RNA-seq experiments and cellular contractility assays  
727 consistently revealed their function in regulating cardiomyocyte contraction,  
728 demonstrating novel heart-specific functions of these genes previously recognized as  
729 brain genes. In addition to studying mature cardiomyocytes, we also performed gene  
730 editing to perturb ASH2L expression in iPSCs and confirmed ASH2L function in  
731 controlling the differentiation process into cardiomyocytes. This observation is  
732 consistent with SALL3 function in the GATA4 cluster (**Figure 4C**), where SALL3 acted  
733 as a switch controlling the developmental trajectory of stem cells towards the neural or

734 cardiac lineage<sup>73</sup>. Taken together, Group-I genes constitute the shared molecular  
735 etiologies between cardiac and neurological conditions given their tissue-context-  
736 dependent functions or their role in controlling cell differentiation processes.

737

738 We analyzed rare mutations from the PCGC cohort on the network, and observed that  
739 Group-II genes were specifically enriched for pathogenic mutations from individuals with  
740 HLHS. This observation was consistent with our transcriptome analysis of the typically  
741 developing and HLHS heart, thereby revealing potential tissue of origins underlying  
742 HLHS. We particularly highlight the role of endothelium cells in HLHS, where its  
743 etiological contribution has just recently begun to be recognized<sup>35</sup>. Although we  
744 demonstrated increased mutational load in Group-II genes in HLHS, this observation did  
745 not fully preclude the role of Group-I genes in HLHS development. In fact, among the  
746 four Group-I genes (RBBP5, ASHL2, TLK1 and TEAD2) we experimentally validated in  
747 iPSC-CMs (**Figure 4**), RBBP5 also displayed significant down-regulation in the  
748 cardiomyocytes of this HLHS sample. Because Group-I genes are also associated with  
749 neurodevelopmental phenotypes, severe phenotypic consequences likely have  
750 suppressed excessive deleterious mutations in the human population, or they likely  
751 underlie many syndromic cases (eg. Kabuki syndrome, **Figure 4C**). Therefore,  
752 screening individual pathogenic mutations in Group-I genes would provide a molecular  
753 basis for clinical evaluation of the comorbidities between CHD and neurodevelopmental  
754 conditions.

755

756 **Acknowledgment**

757 JL acknowledges the startup fund from the Eli and Edythe Broad Center of  
758 Regeneration Medicine and Stem Cell Research, the Bakar Computational Health  
759 Sciences Institute, and the Parker Institute for Cancer Immunotherapy at UCSF. MS  
760 acknowledges grant  
761 award NIH S10OD025212, and NIH/NIDDK P30DK116074.

762 **References**

- 763 1. Hoffman JI, Kaplan S. The incidence of congenital heart disease. *J Am Coll*  
764 *Cardiol.* 2002;39:1890-1900. doi: 10.1016/s0735-1097(02)01886-7
- 765 2. Reller MD, Strickland MJ, Riehle-Colarusso T, Mahle WT, Correa A. Prevalence  
766 of congenital heart defects in metropolitan Atlanta, 1998-2005. *J Pediatr.*  
767 2008;153:807-813. doi: 10.1016/j.jpeds.2008.05.059
- 768 3. Fahed AC, Gelb BD, Seidman JG, Seidman CE. Genetics of congenital heart  
769 disease: the glass half empty. *Circ Res.* 2013;112:707-720. doi:  
770 10.1161/CIRCRESAHA.112.300853
- 771 4. Zaidi S, Brueckner M. Genetics and Genomics of Congenital Heart Disease. *Circ*  
772 *Res.* 2017;120:923-940. doi: 10.1161/CIRCRESAHA.116.309140
- 773 5. Soemedi R, Wilson IJ, Bentham J, Darlay R, Topf A, Zelenika D, Cosgrove C,  
774 Setchfield K, Thornborough C, Granados-Riveron J, et al. Contribution of global  
775 rare copy-number variants to the risk of sporadic congenital heart disease. *Am J*  
776 *Hum Genet.* 2012;91:489-501. doi: 10.1016/j.ajhg.2012.08.003
- 777 6. Demal TJ, Heise M, Reiz B, Dogra D, Braenne I, Reichenspurner H, Manner J,  
778 Aherrahrou Z, Schunkert H, Erdmann J, et al. A familial congenital heart disease  
779 with a possible multigenic origin involving a mutation in BMPR1A. *Sci Rep.*  
780 2019;9:2959. doi: 10.1038/s41598-019-39648-7
- 781 7. Calcagni G, Digilio MC, Sarkozy A, Dallapiccola B, Marino B. Familial recurrence  
782 of congenital heart disease: an overview and review of the literature. *Eur J*  
783 *Pediatr.* 2007;166:111-116. doi: 10.1007/s00431-006-0295-9

- 784 8. Pierpont ME, Brueckner M, Chung WK, Garg V, Lacro RV, McGuire AL, Mital S,  
785 Priest JR, Pu WT, Roberts A, et al. Genetic Basis for Congenital Heart Disease:  
786 Revisited: A Scientific Statement From the American Heart Association.  
787 *Circulation*. 2018;138:e653-e711. doi: 10.1161/CIR.0000000000000606
- 788 9. Pediatric Cardiac Genomics C, Gelb B, Brueckner M, Chung W, Goldmuntz E,  
789 Kaltman J, Kaski JP, Kim R, Kline J, Mercer-Rosa L, et al. The Congenital Heart  
790 Disease Genetic Network Study: rationale, design, and early results. *Circ Res*.  
791 2013;112:698-706. doi: 10.1161/CIRCRESAHA.111.300297
- 792 10. Jin SC, Homsy J, Zaidi S, Lu Q, Morton S, DePalma SR, Zeng X, Qi H, Chang W,  
793 Sierant MC, et al. Contribution of rare inherited and de novo variants in 2,871  
794 congenital heart disease probands. *Nat Genet*. 2017;49:1593-1601. doi:  
795 10.1038/ng.3970
- 796 11. Lahm H, Jia M, Dressen M, Wirth FFM, Puluca N, Gilsbach R, Keavney B,  
797 Cleuziou J, Beck N, Bondareva O, et al. Congenital heart disease risk loci  
798 identified by genome-wide association study in European patients. *J Clin Invest*.  
799 2020. doi: 10.1172/JCI141837
- 800 12. Li J, Li X, Zhang S, Snyder M. Gene-Environment Interaction in the Era of  
801 Precision Medicine. *Cell*. 2019;177:38-44. doi: 10.1016/j.cell.2019.03.004
- 802 13. Li J, Pan C, Zhang S, Spin JM, Deng A, Leung LLK, Dalman RL, Tsao PS,  
803 Snyder M. Decoding the Genomics of Abdominal Aortic Aneurysm. *Cell*.  
804 2018;174:1361-1372.e1310. doi: 10.1016/j.cell.2018.07.021
- 805 14. Li J, Ma Z, Shi M, Maly RH, Aoki H, Minic Z, Phanse S, Jin K, Wall DP, Zhang Z,  
806 et al. Identification of Human Neuronal Protein Complexes Reveals Biochemical



- 807           Activities and Convergent Mechanisms of Action in Autism Spectrum Disorders.  
808           *Cell Syst.* 2015;1:361-374. doi: 10.1016/j.cels.2015.11.002
- 809   15.   Li J, Shi M, Ma Z, Zhao S, Euskirchen G, Ziskin J, Urban A, Hallmayer J, Snyder  
810           M. Integrated systems analysis reveals a molecular network underlying autism  
811           spectrum disorders. *Mol Syst Biol.* 2014;10:774. doi: 10.15252/msb.20145487
- 812   16.   Lage K, Greenway SC, Rosenfeld JA, Wakimoto H, Gorham JM, Segre AV,  
813           Roberts AE, Smoot LB, Pu WT, Pereira AC, et al. Genetic and environmental risk  
814           factors in congenital heart disease functionally converge in protein networks  
815           driving heart development. *Proc Natl Acad Sci U S A.* 2012;109:14035-14040.  
816           doi: 10.1073/pnas.1210730109
- 817   17.   Lage K, Mollgard K, Greenway S, Wakimoto H, Gorham JM, Workman CT,  
818           Bendsen E, Hansen NT, Rigina O, Roque FS, et al. Dissecting spatio-temporal  
819           protein networks driving human heart development and related disorders. *Mol*  
820           *Syst Biol.* 2010;6:381. doi: 10.1038/msb.2010.36
- 821   18.   Krumm N, Turner TN, Baker C, Vives L, Mohajeri K, Witherspoon K, Raja A, Coe  
822           BP, Stessman HA, He ZX, et al. Excess of rare, inherited truncating mutations in  
823           autism. *Nat Genet.* 2015;47:582-588. doi: 10.1038/ng.3303
- 824   19.   Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-  
825           Luria AH, Ware JS, Hill AJ, Cummings BB, et al. Analysis of protein-coding  
826           genetic variation in 60,706 humans. *Nature.* 2016;536:285-291. doi:  
827           10.1038/nature19057
- 828   20.   Kuleshov MV, Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, Koplev  
829           S, Jenkins SL, Jagodnik KM, Lachmann A, et al. Enrichr: a comprehensive gene

