

# 1 PRS-Net: Interpretable polygenic risk scores via geometric learning

2 Han Li<sup>1</sup>, Jianyang Zeng<sup>2,\*</sup>, Michael P. Snyder<sup>3,\*</sup>, and Sai Zhang<sup>4,5,6,\*</sup>

3 <sup>1</sup> Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing, China

4 <sup>2</sup> School of Engineering, Westlake University, Hangzhou, Zhejiang, China

5 <sup>3</sup> Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

6 <sup>4</sup> Department of Epidemiology, University of Florida, Gainesville, FL, USA

7 <sup>5</sup> J. Crayton Pruitt Family Department of Biomedical Engineering, University of Florida, Gainesville, FL, USA

8 <sup>6</sup> The Genetics Institute, University of Florida, Gainesville, FL, USA

9 \* Correspondence: zengjy@westlake.edu.cn; mpsnyder@stanford.edu; sai.zhang@ufl.edu

10 **Abstract.** Polygenic risk score (PRS) serves as a valuable tool for predicting the  
11 genetic risk of complex human diseases for individuals, playing a pivotal role in ad-  
12 vancing precision medicine. Traditional PRS methods, predominantly following a linear  
13 structure, often fall short in capturing the intricate relationships between genotype and  
14 phenotype. We present PRS-Net, an interpretable deep learning-based framework de-  
15 signed to effectively model the nonlinearity of biological systems for enhanced disease  
16 prediction and biological discovery. PRS-Net begins by deconvoluting the genome-  
17 wide PRS at the single-gene resolution, and then it encapsulates gene-gene interac-  
18 tions for genetic risk prediction leveraging a graph neural network, thereby enabling  
19 the characterization of biological nonlinearity underlying complex diseases. An atten-  
20 tive readout module is specifically introduced into the framework to facilitate model in-  
21 terpretation and biological discovery. Through extensive tests across multiple complex  
22 diseases, PRS-Net consistently outperforms baseline PRS methods, showcasing its  
23 superior performance on disease prediction. Moreover, the interpretability of PRS-Net  
24 has been demonstrated by the identification of genes and gene-gene interactions that  
25 significantly influence the risk of Alzheimer's disease and multiple sclerosis. In sum-  
26 mary, PRS-Net provides a potent tool for parallel genetic risk prediction and biological  
27 discovery for complex diseases.

## 28 1 Introduction

29 Complex human diseases display polygenicity in their genetic architectures, characterized by a  
30 multitude of common genetic variants with minor individual effects accumulatively influencing the  
31 disease risk<sup>1-4</sup>. The polygenic risk scores (PRSs) are developed to quantitatively characterize the  
32 genetic susceptibility of individuals to specific traits or complex diseases based on the common  
33 genetic variants<sup>5-7</sup>. This methodology empowers the early deployment of targeted therapeutic  
34 interventions and facilitates the practice of personalized medicine<sup>8-10</sup>.

35 PRS is typically calculated using the summary statistics derived from genome-wide association  
36 studies (GWAS)<sup>11-17</sup>, a widely-used statistical method for identifying disease-associated genetic  
37 variants<sup>18-20</sup>. While GWAS can identify disease risk genetic variants, such as single nucleotide  
38 polymorphisms (SNPs), that exhibit significant differences in frequencies between cases and con-  
39 trols, these variants tend to have modest individual effects on the phenotype, resulting in limited  
40 prediction capability. In an effort to enhance predictive modeling, various statistical methods have  
41 been applied to aggregate the effects of individual SNPs. The widely adopted method for calculat-  
42 ing PRS, exemplified by PLINK<sup>21</sup> and PRSice<sup>12</sup>, is known as clumping and thresholding (C+T)<sup>11</sup>,  
43 which involves summing allele counts weighted by effect sizes estimated from GWAS. More recent  
44 approaches like LDpred2<sup>16</sup> utilize Bayesian modeling to infer the posterior mean effect size of each  
45 marker by incorporating prior information on effect sizes and linkage disequilibrium (LD) data from  
46 an external reference panel. Similarly, lassosum2<sup>17</sup> estimates PRS using summary statistics and  
47 a reference panel within a penalized regression framework. With the notable increase in dataset  
48 sample sizes for GWAS, these methods have achieved enhanced predictive power<sup>22</sup>. Nonethe-  
49 less, these techniques primarily rely on univariate effect sizes derived from linear GWAS models,  
50 thus often overlook potential non-linear associations between genetic factors and phenotypes,  
51 which can undermine their predictive performance.

