# Using network clustering to predict copy number variations associated with health disparities

Substantial health disparities exist between African Americans and Caucasians in the United States. Copy number variations (CNVs) are one form of human genetic variations that have been linked with complex diseases and often occur at different frequencies among African Americans and Caucasian populations. Here, we aimed to investigate whether CNVs with differential frequencies can contribute to health disparities from the perspective of gene networks. We inferred network clusters from human gene/protein networks based on two different data sources. We then evaluated each network cluster for the occurrences of known pathogenic genes and genes located in CNVs with different population frequencies, and used false discovery rates to rank network clusters. This approach let us identify five clusters enriched with known pathogenic genes and with genes located in CNVs with different frequencies between African Americans and Caucasians. These clustering patterns predict two candidate causal genes located in four population-specific CNVs that play potential roles in health disparities

1 **Using network clustering to predict copy number variations**

2 **associated with health disparities.**

3 Yi Jiang[*]

4 Department of Computer Science and Engineering, University of Tennessee at

5 Chattanooga, TN

6 Hong Qin[*§]

7 Department of Biology, Spelman College, Atlanta, GA

8 Li Yang

9 Department of Computer Science and Engineering, University of Tennessee at

10 Chattanooga, TN

11 * Co-first authors.

12 § Corresponding author. Phone: (404) 270-5757, Fax: (404) 270-5725, Email:

13 hqin@spelman.edu

## Abstract

Substantial health disparities exist between African Americans and Caucasians in the United States. Copy number variations (CNVs) are one form of human genetic variations that have been linked with complex diseases and often occur at different frequencies among African Americans and Caucasian populations. Here, we aimed to investigate whether CNVs with differential frequencies can contribute to health disparities from the perspective of gene networks. We inferred network clusters from human gene/protein networks based on two different data sources. We then evaluated each network cluster for the occurrences of known pathogenic genes and genes located in CNVs with different population frequencies, and used false discovery rates to rank network clusters. This approach let us identify five clusters enriched with known pathogenic genes and with genes located in CNVs with different frequencies between African Americans and Caucasians. These clustering patterns predict two candidate causal genes located in four population-specific CNVs that play potential roles in health disparities.

Keywords:

33 **List of Key Abbreviations:**

34 CNV: Copy number variation

35 SNP: Single nucleotide polymorphism

36 PPIN: Protein-protein interaction network

37 HPRD: Human protein reference database

38 PPI: Protein-protein interaction

39 AA: African American

40 MCL: Markov Cluster Algorithm

41 FDR: false discovery rate

42 GO: Gene ontology

43 OMIM: Online Mendelian Inheritance in Man

44 dbSNP: Single Nucleotide Polymorphism Database

45 SERCA1: Sarco/endoplasmic reticulum $Ca^{2+}$-ATPase 1

## 46 Introduction

47 Health disparities refer to differences in the disease distribution and/or health

48 outcomes across racial and ethnic groups. In United States, health disparities

49 in African Americans are found in life expectancy, death rates, and health

50 measures (National Center for Health Statistics 2013). In addition to social

51 determinants, such as socio-economical status, health care access and

52 cultural practices, human genetic variations play a significant role in health

53 disparities. Genetic variations at different frequencies among populations can

54 lead to differences in disease susceptibility. Studies on genetic variations and

55 disease association are greatly advanced by the completion of the

56 International HapMap Project and new genome sequencing techniques

57 (Ramos & Rotimi 2009).

58 Genome-wide association studies (GWAS) are currently an effective approach

59 to identify diseases-associated genetic variations (Hirschhorn & Daly 2005;

60 Wang et al. 2005). Although GWAS have revealed many disease-associated

61 single nucleotide polymorphisms (SNPs), GWAS are often limited to individual

62 genetic variations and often do not address complex gene interactions.

63 Moreover, associated SNPs are often located in haplotype blocks that contain

64 more than one gene.  To address these limitations, human gene networks

65 have been used to improve GWAS detection of genes associated with

66 complex diseases, such as the comorbidity analysis (Sharma et al. 2013), an

67 improved guilt-by-association method (Baranzini et al. 2009; Lee et al. 2011),

68 and a distance-based scoring method using seeded diseases genes (Liu et al.

69 2012).

70 Copy number variations (CNVs) are duplications or deletions of genomic

71 segments that can contain one or more genes (McCarroll & Altshuler 2007).

72 CNVs have been associated with complex diseases such as autism (Gilman et

73 al. 2011; Glessner et al. 2009). Computational tools and methods have been

74 developed, such as the CNV annotator (Zhao & Zhao 2013) and NETBAG

75 (Gilman et al. 2011), to address the potential roles of CNVs in human

76 diseases. Recently, it is reported that CNVs can occur at different frequencies

77 between African Americans and Caucasians (McElroy et al. 2009), and

78 naturally the question about the potential roles of CNVs in health disparity is

79 raised.

80 Here, we aim to investigate the clustering of pathogenic genes and genes in

81 CNVs with different population frequencies in two human gene/protein

82 networks, in order to better understand health disparities between African

83 Americans and Caucasians. The current human gene/protein networks

84 contain thousands of interacting molecules (Barabasi et al. 2011; Vidal et al.

85 2011). We will partition gene networks into clusters and use these clusters to

86 predict potential diseases associated with population-specific CNVs, based on

87 the rationale that interacting genes often share similar functions (Pizzuti et

88 al. 2012).

89 **Materials and Methods**

90   Our overall work flow is shown in Figure 1. To identify potential diseases

91   associated with CNVs, our basic idea is to identify gene interaction clusters

92   that involve genes in population-specific CNVs. The diseases associated with

93   a CNV-gene's interacting genes are potential diseases associated with this

94   CNV. Specifically, we first obtained two human gene/protein networks and

95   partitioned them into gene clusters. We then performed statistical tests on

96   each cluster to estimate its significances of containing pathogenic genes and

97   genes in population-specific CNVs. Finally, we ranked gene clusters based on

98   false discovery rates (FDRs). High-ranked clusters were enriched both for

99   pathogenic genes and for genes in CNVs with differential frequencies

100  between African-Americans and Caucasians. These clusters were then

101  searched for enriched Gene Ontology (GO) terms and related disease

102  phenotypes.

