

Network Assisted Analysis of *De Novo* Variants Using Protein-Protein Interaction Information Identified 46 Candidate Genes for Congenital Heart Disease

Yuhan Xie¹, Wei Jiang¹, Weilai Dong², Hongyu Li¹, Sheng Chih Jin³, Martina Brueckner^{2,4}, Hongyu Zhao^{1,2,5*}

¹ Department of Biostatistics, Yale School of Public Health, New Haven, CT, USA 06520

² Department of Genetics, Yale School of Medicine, New Haven, CT, USA 06520

³ Department of Genetics, Washington University School of Medicine, St Louis, MO, USA 63110

⁴ Department of Pediatrics, Yale University, New Haven, CT, USA 06520

⁵ Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA 06520

* To whom correspondence should be addressed:

Hongyu Zhao, Ph.D.

Department of Biostatistics, Yale School of Public Health, 60 College Street, New Haven, CT, 06520, USA

Email: hongyu.zhao@yale.edu

Abstract

De novo variants (DNVs) with deleterious effects have proved informative in identifying risk genes for early-onset diseases such as congenital heart disease (CHD). A number of statistical methods have been proposed for family-based studies or case/control studies to identify risk genes by screening genes with more DNVs than expected by chance in Whole Exome Sequencing (WES) studies. However, the statistical power is still limited for cohorts with thousands of subjects. Under the hypothesis that connected genes in protein-protein interaction (PPI) networks are more likely to share similar disease association status, we develop a Markov Random Field model that can leverage information from publicly available PPI databases to increase power in identifying risk genes. We identified 46 candidate genes with at least 1 DNV in the CHD study cohort, including 18 known human CHD genes and 35 highly expressed genes in mouse developing heart. Our results may shed new insight on the shared protein functionality among risk genes for CHD.

Keywords

De novo variants, protein-protein interactions, congenital heart disease, network modeling

1 **Background**

2 Congenital heart disease (CHD) is the most common birth defect affecting ~ 1% of live births
3 and accounts for one-third of all major congenital abnormalities [1-3]. There is substantial
4 evidence that CHD has a strong genetic component [4]. Although it is estimated that
5 aneuploidies and copy number variations account for about 23% of CHD cases, few individual
6 disease-causing genes have been identified in published studies [5-8]. Therefore, the limited
7 knowledge on the underlying genetic causes poses an obstacle to the reproductive counseling
8 of CHD patients [9].

9
10 Whole Exome Sequencing (WES) studies have successfully boosted novel causal genes
11 identification for both Mendelian and complex disorders [10, 11]. To narrow down the pool of
12 candidate variants from WES, family-based studies have been conducted to scan for *de novo*
13 variants (DNVs) from parent-offspring trios. DNV studies have been shown to play an important
14 role in risk gene identification for CHD [1, 3, 5, 6, 12-15]. From the analysis of 1,213 CHD
15 parent-offspring trios, Homsy et al. identified greater burden of damaging DNVs, especially in
16 genes with likely functional roles in heart and brain development [12]. Recently, Jin et al.
17 inferred that DNVs in ~440 genes were likely contributors to CHD [5]. Despite these advances, it
18 remains challenging to capture the causal genes with only DNV data as CHD is very genetically
19 heterogeneous [6].

20
21 It is believed that genes interact with each other in biological processes and may jointly affect
22 the disease risk through biological pathways. In recent studies on CHD, researchers reported
23 significant enrichment of genes related to histone modification, chromatin modification, cilia
24 function, transcriptional regulation, neural tube development, and cardiac development and
25 enrichment in pathways including Wnt, Notch, Igf, HDAC, ErbB and NF-kb signaling [1, 3, 12, 14,
26 16]. These results suggest that considering functional relationships of genes through pathway-

27 level analyses may complement gene-level analyses and improve risk gene identification for
28 CHD.

29
30 Recently, Nguyen et al. proposed a framework to conduct pathway-level analysis on
31 schizophrenia DNV data [17]. Their method is built on the extTADA framework and integrates
32 pathway information by multiplying a gene-set related term in the likelihood function. However,
33 this method requires curated pathways related to schizophrenia and treats genes as
34 exchangeable without considering gene-gene interactions. To characterize the functional
35 connectivity between genes, we can consider genes as a network by modeling the topology of
36 pathways.

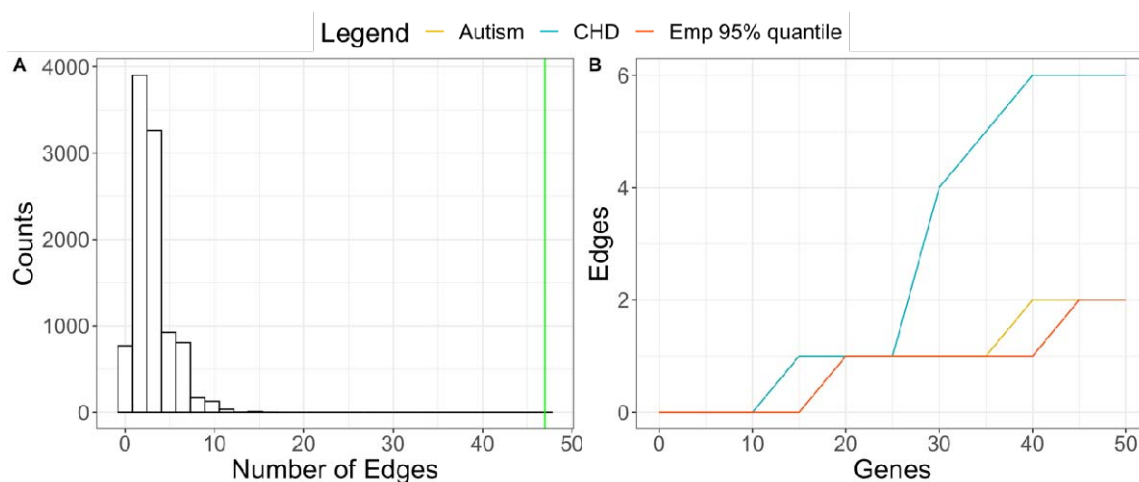
37
38 Network-based approaches have been successful in prioritizing risk genes for downstream
39 analysis of Genome-Wide Association Studies (GWAS) and gene expression studies [18-20].
40 Chen et al. [19] proposed a Markov Random Field (MRF) model to incorporate pathway
41 topology structure for GWAS. They showed that their method is more powerful than single
42 gene-based methods through both simulation and real data analyses. Following the idea of this
43 framework, Hou et al. [20] proposed a method that integrates co-expression networks and
44 GWAS results. By applying their method to Crohn's disease and Parkinson's disease, they
45 showed that their method could lead to more replicable results and find potential disease-
46 associated pathways. More recently, Liu et al. adopted a similar idea as Chen et al. and Hou et
47 al. to analyze DNV data from WES studies [21]. Their framework, namely DAWN, combines
48 TADA p-values with estimated network from gene co-expression data. In their real data analysis
49 for autism, 333 genes were prioritized by integrating DNV summary statistics and expression
50 data from brain tissue. However, all of these above methods require summary statistics (Z
51 scores or p-values) from genetic association analysis as their input, which may not be provided
52 for DNV analysis results [22, 23]. In addition, it is difficult to apply DAWN directly to CHD

53 because of the lack of ideal co-expression data from relevant tissues (such as developmental
54 heart) for CHD.