52 Efforts have also been made to construct models capable of capturing non-linear interac-  
53 tions in PRS calculation. These include tree-based methods like random forests<sup>23,24</sup>, gradient  
54 boosting<sup>25,26</sup>, and AdaBoost<sup>27,28</sup>, as well as deep learning-based techniques such as multiple-  
55 layer perceptrons (MLP)<sup>29</sup> and convolutional neural networks<sup>30</sup>. However, these methods only  
56 take a limited number of variants as their input, and lack the integration of versatile prior biological  
57 knowledge. Indeed, these approaches have demonstrated either comparable or, in many cases,  
58 less effective performance in predicting phenotypes when compared to linear models<sup>31,32</sup>.

59 In this study, we propose PRS-Net, a geometric deep learning-based approach designed to  
60 effectively model the intricate non-linear relationships among genetic factors such as genes in  
61 predicting the disease risk, thus delivering more accurate and robust PRSs. Based on the sum-  
62 mary statistics of GWAS, PRS-Net first maps PRS onto a gene-gene interaction (GGI) network  
63 through the derivation of gene-level PRSs using the C+T method. Subsequently, a graph neural  
64 network is employed to iteratively update the embedding of the genes via performing message  
65 passing on the GGI network, thus capturing the complex GGIs from the network. An attentive  
66 readout module is then introduced to provide interpretable PRS predictions. PRS-Net also inte-  
67 grates a mixture-of-expert module<sup>33</sup> designed to enhance the accuracy of PRS predictions when  
68 dealing with multi-ancestry datasets. Our comprehensive evaluation encompasses six complex  
69 diseases extracted from the UK Biobank database<sup>34</sup>, including Alzheimer's disease, atrial fib-  
70 rillation, rheumatoid arthritis, multiple sclerosis, ulcerative colitis, and asthma. The results con-  
71 sistently demonstrated the superiority of PRS-Net over baseline methods, including PLINK<sup>21</sup>,  
72 PRSice<sup>14</sup>, LDpred-2<sup>16</sup>, and lassosum2<sup>17</sup> in PRS prediction. Notably, through case studies fo-  
73 cused on Alzheimer's disease and multiple sclerosis, we illustrated that PRS-Net provided biolog-  
74 ically meaningful interpretability by identifying specific genes and GGIs that significantly influence

2 H. Li et al.

75 disease risk. In summary, PRS-Net stands as a potent and innovative tool for precise PRS pre-  
 76 diction, addressing the limitations of current linear models and offering a more comprehensive  
 77 approach to unraveling the genetic underpinnings of complex traits and diseases.