- 830 set enrichment analysis web server 2016 update. *Nucleic Acids Res.*  
831 2016;44:W90-97. doi: 10.1093/nar/gkw377
- 832 21. Page L, Brin S, Motwani R, Winograd T. The PageRank citation ranking: Bringing  
833 order to the web. 1999
- 834 22. Maslov S, Sneppen K. Specificity and stability in topology of protein networks.  
835 *Science.* 2002;296:910-913. doi: 10.1126/science.1065103
- 836 23. Li X, Martinez-Fernandez A, Hartjes KA, Kocher JP, Olson TM, Terzic A, Nelson  
837 TJ. Transcriptional atlas of cardiogenesis maps congenital heart disease  
838 interactome. *Physiol Genomics.* 2014;46:482-495. doi:  
839 10.1152/physiolgenomics.00015.2014
- 840 24. Roadmap Epigenomics C, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A,  
841 Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, et al. Integrative analysis  
842 of 111 reference human epigenomes. *Nature.* 2015;518:317-330. doi:  
843 10.1038/nature14248
- 844 25. Consortium EP. An integrated encyclopedia of DNA elements in the human  
845 genome. *Nature.* 2012;489:57-74. doi: 10.1038/nature11247
- 846 26. Asp M, Giacomello S, Larsson L, Wu C, Fürth D, Qian X, Wärdell E, Custodio J,  
847 Reimegård J, Salmén F, et al. A Spatiotemporal Organ-Wide Gene Expression  
848 and Cell Atlas of the Developing Human Heart. *Cell.* 2019;179:1647-1660.e1619.  
849 doi: 10.1016/j.cell.2019.11.025
- 850 27. Hao Y, Hao S, Andersen-Nissen E, Mauck WM, 3rd, Zheng S, Butler A, Lee MJ,  
851 Wilk AJ, Darby C, Zager M, et al. Integrated analysis of multimodal single-cell  
852 data. *Cell.* 2021;184:3573-3587 e3529. doi: 10.1016/j.cell.2021.04.048

- 853 28. Kendig KI, Baheti S, Bockol MA, Drucker TM, Hart SN, Heldenbrand JR,  
854 Hernaez M, Hudson ME, Kalmbach MT, Klee EW, et al. Sentieon DNaseSeq  
855 Variant Calling Workflow Demonstrates Strong Computational Performance and  
856 Accuracy. *Front Genet.* 2019;10:736. doi: 10.3389/fgene.2019.00736
- 857 29. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, Collins  
858 RL, Laricchia KM, Ganna A, Birnbaum DP, et al. The mutational constraint  
859 spectrum quantified from variation in 141,456 humans. *Nature.* 2020;581:434-  
860 443. doi: 10.1038/s41586-020-2308-7
- 861 30. Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, Taliun SAG,  
862 Corvelo A, Gogarten SM, Kang HM, et al. Sequencing of 53,831 diverse  
863 genomes from the NHLBI TOPMed Program. *Nature.* 2021;590:290-299. doi:  
864 10.1038/s41586-021-03205-y
- 865 31. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic  
866 variants from high-throughput sequencing data. *Nucleic Acids Res.* 2010;38:e164.  
867 doi: 10.1093/nar/gkq603
- 868 32. Kircher M, Witten DM, Jain P, O'Roak BJ, Cooper GM, Shendure J. A general  
869 framework for estimating the relative pathogenicity of human genetic variants.  
870 *Nat Genet.* 2014;46:310-315. doi: 10.1038/ng.2892
- 871 33. Rentzsch P, Schubach M, Shendure J, Kircher M. CADD-Splice-improving  
872 genome-wide variant effect prediction using deep learning-derived splice scores.  
873 *Genome Med.* 2021;13:31. doi: 10.1186/s13073-021-00835-9

- 874 34. Xiong M, Heruth DP, Zhang LQ, Ye SQ. Identification of lung-specific genes by  
875 meta-analysis of multiple tissue RNA-seq data. *FEBS Open Bio*. 2016;6:774-781.  
876 doi: 10.1002/2211-5463.12089
- 877 35. Miao Y, Tian L, Martin M, Paige SL, Galdos FX, Li J, Klein A, Zhang H, Ma N,  
878 Wei Y, et al. Intrinsic Endocardial Defects Contribute to Hypoplastic Left Heart  
879 Syndrome. *Cell stem cell*. 2020;27:574-589.e578. doi:  
880 10.1016/j.stem.2020.07.015
- 881 36. Churko JM, Burridge PW, Wu JC. Generation of human iPSCs from human  
882 peripheral blood mononuclear cells using non-integrative Sendai virus in  
883 chemically defined conditions. *Methods Mol Biol*. 2013;1036:81-88. doi:  
884 10.1007/978-1-62703-511-8\_7
- 885 37. Lian X, Zhang J, Azarin SM, Zhu K, Hazeltine LB, Bao X, Hsiao C, Kamp TJ,  
886 Palecek SP. Directed cardiomyocyte differentiation from human pluripotent stem  
887 cells by modulating Wnt/beta-catenin signaling under fully defined conditions. *Nat*  
888 *Protoc*. 2013;8:162-175. doi: 10.1038/nprot.2012.150
- 889 38. Burridge PW, Matsa E, Shukla P, Lin ZC, Churko JM, Ebert AD, Lan F, Diecke S,  
890 Huber B, Mordwinkin NM, et al. Chemically defined generation of human  
891 cardiomyocytes. *Nat Methods*. 2014;11:855-860. doi: 10.1038/nmeth.2999
- 892 39. Wilson KD, Ameen M, Guo H, Abilez OJ, Tian L, Mumbach MR, Diecke S, Qin X,  
893 Liu Y, Yang H, et al. Endogenous Retrovirus-Derived lncRNA BANCR Promotes  
894 Cardiomyocyte Migration in Humans and Non-human Primates. *Developmental*  
895 *cell*. 2020;54:694-709.e699. doi: 10.1016/j.devcel.2020.07.006

- 896 40. Lam CK, Tian L, Belbachir N, Wnorowski A, Shrestha R, Ma N, Kitani T, Rhee  
897 JW, Wu JC. Identifying the Transcriptome Signatures of Calcium Channel  
898 Blockers in Human Induced Pluripotent Stem Cell-Derived Cardiomyocytes. *Circ*  
899 *Res.* 2019;125:212-222. doi: 10.1161/CIRCRESAHA.118.314202
- 900 41. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P,  
901 Chaisson M, Gingeras TR. STAR: ultrafast universal RNA-seq aligner.  
902 *Bioinformatics.* 2013;29:15-21. doi: 10.1093/bioinformatics/bts635
- 903 42. Liao Y, Smyth GK, Shi W. featureCounts: an efficient general purpose program  
904 for assigning sequence reads to genomic features. *Bioinformatics.* 2014;30:923-  
905 930. doi: 10.1093/bioinformatics/btt656
- 906 43. Love MI, Huber W, Anders S. Moderated estimation of fold change and  
907 dispersion for RNA-seq data with DESeq2. *Genome Biol.* 2014;15:550. doi:  
908 10.1186/s13059-014-0550-8
- 909 44. Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing  
910 biological themes among gene clusters. *OMICS.* 2012;16:284-287. doi:  
911 10.1089/omi.2011.0118
- 912 45. Xia F, Liu J, Nie H, Fu Y, Wan L, Kong X. Random Walks: A Review of  
913 Algorithms and Applications. *IEEE Transactions on Emerging Topics in*  
914 *Computational Intelligence.* 2019;PP:1-13. doi: 10.1109/TETCI.2019.2952908
- 915 46. Noh JD, Rieger H. Random walks on complex networks. *Phys Rev Lett.*  
916 2004;92:118701. doi: 10.1103/PhysRevLett.92.118701
- 917 47. Sifrim A, Hitz MP, Wilsdon A, Breckpot J, Turki SH, Thienpont B, McRae J,  
918 Fitzgerald TW, Singh T, Swaminathan GJ, et al. Distinct genetic architectures for

