

1 Identifying Genetic Interactions Associated with Late-Onset

2 Alzheimer's Disease

3 Charalampos S. Floudas, M.D, Ph.D.^{1§}, Nara Um, M.D, M.S.¹, M. Ilyas Kamboh, Ph.D.²,

4 Michael M. Barmada, Ph.D.², Shyam Visweswaran, M.D, Ph.D.^{1,3}

5
6 ¹Department of Biomedical Informatics, University of Pittsburgh, Pittsburgh, PA, USA

7 ²Department of Human Genetics, University of Pittsburgh, Pittsburgh, PA, USA

8 ³The Intelligent Systems Program, University of Pittsburgh, Pittsburgh, PA, USA

9
10 [§]Corresponding author

17 Email addresses:

18 CSF: chfloudas@gmail.com

19 SV: shv3@pitt.edu

20

21 **Abstract**

22 **Background**

23 Identifying genetic interactions in data obtained from genome-wide association studies (GWASs)
24 can help in understanding the genetic basis of complex diseases. The large number of single
25 nucleotide polymorphisms (SNPs) in GWASs however makes the identification of genetic
26 interactions computationally challenging. We developed the Bayesian Combinatorial Method
27 (BCM) that can identify pairs of SNPs that in combination have high statistical association with
28 disease.

29 **Results**

30 We applied BCM to two late-onset Alzheimer's disease (LOAD) GWAS datasets to identify
31 SNP-SNP interactions between a set of known SNP associations and the dataset SNPs. For
32 evaluation we compared our results with those from logistic regression, as implemented in
33 PLINK. Gene Ontology analysis of genes from the top 200 dataset SNPs for both GWAS
34 datasets showed overrepresentation of LOAD-related terms. Four genes were common to both
35 datasets: APOE and APOC1, which have well established associations with LOAD, and
36 CAMK1D and FBXL13, not previously linked to LOAD but having evidence of involvement in
37 LOAD. Supporting evidence was also found for additional genes from the top 30 dataset SNPs.

38 **Conclusion**

39 BCM performed well in identifying several SNPs having evidence of involvement in the
40 pathogenesis of LOAD that would not have been identified by univariate analysis due to small
41 main effect. These results provide support for applying BCM to identify potential genetic
42 variants such as SNPs from high dimensional GWAS datasets.

43 Introduction

44 Elucidating the genetic basis of common diseases will lead to understanding of the biological
45 mechanisms that underlie such diseases and can help in risk assessment, diagnosis, prognosis and
46 development of new therapies. During the past several decades genetic linkage studies have been
47 effective in mapping genetic loci responsible for many Mendelian diseases that are caused by a
48 single genetic variant (Hardy & Singleton, 2009). More recently, genetic studies have indicated
49 that most common diseases are likely to be polygenic where multiple genetic variants acting
50 singly and in combination underlie the expression of disease (Thornton-Wells, Moore & Haines,
51 2004).

52 The commonest type of genetic variation is the single nucleotide polymorphism (SNP)
53 that results when a single nucleotide is replaced by another in the genome sequence. The
54 development of high-throughput genotyping technologies has led to a flurry of genome-wide
55 association studies (GWASs) with the aim of discovering SNPs that are associated with common
56 diseases. GWASs have been moderately successful in identifying SNPs associated with common
57 diseases and traits. However, in most cases the identified SNPs have small effect sizes, and the
58 proportion of heritability explained is quite modest. One view is that SNPs may interact in subtle
59 ways that lead to substantially greater effects than the effect due to any single SNP. Another
60 view is that common diseases may be due to rare and usually deleterious SNPs that cause disease
61 in individual patients and that in different individuals or subpopulations the disease is caused by
62 different deleterious SNPs.

63 This paper addresses the challenge of identifying interacting SNPs that may have small
64 effects and describes a Bayesian combinatorial method (BCM) for identifying such interactions
65 that are associated with disease. This method has been shown empirically to perform well on low

66 dimensional synthetic data (Balding, 2006). However, to our knowledge BCM has not been
67 applied to a disease dataset with a large number of SNPs. In this paper we apply BCM to an
68 Alzheimer's disease GWAS dataset to identify SNPs that interact with known Alzheimer
69 associated SNPs.

70 As background, we provide brief summaries about GWASs, genetic interactions, and
71 Alzheimer's disease in the following sections.

72 **Genome-wide Association Studies**

73 The development of high-throughput genotyping technologies that assay hundreds of thousands
74 of SNPs or more, along with the identification of SNPs in the human genome by the
75 International HapMap Project led to the emergence of GWASs. GWASs are typically case-
76 control studies aimed at discovering SNPs – either as disease causing variants or as markers of
77 disease – that are associated with a common disease or trait. The success of GWASs is based in
78 large part on the common disease - common variant hypothesis. This hypothesis posits that
79 common diseases in most individuals are caused by relatively common genetic variants that have
80 low penetrance and hence have small to moderate influence in causing disease. An alternative
81 hypothesis is the common disease - rare variant hypothesis, which posits that many rare variants
82 underlie common diseases and each variant causes disease in relatively few individuals with high
83 penetrance. Both these hypotheses likely contribute to common diseases with genetic variants
84 may range from rare to the common SNPs.

85 GWAS data is typically analyzed for univariate associations between SNPs and the
86 disease of interest; the statistical tests used include the Pearson's chi-square test, the Fisher's
87 exact test, the Cochran-Armitage trend test, and odds ratios (Cordell, 2009). SNPs identified as

88 significant by univariate analyses may be further examined for interactions among them using
89 methods such as logistic regression.

90 **Genetic Interactions**

91 Genetic interactions, also known as epistasis, can be defined biologically as well as statistically.
92 Biologically, epistasis refers to gene-gene interaction when the action of one gene is modified by
93 one or several other genes. Statistically, epistasis refers to interaction between variants at
94 multiple loci in which the total effect of the combination of variants at the different loci may
95 differ considerably from a linear combination of the effects of individual loci. The detection of
96 statistical epistasis has the potential to indicate genetic loci that have a biological interaction
97 (Hahn, Ritchie & Moore, 2003).

98 Statistical methods for identifying genetic interactions can be broadly divided into
99 exhaustive and non-exhaustive methods. Exhaustive methods examine all possible SNP-subsets
100 and examples include Multifactor Dimensionality Reduction (MDR) (Moore et al, 2006) and the
101 BCM (Visweswaran, Wong & Barmada, 2009) that we describe in the next section. Examples of
102 non-exhaustive methods include Boolean Operation-based Screening (BOOST), SNPHarvester
103 and SNPRuler. We briefly describe these methods below.

104 The software package PLINK that is used widely for the analysis of GWAS datasets also
105 implements logistic regression for the detection of SNP-SNP interactions and offers the option to
106 test either all or specific sets of SNPs in a dataset (Purcell et al, 2007).

107 MDR exhaustively evaluates all 1-,2-,3-,... n -SNP subsets where n is specified by the user.
108 It combines the variables in a SNP subset to construct a single binary variable and uses
109 classification accuracy of the binary variable to evaluate a SNP-subset. Since MDR does not

110 scale up beyond a few hundred SNPs, for high dimensional data a multivariate filtering
111 algorithm called ReliefF is applied to reduce the number of SNPs to a few hundred (M. D.
112 Ritchie et al, 2001; Hahn, Ritchie & Moore, 2003; Moore et al, 2006; Moore & White, 2007).

113 BOOST uses a two-step procedure (Wan et al, 2010a). In the screening step, it uses an
114 approximate likelihood ratio statistic that is computationally efficient and computes it for all
115 pairs of SNPs. Only those SNPs that pass a threshold in the first step are examined for significant
116 interaction effect using the classical likelihood ratio test that is computationally more expensive.

117 SNPHarvester is a stochastic search algorithm that uses a two-step procedure to identify
118 epistatic interactions (Yang et al, 2009). In the first step it identifies 40–50 significant SNP
119 groups using a stochastic search strategy, and in the second step, it fits a penalized logistic
120 regression model to each group.