78 **2 Method**

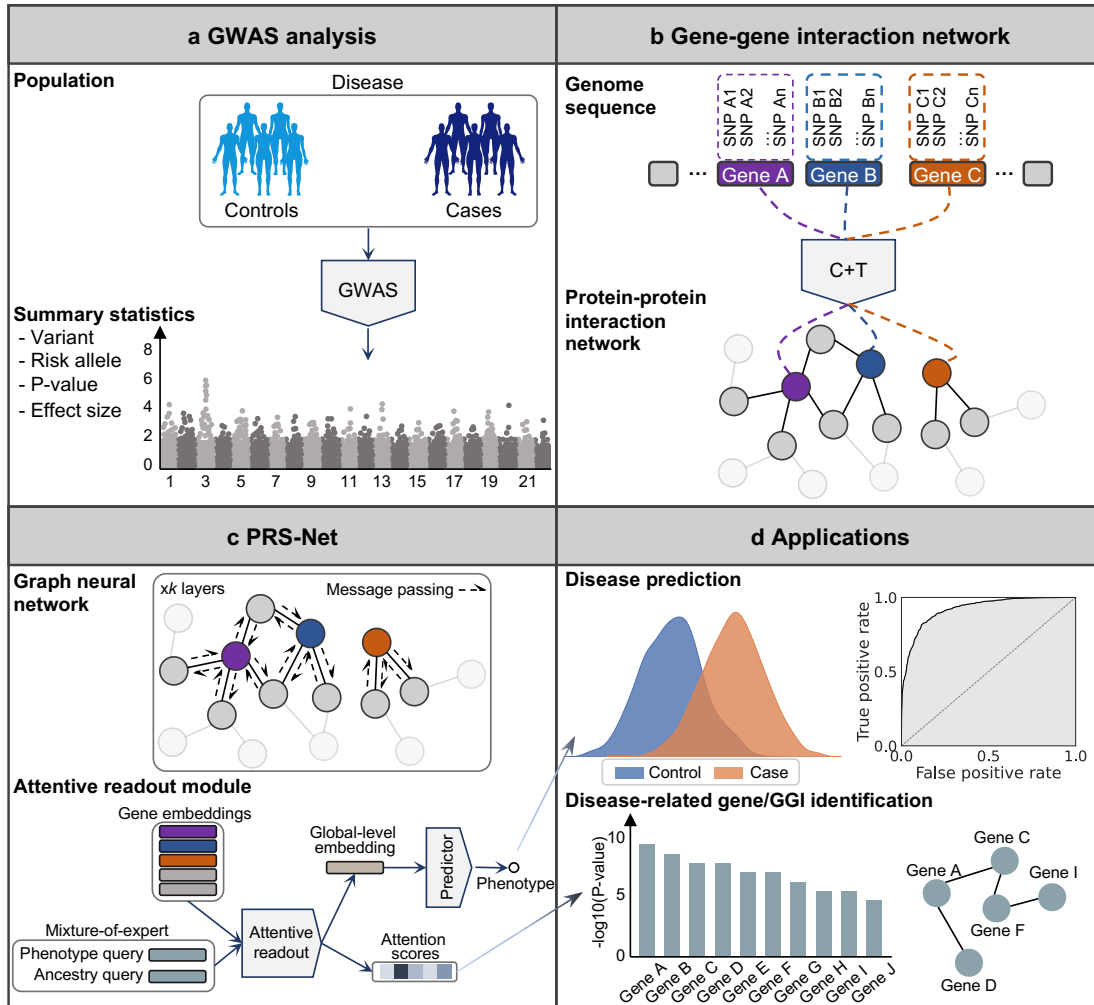


Fig. 1: An illustrative diagram of PRS-Net. **a** The proposed framework is based on summary statistics, including variants, risk alleles, P-values, and effect sizes derived from GWAS. **b** A gene-gene interaction network is constructed based on the protein-protein interaction network. Gene-level PRSs are calculated with the C+T method to serve as the node features for the nodes within the network. **c** A graph neural network is employed to update node features via message passing and subsequently an attentive readout module is applied to provide interpretable PRS predictions. **d** The PRS-Net can be applied for disease prediction and disease-related gene/GGI identification.

79 In this section, we present our proposed framework for PRS estimation (Fig. 1), covering the  
 80 establishment of the GGI network, the derivation of gene-level PRS, and the architecture of PRS-  
 81 Net.

## 82 2.1 GGI network

83 It is widely recognized that the disease phenotype is not solely determined by individual genes  
84 but rather involves the intricate interactions among multiple genes, which can exhibit additive or  
85 non-additive genetic relationships<sup>35–37</sup>. Additive genetic interactions manifest when the cumulative  
86 effects of genes jointly contribute to a specific phenotype. Furthermore, there are increasing stud-  
87 ies highlighting the significance of non-additive genetic interactions<sup>38–40</sup>. Epistasis is a prominent  
88 example of non-additive genetic interaction, which occurs when the impact of a gene mutation  
89 depends on the presence or absence of mutations in one or more other genes<sup>41–43</sup>. We estab-  
90 lish a GGI network that empowers PRS-Net to capture the intricate genetic relationships that are  
91 potentially associated with the target phenotypes (Fig. 1b).