- 919 syndromic and nonsyndromic congenital heart defects identified by exome  
920 sequencing. *Nat Genet.* 2016;48:1060-1065. doi: 10.1038/ng.3627
- 921 48. Zhu Y, Gramolini AO, Walsh MA, Zhou YQ, Slorach C, Friedberg MK, Takeuchi  
922 JK, Sun H, Henkelman RM, Backx PH, et al. Tbx5-dependent pathway regulating  
923 diastolic function in congenital heart disease. *Proc Natl Acad Sci U S A.*  
924 2008;105:5519-5524. doi: 10.1073/pnas.0801779105
- 925 49. Pashmforoush M, Lu JT, Chen H, Amand TS, Kondo R, Pradervand S, Evans  
926 SM, Clark B, Feramisco JR, Giles W, et al. Nkx2-5 pathways and congenital  
927 heart disease; loss of ventricular myocyte lineage specification leads to  
928 progressive cardiomyopathy and complete heart block. *Cell.* 2004;117:373-386.  
929 doi: 10.1016/s0092-8674(04)00405-2
- 930 50. Xu M, Wu X, Li Y, Yang X, Hu J, Zheng M, Tian J. CITED2 mutation and  
931 methylation in children with congenital heart disease. *J Biomed Sci.* 2014;21:7.  
932 doi: 10.1186/1423-0127-21-7
- 933 51. Tomita-Mitchell A, Mahnke DK, Struble CA, Tuffnell ME, Stamm KD, Hidestrand  
934 M, Harris SE, Goetsch MA, Simpson PM, Bick DP, et al. Human gene copy  
935 number spectra analysis in congenital heart malformations. *Physiol Genomics.*  
936 2012;44:518-541. doi: 10.1152/physiolgenomics.00013.2012
- 937 52. Pu T, Liu Y, Xu R, Li F, Chen S, Sun K. Identification of ZFPM2 mutations in  
938 sporadic conotruncal heart defect patients. *Mol Genet Genomics.* 2018;293:217-  
939 223. doi: 10.1007/s00438-017-1373-6
- 940 53. De Luca A, Sarkozy A, Ferese R, Consoli F, Lepri F, Dentici ML, Vergara P, De  
941 Zorzi A, Versacci P, Digilio MC, et al. New mutations in ZFPM2/FOG2 gene in

- 942 tetralogy of Fallot and double outlet right ventricle. *Clin Genet.* 2011;80:184-190.  
943 doi: 10.1111/j.1399-0004.2010.01523.x
- 944 54. Ma L, Selamet Tierney ES, Lee T, Lanzano P, Chung WK. Mutations in ZIC3 and  
945 ACVR2B are a common cause of heterotaxy and associated cardiovascular  
946 anomalies. *Cardiol Young.* 2012;22:194-201. doi: 10.1017/S1047951111001181
- 947 55. Marino BS, Lipkin PH, Newburger JW, Peacock G, Gerdes M, Gaynor JW,  
948 Mussatto KA, Uzark K, Goldberg CS, Johnson WH, Jr., et al.  
949 Neurodevelopmental outcomes in children with congenital heart disease:  
950 evaluation and management: a scientific statement from the American Heart  
951 Association. *Circulation.* 2012;126:1143-1172. doi:  
952 10.1161/CIR.0b013e318265ee8a
- 953 56. Razzaghi H, Oster M, Reefhuis J. Long-term outcomes in children with  
954 congenital heart disease: National Health Interview Survey. *J Pediatr.*  
955 2015;166:119-124. doi: 10.1016/j.jpeds.2014.09.006
- 956 57. Kocak G, Onal C, Kocak A, Karakurt C, Ates O, Cayli SR, Yologlu S. Prevalence  
957 and outcome of congenital heart disease in patients with neural tube defect. *J*  
958 *Child Neurol.* 2008;23:526-530. doi: 10.1177/0883073807309789
- 959 58. Ramakrishna S, Kim IM, Petrovic V, Malin D, Wang IC, Kalin TV, Meliton L, Zhao  
960 YY, Ackerson T, Qin Y, et al. Myocardium defects and ventricular hypoplasia in  
961 mice homozygous null for the Forkhead Box M1 transcription factor. *Dev Dyn.*  
962 2007;236:1000-1013. doi: 10.1002/dvdy.21113
- 963 59. Bolte C, Zhang Y, Wang IC, Kalin TV, Molkentin JD, Kalinichenko VV.  
964 Expression of Foxm1 transcription factor in cardiomyocytes is required for

- 965 myocardial development. *PLoS One*. 2011;6:e22217. doi:  
966 10.1371/journal.pone.0022217
- 967 60. Bult CJ, Blake JA, Smith CL, Kadin JA, Richardson JE, Mouse Genome  
968 Database G. Mouse Genome Database (MGD) 2019. *Nucleic Acids Res*.  
969 2019;47:D801-D806. doi: 10.1093/nar/gky1056
- 970 61. Marcu R, Choi YJ, Xue J, Fortin CL, Wang Y, Nagao RJ, Xu J, MacDonald JW,  
971 Bammler TK, Murry CE, et al. Human Organ-Specific Endothelial Cell  
972 Heterogeneity. *iScience*. 2018;4:20-35. doi: 10.1016/j.isci.2018.05.003
- 973 62. MacGrogan D, Munch J, de la Pompa JL. Notch and interacting signalling  
974 pathways in cardiac development, disease, and regeneration. *Nat Rev Cardiol*.  
975 2018;15:685-704. doi: 10.1038/s41569-018-0100-2
- 976 63. Mukhtar T, Breda J, Grison A, Karimaddini Z, Grobecker P, Iber D, Beisel C, van  
977 Nimwegen E, Taylor V. Tead transcription factors differentially regulate cortical  
978 development. *Sci Rep*. 2020;10:4625. doi: 10.1038/s41598-020-61490-5
- 979 64. Kaneko KJ, Kohn MJ, Liu C, DePamphilis ML. Transcription factor TEAD2 is  
980 involved in neural tube closure. *Genesis*. 2007;45:577-587. doi:  
981 10.1002/dvg.20330
- 982 65. Adam MP, Hudgins L. Kabuki syndrome: a review. *Clin Genet*. 2005;67:209-219.  
983 doi: 10.1111/j.1399-0004.2004.00348.x
- 984 66. Digilio MC, Gnazzo M, Lepri F, Dentici ML, Pisaneschi E, Baban A, Passarelli C,  
985 Capolino R, Angioni A, Novelli A, et al. Congenital heart defects in molecularly  
986 proven Kabuki syndrome patients. *Am J Med Genet A*. 2017;173:2912-2922. doi:  
987 10.1002/ajmg.a.38417



- 988 67. Li L, Ruan X, Wen C, Chen P, Liu W, Zhu L, Xiang P, Zhang X, Wei Q, Hou L, et  
989 al. The COMPASS Family Protein ASH2L Mediates Corticogenesis via  
990 Transcriptional Regulation of Wnt Signaling. *Cell Rep.* 2019;28:698-711 e695.  
991 doi: 10.1016/j.celrep.2019.06.055
- 992 68. Jung Y, Hsieh LS, Lee AM, Zhou Z, Coman D, Heath CJ, Hyder F, Mineur YS,  
993 Yuan Q, Goldman D, et al. An epigenetic mechanism mediates developmental  
994 nicotine effects on neuronal structure and behavior. *Nat Neurosci.* 2016;19:905-  
995 914. doi: 10.1038/nn.4315
- 996 69. Landim-Vieira M, Johnston JR, Ji W, Mis EK, Tijerino J, Spencer-Manzon M,  
997 Jeffries L, Hall EK, Panisello-Manterola D, Khokha MK, et al. Familial Dilated  
998 Cardiomyopathy Associated With a Novel Combination of Compound  
999 Heterozygous TNNC1 Variants. *Front Physiol.* 2019;10:1612. doi:  
1000 10.3389/fphys.2019.01612
- 1001 70. Stein AB, Jones TA, Herron TJ, Patel SR, Day SM, Noujaim SF, Milstein ML,  
1002 Klos M, Furspan PB, Jalife J, et al. Loss of H3K4 methylation destabilizes gene  
1003 expression patterns and physiological functions in adult murine cardiomyocytes.  
1004 *J Clin Invest.* 2011;121:2641-2650. doi: 10.1172/JCI44641
- 1005 71. Stoller JZ, Huang L, Tan CC, Huang F, Zhou DD, Yang J, Gelb BD, Epstein JA.  
1006 Ash2l interacts with Tbx1 and is required during early embryogenesis. *Exp Biol*  
1007 *Med (Maywood).* 2010;235:569-576. doi: 10.1258/ebm.2010.009318
- 1008 72. Barron DJ, Kilby MD, Davies B, Wright JG, Jones TJ, Brawn WJ. Hypoplastic left  
1009 heart syndrome. *Lancet.* 2009;374:551-564. doi: 10.1016/S0140-6736(09)60563-  
1010 8