103  **Network clustering**

104  We obtained two human gene/protein networks, one from Human Protein

105  Reference Database (HPRD) (Mishra et al. 2006; Peri et al. 2003; Prasad et al.

106  2009) and another from MultiNet (Khurana et al. 2013). The HPRD network

107  (referred to as HPRDNet) contains only physical protein-protein interactions

108  (PPIs), whereas MultiNet is a unified network including PPI, phosphorylation,

109  metabolic, signaling, genetic and regulatory networks. The two networks

110  share 8468 genes (89.6% of HPRDNet and 58.6% of MultiNet) but only 8769

111  interactions (23.8% of HPRDNet and 8% of MultiNet). These two networks

112  were both partitioned into gene clusters using the Markov Cluster (MCL)

113　Algorithm (van Dongen 2000). Clustering was done with the inflation

114　parameter I ranging from 1.1 to 2.0 with a step of 0.1. Descriptive statistics

115　of the two networks and their clustering results are summarized in

116　Supporting Table S1.

117　**Mapping of CNVs and SNPs**

118　CNV coordinates were obtained from a CNV map in African Americans and

119　Caucasians (McElroy et al. 2009). There are three types of CNVs in this map:

120　(1) CNVs only occur in African Americans, (2) CNVs only occur in Caucasians,

121　and (3) CNVs occurred in both African Americans and Caucasians. To simplify

122　the analysis, we further partitioned the last type: CNVs that occurred more

123　than 50% in African Americans or in Caucasians were combined with the first

124　and second types of CNVs, respectively. This repartition resulted in two

125　modified CNV sets with differential population frequencies. The coordinates of

126　these CNVs were then searched in the UCSC Genome Database (Karolchik et

127　al. 2014) through its MySQL API to obtain the corresponding gene sets. For

128　simplicity, CNVs that occur more frequently in African Americans were called

129　African-American CNVs or CNV_AA; CNVs that occur more frequently in

130　Caucasians were called Caucasian CNVs or CNV_CA.

131　Disease-associated SNPs were retrieved from a file, OmimVarLocusIdSNP.bcp,

132　from the FTP site of Single Nucleotide Polymorphism Database (dbSNP)

133　(Sherry et al. 2001). Coordinates of these SNPs were then queried against the

134   MySQL API of the UCSC Genome Database to identify genes in which those

135   SNPs are located. This identified gene set was termed as pathogenic genes.

136   Details of gene mapping results are shown in Supporting Table S2.


137   **Cluster Analyses**

138   Clusters were obtained from both HPRDNet and MultiNet using MCL with a

139   range of ten inflation parameters. For each cluster, contingency tables were

140   constructed using the numbers of pathogenic genes and CNVs related genes

141   (Table 1A and 1B). Right-tailed Fisher's exact tests were applied to these

142   contingency tables to calculate enrichment significance of pathogenic genes,

143   and CNV_AA or CNV_CA genes, respectively. Based on obtained $p$-values,

144   false discovery rates (FDRs) were calculated using the Robust FDR Routine

145   (Pounds & Cheng 2006). Fisher's exact tests and Robust FDR Routine were

146   both performed in the R statistical environment (R Development Core Team

147   2013). Ranking were applied to clusters with $p$-value<0.10 and FDR<0.20 in

148   both enrichment tests for pathogenic genes and population-preferred CNVs

149   genes. Assuming both enrichment tests are independent, the FDR values

150   were multiplied to jointly rank the network clusters. The same cluster

151   analysis procedure was applied to clustering results with different MCL

152   inflation parameters.

153  For clarity, we focused our functional analyses on clusters that were

154  consistently ranked at the first place with different MCL inflation parameter

155  values.

156  **Biological Significance Analyses**

157  Biological relevance of selected network clusters were analyzed by GOrilla

158  (Eden et al. 2009) to search for enriched gene ontology (GO) terms. In GOrilla

159  search, genes in the selected clusters were target genes, and all genes in the

160  network were treated as background genes. To investigate the possible links

161  of population-specific CNVs to heath disparities, we first identified

162  significantly enriched GO terms that are associated with CNV_AA or CNV_CA

163  genes. We then focused on the pathogenic genes with the enriched GO

164  terms, and examined their associated disease phenotypes in OMIM database

165  (Online Mendelian Inheritance in Man 2014).

166  **Results and Discussions**

167  **Top-ranked network clusters**

168  We performed cluster analyses with ten MCL inflation parameters values for

169  both HPRDNet and MultiNet (Table S1), and scored the resulted clusters for

170  their potential roles in CNV related health disparities (Table S3). For clarity,

171  we focused on clusters that are consistently top-ranked with different MCL

172  inflation parameters. The graph representations of selected clusters are

173  shown in Figure 2.

174 We found four similar clusters, (AA1, AA2, and AA3 in HPRDNet and AA4 in

175 Multinet), that are enriched both for pathogenic genes and for genes located

176 in African-American CNVs (Table 2). In HPRDNet, cluster AA1, AA2 and AA3

177 together were ranked at first place five times; and cluster AA4 were ranked

178 five times in Multinet (Table S3). Cluster AA1 contains 11 genes, within which

179 eight are pathogenic genes (Figure 2A). Cluster AA2 and AA3 contain one and

180 two more genes than cluster AA1, respectively (Figure S1). In MultiNet,

181 cluster AA4 contains five genes and can be considered as a sub-cluster of

182 cluster AA1, AA2 and AA3 (Figure 2B). In these four clusters, gene *HSPB1* is

183 mainly duplicated in African Americans (Table 2 and Table 3). Since cluster

184 AA1, AA2 and AA3 were selected from the same network and are highly

185 similar to each other, only cluster AA1 and AA4 were studied in biological

186 significance analyses.