121 SNPRuler searches in the space of SNP rules and uses a branch-and-bound strategy to
122 prune the huge number of possible rules in GWAS data (Wan et al, 2010b). An example of a rule
123 is $X_1 = 0 \wedge X_2 = 2 \Rightarrow Z = 1$ (X_1 and X_2 are SNPs, the three genotypes that a SNP can take are
124 coded as 0, 1 and 2 and Z is a binary outcome variable). The quality of a rule is evaluated with
125 the chi-square statistic.

126 **Alzheimer's Disease**

127 Alzheimer's disease (AD) is the commonest neurodegenerative disease associated with aging
128 and the commonest cause of dementia (Goedert & Spillantini, 2006). AD affects about 3% of all
129 people between ages 65 and 74, about 19% of those between 75 and 84, and about 47% of those
130 over 85. AD is characterized by adult onset of progressive dementia that typically begins with

131 subtle memory failure and progresses to a slew of cognitive deficits like confusion, language
132 disturbance and poor judgment (Bertram, Lill & Tanzi, 2010).

133 AD is typically divided into early-onset Alzheimer's disease (EOAD) in which the onset
134 of disease is before 60 years of age and late-onset Alzheimer's disease (LOAD) in which the
135 onset is at or after 60 years of age. EOAD is rare and exhibits an autosomal dominant mode of
136 inheritance. The genetic basis of EOAD is well established, and mutations in one of three genes
137 (amyloid precursor protein gene, presenelin 1, or presenelin 2) account for most cases of EOAD
138 (Avramopoulos, 2009).

139 LOAD is widespread and is estimated to strike almost half of all people over the age of
140 85. LOAD is believed to be a disease with both genetic and environmental influences, and
141 elucidating the role of genetic factors in the pathogenesis and development of LOAD has been a
142 major focus of research for more than a decade. One genetic risk factor for LOAD that has been
143 consistently replicated is the apolipoprotein E (APOE) locus determined by the combined
144 genotypes at the loci rs429358 (APOE*4) at codon 112 and rs7412 (APOE*2) at codon 158
145 (Holtzman, Morris & Goate, 2011). Because the two SNPs are in LD their combination
146 determines six genotypes (2-2, 2-3, 3-3, 3-4, 2-4, 4-4) and the well-established protein
147 polymorphism as one locus with three alleles (E*2, E*3, E*4). In the past few years, GWASs
148 have identified several additional genetic loci associated with LOAD (Reiman et al, 2007;
149 Hollingworth et al, 2011; Hu et al, 2011; Wijsman et al, 2011; Kamboh et al, 2012).

150 **Bayesian Combinatorial Method**

151 BCM uses a Bayesian network (BN) to model a set of SNPs and interactions among them and
152 their association with disease, and the model is evaluated with a Bayesian score. It then
153 exhaustively searches a space of all possible models to identify high scoring models.

154 **Bayesian network model and score.** For a dataset D that contains a set of n SNPs $\{X_1,$
155 $X_2, \dots, X_n\}$ and a binary outcome variable Z (e.g, disease or phenotype) on N individuals, BCM's
156 goal is to identify a set of SNPs that together are most predictive of Z in D . We model the effects
157 of SNPs on Z with a BN that has n SNP-nodes and an additional node for Z . In this BN, which
158 we call a SNP-BN, a subset of the n SNPs is modeled to have an effect on Z and every node in
159 that subset has an arc to Z and every node not in the subset does not have an arc to Z . Also, there
160 are no arcs between the SNP-nodes since we do not model the relations among the SNPs. Figure
161 1 gives an example of a SNP-BN where SNPs X_2 and X_3 are modeled to have a joint effect on Z
162 (as shown by the arcs connecting them to Z) and the remaining SNPs do not have an effect on Z .

163 We evaluate the goodness of fit of a SNP-BN to data using an efficiently computable
164 Bayesian score that computes the posterior probability of the BN given the data. In particular, we
165 compute the BDeu (Bayesian Dirichlet equivalence uniform) score described in (Heckerman,
166 Geiger & Chickering, 1995) which is commonly used in BN learning from data. This score is
167 computed efficiently in closed form as follows:

168
$$P(M | D) = P(M) \prod_{i=1}^{n+1} \prod_{j=1}^{J_i} \frac{\Gamma(\alpha_{ij})}{\Gamma(n_{ij} + \alpha_{ij})} \prod_{k=1}^{K_i} \frac{\Gamma(n_{ijk} + \alpha_{ijk})}{\Gamma(\alpha_{ijk})} \quad (1)$$

169 where, $\Gamma(\cdot)$ is the gamma function, M is a SNP-BN, $P(D | M)$ is the posterior probability of M
170 given D , $P(M)$ is the prior probability of M , K_i is the number of states of variable X_i represented
171 by node i , J_i is the number of joint states of the parents of node i , n_{ijk} is the number of times in

172 the data that node i is in state k given parent state j , α_{ijk} are the parameter priors in a Dirichlet
 173 distribution which define the prior probability over the BN parameters. Also,

174
$$n_{ij} = \sum_{k=1}^{K_i} n_{ijk}, \alpha_{ij} = \sum_{k=1}^{K_i} \alpha_{ijk}, \text{ and } \alpha_{ijk} = \frac{\alpha}{J_i \cdot K_i},$$
 where α is a single user-defined parameter prior.

175 The n_{ijk} are obtained from the data and stored in a counts table that is associated with each node
 176 (an example of a counts table for node Z is shown in Figure 1). We make the following

177 assumptions and simplifications: (1) model the prior probability $P(M)$ as a constant, i.e., *a priori*
 178 we consider all models to be equally plausible, (2) set $\alpha = 1$ which is a commonly used non-
 179 informative parameter prior, (3) use the logarithmic form to simplify computations when dealing
 180 with very small numbers, and (4) assign the score for a SNP-BN model to be the BDeu score
 181 attributable to just node Z (Visweswaran & Wong, 2009). The reason for assumption (4) is as
 182 follows. The BDeu score decomposes over the nodes in the BN and each node makes an
 183 independent contribution to the overall score. In the space of SNP-BNs, the score contributions
 184 of the SNP-nodes is a constant since they have no incoming arcs, and hence variation in the
 185 scores for distinct SNP-BNs is due only to the score attributable to Z . Thus, the score we use for
 186 a SNP-BN is given by the following expression (index i is absent since there is only one node
 187 under consideration, namely, Z , and $K = 2$ since Z is binary):

188
$$\text{score}(M) = \sum_{j=1}^J \log \frac{\Gamma(\alpha_j)}{\Gamma(n_j + \alpha_j)} + \sum_{k=1}^2 \log \frac{\Gamma(n_{jk} + \alpha_{jk})}{\Gamma(\alpha_{jk})} \quad (2)$$

189 We have evaluated the BCM score in low dimensional synthetic data and found that in such data
 190 it has significantly greater power and is computed more efficiently than MDR (Visweswaran,
 191 Wong & Barmada, 2009; Jiang, Barmada & Visweswaran, 2010; Jiang et al, 2011).

192 **Materials & Methods**

193 This section describes the GWAS datasets, the experimental methods, and previously identified
194 LOAD SNPs.

195 **GWAS Dataset**

196 We used two different LOAD GWAS datasets in our experiments. The first dataset was part of
197 the University of Pittsburgh Alzheimer's Disease Research Center (ADRC) that is described
198 elsewhere (Kamboh et al, 2012). This dataset consists of 2245 individuals, of which 1290 had
199 LOAD and 955 did not. For each individual, the genotype data consists of 682,685 SNPs on
200 autosomal chromosomes.

201 The second dataset was collected by the Translational Genomics Research Institute
202 (TGen) (Reiman et al, 2007). This dataset consists of 1411 individuals, of which. 861 had LOAD
203 and 550 did not. For each individual, the imputed genotype data consists of 234,665 SNPs on
204 autosomal chromosomes. For each individual, the genotype data consists of 502,627 SNPs; the
205 original investigators analyzed 312,316 SNPs after applying quality controls. We used those
206 312,316 SNPs, plus two additional APOE SNPs from the same study namely, rs429358 and
207 rs7412.