55
56 As there is a limited number of co-expression data sets for human developmental heart, a
57 natural choice for network information would be human protein-protein interaction (PPI)
58 databases. There are multiple public sources of PPI network databases such as BioGRID [24],
59 IntAct [25], DIP [26], MINT [27], and HPRD [28]. Most of network-based studies apply their real
60 data on two or more of the databases to obtain their results. Nonetheless, it is hard to check the
61 overlapping information between two PPI databases and interpret the divergent results. The
62 STRING [29] database provides a platform to resolve the above problems. It imports protein
63 association knowledge from physical interaction and curated knowledge from the
64 aforementioned PPI databases and other pathway information knowledge such as KEGG and
65 GO. In addition, it provides a score to measure the likelihood of interactions. Some studies have
66 used STRING in their post-association analysis for gene-based DNV studies and showed
67 significant enrichment of candidate CHD risk genes in the STRING PPI network [30, 31]. One
68 recent study [3] constructed a priority score based on canonical pathways and PPI networks
69 and identified 23 novel genes for CHD. These results suggest that incorporating PPI network
70 information from STRING may identify additional risk genes with more biological interpretability.

71
72 In addition, in the post-association analysis of CHD DNV analysis from our previous multi-trait
73 method M-DATA [32], we found that the number of edges formed by candidate CHD genes (47
74 edges, green line in Fig 1A) identified by M-DATA multi-trait model (33 genes with false
75 discovery rate (FDR) q -values < 0.1) is far larger than the upper tail of the empirical distribution
76 sampled from 33 randomly selected genes in the STRING V11.0 database (score threshold:
77 400) for 10,000 times (Fig 1A). This suggests that the candidate CHD genes are highly enriched
78 in terms of their interactions in the STRING database. To further illustrate that PPI information

79 may contribute to CHD gene discovery, we use the CHD result from M-DATA single-trait
80 analysis and use the result of autism as a comparison. Fig 1B shows the number of edges
81 formed by the top genes (ranked by FDR q-values) for CHD and autism from M-DATA single-
82 trait model with a more stringent selection of PPI edges in the STRING database (score
83 threshold: 950), respectively. To compare with the number of edges formed by randomly
84 selected genes, we plot the 95% quantile of the empirical distribution sampled from random
85 genes in the STRING v11.0 database (score threshold: 950) for 10,000 times as a baseline.
86 When more than 25 top CHD genes are selected, the number of edges formed by these genes
87 is significantly more than that from randomly selected genes, whereas the number of edges
88 formed by the autism top genes did not differ much from randomly selected genes. This
89 suggests top genes in CHD tend to be neighbors in the STRING PPI network.



90
91 **Fig 1. CHD top genes are more connected than randomly selected genes in the STRING PPI network** (A) Empirical
92 distribution of the number of edges formed by 33 randomly selected genes. Green line represents the number of edges formed by
93 the 33 CHD top genes from M-DATA. (B) Number of edges formed by CHD top genes, autism top genes from single-trait analyses
94 and randomly selected genes, respectively.

95
96 Motivated by the observation from Fig 1, we develop a **Network assisted model for De novo**
97 **Association Test** using protein-protein interAction information, named N-DATA, to leverage prior

98 information of interactions among genes from the PPI network to boost statistical power in
 99 identifying risk genes for CHD. In the following, we first introduce the inference procedure for
 100 our model, and then demonstrate the performance of our method through simulation studies
 101 and real data applications.

102

103 **Methods**

104 In this section, we introduce the statistical model for the proposed framework. The network
 105 information in the PPI database is represented by an undirected graph $G = (V, E)$, where
 106 $V = \{1, \dots, n\}$ is a set of n genes in the network, and
 107 $E = \{ \langle i, j \rangle : i \text{ and } j \text{ are genes connected by the edges} \}$. The degree of a gene i is defined as the
 108 number of direct neighbors (N_i) for gene i in the network and denoted as d_i . We denote the
 109 latent association status of gene i with a disease of interest, e.g., CHD, as S_i , where $S_i = 1$ if
 110 gene i is associated with the disease, $S_i = -1$ if gene i is not associated with the disease.
 111 $S = \{S_1, \dots, S_n\}$ are the corresponding latent states for genes in $V = \{1, \dots, n\}$. The DNV count of
 112 the cohort is defined as Y_i . To formalize the assumption that genes that have PPIs with risk
 113 genes are more likely to be risk genes, we apply a nearest neighbor Gibbs measure [33] to
 114 arrive at the following model:

$$P(S|\theta_0) \propto \exp \left\{ h \sum_{i \in V} I_1(S_i) + \tau_0 \sum_{\langle i, j \rangle \in E} (w_i + w_j) I_{-1}(S_i) I_{-1}(S_j) + \tau_1 \sum_{\langle i, j \rangle \in E} (w_i + w_j) I_1(S_i) I_1(S_j) \right\}$$

$$Y_i | S_i = -1 \sim \text{Poisson}(2N\mu_i)$$

$$Y_i | S_i = 1 \sim \text{Poisson}(2N\mu_i\gamma)$$

$$\theta_0 = (h, \tau_0, \tau_1); \theta_1 = \gamma,$$

115

116 where w_i is the weight for gene i and will be chosen based on the characteristics of the
 117 network, $\theta_0 = (h, \tau_0, \tau_1)$ are hyperparameters related to the network. Specifically, h determines

118 the marginal distribution of S_i when all genes are independent i.e., $P(S_i = 1|h, \tau_0 = \tau_1 = 0) =$
119 $\frac{\exp(h)}{1+\exp(h)}$. τ_0 and τ_1 characterize the prior weights of edges between non-associated genes and
120 associated genes, respectively. N is the sample size of the case cohort, μ_i is the mutability of
121 gene i estimated using the framework in Samocha et al. [34], and θ_1 (γ) is the relative risk of
122 the DNVs in the risk gene.

123
124 To reduce the computational burden from a fully Bayesian solution for maximizing the marginal
125 likelihood, we propose an empirical Bayes method to estimate the parameters θ_0 and θ_1 , and
126 the latent association status S by maximizing the pseudo conditional likelihood (PCLK) for n
127 genes as follows

$$\text{PCLK} = \prod_{i=1}^n f(Y_i|S_i, \theta_1) \Pr(S_i|S_{N_i}, \theta_0).$$

128 It has been shown that the estimator from the PCLK in a general Markov random field setting is
129 consistent under mild regularity conditions [19] [35]. When maximizing the PCLK, we can
130 estimate the hyperparameters θ_0 , θ_1 and latent status S iteratively.