187 In both HPRDNet and MultiNet, the same cluster, named as CA1, was

188 identified to be enriched with both pathogenic genes and genes located in

189 Caucasian CNVs (Table 2). Cluster CA1 was ranked at first place four times in

190 HPRDNet and seven times in MultiNet (Table S3). This cluster contains five

191 genes, and four of them are associated with diseases (Figure 2C). Cluster CA1

192 contains gene *ATP2A1* that is duplicated only in Caucasians (Table 3).

193 **Duplication of *HSPB1* and health disparities in African Americans.**

194 Gene *HSPB1* is located in genomic duplication regions occurring more

195 frequently in African Americans (Table 3), and is found in the cluster family of

196 AA1, AA2, AA3, and AA4 (Table 2). For cluster AA1, only one GO molecular

197  function term related to gene *HSPB1* is significantly enriched (Cluster AA1 in

198  Table 4). For cluster AA4, in addition to the same enriched GO molecular

199  functions term, three GO biological process terms and one GO cellular

200  component term are found significantly enriched (Cluster AA4 in Table 4). In

201  the genes with the enriched GO terms, four of them are known to be

202  associated with diseases (Cluster AA1/AA4 in Table 5). Among these four

203  genes, three of them are implicated in health disparities of African

204  Americans. Specifically, gene *CRYAB* is related to dilated cardiomyopathy and

205  myofibrillar myopathy. African Americans were found at higher risk for

206  idiopathic dilated cardiomyopathy compared with Caucasian, and this could

207  not be explained by income, education, alcohol use, smoking, or history of

208  some other diseases (Coughlin et al. 1993). Moreover, gene *CRYAA*, *CRYAB*

209  and *CRYBB2* are all related to various types of cataract. It was reported that

210  age-specific blindness prevalence was higher for African Americans compared

211  with Caucasian, and cataract accounts for 36.8% of all blindness in African

212  American, but for only 8.7% in Caucasian (Congdon et al. 2004).

213  How could *HSPB1* duplication contribute to health disparities? Based on the

214  direct interaction between *HSPB1* and *CRYAB* and the fact that both genes

215  are expressed in Z-disc (Table 4), it is plausible that *HSPB1* may play an

216  unknown role in cardiomyopathy. Alternatively, *HSPB1* might be involved in

217  cataract, because *HSPB1*, *CRYAA* and *CRYAB* interact with each other and all

218  can negatively regulate apoptotic process (Table 4). Studies suggested that

219  lens epithelial cell apoptosis may be a common cellular basis for initiation of

220  non-congenital cataract formation (Li et al. 1995), and inhibition of epithelial

221 cell apoptosis may be one possible mechanism that inhibits cataract

222 development (Nahomi et al. 2013). Our results here argue for further

223 experimental studies to test the possible role of *HSPB1* CNVs in

224 cardiomyopathy or cataract/blindness in African Americans.

225 **Duplication of *ATP2A1* and cardiomyopathy.**

226 Gene *ATP2A1* in cluster CA1 is located in a genomic duplication region that

227 occurs only in Caucasians (Table 3). We found that three genes in cluster CA1

228 are enriched with various GO biological process terms that involve *ATP2A1*

229 (Cluster CA1 in Table 4). All of the three genes are related to diseases when

230 they are mutated (Cluster CA1 in Table 5).

231 How would *ATP2A1* influence health disparities? Among the diseases related

232 to the pathogenic genes in cluster CA1, idiopathic dilated cardiomyopathy

233 occurs less often in Caucasians than in African Americans (Coughlin et al.

234 1993). One possibility is that higher copies of *ATP2A1* may offer some

235 benefits to Caucasians. Studies have shown that increased activity of

236 sarco/endoplasmic reticulum $Ca^{2+}$-ATPase 1 (SERCA1), which is encoded by

237 *ATP2A1*, can partially rescue the heart from ·OH-induced injury (Hiranandani

238 et al. 2006), and protect the heart from ischemia-reperfusion (I/R) injury

239 (Talukder et al. 2007). Another possibility is that higher copies of *ATP2A1* only

240 lead to moderate risk of cardiomyopathy in Caucasians, and this moderate

241 effect is overshadowed by other genetics factors not covered by our CNV

242 dataset.

## Remarks and future directions

244 Although genetic factors play a crucial role in health disparities, only a few

245 association studies have been reported in health disparities in common

246 complex diseases, such as breast cancer (Long et al. 2013), prostate cancer

247 (Bensen et al. 2014; Bensen et al. 2013; Xu et al. 2011), type 2 diabetes (Ng

248 et al. 2014) and vascular diseases (Wei et al. 2011).

249 Our study here is closely related to network-based meta-analyses of GWAS

250 results (Atias et al. 2013; Leiserson et al. 2013). One important aim of

251 network-based meta-analysis of GWAS data is to distinguish the bona fide

252 causal gene from others in the same haplotype block associated with the

253 significant SNP. Likewise, our network approach aims to predict a potential

254 causal gene from a population-specific CNV that can be associated with

255 pathogenic genes.

256 Noticeably, our method does not require network permutations, whereas

257 many existing methods of network/pathway based meta-analyses of GWAS

258 data do. This difference is because we first partitioned the network into

259 clusters and then perform association tests. In comparison, many network

260 based GWAS meta-analysis methods use traversal distances to seed genes to

261 evaluate candidate genes. This kind of traversal distance based method

262 generally prohibits pre-partition of network into clusters and require network

263 permutations for estimation of p-values.  It can be seen that our cluster-

264 based method naturally accommodate multiple candidate genes in the

265 association analysis, whereas traversal distance in a network is by definition

266 often limited to single candidate gene evaluation.