92 We construct our GGI network based on the protein-protein interactions derived from the  
93 STRING database<sup>44</sup>, as protein-protein interactions represent potent indicators of functional rela-  
94 tionships between genes. Formally, we construct a GGI network, denoted as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  
95  $\mathcal{V}$  stands for the set of nodes and  $\mathcal{E}$  stands for the set of edges. Each node  $v_i \in \mathcal{V}$  stands for a  
96 coding gene and each edge  $(v_i, v_j) \in \mathcal{E}$  stands for an interaction between nodes  $v_i$  and  $v_j$  derived  
97 from the STRING database<sup>44</sup>. Note that, we add a self-loop  $(v_i, v_i)$  for each node  $v_i \in \mathcal{V}$ . This  
98 network construction results in a GGI network encompassing 19,836 coding genes and 250,236  
99 interactions.

100 Upon deriving the GGI network, we proceed to compute gene-level PRSs for the genes within  
101 the network using a C+T approach<sup>11,21</sup>. More precisely, for each gene in the network, we focus  
102 on the SNPs falling within a designated range, spanning from its transcription start site  $-L$  to  
103 its transcription end site  $+L$ . In our tests, we set  $L$  to 10 kilobases (KB), thereby encompassing  
104 the SNPs situated in non-coding regions, such as the promoters of the genes. Subsequently, for  
105 each gene, we perform LD clumping on the associated SNPs from the GWAS data, utilizing the  
106 LD information estimated in the target data. Following this, we filter the SNPs based on a specific  
107 P-value threshold. The gene-level PRSs are then derived by multiplying the genotype matrix by  
108 the effect sizes obtained from the GWAS data, and then dividing this by the number of allele  
109 observations for each gene. For the LD clumping process, we set the LD threshold  $R^2$  to 0.5 and  
110 the physical distance threshold to 250 KB. As for the thresholding step, we set the P-values to  
111  $1e^{-5}$ ,  $1e^{-4}$ ,  $1e^{-3}$ ,  $1e^{-2}$ ,  $5e^{-2}$ , 0.1, 0.2, 0.3, 0.5, and 1, respectively. This process results in the  
112 computation of eleven PRSs for each gene, which serves as their initial features. We denote the  
113 initial feature of  $v_i \in \mathcal{V}$  as  $\mathbf{h}_i \in \mathbf{H}$ , where  $\mathbf{H} \in \mathbb{R}^{|\mathcal{V}| \times 11}$  and  $|\mathcal{V}|$  stands for the number of genes in  
114  $\mathcal{G}$ .

## 115 2.2 PRS-Net

### 116 Graph neural network

117 We harness the power of a graph neural network to capture the complex interactions among genes  
118 within our established GGI network (Fig. 1c). In our tests, we specifically opt for a graph isomor-  
119 phism network (GIN)<sup>45</sup> due to its proven theoretical and experimental expressiveness. Formally,  
120 we first encode the initial feature of nodes, denoted as  $\mathbf{H}$ , by employing an MLP in the following  
121 manner:

$$\mathbf{H}^0 = \text{MLP}^0(\mathbf{H}), \quad (1)$$

122 where  $\mathbf{H}^0 \in \mathbb{R}^{|\mathcal{V}| \times D}$  and  $D$  is the dimension of hidden states. Subsequently, we apply multiple GIN  
123 layers to iteratively update the representation of each node by aggregating the representations of

124 its neighbors, as depicted below:

$$\mathbf{h}_i^k = \text{MLP}^k((1 + \epsilon^k) \cdot \mathbf{h}_i^{k-1} + \sum_{v_j \in \mathcal{N}(v_i)} \mathbf{h}_j^{k-1}), \quad (2)$$

125 where  $\mathbf{h}_i^{k-1}$  is the hidden states of  $v_i$  at the  $(k-1)$ -th layer,  $\mathcal{N}(v_i)$  stands for the neighbors of  $v_i$  in  
 126 the GGI network,  $\mathbf{h}_i^k$  stands for the updated hidden states of  $v_i$  at the  $k$ -th layer,  $\text{MLP}^k$  is the MLP  
 127 at the  $k$ -th layer, and  $\epsilon$  stands for a learnable variable. Following  $k$  iterations of aggregation, each  
 128 gene effectively encapsulates the interaction information within its  $k$ -hop neighborhood.