131
132 We can obtain an empirical estimate for θ_0 by maximizing $\prod_{i=1}^n \Pr(S_i|S_{N_i}, \theta_0)$, which is
133 equivalent to maximizing the parameters in the following logistic regression model:

$$\text{logit} \Pr(S_i|S_{N_i}, \theta_0) = h + \tau_1 X_{i1} - \tau_0 X_{i0},$$

134 where $X_{i1} = w_i \sum_{k \in N_i} I_1(S_k) + \sum_{k \in N_i} w_k I_1(S_k)$ and $X_{i0} = w_i \sum_{k \in N_i} I_{-1}(S_k) + \sum_{k \in N_i} w_k I_{-1}(S_k)$, $i =$
135 $1, \dots, n$. To make sure the estimated θ_0 is finite, we can add a ridge penalty term $\lambda(h^2 + \tau_0^2 + \tau_1^2)$
136 to the likelihood function to solve the maximization problem by the Newton-Raphson's method
137 [36].

138

139 We then update the latent status S by maximizing the PCLK using the iterative conditional mode
 140 method [35]. After we obtain the updated values θ_0 and S , we can estimate the hyperparameter
 141 θ_1 by maximizing $\prod_{i=1}^n f(Y_i | S_i, \theta_1)$ by using the following closed-form expression:

$$\log L(\theta_1 | Y) \propto \log \prod_{S_i=1} \exp(-2\mu N \gamma + Y_i \log \gamma)$$

$$\frac{\partial \log L(\theta_1 | Y)}{\partial \gamma} = - \sum_{S_i=1} 2\mu N + \frac{\sum_{S_i=1} Y_i}{\gamma}$$

$$\hat{\gamma} = \frac{\sum_{S_i=1} Y_i}{\sum_{S_i=1} 2\mu N}$$

142

Algorithm 1: Procedure for Parameter Estimation

1. Set initial configuration S^0
2. In the j th iteration, for given $s^{(j-1)}$, obtain $\hat{\theta}_0^j$ from

$$\text{logit Pr}(S_i^{(j-1)} | S_{N_i}^{(j-1)}, \theta_0) = h + \tau_1 X_{i1} - \tau_0 X_{i0}, \quad i = 1, \dots, n$$

3. Sequentially update the labels of nodes to obtain $S^{(j)}$ (ICM)

$$S_i^{(j)} = \arg \max_{S_i} f(Y_i | S_i, \hat{\theta}_1^{(j-1)}) \Pr(S_i | S_{N_i}^{(j-1)}, \hat{\theta}_0^{(j)}) \prod_{k \in S_N} \Pr(S_k^{(j-1)} | S_i, S_{N_k-i}^{(j-1)}, \hat{\theta}_0^{(j)})$$

4. Obtain $\hat{\theta}_1^j$ ($\hat{\gamma}^{(j)}$) from

$$\hat{\theta}_1^{(j)} = \operatorname{argmax}_{\theta_1} \log L(\theta_1 | \theta_0^{(j)}, S^{(j)}, Y)$$

5. Repeat steps 2, 3, and 4 until convergence
-

143

144

145 Finally, after we obtain the estimated $\hat{\theta}_0$ and $\hat{\theta}_1$, we use Gibbs sampling based on the
 146 conditional distribution $P(S_i | S_{N_i}, \hat{\theta}_0, \hat{\theta}_1)$. This method has been proved to be valid for multiple
 147 testing under dependence in a compound decision theoretic framework [37, 38]. Then, we can

148 estimate the marginal posterior probability $q_i = P(S_i = -1|Y)$. Let $q_{(i)}$ be the sorted values of q_i
149 in descending order. For each gene i , the null hypothesis and alternative hypothesis are

150 H_{i0} : Gene i is not associated with the trait of interest

151 H_{i1} : Gene i is associated with the trait of interest

152 As shown by Jiang and Yu [39], the relationship between global FDR and local FDR (lfd_r) is
153 $FDR = E(lfd_r|Y \in \mathcal{R})$, where the rejection region \mathcal{R} is the set of Y such that the null hypothesis
154 can be rejected based on a specific rejection criterion. To control the expected global FDR less
155 than α , we propose the following procedure: let $m = \max \{s: \frac{1}{s} \sum_{i=1}^s q_{(i)}\}$, we reject all the null
156 hypotheses corresponding to $H_{(1)}, \dots, H_{(m)}$.

157

158 **Results**

159 **Simulation Studies**

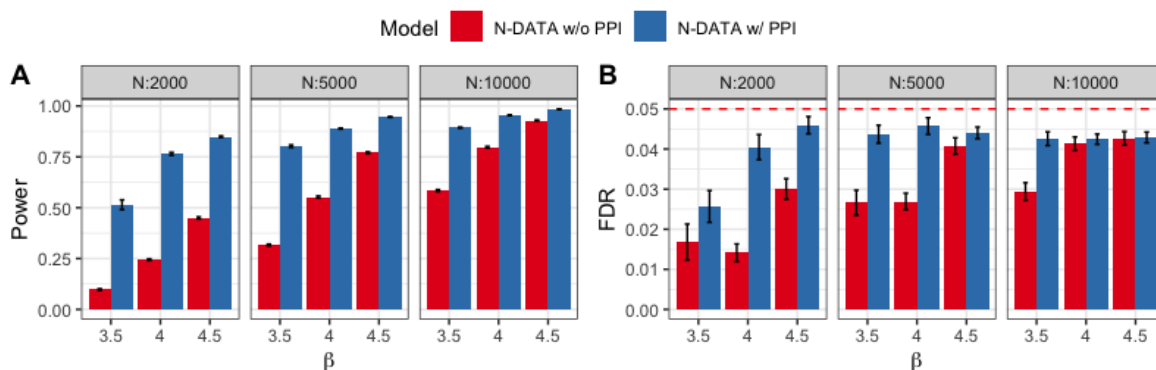
160 For DNV and network data, we considered similar settings as our real data. We randomly
161 selected 2,000 genes, retrieved their mutability from the real data, and extracted the
162 corresponding PPI network formed by these 2,000 genes. First, we simulated the true model
163 with parameter values similar as those from real data analysis to see whether N-DATA can
164 control FDR. Then, we evaluated the power under various settings of sample sizes N and
165 relative risk parameter γ .

166

167 We set true network parameter θ_0 as (-4, 0.2, 0) to make the total number of risk genes in the
168 network of 2,000 random genes to be around 100. We varied the sample size N at 2,000, 5,000
169 and 10,000 to evaluate the performance of N-DATA in small, medium, and large WES cohorts,
170 respectively. In addition, we varied β (log relative risk parameter γ) at 3.5, 4 and 4.5 to
171 investigate the performance of N-DATA around the burden estimated results from real data (In
172 real data analysis, $\hat{\beta} = 3.60$). Each simulation setting was replicated 100 times. For Gibbs

173 sampling-based inference, we used 2,000 MCMC iterations, and set the first 1,000 iterations as
174 burn-ins. These numbers were chosen empirically based on the diagnostic plots for
175 convergence. We report the performance under FDR threshold 0.05 in the main text (Fig 2) and
176 FDR threshold 0.01 and 0.1 in Additional File 1.

177
178 First, we compared the performance of N-DATA model with and without the PPI network as
179 input. For N-DATA model without the PPI network, we assigned the weight of gene for
180 inference. Both models controlled FDR well under all the settings. N-DATA model with PPI
181 network had much better power than the model without PPI network when the sample size and
182 were both low. When the sample size and were larger, the power of N-DATA without PPI
183 network improved as expected.