267 In future studies, we plan to improve network clustering results by integrating

268 functional genomics data sets, such as gene expressions, into gene networks

269 to generate weighted interactions.

## Conclusions

271 In this study, gene clusters were inferred from two human gene/protein

272 networks, HPRDNet and MultiNet, by MCL clustering algorithm with different

273 parameters. Each cluster was ranked based on products of FDR values based

274 on the right-tailed Fisher's exact tests for enrichment of pathogenic or CNV-

275 genes. Five clusters were consistently found to be enriched with both

276 pathogenic genes and genes located in African-American or Caucasian CNVs.

277 In cluster AA1, AA2, AA3 and AA4, gene *HSPB1* is duplicated more frequently

278 in African-Americans. In clusters CA1, gene *ATP2A1* is duplicated only in

279 Caucasians. All gene clusters are associated with certain diseases that occur

280 more often in one population than in the other. Although we only studied

281 population-preferred CNVs and did not consider the roles of other genetic

282 factors, our computational studies have generated some interesting

283 hypotheses for further experimental studies to understand health disparities

284 in these diseases.

## Author contributions

286 HQ initiated this study. HQ and LY designed the overall project. YJ

287 implemented the methods and performed data analyses. All authors

288 participated in writing.

289 **Acknowledgements**

294 **Reference:**

295 Atias N, Istrail S, and Sharan R. 2013. Pathway-based analysis of genomic variation
296       data. *Curr Opin Genet Dev* 23:622-626.
297 Barabasi AL, Gulbahce N, and Loscalzo J. 2011. Network medicine: a network-based
298       approach to human disease. *Nat Rev Genet* 12:56-68.
299 Baranzini SE, Galwey NW, Wang J, Khankhanian P, Lindberg R, Pelletier D, Wu W,
300       Uitdehaag BM, Kappos L, Gene MSAC, Polman CH, Matthews PM, Hauser SL,
301       Gibson RA, Oksenberg JR, and Barnes MR. 2009. Pathway and network-based
302       analysis of genome-wide association studies in multiple sclerosis. *Hum Mol*
303       *Genet* 18:2078-2090.
304 Bensen JT, Xu Z, McKeigue PM, Smith GJ, Fontham ET, Mohler JL, and Taylor JA. 2014.
305       Admixture mapping of prostate cancer in African Americans participating in
306       the North Carolina-Louisiana Prostate Cancer Project (PCaP). *Prostate* 74:1-9.
307 Bensen JT, Xu Z, Smith GJ, Mohler JL, Fontham ET, and Taylor JA. 2013. Genetic
308       polymorphism and prostate cancer aggressiveness: a case-only study of
309       1,536 GWAS and candidate SNPs in African-Americans and European-
310       Americans. *Prostate* 73:11-22.
311 Congdon N, O'Colmain B, Klaver CC, Klein R, Munoz B, Friedman DS, Kempen J,
312       Taylor HR, and Mitchell P. 2004. Causes and prevalence of visual impairment
313       among adults in the United States. *Arch Ophthalmol* 122:477-485.
314 Coughlin SS, Labenberg JR, and Tefft MC. 1993. Black-white differences in idiopathic
315       dilated cardiomyopathy: the Washington DC dilated Cardiomyopathy Study.
316       *Epidemiology* 4:165-172.
317 Eden E, Navon R, Steinfeld I, Lipson D, and Yakhini Z. 2009. GOrilla: a tool for
318       discovery and visualization of enriched GO terms in ranked gene lists. *BMC*
319       *Bioinformatics* 10:48.
320 Gilman SR, Iossifov I, Levy D, Ronemus M, Wigler M, and Vitkup D. 2011. Rare de
321       novo variants associated with autism implicate a large functional network of
322       genes involved in formation and function of synapses. *Neuron* 70:898-907.

Glessner JT, Wang K, Cai G, Korvatska O, Kim CE, Wood S, Zhang H, Estes A, Brune CW, Bradfield JP, Imielinski M, Frackelton EC, Reichert J, Crawford EL, Munson J, Sleiman PM, Chiavacci R, Annaiah K, Thomas K, Hou C, Glaberson W, Flory J, Otieno F, Garris M, Soorya L, Klei L, Piven J, Meyer KJ, Anagnostou E, Sakurai T, Game RM, Rudd DS, Zurawiecki D, McDougle CJ, Davis LK, Miller J, Posey DJ, Michaels S, Kolevzon A, Silverman JM, Bernier R, Levy SE, Schultz RT, Dawson G, Owley T, McMahon WM, Wassink TH, Sweeney JA, Nurnberger JI, Coon H, Sutcliffe JS, Minshew NJ, Grant SF, Bucan M, Cook EH, Buxbaum JD, Devlin B, Schellenberg GD, and Hakonarson H. 2009. Autism genome-wide copy number variation reveals ubiquitin and neuronal genes. *Nature* 459:569-573.

Hiranandani N, Bupha-Intr T, and Janssen PM. 2006. SERCA overexpression reduces hydroxyl radical injury in murine myocardium. *Am J Physiol Heart Circ Physiol* 291:H3130-3135.

Karolchik D, Barber GP, Casper J, Clawson H, Cline MS, Diekhans M, Dreszer TR, Fujita PA, Guruvadoo L, Haeussler M, Harte RA, Heitner S, Hinrichs AS, Learned K, Lee BT, Li CH, Raney BJ, Rhead B, Rosenbloom KR, Sloan CA, Speir ML, Zweig AS, Haussler D, Kuhn RM, and Kent WJ. 2014. The UCSC Genome Browser database: 2014 update. *Nucleic Acids Res* 42:D764-770.

Khurana E, Fu Y, Chen J, and Gerstein M. 2013. Interpretation of genomic variants using a unified biological network approach. *PLoS Comput Biol* 9:e1002886.

Lee I, Blom UM, Wang PI, Shim JE, and Marcotte EM. 2011. Prioritizing candidate disease genes by network-based boosting of genome-wide association data. *Genome Res* 21:1109-1121.