208 **Experimental Methods**

209 BCM searches exhaustively over all possible SNP-BN models in a dataset. For a GWAS dataset
210 with half a million SNPs, the number of SNP-BN models is $2^n = 9.95 \times 10^{150514}$ and the number

211 of SNP-BN models with just 2 SNPs is $\binom{500000}{2} = 1.25 \times 10^{11}$. Thus, the search space is very

212 large and it is computationally infeasible to evaluate every model in the space (Ritchie, 2011).

213 We addressed this challenge by applying BCM to a restricted space of SNP-BN models
214 that consisted of a subset of all possible 2-SNP models. We considered only those 2-SNP models
215 where one of the SNPs in a model is a member of a set of SNPs previously known to be
216 associated with LOAD and the second SNP is any SNP (excluding the first SNP) in the dataset
217 of interest. Since the number of known LOAD associated SNPs is much smaller than the number
218 of SNPs in a dataset, it was computationally tractable to search this space of SNP-BN models.
219 The selection of the previously identified LOAD SNPs that we used is described in the next
220 section.

221 We applied BCM to each of the two GWAS datasets separately and analyzed in detail the
222 top scoring 200 SNP-BN models. From each SNP-BN model, we extracted the SNP that was not
223 in the set of previously identified LOAD SNPs. We mapped these SNPs to genes and considered
224 only intragenic SNPs for further analyses. We performed the SNP to gene mapping with BioQ, a
225 web-service which uses dbSNP build 135 and Genome Assembly GRCh37.p5 (Saccone, Quan &
226 Jones, 2012). We performed enrichment analysis of the annotations of the associated genes in the
227 Gene Ontology (GO) with the web-based tool GeneCoDis. For a set of genes GeneCoDis
228 retrieves the associated GO terms, and identifies and ranks those GO terms that are significantly
229 enriched in the set of genes (Carmona-Saez et al, 2007; Nogales-Cadenas et al, 2009). Enriched
230 functional descriptors facilitate the interpretation of the gene set. The hierarchical nature of the
231 GO annotations however means that the set of enriched GO terms may contain terms closely
232 related in a parent-child relationship (Khatri & Drăghici, 2005). Such redundant terms confound

233 the interpretation. Therefore, we further examined the GO terms associated with the intragenic
234 SNPs using the REViGo webserver. The REViGo software evaluates the semantic similarity
235 between the enriched terms, identifies the most informative common ancestors and the related
236 redundant GO terms and groups the latter under their ancestors (Supek et al, 2011). The resulting
237 set facilitates simultaneous examination of the enriched GO terms at two levels: a detailed one, at
238 the lowest level overrepresented term and a more abstract one at the highest level common
239 ancestor of overrepresented terms. The detailed level can reveal specific genes of interest
240 whereas the abstract level serves a compact overview of the processes, functions and cellular
241 compartments associated with the genes in the set.

242 In addition to the analysis of the top scoring 200 SNP-BN models, we performed
243 additional analyses of the top scoring 30 SNP-BN models. We analyzed the genes associated
244 with the intragenic SNPs for differential expression in AD, through the ArrayExpress web server
245 (Parkinson et al, 2011) and biological function analysis. Differential gene expression in relation
246 to AD aims to integrate experimental evidence from transcriptomic analysis with those of
247 genomic analysis. Up-regulation or down-regulation in AD of a gene in our results indicates
248 increased biological plausibility for the reported genetic interaction. Finally, elements from the
249 functional description of a gene (expression site, function related to the nervous system or
250 pathways of LOAD, previous literature) were considered as supporting the biological relevance
251 of an identified interaction.

252 We also compared BCM with logistic regression as implemented in PLINK for the
253 identification of statistical genetic interactions. For this comparison we used as previous
254 knowledge SNPs for all methods the rs429358 (APOE *4), a known LOAD risk SNP. We
255 applied the methods to the ADRC LOAD dataset.

256 **Previously Identified LOAD SNPs**

257 We obtained a set of SNPs that are known to be associated with LOAD from the AlzGene
258 website. The AlzGene website contains a regularly updated database of SNPs that have been
259 shown to be associated with LOAD mostly in GWAS studies (Bertram et al, 2007). The curators
260 of the AlzGene website use criteria established by the Human Genome Epidemiology Network
261 (HuGENet) for assessing the cumulative evidence of associations of SNPs with disease
262 (Ioannidis et al, 2008). We obtained 10 SNPs that were assessed to have sufficiently strong
263 evidence of being associated with LOAD from the AlzGene website in March 2012. If a
264 previously identified LOAD SNP was not present in our datasets, we selected a replacement
265 SNP. The replacement SNP was within 500 kb, in the same gene, as the original SNP with
266 pairwise linkage disequilibrium threshold of $r^2 \geq 0.8$, using the SNAP web-based tool (Johnson
267 et al, 2008). Using this protocol, we were unable to identify replacement SNPs in the TGen
268 dataset for three previously identified LOAD SNPs; therefore we replaced them with SNPs from
269 other genes, also reported as significantly associated with LOAD in the AlzGene website. Table
270 1 gives the list of previously identified LOAD SNPs that we used in the experiments.

271 **Results and Discussion**

272 This section describes the results that were obtained from applying BCM to the ADRC LOAD
273 dataset and from applying BCM to the ADRC and the TGen GWAS datasets.

274 **Top Scoring SNP-BN Models**

275 Each SNP-BN model includes two SNPs of which one SNP is a previously identified LOAD
276 SNP and the other SNP is not. We call the former SNP a *known SNP* and the latter SNP a *dataset*

277 *SNP*. The known and dataset SNPs from the top scoring 200 SNP-BN models are given in Table
278 S1 (for ADRC) and Table S2 (for TGen) in the Supplemental Tables. A plot of the scores of the
279 top scoring 200 SNP-BN models for the two datasets is shown in Figure 2.

280 In the ADRC dataset, the known SNP in each of the top scoring 200 SNP-BN models is
281 rs429358 (APOE*4). In the TGen dataset, rs429358 is the known SNP in 192 of the top scoring
282 200 SNP-BN models, specifically models ranked 1 and 10-200, and in the 8 remaining models
283 (ranked 2-9) the known SNP belongs to genes GAB2, MS4A6A, MS4A4E, CR1, PICALM,
284 SORL1, TF whereas the dataset SNP is rs7412 for all 8 models. SNPs rs429358 and rs7412 are
285 located on the APOE gene and their combined genotypes determine the APOE allelic status
286 which is known to be the strongest genetic variant that is predictive of LOAD.

287 In the ADRC dataset, the dataset SNPs from the top scoring 200 models included 92
288 intragenic SNPs that mapped to 77 distinct genes, and the dataset SNPs from the top scoring 30
289 models included 18 intragenic SNPs that mapped to 15 distinct genes. In the TGen dataset, the
290 dataset SNPs from the top scoring 200 models included 82 intragenic SNPs that mapped to 69
291 genes, and the dataset SNPs from the top scoring 30 models included 19 intragenic SNPs that
292 mapped to 11 genes.

293 In the top scoring 200 SNP-BN models the two datasets have in common two intragenic
294 SNPs, rs7412 (APOE gene) and rs4420638 (APOC1 gene) as well as two genes mapped from
295 intragenic SNPs, CAMK1D (rs11257738 in ADRC and rs17151584 in TGen) and FBXL13
296 (rs7779121 in ADRC and rs17475512 in TGen).

297 **GO Term Analysis**

298 The most informative common ancestors of the overrepresented GO terms obtained from
299 GeneCoDis for the ADRC dataset are given in Table S3 and for the TGen dataset are given in
300 Table S4 in the Supplemental Tables. In both sets nervous system-related terms are enriched
301 (e.g, *regulation of dendrite development, nervous system development, regulation of axon*
302 *extension, short term memory*), as well as terms related to cholesterol and lipid metabolism (e.g,
303 *lipid metabolic process, chylomicron*), beta amyloid (*beta amyloid binding*) cell membranes (e.g,
304 *integral to membrane, plasma membrane, postsynaptic , clathrin-coated endocytic vesicle*),
305 calmodulin and intracellular calcium homeostasis (e.g, *calmodulin binding, cytosolic calcium ion*
306 *transport*) and the immune system (*immunoglobulin binding*). Overrepresentation of these terms
307 shows that the identified genes from both datasets include genes that are members of
308 biochemical pathways involved in LOAD pathophysiology (Holtzman, Morris & Goate, 2011;
309 Morgan, 2011).