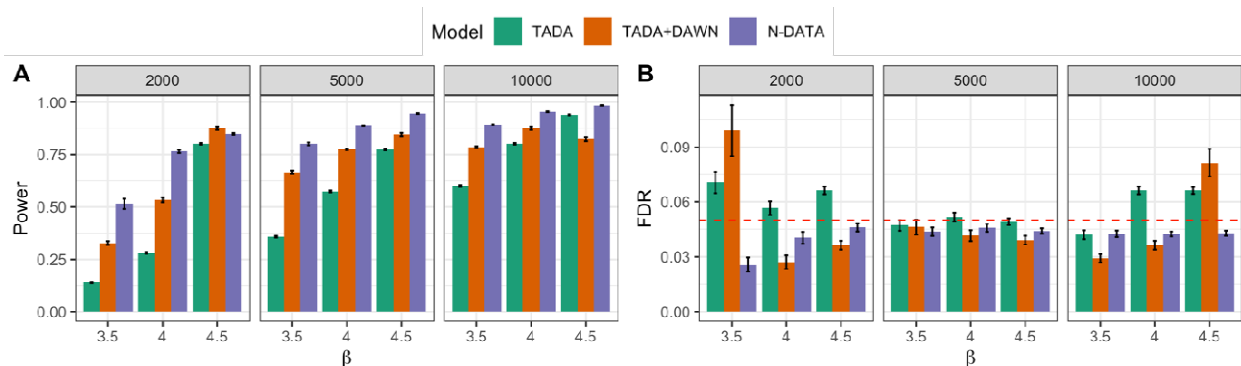


184
185 **Fig 2. Comparison of N-DATA w/o and w/ PPI network models.** (A) Power comparison of the two models. Three
186 panels from left to right represent cohorts with small, medium, and large sample sizes, respectively. (B) FDR control
187 for the two models. Red dash line represents the preset FDR threshold 0.05. Three panels from left to right represent
188 cohorts with small, medium, and large sample sizes, respectively.

189
190 Then, we compared the power of TADA-*De novo*, DAWN, and N-DATA using the same
191 parameter settings. Hyperprior of TADA-*De novo* was estimated from the function *denovo.MOM*
192 based on the recommendation from the authors [40]. Power of TADA was calculated based on
193 TADA q-values. DAWN v1.0 was downloaded from

194 http://www.compgen.pitt.edu/DAWN/DAWN_homepage.htm. We adapted the code of DAWN to
 195 use adjacency matrices of networks as its input. We used TADA-*De novo* p-values and PPI
 196 network as the input of DAWN. The parameter trim threshold was set to 4 based on the
 197 observation in our real data that, which corresponds to select the top 5 genes in TADA-*De novo*
 198 as fixed risk genes.

199
 200 We compared the performance of TADA-*De novo*, TADA-*De novo* p-values + DAWN and N-
 201 DATA under different simulation settings. We reported the performance under FDR threshold
 202 0.05 in the main text (Fig 3) and FDR threshold 0.01 and 0.1 in Additional File 1. We first
 203 checked if all three methods could control the global FDR when the threshold is 0.05. As
 204 discussed above, N-DATA controlled FDR well under all settings, while there was slight FDR
 205 inflation of TADA-*De novo* and DAWN under a couple of settings, especially when both α and β
 206 were very small or very large. N-DATA had the best power among the three methods under all
 207 settings. With an increase of α and β , the performance of TADA-*De novo* and DAWN also
 208 became better. For settings with FDR well controlled, DAWN performed better than TADA-*De*
 209 *novo*.



210
 211 **Fig 3. Comparison of TADA-*De novo*, TADA-*De novo* p-values + DAWN and N-DATA.** Error bars represent
 212 standard errors estimated from 100 replications of simulation. (A) Power comparison of the three methods. Three
 213 panels from left to right represent cohorts with small, medium, and large sample sizes, respectively. (B) FDR control

214 for the three methods. Red dash line represents the preset FDR threshold 0.05. Three panels from left to right
215 represent cohorts with small, medium, and large sample sizes, respectively.

216

217 **Real Data Applications**

218 We applied N-DATA to DNV data from 2,645 CHD trios reported in Jin et al [5]. We only
219 considered damaging variants (loss of function (LoF) and deleterious missense (Dmis) variants
220 by the MetaSVM algorithm) in our analysis as the number of non-deleterious variants is not
221 expected to provide information to differentiate cases from controls biologically [41].

222

223 For network information, we first downloaded STRING v11.0 with medium edge likelihood via
224 interface from STRINGdb package in R, and call this original network from STRING \mathcal{G}_0 . We
225 obtained the curated list of known human CHD genes from Jin et al [5] and expanded the gene
226 list by including additional candidate genes (FDR<0.1) from the single-trait analysis in our
227 previous work [32]. Then, we extracted the subnetwork including the aforementioned gene list
228 and the direct neighbors with likelihood score larger than 950 of those genes, and call this
229 subnetwork \mathcal{G}_1 . We only kept overlapping genes with our DNV data in \mathcal{G}_1 and called the final
230 network used in our real application as \mathcal{G}_2 . There were in total 1,818 genes and 21,534 edges in
231 \mathcal{G}_2 .

232

233 To show that our method can leverage network information to boost risk gene identification, we
234 applied our algorithm without using the network as an input. When there was no prior
235 information from the network, we identified 18 significant genes with FDR<0.05. To include the
236 network information from \mathcal{G}_2 , we denote the degree of gene i in network \mathcal{G}_2 as d_i , and let the
237 weight in the prior as $w_i = \sqrt{d_i}$ following Chen et al. [19]. After adding the network information
238 from \mathcal{G}_2 , we identified 46 genes with at least 1 DNV, and 26 genes harboring at least 2 DNVs
239 with FDR<0.05 in the CHD cohort.