Leiserson MD, Eldridge JV, Ramachandran S, and Raphael BJ. 2013. Network analysis of GWAS data. *Curr Opin Genet Dev* 23:602-610.

Li WC, Kuszak JR, Dunn K, Wang RR, Ma W, Wang GM, Spector A, Leib M, Cotliar AM, Weiss M, and et al. 1995. Lens epithelial cell apoptosis appears to be a common cellular basis for non-congenital cataract development in humans and animals. *J Cell Biol* 130:169-181.

Liu ZP, Wang Y, Zhang XS, and Chen L. 2012. Network-based analysis of complex diseases. *IET Syst Biol* 6:22-33.

Long J, Zhang B, Signorello LB, Cai Q, Deming-Halverson S, Shrubsole MJ, Sanderson M, Dennis J, Michailiou K, Easton DF, Shu XO, Blot WJ, and Zheng W. 2013. Evaluating genome-wide association study-identified breast cancer risk variants in African-American women. *PLoS One* 8:e58350.

McCarroll SA, and Altshuler DM. 2007. Copy-number variation and association studies of human disease. *Nat Genet* 39:S37-42.

McElroy JP, Nelson MR, Caillier SJ, and Oksenberg JR. 2009. Copy number variation in African Americans. *BMC Genet* 10:15.

Mishra GR, Suresh M, Kumaran K, Kannabiran N, Suresh S, Bala P, Shivakumar K, Anuradha N, Reddy R, Raghavan TM, Menon S, Hanumanthu G, Gupta M, Upendran S, Gupta S, Mahesh M, Jacob B, Mathew P, Chatterjee P, Arun KS, Sharma S, Chandrika KN, Deshpande N, Palvankar K, Raghavnath R, Krishnakanth R, Karathia H, Rekha B, Nayak R, Vishnupriya G, Kumar HG, Nagini M, Kumar GS, Jose R, Deepthi P, Mohan SS, Gandhi TK, Harsha HC, Deshpande KS, Sarker M, Prasad TS, and Pandey A. 2006. Human protein reference database--2006 update. *Nucleic Acids Res* 34:D411-414.

Nahomi RB, Wang B, Raghavan CT, Voss O, Doseff AI, Santhoshkumar P, and Nagaraj RH. 2013. Chaperone peptides of alpha-crystallin inhibit epithelial cell apoptosis, protein insolubilization, and opacification in experimental cataracts. *J Biol Chem* 288:13022-13035.

National Center for Health Statistics. 2013. Health, United States, 2012: With Special Feature on Emergency Care. Hyattsville, MD.

Ng MC, Shriner D, Chen BH, Li J, Chen WM, Guo X, Liu J, Bielinski SJ, Yanek LR, Nalls
    MA, Comeau ME, Rasmussen-Torvik LJ, Jensen RA, Evans DS, Sun YV, An P,
    Patel SR, Lu Y, Long J, Armstrong LL, Wagenknecht L, Yang L, Snively BM,
    Palmer ND, Mudgal P, Langefeld CD, Keene KL, Freedman BI, Mychaleckyj JC,
    Nayak U, Raffel LJ, Goodarzi MO, Chen YD, Taylor HA, Jr., Correa A, Sims M,
    Couper D, Pankow JS, Boerwinkle E, Adeyemo A, Doumatey A, Chen G,
    Mathias RA, Vaidya D, Singleton AB, Zonderman AB, Igo RP, Jr., Sedor JR,
    Kabagambe EK, Siscovick DS, McKnight B, Rice K, Liu Y, Hsueh WC, Zhao W,
    Bielak LF, Kraja A, Province MA, Bottinger EP, Gottesman O, Cai Q, Zheng W,
    Blot WJ, Lowe WL, Pacheco JA, Crawford DC, Grundberg E, Rich SS, Hayes MG,
    Shu XO, Loos RJ, Borecki IB, Peyser PA, Cummings SR, Psaty BM, Fornage M,
    Iyengar SK, Evans MK, Becker DM, Kao WH, Wilson JG, Rotter JI, Sale MM, Liu
    S, Rotimi CN, and Bowden DW. 2014. Meta-analysis of genome-wide
    association studies in African Americans provides insights into the genetic
    architecture of type 2 diabetes. *PLoS Genet* 10:e1004517.
Online Mendelian Inheritance in Man O. 2014. Baltimore, MD: McKusick-Nathans
    Institute of Genetic Medicine, Johns Hopkins University.
Peri S, Navarro JD, Amanchy R, Kristiansen TZ, Jonnalagadda CK, Surendranath V,
    Niranjan V, Muthusamy B, Gandhi TK, Gronborg M, Ibarrola N, Deshpande N,
    Shanker K, Shivashankar HN, Rashmi BP, Ramya MA, Zhao Z, Chandrika KN,
    Padma N, Harsha HC, Yatish AJ, Kavitha MP, Menezes M, Choudhury DR,
    Suresh S, Ghosh N, Saravana R, Chandran S, Krishna S, Joy M, Anand SK,
    Madavan V, Joseph A, Wong GW, Schiemann WP, Constantinescu SN, Huang L,
    Khosravi-Far R, Steen H, Tewari M, Ghaffari S, Blobe GC, Dang CV, Garcia JG,
    Pevsner J, Jensen ON, Roepstorff P, Deshpande KS, Chinnaiyan AM, Hamosh A,
    Chakravarti A, and Pandey A. 2003. Development of human protein reference
    database as an initial platform for approaching systems biology in humans.
    *Genome Res* 13:2363-2371.
Pizzuti C, Rombo S, and Marchiori E. 2012. Complex Detection in Protein-Protein
    Interaction Networks: A Compact Overview for Researchers and Practitioners.
    In: Giacobini M, Vanneschi L, and Bush W, eds. *Evolutionary Computation,
    Machine Learning and Data Mining in Bioinformatics*: Springer Berlin
    Heidelberg, 211-223.
Pounds S, and Cheng C. 2006. Robust estimation of the false discovery rate.
    *Bioinformatics* 22:1979-1987.
Prasad TSK, Goel R, Kandasamy K, Keerthikumar S, Kumar S, Mathivanan S,
    Telikicherla D, Raju R, Shafreen B, Venugopal A, Balakrishnan L, Marimuthu A,
    Banerjee S, Somanathan DS, Sebastian A, Rani S, Ray S, Harrys Kishore CJ,
    Kanth S, Ahmed M, Kashyap MK, Mohmood R, Ramachandra YL, Krishna V,
    Rahiman BA, Mohan S, Ranganathan P, Ramabadran S, Chaerkady R, and
    Pandey A. 2009. Human Protein Reference Database--2009 update. *Nucleic
    Acids Res* 37:D767-772.
R Development Core Team. 2013. R: A language and environment for statistical
    computing. Vienna, Austria: R Foundation for Statistical Computing.
Ramos E, and Rotimi C. 2009. The A's, G's, C's, and T's of health disparities. *BMC
    Med Genomics* 2:29.
Sharma A, Gulbahce N, Pevzner SJ, Menche J, Ladenvall C, Folkersen L, Eriksson P,
    Orho-Melander M, and Barabasi AL. 2013. Network-based analysis of genome
    wide association data provides novel candidate genes for lipid and lipoprotein
    traits. *Mol Cell Proteomics* 12:3398-3408.
Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, and Sirotkin K.
    2001. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res* 29:4.