310 **Expression Analysis**

311 The 15 genes corresponding to the 18 dataset intragenic SNPs from the top scoring 30 models in
312 the ADRC dataset and the 19 genes corresponding to the 19 dataset intragenic SNPs from the top
313 scoring 30 models in the TGen dataset were examined for relative expression in AD (Table 2 for
314 the ADRC dataset and in Table 3 for the TGen dataset). In these tables, the second to last column
315 gives the rank of the corresponding SNP based on the model score obtained by applying BCM to
316 1-SNP models. For some of the SNPs the rank based on the 1-SNP model is very low compared
317 to the score of the corresponding 2-SNP model which implies that these SNPS would not have

318 been identified by univariate analysis. The last column gives the p values for the pairs of SNPs
319 that were obtained from using logistic regression in PLINK.

320 **Comparison of BCM with PLINK**

321 Among the top 50 models ranked by BCM and PLINK respectively, there are 4 models in
322 common, and among the top 200 models, there were 25 models in common between BCM and
323 PLINK.

324

325 **Discussion**

326 Examining all pairs of SNPs in a GWAS dataset for identifying interacting SNP pairs is usually
327 not computationally tractable due to the large number of SNP pairs. We addressed this challenge
328 by examining only a subset of all pairs of SNPs where one member of the pair is drawn from a
329 small set of previously known disease-associated SNPs and by using a Bayesian score to
330 evaluate that statistical association of a SNP pair with the disease. We applied this strategy to
331 two LOAD GWAS datasets and our results show that it can identify interacting SNPs of
332 plausible biological significance. Moreover, this strategy finds SNPs that would be overlooked in
333 a univariate analysis because they exhibit small main effects; however, they are detected when
334 paired with another SNP due to interaction effects.

335 In both LOAD GWAS datasets that we examined, the previously known disease-
336 associated SNP that was identified is either rs429358 (APOE*4) or rs7412 (APOE*2); these
337 SNPs reside in the APOE gene which is known to be the strongest genetic determinant for
338 LOAD. GO term enrichment analysis of the dataset SNPs identified terms that are relevant to

339 biochemical pathways implicated in the pathogenesis of LOAD such as *lipid metabolic process*,
340 *calmodulin binding*, *nervous system development* and multiple membrane-related terms.

341 Gene expression analysis of the dataset SNPs showed that for each dataset studied a
342 majority of the genes corresponding to the top 30 dataset SNPs are differentially expressed in
343 LOAD. Functional annotations and literature evidence that are presented with the expression
344 data in the relevant tables further support the role of these genes in the pathogenesis of LOAD.

345 Among the genes corresponding to the top 200 dataset SNPs, besides the APOE gene,
346 three other genes are common to both datasets: APOC1, CAMK1D and FBXL13. Evidence
347 supporting the interaction of APOC1 with APOE are presented in the analysis for the top 30
348 dataset SNPs. CAMK1D (calcium/calmodulin-dependent protein kinase ID) belongs to the
349 family of calmodulin kinases which modulate neuronal development and plasticity (Wayman et
350 al, 2008). It has been found to be overexpressed in AD and is expressed in the brain especially
351 during hippocampal formation with high expression in the pyramidal cell layers (Lukk et al,
352 2010; Pugazhenthii et al, 2011). It encodes a member of the Ca²⁺/calmodulin-dependent protein
353 kinase 1 subfamily of serine/threonine kinases (Maglott et al, 2011). CAMK1D interacts with
354 CALM1 (calmodulin), which has been associated with AD risk (Lambert et al, 2010). The
355 encoded protein may regulate calcium-mediated granulocyte function and activates MAPK3
356 (Mitogen-activated protein kinase 3 - inferred function by similarity). In vitro, it phosphorylates
357 transcription factor CREM (cAMP responsive element binding) isoform Beta and probably
358 CREB1 (Pugazhenthii et al, 2011). The CREB pathway has a role in memory formation and
359 CREB phosphorylation has been proposed as a signalling pathway involved in the pathogenesis
360 of AD (Müller et al, 2011) and also its down-regulation may have a role in exacerbations of AD
361 (Pugazhenthii et al, 2011). Another member of the same family, neuronal CaM kinase II

362 phosphorylates tau protein on ser262, an important step in the formation of neurofibrillary
363 tangles in AD (Yamauchi, 2005). FBXL13 (F-box and leucine-rich repeat protein 13) belongs to
364 the F-box protein family. Members of this family have a characteristic approximately 40-amino
365 acid F-box motif and take part in SCF (SKP1-CUL1-F-box protein) complexes that act as
366 protein-ubiquitin ligases (Maglott et al, 2011). The ubiquitin-proteasome system is involved in
367 protein turnover and degradation and is perturbed in AD (Riederer et al, 2011). An SCF complex
368 of another F box protein (FBXW7) is involved in the degradation of NICD (NOTCH1 released
369 notch intracellular domain) and probably of PSEN1 (The UniProt Consortium, 2013).

370 In addition to the genes corresponding to the top 30 top scoring SNP-BN models in the
371 ADRC dataset, we found other genes in lower scoring SNP-BN models with plausible
372 associations with LOAD. In the 80th scoring model (dataset SNP rs7793977), gene PION [pigeon
373 homolog (*Drosophila*)], also known as GSAP (gamma-secretase-activating protein), is known to
374 increase amyloid beta production (Maglott et al, 2011). In the 196th scoring model, (dataset SNP
375 rs6534145), gene PDE5A (phosphodiesterase 5A, cGMP-specific) could be implicated to LOAD
376 pathogenesis via two different mechanisms. PDE5A is a substrate of CASP3 (caspase 3) (Frame
377 et al, 2001), which in turn has been shown to be involved in the early synaptic dysfunction in a
378 mouse model of AD (D'Amelio et al, 2011). It has also been shown that inhibition of PDE5A
379 results in a decrease in the transcription of Wnt/ β -catenin (Tinsley et al, 2011). A reduction in
380 Wnt signalling has been implicated in the amyloid beta-dependent neurodegeneration in LOAD
381 (Inestrosa & Toledo, 2008).

382 While BCM has been applied to low dimensional synthetic data with good results
383 (Visweswaran & Wong, 2009), in this paper we have applied it to GWAS datasets. BCM has
384 several advantages. It is computationally more efficient than the widely used MDR

385 (Visweswaran, Wong & Barmada, 2009). Since BCM uses the Bayesian paradigm, the BCM
386 score represents a coherent way to combine knowledge with data. Biological knowledge or
387 results from analyses of earlier studies can be encoded as a prior distribution over the models that
388 can then be used in Equation 1. Use of informative priors is becoming common in the analysis of
389 microarray expression studies, and a similar strategy can be employed for genomic data.

390 A limitation of our study is the use of GWAS datasets related to a single disease,
391 although it is an important disease. In future research, we plan to apply and investigate the utility
392 of BCM on GWAS datasets related to additional diseases. Another limitation is the use just 10
393 previously known LOAD-associated SNPs. In future work, we plan to explore the use of a larger
394 set of known LOAD associated SNPs that will include SNPs with weaker evidence of being
395 associated with LOAD. In addition, we plan to study the effect of excluding the APOE SNPs
396 rs429358 and rs7412 which are present in every SNP pair we examined for biological
397 plausibility. Another limitation is that we did not use informative prior probabilities for encoding
398 prior knowledge from the literature and previous GWASs. BCM can be extended easily to allow
399 the incorporation of informative priors and inclusion of informative priors in the analysis is an
400 interesting area for study.

401 **Conclusion**

402 We applied BCM to two LOAD GWAS datasets to identify pairs of SNPs that in combination
403 have high statistical association with development of LOAD. To reduce the large search space of
404 all possible parts of SNPs in a GWAS dataset we restricted BCM to evaluate those SNP pairs
405 where one of the SNP was drawn from a set of 10 previously known LOAD associated SNPs.
406 Our results identified several SNPs that have biological evidence of being involved in the

407 pathogenesis of LOAD that would not have been identified by univariate analysis alone due to
408 small main effect but were identified in conjunction with another SNP. These results provide
409 support for applying BCM to identify potential genetic variants such as SNPs from high
410 dimensional GWAs datasets.