- 1011 73. Kuroda T, Yasuda S, Tachi S, Matsuyama S, Kusakawa S, Tano K, Miura T,  
1012 Matsuyama A, Sato Y. SALL3 expression balance underlies lineage biases in  
1013 human induced pluripotent stem cell differentiation. *Nat Commun.* 2019;10:2175.  
1014 doi: 10.1038/s41467-019-09511-4  
1015  
1016

1017 **Table 1. Enriched Gene Ontology Terms for Group-I and II Genes on the Network.**

	Sources	Function Terms	FDR
<b>G-I</b>	Gene Ontology	Histone lysine methylation	7.28E-05
		Neuronal ion channel clustering	5.01E-03
		Neuron differentiation	1.31E-02
	MGI Phenotypes	Exencephaly	1.17E-04
		Open neural tube	1.28E-04
		Right atrial isomerism	8.35E-03
		Abnormal cardiovascular development	2.25E-02
<b>G-II</b>	Gene Ontology	Ventricular septum development	6.58E-06
		Regulation of cardiac muscle cell proliferation	2.03E-05
		Outflow tract septum morphogenesis	2.03E-05
		BMP signaling pathway	2.03E-05
	MGI Phenotypes	Atrioventricular block	3.60E-03
		Abnormal heart right ventricle morphology	3.60E-03
		Abnormal heart atrium morphology	7.62E-03

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030 **Figure Legends**

1031 **Figure 1. NetWalker Identified a highly connected network in CHD.**

1032 **A.** PCGC candidate genes tended to occupy central positions on the protein interaction  
1033 network. PCGC candidate genes (red) were identified by dosage-sensitive genes  
1034 affected by de novo loss-of-function (LoF) mutations from the PCGC CHD cohorts. The  
1035 same procedure also identified affected genes (blue) from the matched unaffected  
1036 sibling cohort. Network connectivity indicates the number of interacting partners for  
1037 each protein on the network. As an independent control experiment, the same  
1038 comparison was also performed on genes affected by de novo synonymous mutations,  
1039 whose functions were presumably neutral. P-values were derived from Wilcoxon rank-  
1040 sum test.

1041 **B.** PCGC candidate genes were more likely to maintain mutual interactions on the  
1042 network. The fractions of interacting proteins were computed among the candidate and  
1043 control genes identified from proband and sibling cohorts, respectively. Genes affected  
1044 by de novo synonymous mutations were used as an independent control experiment. P-  
1045 values were derived from the Fisher's exact test.

1046 **C.** The schematic presentation of the random walk algorithm on the protein interaction  
1047 network. The random walk scheme starts from every node on the network following  
1048 stochastic flow till convergence. For each protein, the probabilities of visiting all other  
1049 proteins on the network will be calculated, which defines the reachability of this node to  
1050 any other nodes on the network.

1051 **D.** The identified network component has substantially increased reachability to PCGC  
1052 CHD candidate proteins relative to all other proteins on the network. P-values were  
1053 derived from Wilcoxon rank-sum test.

1054

1055 **Figure 2. Functional characterization of the CHD network.**

1056 **A.** An overview of the identified network seeded with PCGC candidate proteins  
1057 (grouped by CHD subtypes that were color coded), and the orange nodes were novel  
1058 proteins identified by the NetWalker algorithm. The subtype annotations were derived  
1059 from the original publication<sup>10</sup>, where CTD, HTX and LVO stand for conotruncal defects,  
1060 heterotaxy and left ventricular outflow tract obstruction, respectively. Other indicates  
1061 more than one subtype was associated with the corresponding proteins.

1062 **B.** Temporal expression of the network genes across heart developmental stages.  
1063 Human genes were mapped onto mouse orthologs, and hierarchical clustering revealed  
1064 two expression components of the network, where Group-I (G-I) genes displayed  
1065 preferential expression from embryonic stem cells (ESC) to E7.3, whereas Group-II (G-  
1066 II) genes exhibited substantial expression enrichment from E7.3 to postnatal and adult  
1067 stages. Close examination of Group-II genes further revealed two subcluster structure,  
1068 where Group-II-A (G-II-A) was preferentially expressed across fetal developmental  
1069 stages and Group-II-B (G-II-B) was more specific in postnatal stages, particularly strong  
1070 in the adult heart.

1071

1072 **Figure 3. Spatiotemporal expression analysis of the identified CHD network.**

1073 **A-B.** Group-I (G-I) genes and Group-II-B (G-II-B) genes did not show significance  
1074 across all time points. Group-II-A (G-II-A) genes showed pervasive expression activities  
1075 across most spatial spots in the fetal heart in both PCW 6.5 (**A**) and 9 (**B**).

1076 **C.** Group-II-A (G-II-A) genes showed significantly increased expression in the fetal heart  
1077 from postconceptional day 96 to day 147, whereas Group-I (G-I) genes and Group-II-B  
1078 (G-II-B) genes did not display statistical significance ( $p > 0.05$ , Wilcoxon rank-sum test).

1079 **D.** Group-II-A (G-II-A) genes showed significantly increased expression in the fetal heart  
1080 in the gestational weeks 19 and 28 based on RNA-seq data. The statistical significance  
1081 was not observed from other gene groups.

1082

1083 **Figure 4. Validating novel functions of the identified genes in regulating fetal**  
1084 **heart development.**

1085 **A-C.** Network clustering identified 33 local clustering structures on the identified network,  
1086 where cluster #4 (**A**), #2 (**B**), #3 (**C**) were presented as representative pathways  
1087 regulating heart development.

1088 **D.** Gene ontology enrichment of the differentially expressed genes in iPSC-CMs upon  
1089 siRNA knockdown of TEAD2, TLK1 and RBBP5, respectively. These differentially  
1090 expressed genes consistently displayed strong functional enrichment for heart  
1091 development and cardiac muscle contraction. The color intensities of the circles  
1092 represent false discovery rates (FDRs). Sizes of the circles represent the enrichment  
1093 scores.

1094 **E,F,L,M.** RNA-seq identified differentially expressed genes in iPSC-CMs upon siRNA  
1095 knockdown of TEAD2 (**E**), TLK1 (**F**), RBBP5 (**L**) and ASH2L (**M**), respectively. X-axis is

1096 the mean expression of each gene in iPSC-CMs, and Y-axis indicates their respective  
1097 fold changes upon siRNA knockdown (siRNA treatment vs. siRNA control). Genes with  
1098 false discovery rates (FDRs) less than 0.05 were highlighted in red.

1099 **G-K.** Cellular contractility assay in the iPSC-CMs. siRNA knockdown against TEAD2 in  
1100 iPSC-CMs resulted in a marked reduction of the cardiomyocyte beating rate (**G**),  
1101 increased contraction velocity (**H**), contraction deformation distance(**I**), relaxation  
1102 velocity (**J**) and relaxation deformation distance (**K**) relative to the siRNA control. P-  
1103 values were derived from t-test.

1104 **N.** Cellular contractility assay in the iPSC-CMs. RBBP5 knockdown in iPSC-CMs  
1105 displayed an increased beating rate. P-values were derived from t-test.

1106 **O-R.** Cellular contractility assay in the iPSC-CMs. ASH2L knockdown in iPSC-CMs  
1107 showed increased contraction velocity (**O**), contraction deformation distance (**P**),  
1108 relaxation velocity (**Q**) and relaxation deformation distance (**R**). P-values were derived  
1109 from t-test.