240

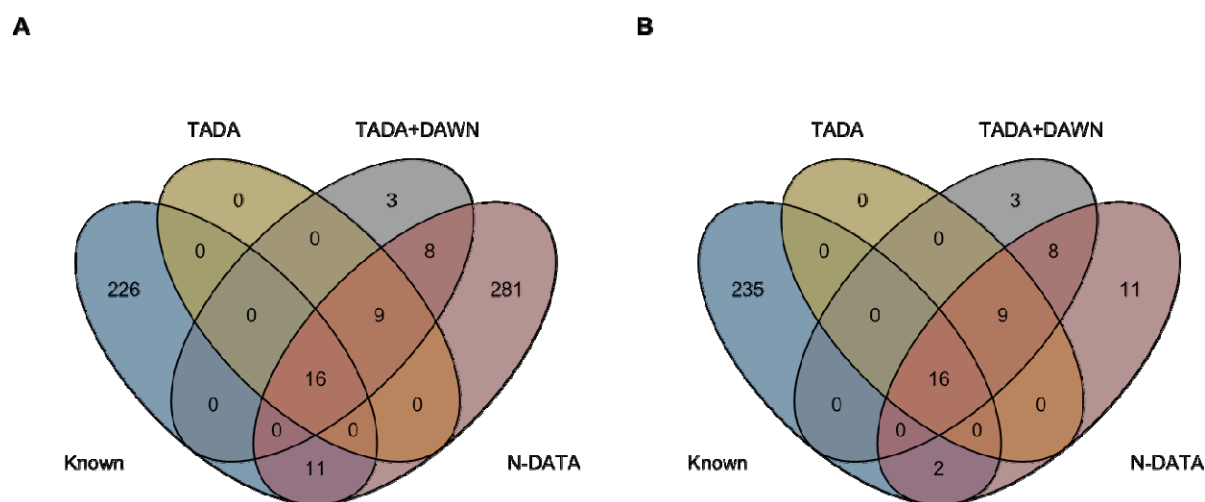
241 We also compared our results with TADA-*De novo* test [40] and DAWN [21, 42]. We
 242 downloaded the software of DAWN and adapted its code by substituting the adjacency matrix
 243 inferred from its Partial Neighborhood Selection algorithm to the adjacency matrix from network
 244 \mathcal{G}_2 . The parameter settings were discussed in the simulation section. Specifically, we also
 245 compared two approaches to control the global FDR for TADA-*De novo* test (FDR q-values and
 246 p-values with FDR adjustment). FDR q-values given by TADA-*De novo* identified 25 significant
 247 genes, and p-values with FDR adjustment identified the same 25 genes. After integrating the p-
 248 values with the \mathcal{G}_2 network, 36 genes were identified by the adapted DAWN method. In total, N-
 249 DATA identified 325 genes with FDR<0.05. As some of the genes may be prioritized due to high
 250 prior probability, but did not have DNV count in the study cohort, we further filtered out genes
 251 without DNV and considered the 46 genes identified with FDR<0.05 and at least 1 DNV as the
 252 candidate genes. (Table 1)

Method	Criteria	Number of Identified Genes
TADA- <i>De novo</i> q-values	FDR<0.05	25
TADA- <i>De novo</i> p-values	FDR<0.05	25
DAWN (Network \mathcal{G}_2 + TADA- <i>De novo</i> p-values)	FDR<0.05	36
N-DATA (Network \mathcal{G}_2 network + DNV counts)	FDR<0.05 ▪ DNVs \geq 1	325 46

253 **Table 1.** Comparison of N-DATA with other methods

254

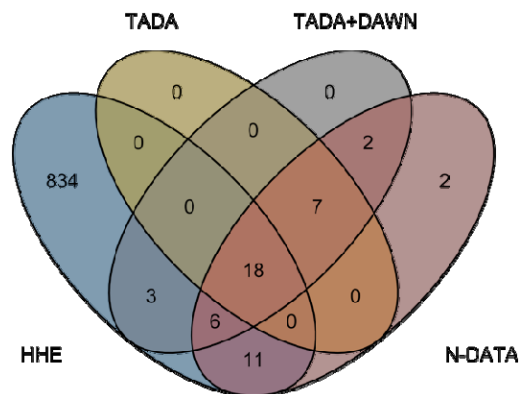
255 We visualized the overlap of 253 known human CHD genes [5], genes that were identified by
256 TADA-*De novo* q-values, DAWN, and N-DATA in Fig 4. Fig 4A shows the 325 genes identified
257 by N-DATA, while Fig 4B shows the 46 genes with at least 1 DNV. From Fig 4B, N-DATA found
258 all genes that can be identified by TADA, and identified 11 different genes compared with
259 DAWN.
260



261
262 **Fig 4. Venn diagram of known human CHD genes, TADA genes, DAWN genes and N-DATA genes. (A)**
263 Overlapping genes between known human CHD genes, TADA- *De novo*, TADA-*De novo* +DAWN and all 325 genes
264 identified by N-DATA. (B) Overlapping genes between known human CHD genes, TADA- *De novo*, TADA-*De novo*
265 +DAWN and 46 candidate genes identified by N-DATA.

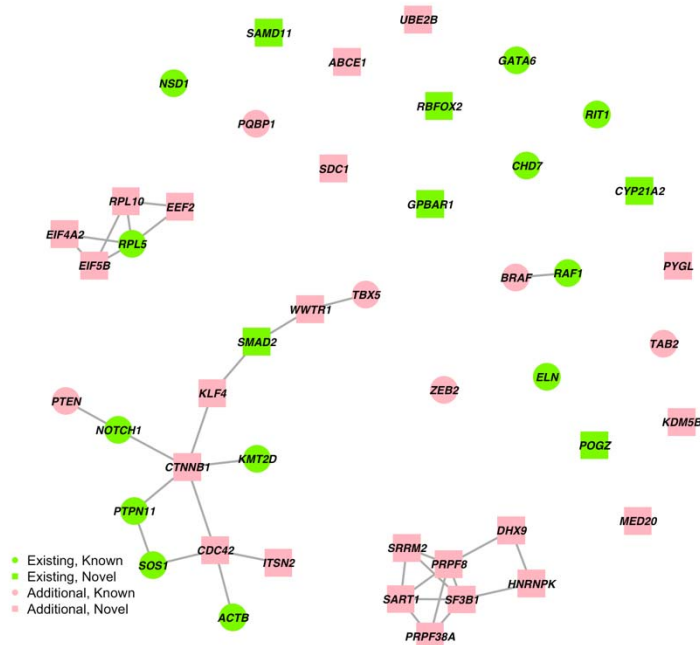
266
267 Among the 46 genes, 18 are known human CHD genes and 35 are in the top 25% in mouse
268 developing heart (HHE) at E14.5 [12]. Further, we calculated the overlap of the 46 N-DATA
269 candidate genes, TADA genes, DAWN genes and 872 HHE genes that were analyzed in the
270 1,818 gene network. (Fig 5).

271



272
 273 **Fig 5. Venn diagram of HHE genes, TADA genes, DAWN genes and the 46 N-DATA candidate genes.** N-DATA
 274 candidate genes identified 11 additional HHE genes compared with other methods.

275
 276 We also visualized these 46 genes in the network to demonstrate that the PPI network
 277 information can help boost statistical power and provide biological interpretation.



278
 279 **Fig 6. N-DATA model identified 46 candidate genes with at least 1 DNV.** Green labels indicate the 18 genes
 280 identified when no network information was provided for N-DATA, and red labels indicate the additional 28 genes

281 identified when the G_2 network was integrated. Circles indicate the 18 known human CHD genes, and squares
282 indicate the 28 novel genes identified by N-DATA.

283

284 Among the 46 candidate genes, *PTPN11*, *RAF1* and *RIT1* had 2 recurrent DNVs, and *CHD7*,
285 *NOTCH1*, *NSD1* and *PYGL* also had recessive genotypes in the CHD cohort [5]. The 46
286 candidate genes form 4 clusters (Fig 6) in the G_2 network. The biggest cluster includes seven
287 known CHD genes *TBX5*, *KMT2D*, *PTPN11*, *SOS1*, *ACTB*, *NOTCH1*, and *PTEN*, which are
288 involved in transcriptional regulation and the early cell growth or differentiation processes. The
289 six new genes *SMAD2*, *KLF4*, *CTNNB1*, *CDC42*, *ITSN2*, and *WWTR1* also function in similar
290 pathways and have varied implications in cardiac development. For instance, *KLF4* and
291 *CTNNB1* have been implicated in cardiac cell differentiation [43]. *Cdc42* cardiomyocyte knock-
292 out mice presented heart defects such as ventricular septum defects and thin ventricular walls
293 [44]. *WWTR1* encodes a transcription regulator, which serves as an effector of Hippo pathway
294 and regulates cardiac wall maturation in zebrafish [45].