411

412 **Funding**

413 CSF was supported in part by an award from the Gerondelis Foundation, SV was supported in
414 part by NLM grant HHSN276201000030C, and M. Ilyas Kamboh was supported by National
415 Institutes of Health grants AG030653, AG005133 and AG041718.

416 **Acknowledgements**

417 We thank Mr. Kevin Bui for his help in data preparation and in implementing the BCM
418 algorithm in software.

419

420 **References**

- 421
- 422 Avramopoulos D. 2009. Genetics of Alzheimer's disease: recent advances. *Genome Medicine*
423 1:34.
- 424 Balding DJ. 2006. A tutorial on statistical methods for population association studies. *Nature*
425 *Reviews Genetics* 7:781–791.
- 426 Bernstein HG & Müller M. 1999. The cellular localization of the L-ornithine
427 decarboxylase/polyamine system in normal and diseased central nervous systems. *Progress in*
428 *Neurobiology* 57:485–505.
- 429 Bertram L, Lill CM & Tanzi RE. 2010. The genetics of Alzheimer disease: back to the future.
430 *Neuron* 68:270–281.
- 431 Bertram L, McQueen MB, Mullin K, Blacker D & Tanzi RE. 2007. Systematic meta-analyses of
432 Alzheimer disease genetic association studies: the AlzGene database. *Nature Genetics* 39:17–23.
- 433 Blacker D, Bertram L, Saunders AJ, Moscarillo TJ, Albert MS, Wiener H, Perry RT, Collins JS,
434 Harrell LE, Go RCP et al. 2003. Results of a high-resolution genome screen of 437 Alzheimer's
435 disease families. *Human Molecular Genetics* 12:23–32.
- 436 Carmona-Saez P, Chagoyen M, Tirado F, Carazo JM & Pascual-Montano A. 2007.
437 GENECODIS: a web-based tool for finding significant concurrent annotations in gene lists.
438 *Genome Biology* 8:R3.
- 439 Cordell HJ. 2009. Detecting gene-gene interactions that underlie human diseases. *Nature*
440 *Reviews Genetics* 10:392–404.
- 441 Cramer PE, Cirrito JR, Wesson DW, Lee CYD, Karlo JC, Zinn AE, Casali BT, Restivo JL,
442 Goebel WD, James MJ et al. 2012. ApoE-directed therapeutics rapidly clear β -amyloid and
443 reverse deficits in AD mouse models. *Science* 335:1503–1506.
- 444 Crews L, Adame A, Patrick C, Delaney A, Pham E, Rockenstein E, Hansen L & Masliah E.
445 2010. Increased BMP6 levels in the brains of Alzheimer's disease patients and APP transgenic
446 mice are accompanied by impaired neurogenesis. *The Journal of Neuroscience* 30:12252–12262.
- 447 D'Amelio M, Cavallucci V, Middei S, Marchetti C, Pacioni S, Ferri A, Diamantini A, De Zio D,
448 Carrara P, Battistini L et al. 2011. Caspase-3 triggers early synaptic dysfunction in a mouse
449 model of Alzheimer's disease. *Nature Neuroscience* 14:69–76.
- 450 Dunn CD, Sulis ML, Ferrando AA & Greenwald I. 2010. A conserved tetraspanin subfamily
451 promotes Notch signaling in *Caenorhabditis elegans* and in human cells. *Proceedings of the*
452 *National Academy of Sciences of the United States of America* 107:5907–5912.

- 453 Frame M, Wan KF, Tate R, Vandenabeele P & Pyne NJ. 2001. The gamma subunit of the rod
454 photoreceptor cGMP phosphodiesterase can modulate the proteolysis of two cGMP binding
455 cGMP-specific phosphodiesterases (PDE6 and PDE5) by caspase-3. *Cellular Signalling* 13:735–
456 741.
- 457 Frykman S, Teranishi Y, Hur J.-Y, Sandebring A, Goto Yamamoto N, Ancarcrona M, Nishimura
458 T, Winblad B, Bogdanovic N, Schedin-Weiss S et al. 2012. Identification of two novel synaptic
459 γ -secretase associated proteins that affect amyloid β -peptide levels without altering Notch
460 processing. *Neurochemistry International* 61:108–118.
- 461 Goedert M & Spillantini MG. 2006. A century of Alzheimer's disease. *Science* 314:777–781.
- 462 Hahn LW, Ritchie MD & Moore JH. 2003. Multifactor dimensionality reduction software for
463 detecting gene-gene and gene-environment interactions. *Bioinformatics* 19:376–382.
- 464 Hardy J & Singleton A. 2009. Genomewide association studies and human disease. *New*
465 *England Journal of Medicine* 360:1759–1768.
- 466 Heckerman D, Geiger D & Chickering DM. 1995. Learning Bayesian Networks: The
467 Combination of Knowledge and Statistical Data. *MACHINE LEARNING* 20:197–243.
- 468 Hollingworth P, Harold D, Sims R, Gerrish A, Lambert J.-C, Carrasquillo MM, Abraham R,
469 Hamshere ML, Pahwa JS, Moskvina V et al. 2011. Common variants at ABCA7,
470 MS4A6A/MS4A4E, EPHA1, CD33 and CD2AP are associated with Alzheimer's disease.
471 *Nature Genetics* 43:429–435.
- 472 Holtzman DM, Morris JC & Goate AM. 2011. Alzheimer's disease: the challenge of the second
473 century. *Science Translational Medicine* 3:77sr1.
- 474 Hu X, Pickering E, Liu YC, Hall S, Fournier H, Katz E, Dechairo B, John S, Van Eerdewegh P
475 & Soares H. 2011. Meta-analysis for genome-wide association study identifies multiple variants
476 at the BIN1 locus associated with late-onset Alzheimer's disease. *PLoS One* 6:e16616.
- 477 Inestrosa NC & Toledo EM. 2008. The role of Wnt signaling in neuronal dysfunction in
478 Alzheimer's Disease. *Molecular neurodegeneration* 3:9.
- 479 Ioannidis JPA, Boffetta P, Little J, O'Brien TR, Uitterlinden AG, Vineis P, Balding DJ,
480 Chokkalingam A, Dolan SM, Flanders WD et al. 2008. Assessment of cumulative evidence on
481 genetic associations: interim guidelines. *International Journal of Epidemiology* 37:120–132.
- 482 Jiang X, Barmada MM & Visweswaran S. 2010. Identifying genetic interactions in genome-wide
483 data using Bayesian networks. *Genetic Epidemiology* 34:575–581.
- 484 Jiang X, Neapolitan RE, Barmada MM & Visweswaran S. 2011. Learning genetic epistasis using
485 Bayesian network scoring criteria. *BMC Bioinformatics* 12:89.