1110 **S-T.** ASH2L<sup>+/-</sup> knockout lines clone 1 (**S**) and clone 2 (**T**) significantly reduced the  
1111 differentiation efficiencies into cardiomyocytes (TNNT2 positive cells) from iPSCs. P-  
1112 values were derived from t-test.

1113

1114 **Figure 5. Group-II genes on the network were enriched for pathogenic mutations**  
1115 **in probands with HLHS.**

1116 **A.** Rare non-synonymous variants displayed a significant increase in mutational  
1117 deleteriousness in HLHS probands, but not in individuals with ASD, VSD, TGA or TOF,  
1118 relative to unaffected siblings.

1119 **B.** Significant differences in mutational deleteriousness of rare synonymous mutations  
1120 were not observed in any CHD subtypes relative to unaffected siblings. HLHS, ASD,  
1121 VSD, TGA and TOF stand for hypoplastic left heart syndrome, atrial septal defects,  
1122 ventricular septal defects, transposition of the great arteries, and tetralogy of fallot.  
1123 Mutational pathogenicity was measured by CADD scores. We computed the mean  
1124 CADD scores for non-synonymous in Group-II genes in each personal exome, and we  
1125 compared the mean CADD score distribution among probands in each CHD subtype  
1126 against the distribution among the unaffected siblings. The same comparison was  
1127 performed on synonymous variants as a set of negative controls. Controls were the  
1128 unaffected siblings from the original publication<sup>10</sup>. P-values were derived from Wilcoxon  
1129 rank-sum test.

1130 **C.** Excluding individuals comorbid with ASD further boosted statistical significance  
1131 (Group-II genes).

1132 **D.** The statistical significance was absent on lung-specific genes for both non-  
1133 synonymous and synonymous variants.

1134 **E.** Permutation analysis confirmed that CADD scores had a similar distribution for rare  
1135 non-synonymous variants in the background exomes in the HLHS cohort relative to the  
1136 control cohort. In each permutation, rare non-synonymous variants were randomly  
1137 sampled from the exome backgrounds in the HLHS and control cohorts respectively,  
1138 matching the number of rare non-synonymous variants in Group-II genes in cases or in  
1139 controls, followed a comparison between their CADD scores using Wilcoxon rank-sum  
1140 test. Among 100 permutations, 98 were statistically insignificant ( $p = 0.98$ ), suggesting



1141 that our observation cannot be merely explained by exome background differences in  
1142 CADD scores. Error bars represent standard error of the mean.

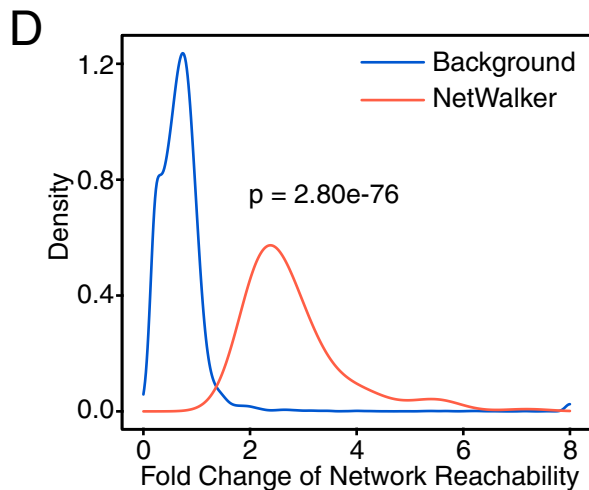
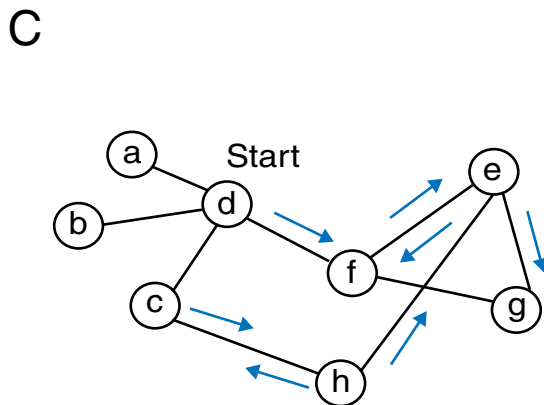
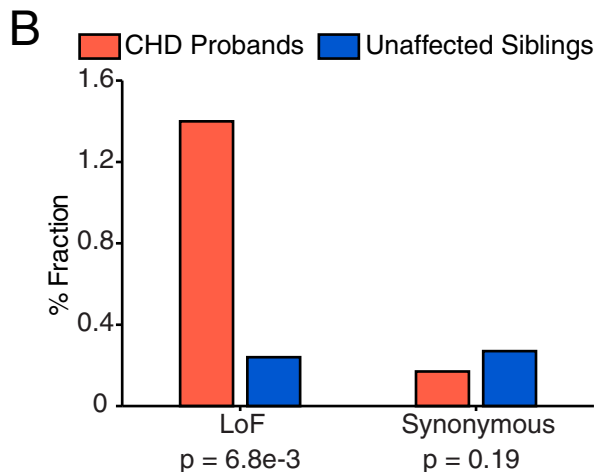
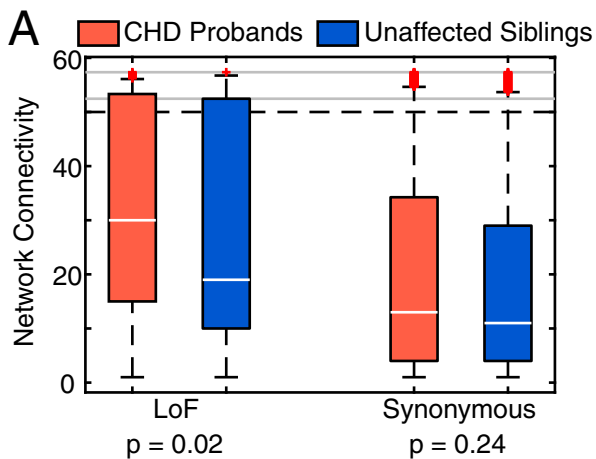
1143