295

296 The second biggest cluster is constituted of 7 new genes, all of which are involved in mRNA
297 splicing. Specifically, *SART1*, *SRRM2*, *PRPF38A*, *PRPF8*, and *SF3B1* are associated factors or
298 components of spliceosome; *HNRNPK* encodes a pre-mRNA-binding protein; *DHX9* encodes
299 an RNA helicase which promotes R-loop formation while RNA splicing is perturbed [46].

300 Alternative splicing plays an essential role in heart development, homeostasis, and disease
301 pathogenesis. Mouse knockouts of multiple splice factors had impaired cardiogenesis [47].

302 *SF3B1*, specifically, has been shown to upregulate to induce heart disease in both human and
303 mice [48]. Thus, though not fully investigated, DNVs in those mRNA splicing-related genes may
304 contribute to CHD pathogenesis.

305

306 The third cluster contains genes involved in protein synthesis, includes the known gene *RPL5*
307 and genes not previously associated with CHD (*EIF4*, *EIF5*, *EEF2*, and *RPL10*). *RPL5* and
308 *RPL10* encode the ribosome subunits. Mutations in *RPL5* and other ribosomal genes can lead
309 to multiple congenital anomalies, including CHD [49]. *EIF4* and *EIF5* encode translation initiation
310 factors while *EEF2* encodes the elongation factor that regulate peptide chain elongation during
311 protein synthesis. A recent study reported that the deficiency in ribosome associated NatA
312 complex reduces ribosomal protein and subsequently impact cell development as a mechanism
313 to cause CHD [50]. Thus, DNVs in the above genes may lead to CHD via impairment of protein
314 synthesis.

315
316 The last cluster contains the known CHD genes *BRAF* and *RAF1*, both of which encode key
317 kinases in Ras signaling and are related to Noonan syndrome with CHD as a common feature.

318
319 Among the un-clustered genes, six are identified after using the network information: *ABCE1*,
320 *UBE2B*, *SDC1*, *PYGL*, *KDM5B*, *MED20*. *UBE2B* and *KDM5B*, encoding epigenetic modifiers,
321 have shown suggestive evidence in cardiac development or CHD [51] [52] and might be
322 potential CHD genes.

323

324 **Discussion**

325 In this article, we have introduced a Bayesian framework to integrate PPI network information as
326 the prior knowledge into DNV analysis for CHD. This approach adopts MRF to model the
327 interactions among genes. We apply an empirical Bayes strategy to estimate parameters in the
328 model and conduct statistical inference based on the posterior distribution sampled from a
329 Gibbs sampler. The simulation studies and real data analysis on CHD suggest that the
330 proposed method has improved power to identify risk genes over methods without integrating
331 network information.

332

333 Our proposed framework is innovative for the following aspects. First, it can directly infer
334 disease-associated genes using a network-based model without relying on summary statistics
335 from other DNV association software. Second, it does not need to estimate hyperprior based on
336 other sources compared to the existing pathway-based test for DNV data [17, 31]. Third, it does
337 not require external expression data for the DNV cohort and uses the publicly available PPI
338 database instead, which makes it more applicable to different diseases. This method will not
339 only increase power in risk gene identification, but will also assist in biological interpretation by
340 visualizing clusters of risk genes with functional relevance in the network.

341

342 However, there are some limitations in the current N-DATA model. In real application, it is
343 important to conduct an initial analysis on the enrichment of top genes identified from *de novo*
344 association test in the network like our motivating example. Another limitation is that the
345 likelihood-based inference may suffer from local maxima [19]. Thus, we recommend to initiate
346 the labels of genes from a known risk gene set or run with multiple starts. Also, we observe the
347 Gibbs sampler tends to move around local maxima for some time before convergence. We
348 suggest running at least 2,000 times of iterations and discard the first 1,000 iterations as burn-
349 ins. In addition, we only considered damaging DNVs and assumed the relative risk parameter γ
350 is the same across all genes in N-DATA, which may cause our model to lose power if it varies
351 across variants with different functions (e.g., LoF and Dmis). Future studies may explore adding
352 functional annotation of variants as a layer in the model to further improve statistical power.

353

354 **Conclusions**

355 The topologic information in a pathway may be informative to identify functionally interrelated
356 genes and help improve statistical power in DNV studies. Under the hypothesis that connected
357 genes in PPI networks are more likely to share similar disease association status, we developed

358 a novel statistical model that can leverage information from publicly available PPI databases.
359 Through simulation studies under multiple settings, we proved our method can increase
360 statistical power in identifying additional risk genes compared to methods without using the PPI
361 network information. We then applied our method to the CHD DNV data, and then visualized the
362 subnetwork of candidate genes to find potential functional gene clusters for CHD. Our results
363 may shed new insight on the shared protein functionality among risk genes for CHD.

364

365 **List of abbreviations**

366 DNV: *de novo* variant

367 CHD: congenital heart disease

368 WES: Whole Exome Sequencing

369 PPI: protein-protein Interaction

370 FDR: false discovery rate

371 MRF: Markov Random Field

372 GWAS: Genome-Wide Association Studies

373

374 **References**

- 375 1. Zaidi S, Choi M, Wakimoto H, Ma L, Jiang J, Overton JD, Romano-Adesman A, Bjornson
376 RD, Breitbart RE, Brown KK, et al: **De novo mutations in histone-modifying genes in**
377 **congenital heart disease.** *Nature* 2013, **498**:220-223.
- 378 2. Postma AV, Bezzina CR, Christoffels VM: **Genetics of congenital heart disease: the**
379 **contribution of the noncoding regulatory genome.** *Journal of Human Genetics* 2016,
380 **61**:13-19.
- 381 3. Sevim Bayrak C, Zhang P, Tristani-Firouzi M, Gelb BD, Itan Y: **De novo variants in exomes**
382 **of congenital heart disease patients identify risk genes and pathways.** *Genome Med*
383 2020, **12**:9.
- 384 4. Diab NS, Barish S, Dong W, Zhao S, Allington G, Yu X, Kahle KT, Brueckner M, Jin SC:
385 **Molecular Genetics and Complex Inheritance of Congenital Heart Disease.** *Genes*
386 (*Basel*) 2021, **12**.