- 486 Johnson AD, Handsaker RE, Pulit SL, Nizzari MM, O'Donnell CJ & de Bakker PIW. 2008.
487 SNAP: a web-based tool for identification and annotation of proxy SNPs using HapMap.
488 *Bioinformatics* 24:2938–2939.
- 489 Kamboh MI, Demirci FY, Wang X, Minster RL, Carrasquillo MM, Pankratz VS, Younkin SG,
490 Saykin AJ, Jun G, Baldwin C et al. 2012. Genome-wide association study of Alzheimer's
491 disease. *Translational Psychiatry* 2:e117.
- 492 Khatri P & Drăghici S. 2005. Ontological analysis of gene expression data: current tools,
493 limitations, and open problems. *Bioinformatics* 21:3587–3595.
- 494 Lambert J-C, Grenier-Boley B, Chouraki V, Heath S, Zelenika D, Fievet N, Hannequin D,
495 Pasquier F, Hanon O, Brice A et al. 2010. Implication of the immune system in Alzheimer's
496 disease: evidence from genome-wide pathway analysis. *Journal of Alzheimer's Disease*
497 20:1107–1118.
- 498 Lipinski MM, Zheng B, Lu T, Yan Z, Py BF, Ng A, Xavier RJ, Li C, Yankner BA, Scherzer CR
499 et al. 2010. Genome-wide analysis reveals mechanisms modulating autophagy in normal brain
500 aging and in Alzheimer's disease. *Proceedings of the National Academy of Sciences of the*
501 *United States of America* 107:14164–14169.
- 502 Liu F, Arias-Vásquez A, Sleegers K, Aulchenko YS, Kayser M, Sanchez-Juan P, Feng B-J,
503 Bertoli-Avella AM, van Swieten J, Axenovich TI et al. 2007. A genomewide screen for late-
504 onset Alzheimer disease in a genetically isolated Dutch population. *American Journal of Human*
505 *Genetics* 81:17–31.
- 506 Lourenço FC, Galvan V, Fombonne J, Corset V, Llambi F, Müller U, Bredesen DE & Mehlen P.
507 2009. Netrin-1 interacts with amyloid precursor protein and regulates amyloid-beta production.
508 *Cell Death & Differentiation* 16:655–663.
- 509 Lukk M, Kapushesky M, Nikkilä J, Parkinson H, Goncalves A, Huber W, Ukkonen E & Brazma
510 A. 2010. A global map of human gene expression. *Nature Biotechnology* 28:322–324.
- 511 Maglott D, Ostell J, Pruitt KD & Tatusova T. 2011. Entrez Gene: gene-centered information at
512 NCBI. *Nucleic Acids Research* 39(Database issue):D52–57.
- 513 Moore JH, Gilbert JC, Tsai C-T, Chiang F-T, Holden T, Barney N & White B. C. 2006. A
514 flexible computational framework for detecting, characterizing, and interpreting statistical
515 patterns of epistasis in genetic studies of human disease susceptibility. *Journal of Theoretical*
516 *Biology* 241:252–261.
- 517 Moore J & White B. 2007. Tuning ReliefF for genome-wide genetic analysis. *Evolutionary*
518 *computation, machine learning and data mining in bioinformatics* 4447:166–175.
- 519 Morgan K. 2011. The three new pathways leading to Alzheimer's disease. *Neuropathology and*
520 *Applied Neurobiology* 37:353–357.

- 521 Müller M, Cárdenas C, Mei L, Cheung K-H & Foskett J. K. 2011. Constitutive cAMP response
522 element binding protein (CREB) activation by Alzheimer's disease presenilin-driven inositol
523 trisphosphate receptor (InsP3R) Ca²⁺ signaling. *Proceedings of the National Academy of*
524 *Sciences of the United States of America* 108:13293–13298.
- 525 Nilsson T, Bogdanovic N, Volkman I, Winblad B, Folkesson R & Benedikz E. 2006. Altered
526 subcellular localization of ornithine decarboxylase in Alzheimer's disease brain. *Biochemical*
527 *and Biophysical Research Communications* 344:640–646.
- 528 Nogales-Cadenas R, Carmona-Saez P, Vazquez M, Vicente C, Yang X, Tirado F, Carazo JM &
529 Pascual-Montano A. 2009. GeneCodis: interpreting gene lists through enrichment analysis and
530 integration of diverse biological information. *Nucleic Acids Research* 37(Web Server
531 issue):W317–322.
- 532 Parkinson H, Sarkans U, Kolesnikov N, Abeygunawardena N, Burdett T, Dylag M, Emam I,
533 Farne A, Hastings E, Holloway E et al. 2011. ArrayExpress update--an archive of microarray and
534 high-throughput sequencing-based functional genomics experiments. *Nucleic Acids Research*
535 39(Database issue):D1002–1004.
- 536 Perry T & Greig NH. 2005. Enhancing central nervous system endogenous GLP-1 receptor
537 pathways for intervention in Alzheimer's disease. *Current Alzheimer Research* 2:377–385.
- 538 Pugazhenti S, Wang M, Pham S, Sze C-I & Eckman CB. 2011. Downregulation of CREB
539 expression in Alzheimer's brain and in A β -treated rat hippocampal neurons. *Molecular*
540 *Neurodegeneration* 6:60.
- 541 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar P, de
542 Bakker PIW, Daly MJ et al. 2007. PLINK: A Tool Set for Whole-Genome Association and
543 Population-Based Linkage Analyses. *American Journal of Human Genetics* 81:559–575.
- 544 Reiman EM, Webster JA, Myers AJ, Hardy J, Dunckley T, Zismann VL, Joshipura KD, Pearson
545 JV, Hu-Lince D, Huentelman MJ et al. 2007. GAB2 alleles modify Alzheimer's risk in APOE
546 epsilon4 carriers. *Neuron* 54:713–720.
- 547 Riederer BM, Leuba G, Vernay A & Riederer IM. 2011. The role of the ubiquitin proteasome
548 system in Alzheimer's disease. *Experimental Biology and Medicine (Maywood)* 236:268–276.
- 549 Ritchie MD, Hahn LW, Roodi N, Bailey LR, Dupont WD, Parl FF & Moore JH. 2001.
550 Multifactor-dimensionality reduction reveals high-order interactions among estrogen-metabolism
551 genes in sporadic breast cancer. *American Journal of Human Genetics* 69:138–417.
- 552 Ritchie MD. 2011. Using biological knowledge to uncover the mystery in the search for epistasis
553 in genome-wide association studies. *Annals of Human Genetics* 75:172–182.
- 554 Saccone SF, Quan J & Jones PL. 2012. BioQ: tracing experimental origins in public genomic
555 databases using a novel data provenance model. *Bioinformatics* 28:1189–1191.

- 556 Supek F, Bošnjak M, Škunca N & Šmuc T. 2011. REVIGO summarizes and visualizes long lists
557 of gene ontology terms. *PLoS One* 6:e21800.
- 558 Tesseur I, Zou K, Esposito L, Bard F, Berber E, Can JV, Lin AH, Crews L, Tremblay P,
559 Mathews P et al. 2006. Deficiency in neuronal TGF-beta signaling promotes neurodegeneration
560 and Alzheimer's pathology. *Journal of Clinical Investigation* 116:3060–3069.
- 561 The UniProt Consortium. 2013. Update on activities at the Universal Protein Resource (UniProt)
562 in 2013. *Nucleic Acids Research* 41:D43–47.
- 563 Thornton-Wells TA, Moore JH & Haines JL. 2004. Genetics, statistics and human disease:
564 analytical retooling for complexity. *Trends in Genetics* 20:640–647.
- 565 Tinsley HN, Gary BD, Keeton AB, Lu W, Li Y & Piazza GA. 2011. Inhibition of PDE5 by
566 sulindac sulfide selectively induces apoptosis and attenuates oncogenic Wnt/ β -catenin-mediated
567 transcription in human breast tumor cells. *Cancer Prevention Research (Philadelphia)* 4:1275–
568 1284.
- 569 Tzschach A, Bisgaard A-M, Kirchhoff M, Graul-Neumann LM, Neitzel H, Page S, Ahmed A,
570 Müller I, Erdogan F, Ropers H-H et al. 2010. Chromosome aberrations involving 10q22: report
571 of three overlapping interstitial deletions and a balanced translocation disrupting C10orf11.
572 *European Journal of Human Genetics* 18:291–295.
- 573 Visweswaran S & Wong A-KI. 2009. Bayesian combinatorial partitioning for detecting
574 interactions among genetic variants. *Summit on Translational Bioinformatics* 2009:133.
- 575 Visweswaran S, Wong A-KI & Barmada MM. 2009. A Bayesian method for identifying genetic
576 interactions. *AMIA Annual Symposium Proceedings* 2009:673–677.
- 577 Wan X, Yang C, Yang Q, Xue H, Fan X, Tang NLS & Yu W. 2010. BOOST: A fast approach to
578 detecting gene-gene interactions in genome-wide case-control studies. *American Journal of*
579 *Human Genetics* 87:325–340.
- 580 Wan X, Yang C, Yang Q, Xue H, Tang NLS & Yu W. 2010. Predictive rule inference for
581 epistatic interaction detection in genome-wide association studies. *Bioinformatics* 26:30–37.
- 582 Wayman GA, Lee Y-S, Tokumitsu H, Silva AJ & Soderling TR. 2008. Calmodulin-kinases:
583 modulators of neuronal development and plasticity. *Neuron* 59:914–931.
- 584 Webster JA, Gibbs JR, Clarke J, Ray M, Zhang W, Holmans P, Rohrer K, Zhao A, Marlowe L,
585 Kaleem M et al. 2009. Genetic control of human brain transcript expression in Alzheimer
586 disease. *American Journal of Human Genetics* 84:445–458.
- 587 Wijsman E. M, Pankratz ND, Choi Y, Rothstein JH, Faber KM, Cheng R, Lee JH, Bird TD,
588 Bennett DA, Diaz-Arrastia R et al. 2011. Genome-wide association of familial late-onset
589 Alzheimer's disease replicates BIN1 and CLU and nominates CUGBP2 in interaction with
590 APOE. *PLoS Genetics* 7:e1001308.