428  Talukder MA, Kalyanasundaram A, Zhao X, Zuo L, Bhupathy P, Babu GJ, Cardounel AJ,
429      Periasamy M, and Zweier JL. 2007. Expression of SERCA isoform with faster
430      Ca2+ transport properties improves postischemic cardiac function and Ca2+
431      handling and decreases myocardial infarction. *Am J Physiol Heart Circ Physiol*
432      293:H2418-2428.
433  van Dongen S. 2000. Graph Clustering by Flow Simulation PhD. University of Utrecht.
434  Vidal M, Cusick ME, and Barabasi AL. 2011. Interactome networks and human
435      disease. *Cell* 144:986-998.
436  Wei P, Milbauer LC, Enenstein J, Nguyen J, Pan W, and Hebbel RP. 2011. Differential
437      endothelial cell gene expression by African Americans versus Caucasian
438      Americans: a possible contribution to health disparity in vascular disease and
439      cancer. *BMC Med* 9:2.
440  Xu Z, Bensen JT, Smith GJ, Mohler JL, and Taylor JA. 2011. GWAS SNP Replication
441      among African American and European American men in the North Carolina-
442      Louisiana prostate cancer project (PCaP). *Prostate* 71:881-891.
443  Zhao M, and Zhao Z. 2013. CNVannotator: a comprehensive annotation server for
444      copy number variation in the human genome. *PLoS One* 8:e80170.

# Figure 1

Overview of our approach to identify CNVs associated with health disparities
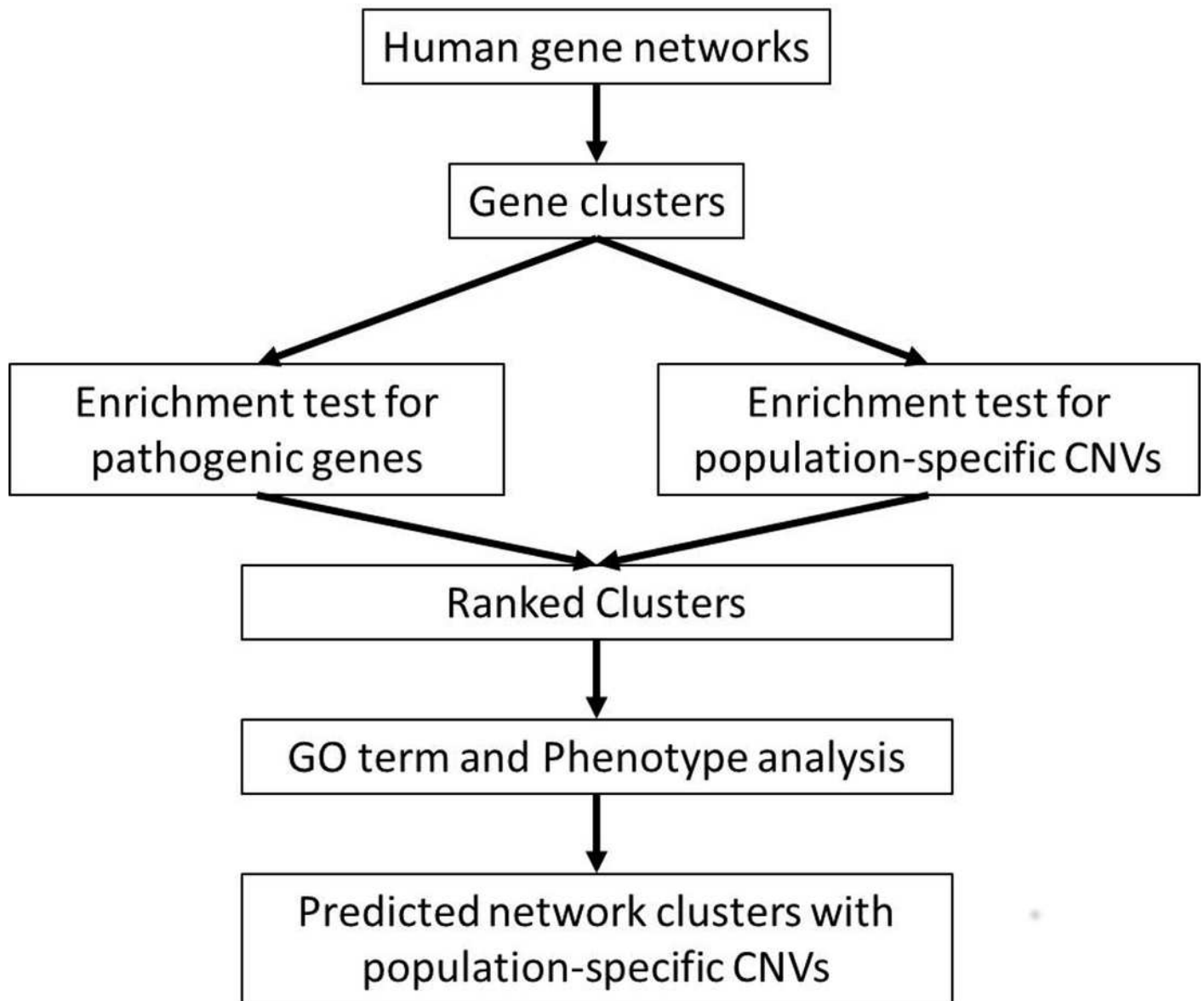
# Figure 2