- 387 5. Jin SC, Homsy J, Zaidi S, Lu Q, Morton S, DePalma SR, Zeng X, Qi H, Chang W, Sierant MC,
388 et al: **Contribution of rare inherited and de novo variants in 2,871 congenital heart**
389 **disease probands.** *Nat Genet* 2017, **49**:1593-1601.
- 390 6. Zaidi S, Brueckner M: **Genetics and Genomics of Congenital Heart Disease.** *Circ Res*
391 2017, **120**:923-940.
- 392 7. Glessner JT, Bick AG, Ito K, Homsy J, Rodriguez-Murillo L, Fromer M, Mazaika E,
393 Vardarajan B, Italia M, Leipzig J, et al: **Increased frequency of de novo copy number**
394 **variants in congenital heart disease by integrative analysis of single nucleotide**
395 **polymorphism array and exome sequence data.** *Circ Res* 2014, **115**:884-896.
- 396 8. Soemedi R, Wilson IJ, Bentham J, Darlay R, Töpf A, Zelenika D, Cosgrove C, Setchfield K,
397 Thornborough C, Granados-Riveron J, et al: **Contribution of global rare copy-number**
398 **variants to the risk of sporadic congenital heart disease.** *Am J Hum Genet* 2012,
399 **91**:489-501.
- 400 9. Pierpont ME, Brueckner M, Chung WK, Garg V, Lacro RV, McGuire AL, Mital S, Priest JR,
401 Pu WT, Roberts A, et al: **Genetic Basis for Congenital Heart Disease: Revisited: A**
402 **Scientific Statement From the American Heart Association.** *Circulation* 2018, **138**:e653-
403 e711.
- 404 10. Teer JK, Mullikin JC: **Exome sequencing: the sweet spot before whole genomes.** *Human*
405 *Molecular Genetics* 2010, **19**:R145-R151.
- 406 11. Rabbani B, Tekin M, Mahdieh N: **The promise of whole-exome sequencing in medical**
407 **genetics.** *Journal of Human Genetics* 2014, **59**:5-15.
- 408 12. Homsy J, Zaidi S, Shen Y, Ware JS, Samocha KE, Karczewski KJ, DePalma SR, McKean D,
409 Wakimoto H, Gorham J, et al: **De novo mutations in congenital heart disease with**
410 **neurodevelopmental and other congenital anomalies.** *Science* 2015, **350**:1262-1266.
- 411 13. Richter F, Morton SU, Kim SW, Kitaygorodsky A, Wasson LK, Chen KM, Zhou J, Qi H, Patel
412 N, DePalma SR: **Genomic analyses implicate noncoding de novo variants in congenital**
413 **heart disease.** *Nature genetics* 2020, **52**:769-777.
- 414 14. Watkins WS, Hernandez EJ, Wesolowski S, Bisgrove BW, Sunderland RT, Lin E, Lemmon
415 G, Demarest BL, Miller TA, Bernstein D: **De novo and recessive forms of congenital**
416 **heart disease have distinct genetic and phenotypic landscapes.** *Nature*
417 *communications* 2019, **10**:1-12.
- 418 15. Sifrim A, Hitz M-P, Wilsdon A, Breckpot J, Turki SHA, Thienpont B, McRae J, Fitzgerald
419 TW, Singh T, Swaminathan GJ, et al: **Distinct genetic architectures for syndromic and**
420 **nonsyndromic congenital heart defects identified by exome sequencing.** *Nature*
421 *Genetics* 2016, **48**:1060-1065.
- 422 16. Sifrim A, Hitz MP, Wilsdon A, Breckpot J, Turki SH, Thienpont B, McRae J, Fitzgerald TW,
423 Singh T, Swaminathan GJ, et al: **Distinct genetic architectures for syndromic and**
424 **nonsyndromic congenital heart defects identified by exome sequencing.** *Nat Genet*
425 2016, **48**:1060-1065.
- 426 17. Nguyen TH, He X, Brown RC, Webb BT, Kendler KS, Vladimirov VI, Riley BP, Bacanu SA:
427 **DECO: a framework for jointly analyzing de novo and rare case/control variants, and**
428 **biological pathways.** *Brief Bioinform* 2021.

- 429 18. Lee I, Blom UM, Wang PI, Shim JE, Marcotte EM: **Prioritizing candidate disease genes by**
430 **network-based boosting of genome-wide association data.** *Genome Res* 2011, **21**:1109-
431 1121.
- 432 19. Chen M, Cho J, Zhao H: **Incorporating biological pathways via a Markov random field**
433 **model in genome-wide association studies.** *PLoS Genet* 2011, **7**:e1001353.
- 434 20. Hou L, Chen M, Zhang CK, Cho J, Zhao H: **Guilt by rewiring: gene prioritization through**
435 **network rewiring in genome wide association studies.** *Hum Mol Genet* 2014, **23**:2780-
436 2790.
- 437 21. Liu L, Lei J, Roeder K: **Network assisted analysis to reveal the genetic basis of autism.**
438 *The Annals of Applied Statistics* 2015, **9**:1571-1600, 1530.
- 439 22. Nguyen HT, Bryois J, Kim A, Dobbyn A, Huckins LM, Munoz-Manchado AB, Ruderfer DM,
440 Genovese G, Fromer M, Xu X, et al: **Integrated Bayesian analysis of rare exonic variants**
441 **to identify risk genes for schizophrenia and neurodevelopmental disorders.** *Genome*
442 *Med* 2017, **9**:114.
- 443 23. Nguyen T-H, Dobbyn A, Brown RC, Riley BP, Buxbaum JD, Pinto D, Purcell SM, Sullivan PF,
444 He X, Stahl EA: **mTADA is a framework for identifying risk genes from de novo**
445 **mutations in multiple traits.** *Nature Communications* 2020, **11**:2929.
- 446 24. Oughtred R, Stark C, Breitkreutz BJ, Rust J, Boucher L, Chang C, Kolas N, O'Donnell L,
447 Leung G, McAdam R, et al: **The BioGRID interaction database: 2019 update.** *Nucleic*
448 *Acids Res* 2019, **47**:D529-d541.
- 449 25. Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, Campbell NH,
450 Chavali G, Chen C, del-Toro N, et al: **The MIntAct project--IntAct as a common curation**
451 **platform for 11 molecular interaction databases.** *Nucleic Acids Res* 2014, **42**:D358-363.
- 452 26. Salwinski L, Miller CS, Smith AJ, Pettit FK, Bowie JU, Eisenberg D: **The Database of**
453 **Interacting Proteins: 2004 update.** *Nucleic Acids Res* 2004, **32**:D449-451.
- 454 27. Licata L, Briganti L, Peluso D, Perfetto L, Iannuccelli M, Galeota E, Sacco F, Palma A,
455 Nardoza AP, Santonico E, et al: **MINT, the molecular interaction database: 2012**
456 **update.** *Nucleic Acids Res* 2012, **40**:D857-861.
- 457 28. Mishra GR, Suresh M, Kumaran K, Kannabiran N, Suresh S, Bala P, Shivakumar K,
458 Anuradha N, Reddy R, Raghavan TM, et al: **Human protein reference database--2006**
459 **update.** *Nucleic Acids Res* 2006, **34**:D411-414.
- 460 29. Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-Cepas J, Simonovic M,
461 Doncheva NT, Morris JH, Bork P, et al: **STRING v11: protein-protein association**
462 **networks with increased coverage, supporting functional discovery in genome-wide**
463 **experimental datasets.** *Nucleic Acids Res* 2019, **47**:D607-d613.
- 464 30. Nguyen TH, Dobbyn A, Brown RC, Riley BP, Buxbaum JD, Pinto D, Purcell SM, Sullivan PF,
465 He X, Stahl EA: **mTADA is a framework for identifying risk genes from de novo**
466 **mutations in multiple traits.** *Nat Commun* 2020, **11**:2929.
- 467 31. Nguyen HT, Dobbyn A, Charney AW, Bryois J, Kim A, Mcfadden W, Skene NG, Huckins
468 LM, Wang W, Ruderfer DM, et al: **Integrative analysis of rare variants and pathway**
469 **information shows convergent results between immune pathways, drug targets and**
470 **epilepsy genes.** *bioRxiv* 2018:410100.
- 471 32. Xie Y, Li M, Dong W, Jiang W, Zhao H: **M-DATA: A statistical approach to jointly**
472 **analyzing de novo mutations for multiple traits.** *PLoS Genet* 2021, **17**:e1009849.