591 Yamauchi T. 2005. Neuronal Ca²⁺/calmodulin-dependent protein kinase II--discovery, progress
592 in a quarter of a century, and perspective: implication for learning and memory. *Biological &*
593 *Pharmaceutical Bulletin* 28:1342–1354.

594 Yang C, He Z, Wan X, Yang Q, Xue H & Yu, W. 2009. SNPHarvester: a filtering-based
595 approach for detecting epistatic interactions in genome-wide association studies. *Bioinformatics*
596 25:504–511.

597

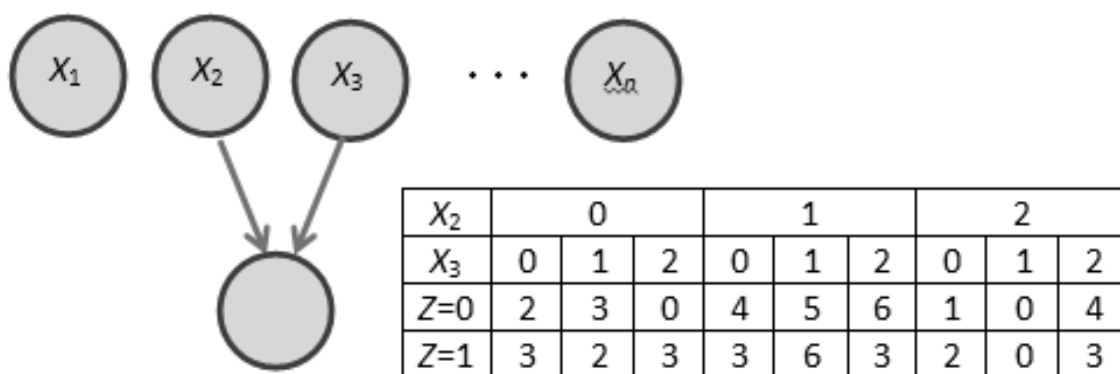
598

599 **Figures**

600

601 **Figure 1.** A SNP-BN model where SNPs X_2 and X_3 have an effect on Z and the remaining SNPs
 602 do not have an effect on Z . The table gives counts for the states of Z conditioned on the joint
 603 states of X_2 and X_3 .

604



605

606

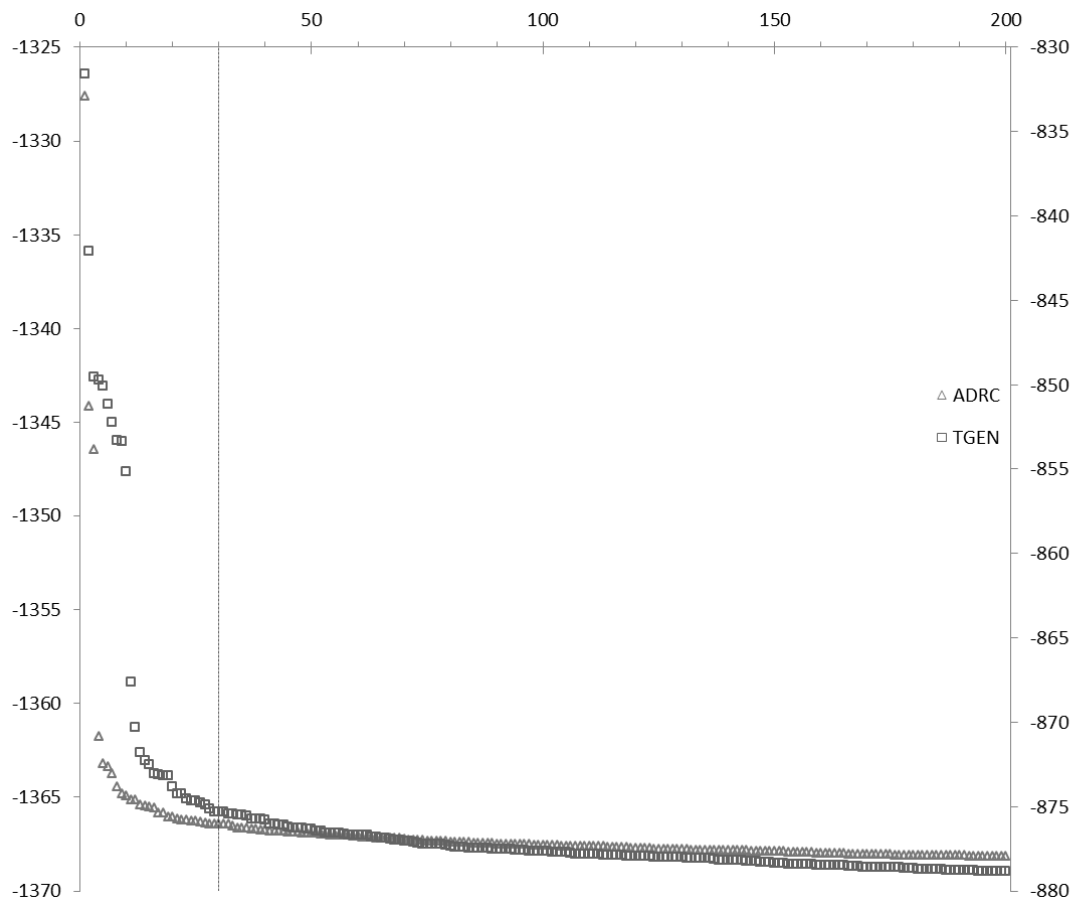
607 **Figure 2.** Plots of the distribution of BCM model scores for the top ranked 200 SNP-BN models
 608 for the two datasets, ADRC and TGen. The scores for the ADRC dataset (blue points)
 609 correspond to the left hand Y axis, while those for the TGen dataset correspond to the right hand
 610 Y axis. The dotted vertical line marks the top ranked 200 SNP-BN models.

611

612

613

614



615 **Tables**

616

617 **Table 1** Previously identified LOAD SNPs.

618

#	Gene	AlzGene SNP	Odds ratio (95% CI)	p value	ADRC SNP	r ²	TGen SNP	r ²
1	APOE 4	rs429358	3.685 (3.30-4.12)	<1E-50	Same	-	Same	-
2	CR1	rs3818361	1.174 (1.14-1.21)	4.72E-21	Same	-	rs6656401	0.840
3	PICALM	rs3851179	0.879 (0.86-0.9)	2.85E-20	Same	-	rs7110631	0.841
4	MS4A6A	rs610932	0.904 (0.88-0.93)	1.81E-11	Same	-	rs574695	0.935
5	CD33	rs3865444	0.893 (0.86-0.93)	2.04E-10	Same	-	Same	-
6	MS4A4E	rs670139	1.079 (1.05-1.11)	9.51E-10	rs600550	1	rs676309	1
7	CD2AP	rs9349407	1.117 (1.08-1.16)	2.75E-09	rs9296559	1	rs9296558	1
8	GAB2	rs2373115	0.85 (0.76-0.94)		Same	-	Same	-
9	SORL1	rs2282649	1.10 (1.03-1.17)		rs726601	0.922	rs726601	0.922
10	TF	rs1049296	1.18 (1.06-1.31)		Same	-	Same	-

619

620

621

622

623

624

AlzGene SNP: the SNP in the AlzGene meta-analysis, along with the relevant odds ratios and p values (the latter for those SNPs with p values <0.00001); *ADRC SNP*: the corresponding SNP in the ADRC dataset, along with the r² scores; *TGen SNP*: the corresponding SNP in the TGen dataset, along with the r² scores for linkage disequilibrium.