1144 **Figure 6. Single-cell analysis of the network in the HLHS heart.**

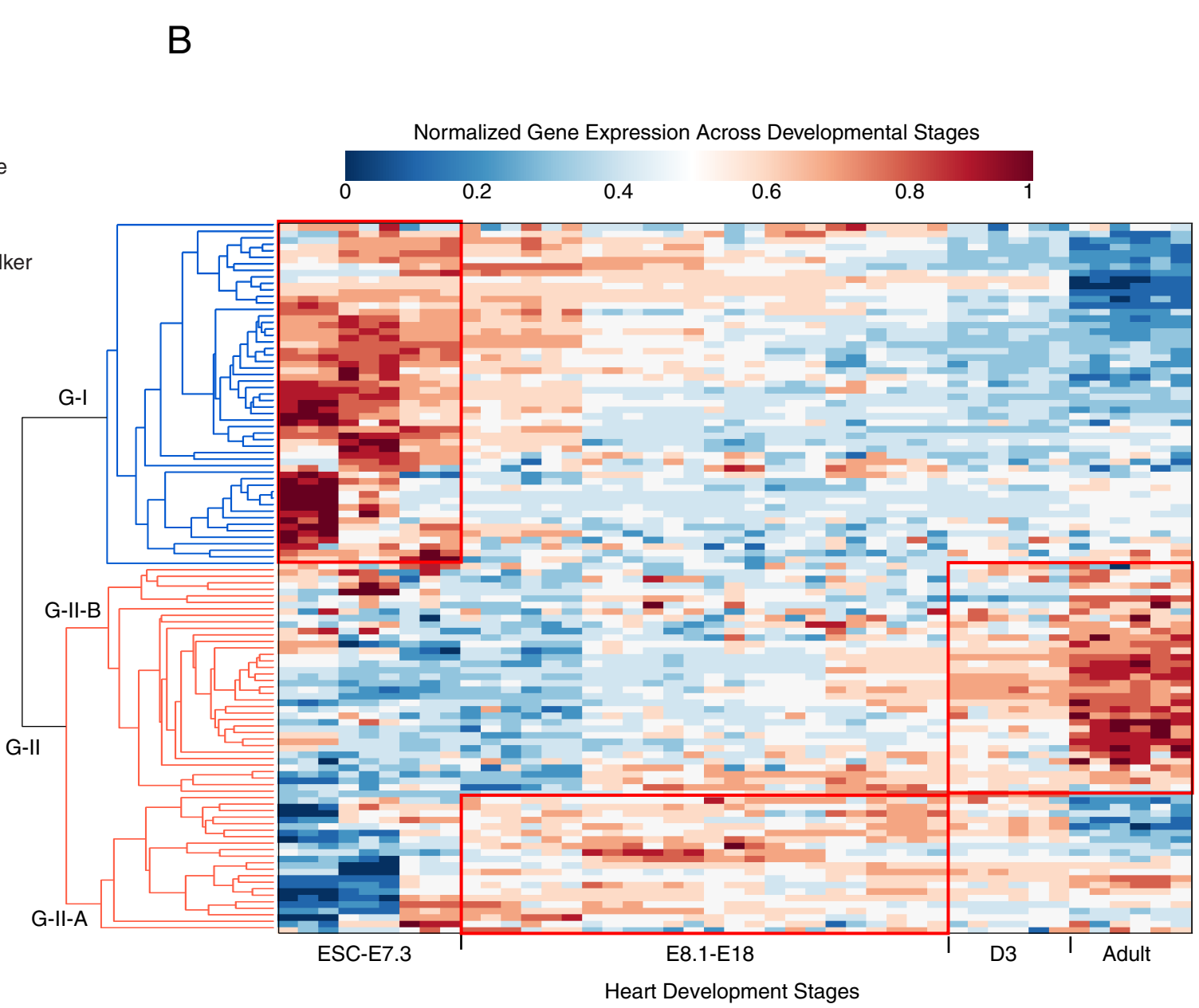
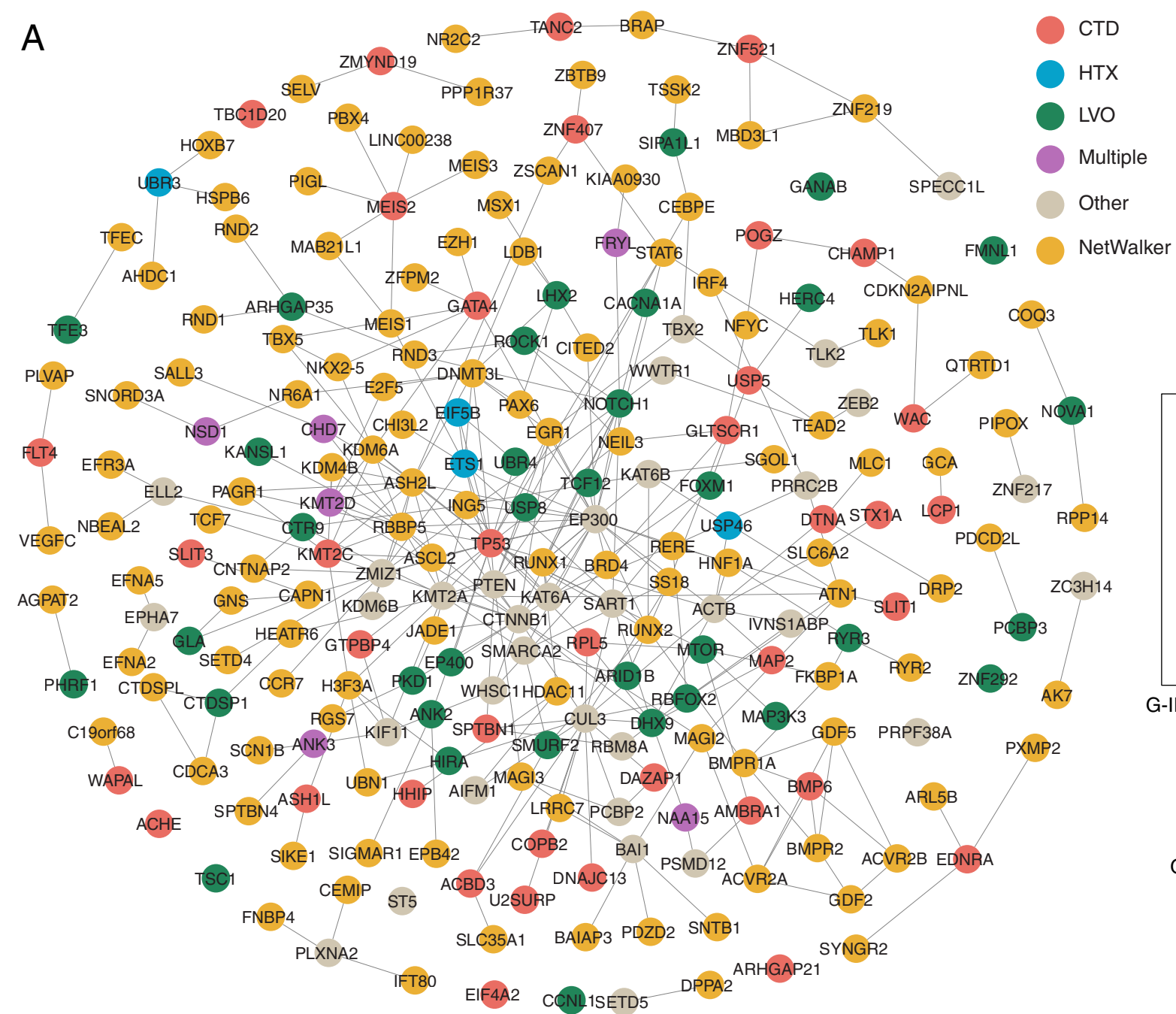
1145 **A.** Across all the cell types in the HLHS left ventricle, Group-II (G-II) genes showed the  
1146 strongest expression reduction in the endothelium cells. P-values were derived from  
1147 Wilcoxon rank-sum test,

1148 **B.** Examining two subgroups of Group-II (G-II) genes revealed significant  
1149 downregulation of both subgroups (G-II-A and G-II-B) in the endothelium cells.

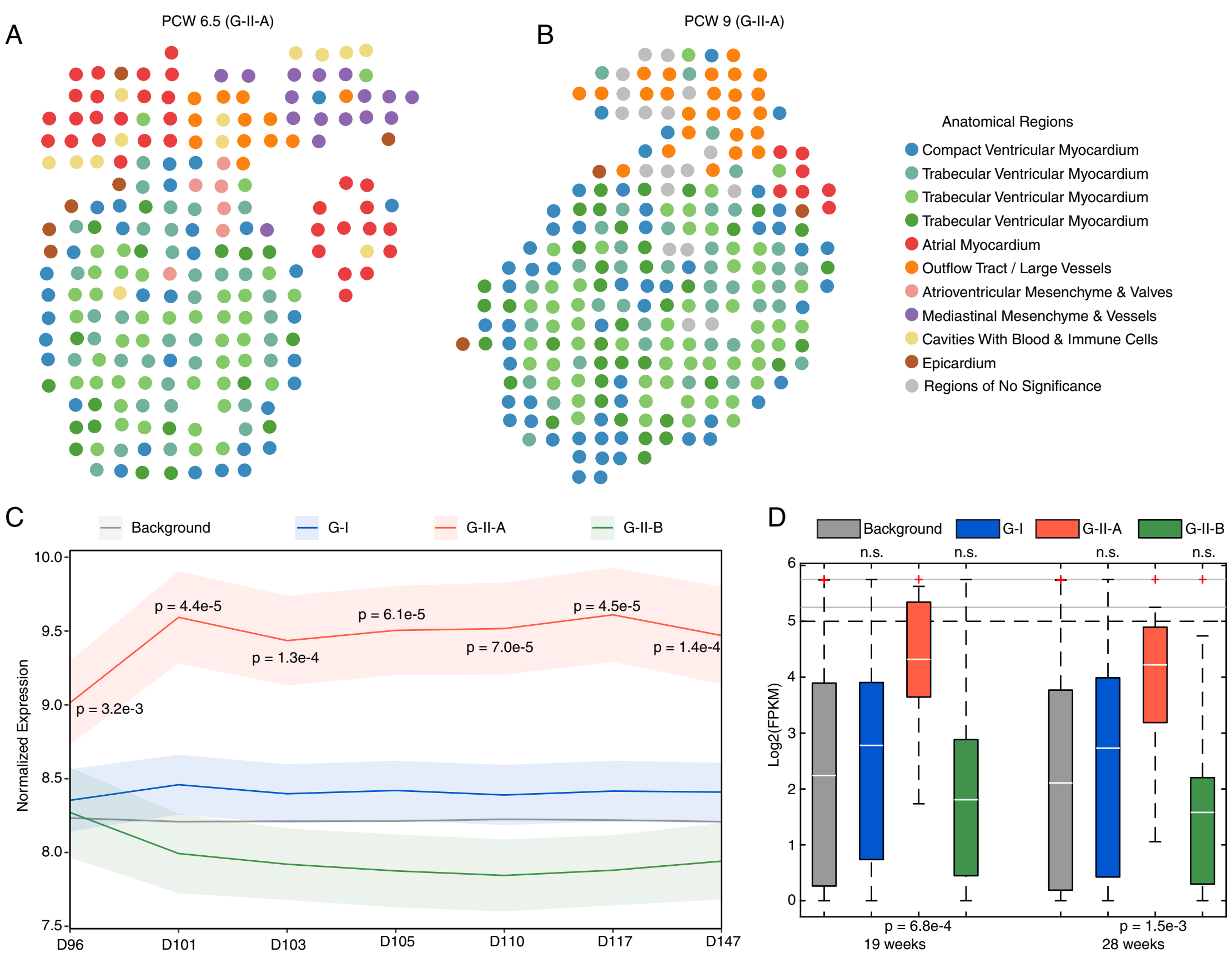
1150 **C.** Only Group-II-A (G-II-A) genes displayed significant expression reduction in the  
1151 conduction system in the HLHS left ventricle.