Graph representations of selected clusters for biological significance analysis.

Each rounded rectangle represents a gene and each gray line represents a gene-gene interaction. Black rounded rectangles represent non-pathogenic genes and orange rounded rectangles represent pathogenic genes. Genes labeled with red or blue ovals are located in African American CNVs or in Caucasian CNVs. Genes with Green lines share the same GO terms. In each cluster, different line types represent the enrichment of different GO terms. Line types shown in different clusters refer to the enrichment of different GO terms.
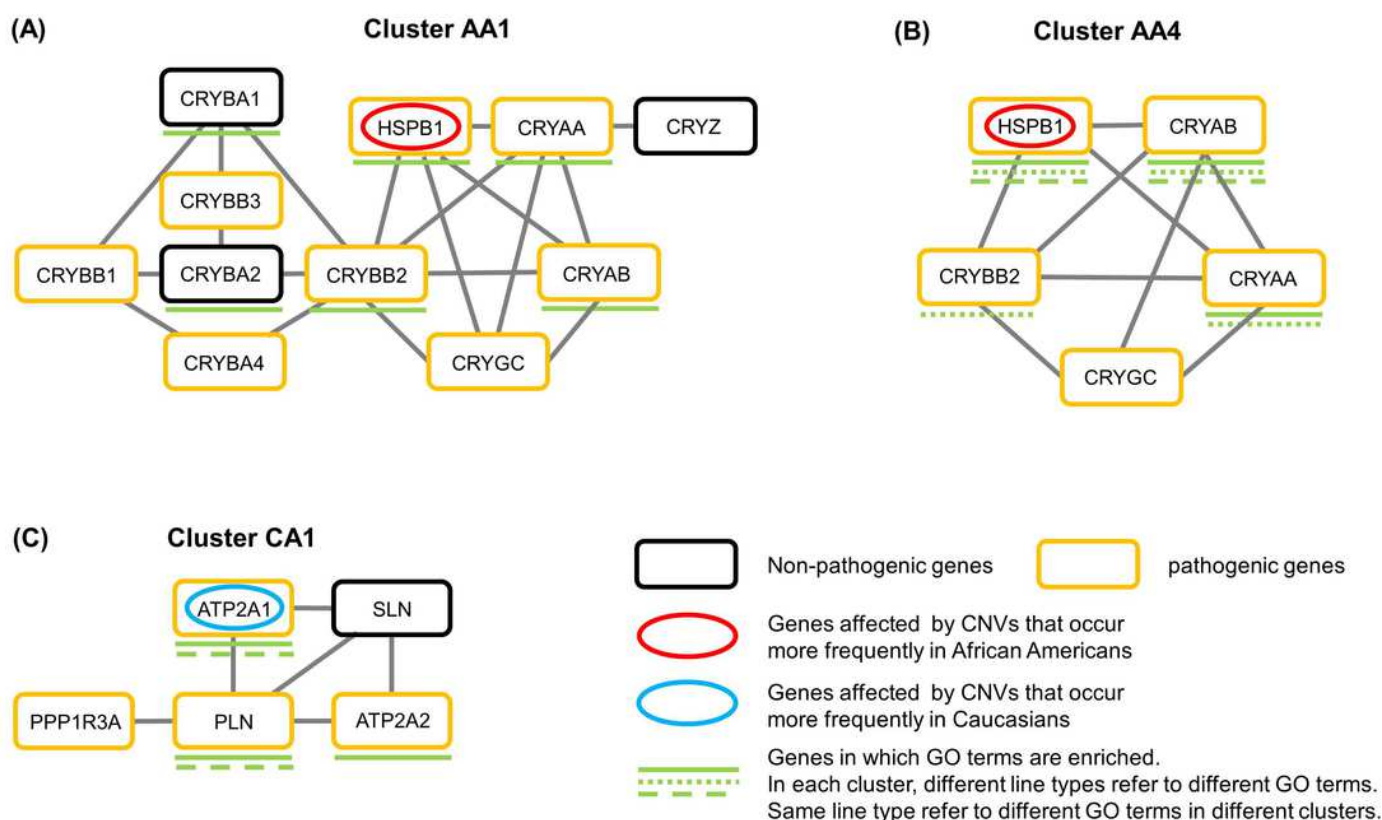
# Table 1(on next page)

Contingency tables

Table 1A. Contingency Table for Fisher's exact Test on Pathogenic Genes. Table 1B. Contingency Table for Fisher's exact Test on CNV genes. For each cluster, contingency tables were constructed for right-tailed Fisher's exact Tests. Table 1A is for pathogenic significance test, and Table 1B is for tests of enrichment significance of CNV genes (CNV_AA or CNV_CA genes). Q and q are the number of pathogenic genes in the whole networks and that in current cluster, respectively. N and m are the number of genes in whole networks and that in current cluster, respectively. S and s are the number of CNV_AA or CNV_CA genes in the whole networks and that in current cluster, respectively.