625 **Table 2.** Functional description and expression of genes associated with the top 30 dataset SNPs
 626 in the ADRC dataset.

Gene Symbol (SNP)	Name	Description	Expression in AD	1-SNP model rank	p value of pair from PLINK
APOC1 (rs4420638)	Apolipoprotein C-I	Appears to modulate the interaction of APOE with beta-migrating VLDL. Binds free fatty acids.	Overexpressed (Lukk et al, 2010)	2	0.3326
TOMM40 (rs157582)	Translocase of outer mitochondrial membrane 40 homolog	Channel-forming subunit of the translocase of the mitochondrial outer membrane (TOM) complex, essential for protein import into mitochondria.	Underexpressed (Lukk et al, 2010)	3	0.4139
APOE (rs7412)	Apolipoprotein E	ApoE is essential for the normal catabolism of triglyceride-rich lipoprotein constituents. Known risk factor for LOAD.	Overexpressed (Lukk et al, 2010)	5	0.8172
SNTG1 (rs16914489)	Gamma-1-syntrophin	Specifically expressed in the brain, highly expressed in the cortex. Organizes the subcellular localization of a variety of proteins.	Overexpressed (Lukk et al, 2010)	24906	0.004546
TMEM217 (rs9470543)	Transmembrane protein 217	Expressed in the brain	-	4584	0.004643
SMAD6 (rs3934907)	Mothers against DPP homolog 6	Negative regulation of BMP and TGF-beta/activin-signaling. BMP-6 is increased in AD brains and leads to impaired neurogenesis (Crews et al, 2010). Reduced TGF-beta signaling is involved in neurodegeneration and promotes AD like changes in mice (Tesseur et al, 2006).	Underexpressed (Lukk et al, 2010)	41282	0.0000998
NPAS3 (rs4981180)	Neuronal PAS domain protein 3.	Transcription factor. May regulate genes involved in neurogenesis. Associated with schizophrenia and mental retardation	Overexpressed (Lukk et al, 2010)	1086	0.1225
NTM (rs11222692)	Neurotrimin	May promote neurite outgrowth and adhesion. NTM lies at locus 11q25, which has been associated	Overexpressed (Lukk et al, 2010)	12209	0.1422

		with AD (Blacker et al, 2003; Liu et al, 2007).			
PPAPDC1A (rs4752432)	Phosphatidic acid phosphatase type 2 domain containing 1A	-	-	6852	0.3963
NPFF (rs8192593)	Neuropeptide FF-amide peptide precursor	Modulation of morphine-induced antinociception.	-	3981	0.1251
SLC25A21 (rs7140725)	Solute carrier family 25	Known also as ornithine decarboxylase (ODC). Mitochondrial oxoadipate carrier, part of polyamine synthesis pathway.	Overexpressed (Bernstein & Müller, 1999; Nilsson et al, 2006)	444	0.06767
RAB23 (rs182662)	Member RAS oncogene family	Intracellular protein transportation. Regulated by miRNA155, which also regulates PICALM (a known AD association).	Underexpressed (Lukk et al, 2010)	96	0.08251
UNC5D (rs4577954)	unc-5 homolog D (C. elegans)	Netrin receptor: netrins are secreted proteins that direct axon extension and cell migration during neural development. APP also binds Netrin-1 and in transgenic mice this suppresses amyloid beta peptide production (Lourenço et al, 2009).	-	63972	0.6731
CHD9 (rs3852742)	Chromodomain helicase DNA binding protein 9, PPARA -interacting complex 320 kDa protein	Transcriptional co-activator for PPARA. The APOE gene promoter has a binding site for PPAR alpha. Low CHD9 activity could reduce apoE levels. Increase in APOE transcription has been shown to clear amyloid beta in AD mouse models (Cramer et al, 2012).	Overexpressed (Lukk et al, 2010)	1061	0.04696
CNTN4 (rs9819935)	Contactin 4, Brain-derived immunoglobulin superfamily protein 2	Mainly expressed in brain. Neuronal membrane protein that may play a role in the formation of axon connections in the developing nervous system. Associated with	-	2149	0.002386

		Spinocerebellar Ataxia, Amyotrophic Lateral Sclerosis, 3p deletion syndrome.			
--	--	---	--	--	--

627

628 *1-SNP model rank*: rank of the corresponding SNP in terms of univariate 1-SNP model score

629

630 **Table 3.** Functional description and expression of genes associated with the top 30 dataset SNPs

631 in the TGen dataset.

Gene Symbol (SNP)	Name	Description	Differential Expression in AD	1-SNP model rank	p value of pair from PLINK
APOE 2 (rs7412)	Apolipoprotein E	ApoE is essential for the normal catabolism of triglyceride-rich lipoprotein constituents. Known risk factor for LOAD.	Overexpressed (Lukk et al, 2010)	1	0.05993
APOC1 (rs4420638)	Apolipoprotein C-1	Appears to modulate the interaction of APOE with beta-migrating VLDL. Binds free fatty acids.	Overexpressed (Lukk et al, 2010)	3	0.705
C10orf11 (rs7079348)	Chromosome 10 open reading frame 11	A brain-expressed gene. Haploinsufficiency of C10orf11 contributes to the cognitive defects in 10q22 syndrome (Tzschach et al, 2010).	-	4	0.009623
VWC2 (rs10499687)	von Willebrand factor C domain-containing protein 2 (Brorin, Brain-specific chordin-like protein)	Encodes a secreted bone morphogenic protein (BMP) antagonist. The encoded protein is possibly involved in neural function and development and may have a role in cell adhesion. BMP-6 is increased in AD brains and leads to impaired neurogenesis (Crews et al, 2010).	Underexpressed (Webster et al, 2009)	12	0.7698
PSD3 (rs17126808)	Pleckstrin and Sec7 domain containing 3	Guanine nucleotide exchange factor for ARF6 that contributes to the regulation of dendritic branching (The UniProt Consortium, 2013).	Overexpressed (Lukk et al, 2010)	34	0.001623
GXYLT2 (rs3732443)	Glucoside xylosyltransferase 2	Elongates the O-linked glucose attached to EGF-like repeats in the extracellular domain of Notch proteins (The UniProt Consortium, 2013), which are substrates of γ -secretase, the enzyme involved in amyloid beta production (Frykman et al,	Underexpressed in a murine AD model (D'Amelio et al, 2011)	6	0.211

		2012).			
GABBR2 (rs2779550)	Gamma-aminobutyric acid (GABA) B receptor, 2	Target for autophagy regulation in neurodegenerative diseases (Lipinski et al, 2010).	Overexpressed (Lukk et al, 2010)	391	0.0002945
ENPP2 (rs16892852)	Ectonucleotide pyrophosphatase/phosphodiesterase 2	Hydrolyzes lysophospholipids to produce lysophosphatidic acid (LPA) in extracellular fluids. Predominantly expressed in brain, placenta, ovary, and small intestine. Secreted by most body fluids including serum and cerebrospinal fluid (The UniProt Consortium, 2013).	Overexpressed (Lukk et al, 2010)	92	0.04851
GLP1R (rs910171)	Glucagon-like peptide 1 receptor	Member of the glucagon receptor family (also includes glucagon, GLP-2, secretin, GHRH and GIP receptors). In the brain located in hypothalamus and brainstem. Protective against amyloid beta accumulation in rats (Perry & Greig, 2005).	Overexpressed (Lukk et al, 2010)	193	0.01462
MOSC1 (rs746767)	MOCO sulphurase C-terminal domain containing 1	A mitochondrial oxidoreductase, cofactor: molybdenum, is expressed in the brain. MOSC1 is a target for miR-129-5p, like GABBR2, and miR-155, like PICALM.	-	66	0.04507
TM4SF20 (rs4408717)	Transmembrane 4 L six family member 20	Tetraspanin superfamily member. Tetraspanins are often thought to act as scaffolding proteins, anchoring multiple proteins to one area of the cell membrane. Other tetraspanin superfamily members have been implicated in Notch signaling and g-secretase activity modulation (Dunn et al, 2010).	-	95	0.004495

632

633 *1-SNP model rank*: rank of the corresponding SNP in terms of univariate 1-SNP model score