**Figure 1**



**Figure 2**



**Figure 3**

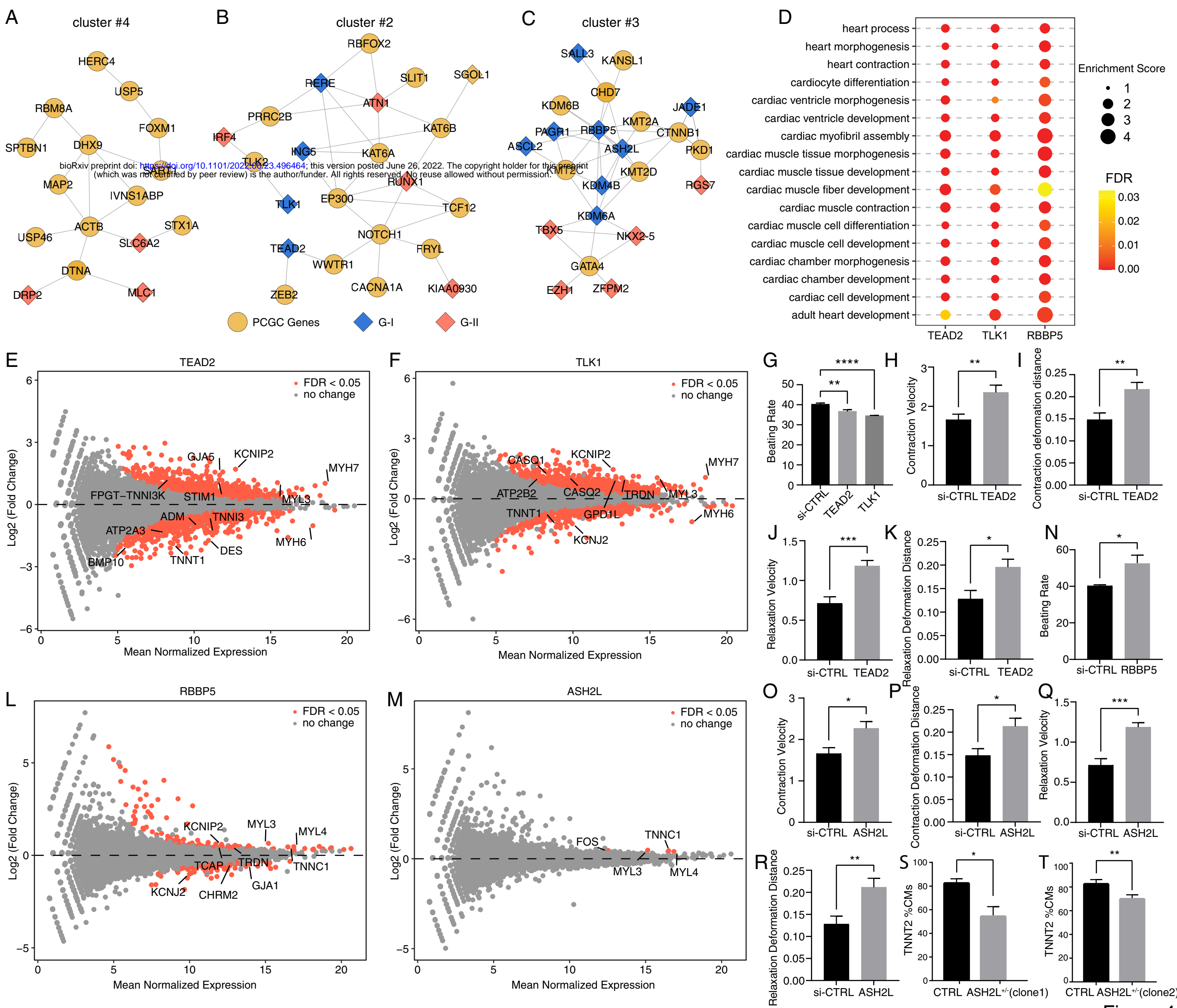
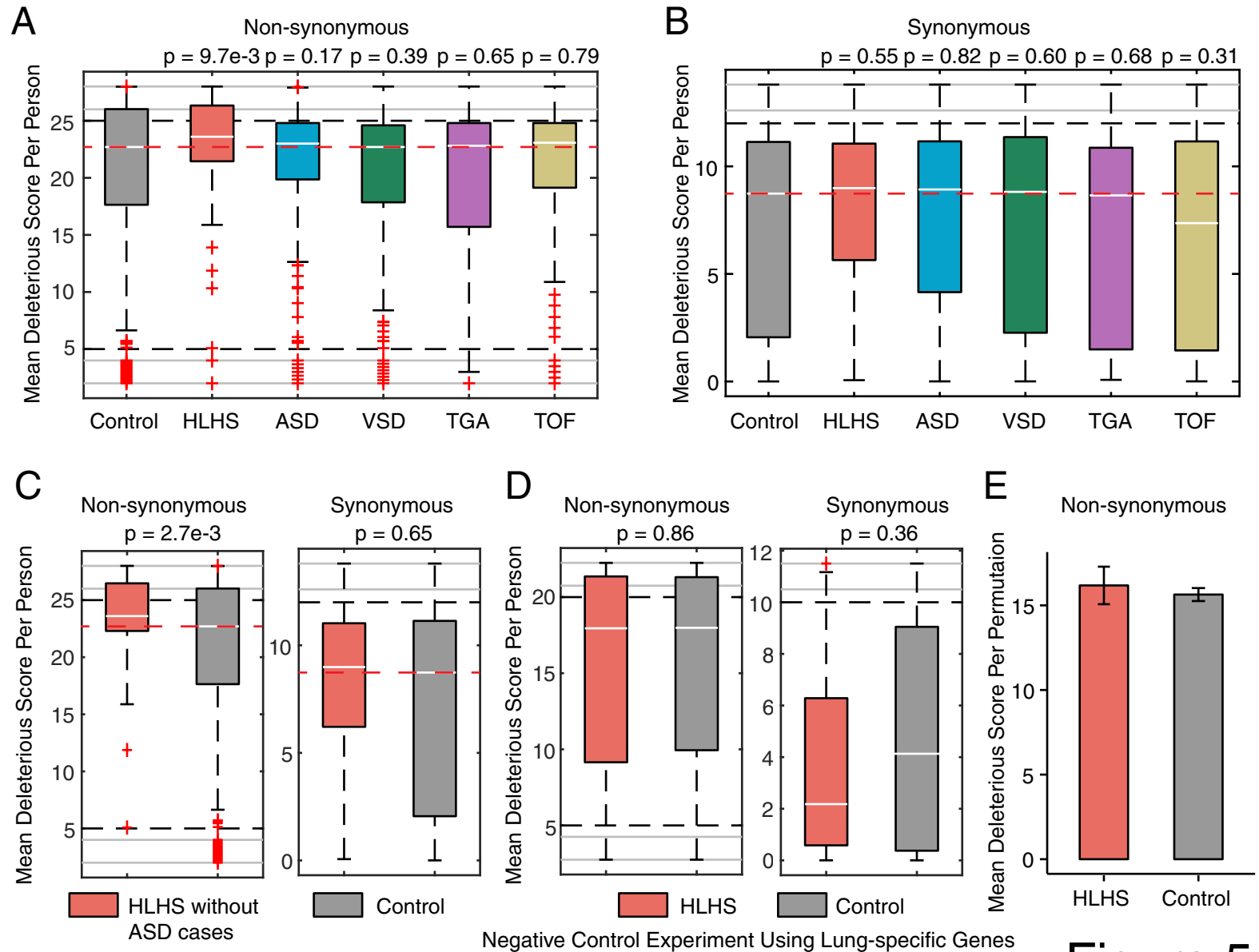