Table 1A. Contingency Table for Fisher's exact Test on Pathogenic Genes

|  | Pathogenic Genes | Non-pathogenic Genes | Total |
|---|---|---|---|
| **Genes in this cluster** | q | m-q | m |
| **Genes in other clusters** | Q-q | N-Q-m+q | N-m |
| **Total** | Q | N-Q | N |

Table 1B. Contingency Table for Fisher's exact Test on CNV genes

|  | CNV Genes | Non-CNV Genes | Total |
|---|---|---|---|
| **Genes in this cluster** | s | m-s | m |
| **Genes in other clusters** | S-s | N-S-m+s | N-m |
| **Total** | S | N-S | N |

For each cluster, contingency tables were constructed for right-tailed Fisher's exact Tests. Table 1A is for pathogenic significance test, and Table 1B is for tests of enrichment significance of CNV genes (CNV_AA or CNV_CA genes). Q and q are the number of pathogenic genes in the whole networks and that in current cluster, respectively. N and m are the number of genes in whole networks and that in current cluster, respectively. S and s are the number of CNV_AA or CNV_CA genes in the whole networks and that in current cluster, respectively.

# Table 2<sup>(on next page)</sup>

Cluster analysis results for HPRDNet and MultiNet

**Table 2**. Cluster analysis results for HPRDNet and MultiNet

| Network | Cluster Name | CNV_AA | CNV_CA | Pathogenic gene number | Cluster Size |
|---------|-------------|--------|--------|------------------------|--------------|
| **HPRDNet** | AA1 | *HSPB1* | - | 8 | 11 |
| | AA2 | *HSPB1* | - | 8 | 12 |
| | AA3 | *HSPB1* | - | 8 | 13 |
| | CA1 | - | *ATP2A1* | 4 | 5 |
| **MultiNet** | AA4 | *HSPB1* | - | 5 | 5 |
| | CA1 | - | *ATP2A1* | 4 | 5 |

Selected clusters were listed. CNV_AA and CNV_CA are CNV-related genes.

# Table 3<span>(on next page)</span>

Detected genes with potential roles in health disparity and their located CNVs

Table 3. Detected genes with potential roles in health disparity and their located CNVs

| Gene | Chr | Gene Coordinates | CNV Region | CNV Type | CNV Occurrence preference |
|------|-----|------------------|------------|----------|---------------------------|
| *HSPB1* | 7 | 75,931,861-75,933,614 | 75,867,431-76,481,102 | Duplication | Only in African American |
| | | | 75,929,740-76,481,102 | Duplication | Only in African American |
| | | | 75,929,740-76,568,388 | Duplication | More in African American than in Caucasian |
| *ATP2A1* | 16 | 28,889,726-28,915,830 | 28,306,730-28,936,772 | Duplication | Only in Caucasian |

Chr represents chromosomes. CNV Regions are regions of CNVs identified in more than a single individual; all CNVs listed have a type of Duplication, referring to one copy increase. CNV Regions and Types are from the CNV map (McElroy et al. 2009). CNV Occurrence preference describes in which population those CNVs have higher occurrence frequency.

# Table 4<sub>(on next page)</sub>

Enriched GO terms with CNV-genes in the identified network clusters

Table 4. Enriched GO terms with CNV-genes in the identified network clusters

| Clusters | Involved Genes | GO Domain | GO ID | GO term |
|---|---|---|---|---|
| AA1 | *HSPB1, CRYAA, CRYAB, CRYBB2, CRYBA1, CRYBA2* | Molecular Function | GO:0042802 | Identical protein binding |
| AA4 | *HSPB1, CRYAA, CRYAB* | Biological Process | GO:0043086 | negative regulation of catalytic activity |
| | | | GO:0043066 | negative regulation of apoptotic process |
| | | | GO:0043069 | negative regulation of programmed cell death |
| | *HSPB1, CRYAA, CRYAB, CRYBB2* | Molecular Function | GO:0042802 | Identical protein binding |
| | *HSPB1, CRYAB* | Cellular Component | GO:0030018 | Z disc |
| CA1 | *ATP2A1, ATP2A2, PLN* | Biological Process | GO:0048878 | chemical homeostasis |
| | *ATP2A1, PLN* | Biological Process | GO:0006937 | regulation of muscle contraction |
| | | | GO:0008016 | regulation of heart contraction |

Biological relevance of network clusters was analyzed by GOrilla (Eden et al. 2009) to search for enriched gene ontology (GO) terms. Genes in the selected clusters were used as target genes, and all genes in the networks were treated as background genes. Three types of GO terms were analyzed: biological process, molecular function and cellular component. The default *p*-value threshold ($1\times10^{-3}$) was used. In the results, enriched GO terms that are associated with CNV_AA gene HSPB1 and CNV_CA gene ATP2A1 were selected and listed in the table.

**Table 5**(on next page)

Associated diseases of genes with enriched GO terms.

Table 5. Associated diseases of genes with enriched GO terms.

| Cluster | Gene | Associated Disease |
|---------|------|--------------------|
| **AA1 and AA4** | HSPB1 | Axonal Charcot-Marie-Tooth disease type 2F |
| | | Distal hereditary motor neuronopathy type 2B |
| | CRYAA | Multiple types of cataract 9 |
| | CRYAB | Multiple types of cataract 16 |
| | | Dilated cardiomyopathy-1II |
| | | Myofibrillar myopathy-2 |
| | | CRYAB-related fatal infantile hypertonic myofibrillar myopathy |
| | CRYBB2 | Multiple types of Cataract 3 |
| **CA1** | ATP2A1 | Brody myopathy |
| | ATP2A2 | Acrokeratosis verruciformis |
| | | Darier disease |
| | PLN | Dilated cardiomyopathy-1P |
| | | Familial hypertrophic cardiomyopathy-18 |

Only GO terms that contain CNV-genes are studied due to our focus on the role of CNV-genes in health disparity.