- 473 33. Kindermann R: **Markov random fields and their applications**. *American mathematical*
474 *society* 1980.
- 475 34. Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, Kosmicki JA,
476 Rehnröm K, Mallick S, Kirby A, et al: **A framework for the interpretation of de novo**
477 **mutation in human disease**. *Nat Genet* 2014, **46**:944-950.
- 478 35. Besag J: **On the statistical analysis of dirty pictures**. *Journal of the Royal Statistical*
479 *Society: Series B (Methodological)* 1986, **48**:259-279.
- 480 36. Le Cessie S, Van Houwelingen JC: **Ridge estimators in logistic regression**. *Journal of the*
481 *Royal Statistical Society: Series C (Applied Statistics)* 1992, **41**:191-201.
- 482 37. Sun W, Tony Cai T: **Large - scale multiple testing under dependence**. *Journal of the*
483 *Royal Statistical Society: Series B (Statistical Methodology)* 2009, **71**:393-424.
- 484 38. Li H, Wei Z, Maris J: **A hidden Markov random field model for genome-wide association**
485 **studies**. *Biostatistics* 2010, **11**:139-150.
- 486 39. Jiang W, Yu W: **Controlling the joint local false discovery rate is more powerful than**
487 **meta-analysis methods in joint analysis of summary statistics from multiple genome-**
488 **wide association studies**. *Bioinformatics* 2016, **33**:500-507.
- 489 40. He X, Sanders SJ, Liu L, De Rubeis S, Lim ET, Sutcliffe JS, Schellenberg GD, Gibbs RA, Daly
490 MJ, Buxbaum JD, et al: **Integrated model of de novo and inherited genetic variants**
491 **yields greater power to identify risk genes**. *PLoS Genet* 2013, **9**:e1003671.
- 492 41. Li M: **Gene-based Association Analysis for Genome-wide Association and Whole-**
493 **exome Sequencing Studies**. Yale University, Biostatistics; 2020.
- 494 42. Liu L, Lei J, Sanders SJ, Willsey AJ, Kou Y, Cicek AE, Klei L, Lu C, He X, Li M, et al: **DAWN: a**
495 **framework to identify autism genes and subnetworks using gene expression and**
496 **genetics**. *Mol Autism* 2014, **5**:22.
- 497 43. Kami D, Kitani T, Kawasaki T, Gojo S: **Cardiac mesenchymal progenitors differentiate**
498 **into adipocytes via Klf4 and c-Myc**. *Cell death & disease* 2016, **7**:e2190-e2190.
- 499 44. Liu Y, Wang J, Li J, Wang R, Tharakan B, Zhang SL, Tong CW, Peng X: **Deletion of Cdc42 in**
500 **embryonic cardiomyocytes results in right ventricle hypoplasia**. *Clinical and*
501 *translational medicine* 2017, **6**:40-40.
- 502 45. Lai JKH, Collins MM, Uribe V, Jiménez-Amilburu V, Günther S, Maischein HM, Stainier
503 DYR: **The Hippo pathway effector Wwtr1 regulates cardiac wall maturation in**
504 **zebrafish**. *Development* 2018, **145**.
- 505 46. Chakraborty P, Huang JTJ, Hiom K: **DHX9 helicase promotes R-loop formation in cells**
506 **with impaired RNA splicing**. *Nat Commun* 2018, **9**:4346.
- 507 47. Zahr HC, Jaalouk DE: **Exploring the Crosstalk Between LMNA and Splicing Machinery**
508 **Gene Mutations in Dilated Cardiomyopathy**. *Front Genet* 2018, **9**:231.
- 509 48. van den Hoogenhof MM, Pinto YM, Creemers EE: **RNA Splicing: Regulation and**
510 **Dysregulation in the Heart**. *Circ Res* 2016, **118**:454-468.
- 511 49. Gazda HT, Sheen MR, Vlachos A, Choessel V, O'Donohue MF, Schneider H, Darras N,
512 Hasman C, Sieff CA, Newburger PE, et al: **Ribosomal protein L5 and L11 mutations are**
513 **associated with cleft palate and abnormal thumbs in Diamond-Blackfan anemia**
514 **patients**. *Am J Hum Genet* 2008, **83**:769-780.

- 515 50. Ward T, Tai W, Morton S, Impens F, Van Damme P, Van Haver D, Timmerman E,
516 Venturini G, Zhang K, Jang MY, et al: **Mechanisms of Congenital Heart Disease Caused**
517 **by NAA15 Haploinsufficiency**. *Circ Res* 2021, **128**:1156-1169.
- 518 51. Robson A, Makova SZ, Barish S, Zaidi S, Mehta S, Drozd J, Jin SC, Gelb BD, Seidman CE,
519 Chung WK, et al: **Histone H2B monoubiquitination regulates heart development via**
520 **epigenetic control of cilia motility**. *Proc Natl Acad Sci U S A* 2019, **116**:14049-14054.
- 521 52. Audain E, Wilsdon A, Breckpot J, Izarzugaza JMG, Fitzgerald TW, Kahlert AK, Sifrim A,
522 Wünnemann F, Perez-Riverol Y, Abdul-Khaliq H, et al: **Integrative analysis of genomic**
523 **variants reveals new associations of candidate haploinsufficient genes with congenital**
524 **heart disease**. *PLoS Genet* 2021, **17**:e1009679.
525

526 **Acknowledgements**

527 We thank Jin et al. [5] for sharing the *de novo* variant data of CHD. We thank Andrew Xu and Dr.
528 Min Chen for discussions on coding, and Ziyu Jiang for discussions on diagnostic for Bayesian
529 inference.

530

531 **Funding**

532 This work was supported in part by NIH grant R03HD100883-01A1 (Y.X. and H.Z.) and
533 R01GM134005-01A1 (W.J., H.L., and H.Z.). The funders had no role in study design, data
534 collection and analysis, decision to publish, or preparation of the manuscript.

535

536 **Ethics Declaration**

537 **Ethics approval and consent to participate**

538 This manuscript is a mainly methodology paper. The data used in this manuscript is all public
539 available. The research has no procedure related to collecting human subject data.

540 **Consent for publication**

541 Not applicable.

542 **Competing interests**

543 The authors declare that they have no competing interests.

544 **Supplementary Information**

545 Additional File 1. Supplementary figures 1-4

546