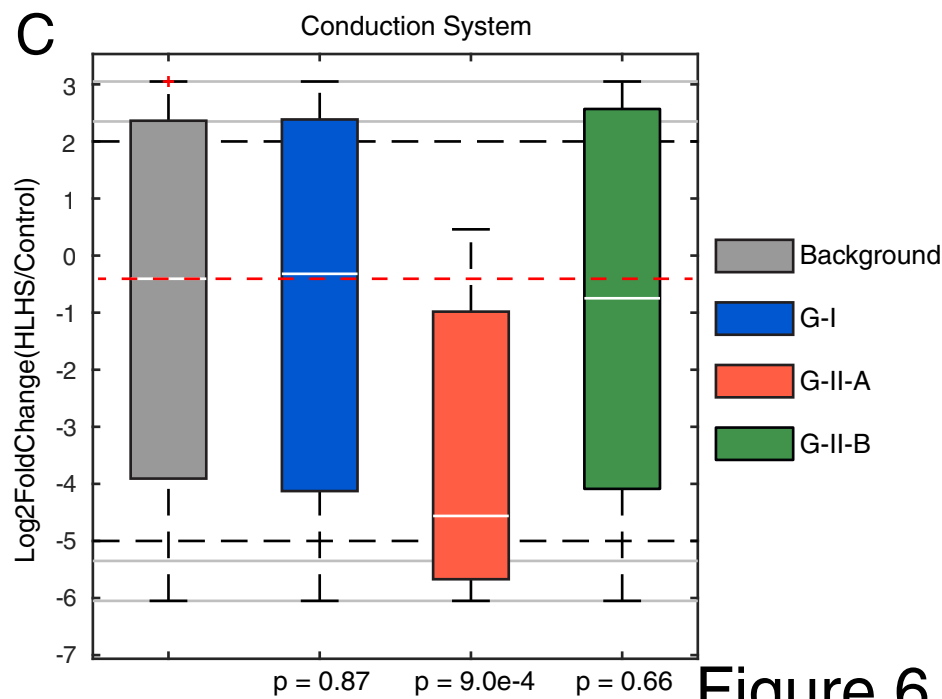
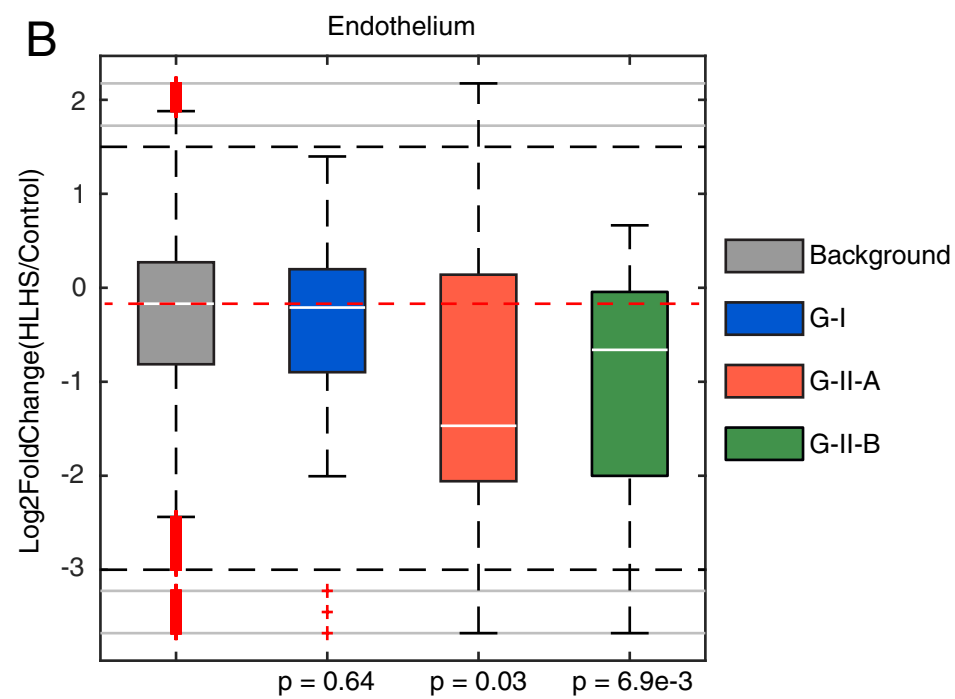
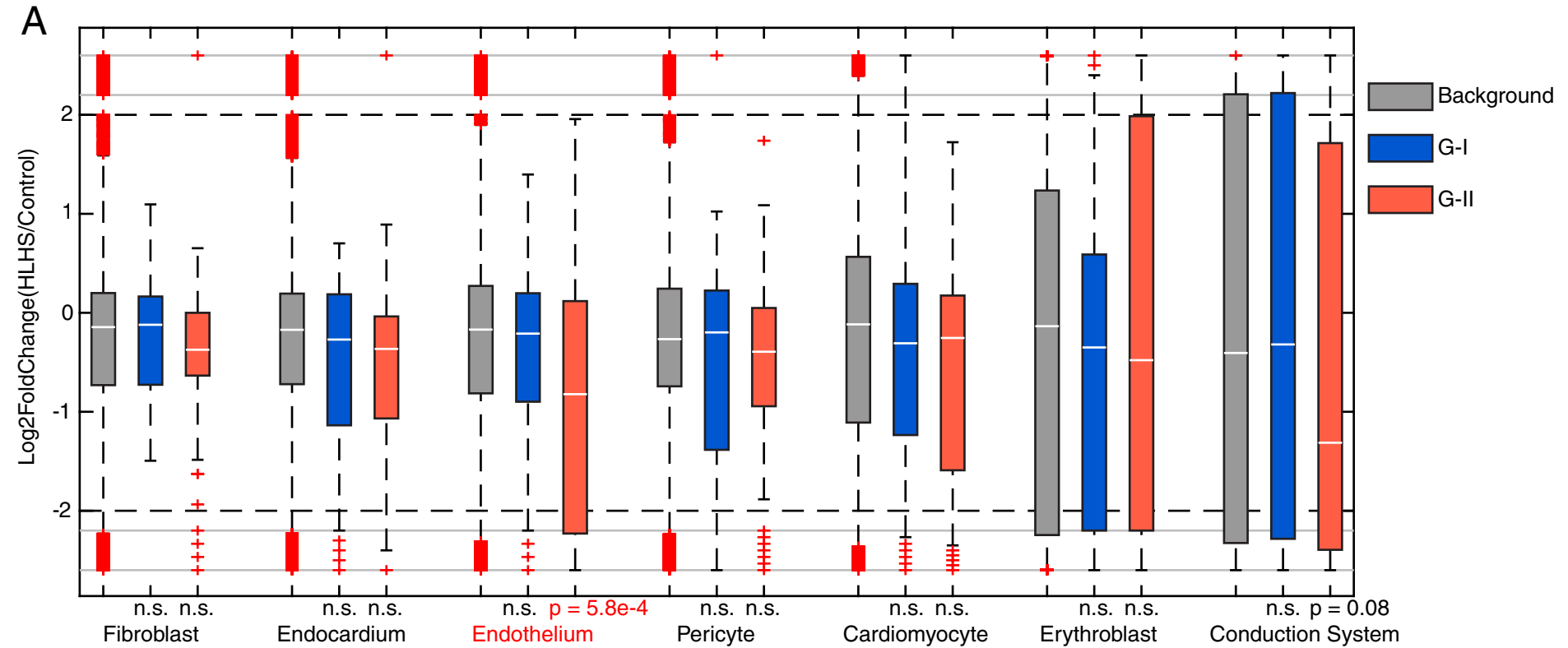


Figure 4



**Figure 5**



**Figure 6**