### 129 Attentive readout module

130 To make predictions for each data sample, we derive the global-level representation for each  
 131 sample through an attentive readout module, illustrated as follows:

$$\begin{aligned} \mathbf{h}_G &= \text{Attentive readout}(\mathbf{Q}, \mathbf{K}, \mathbf{V}), \\ \mathbf{h}_G &= \mathbf{A} \cdot \mathbf{V}, \\ \mathbf{A} &= \text{Sigmoid}(\mathbf{Q} \cdot \mathbf{K}), \\ \mathbf{K} &= \mathbf{H}^k \cdot \mathbf{W}_K, \mathbf{V} = \mathbf{H}^k \cdot \mathbf{W}_V, \end{aligned} \quad (3)$$

132 where  $\mathbf{W}_K \in \mathbb{R}^{D \times D}$  and  $\mathbf{W}_V \in \mathbb{R}^{D \times D}$  stand for trainable projection matrices to derive the key  
 133 (i.e.,  $\mathbf{K}$ ) and value (i.e.,  $\mathbf{V}$ ) matrices, respectively.  $\mathbf{Q} \in \mathbb{R}^{1 \times D}$  stands for a trainable query vector.  
 134 Sigmoid stands for the sigmoid function.  $\mathbf{A} \in \mathbb{R}^{1 \times |\mathcal{V}|}$  stands for the attention scores, with elevated  
 135 scores signifying a greater significance of the associated genes.  $\mathbf{h}_G \in \mathbb{R}^{1 \times D}$  stands for the global-  
 136 level representation.

137 After deriving the global-level representation, we employ an MLP to derive the final prediction,  
 138 denoted as  $\hat{\text{PRS}}$ , as follows:

$$\hat{\text{PRS}} = \text{MLP}(\mathbf{h}_G). \quad (4)$$

139 Additionally, we implement a mixture-of-expert module<sup>33</sup> to effectively handle datasets that  
 140 encompass data samples from multiple ancestries. More specifically, we introduce a specialized  
 141 attentive readout module for each distinct ancestry. These dedicated attentive readout modules  
 142 are activated when processing data from individuals with specific ancestral origins. To illustrate,  
 143 when dealing with input samples of Western European ancestry, we derive the ancestry-specific  
 144 global-level representation as follows:

$$\mathbf{h}_G^{\text{EUR}} = \text{Attentive readout}(\mathbf{Q}^{\text{EUR}}, \mathbf{K}^{\text{EUR}}, \mathbf{V}^{\text{EUR}}). \quad (5)$$

145 The ancestry-specific readout module is designed to capture the unique knowledge pertaining to  
 146 each ancestry in relation to the disease. In addition, we introduce another shared readout module  
 147 to capture disease-related knowledge that holds general applicability across all ancestries:

$$\mathbf{h}_G^{\text{PH}} = \text{Attentive readout}(\mathbf{Q}^{\text{PH}}, \mathbf{K}^{\text{PH}}, \mathbf{V}^{\text{PH}}). \quad (6)$$

148 Then, we derive the final global-level representation by combining the aforementioned two repre-  
 149 sentations:

$$\mathbf{h}_G = \mathbf{h}_G^{\text{EUR}} + \mathbf{h}_G^{\text{PH}}. \quad (7)$$

150 The process for deriving global-level representations of individuals from other ancestries follows a  
 151 similar approach. The final PRS prediction can be computed with equation 4, utilizing the derived  
 152 global-level representation. We refer to the single-ancestry variation as PRS-Net and the multiple-  
 153 ancestry variation as PRS-Net<sub>MA</sub>.

### 154 3 Results

#### 155 3.1 PRS-Net outperforms baseline methods in PRS prediction

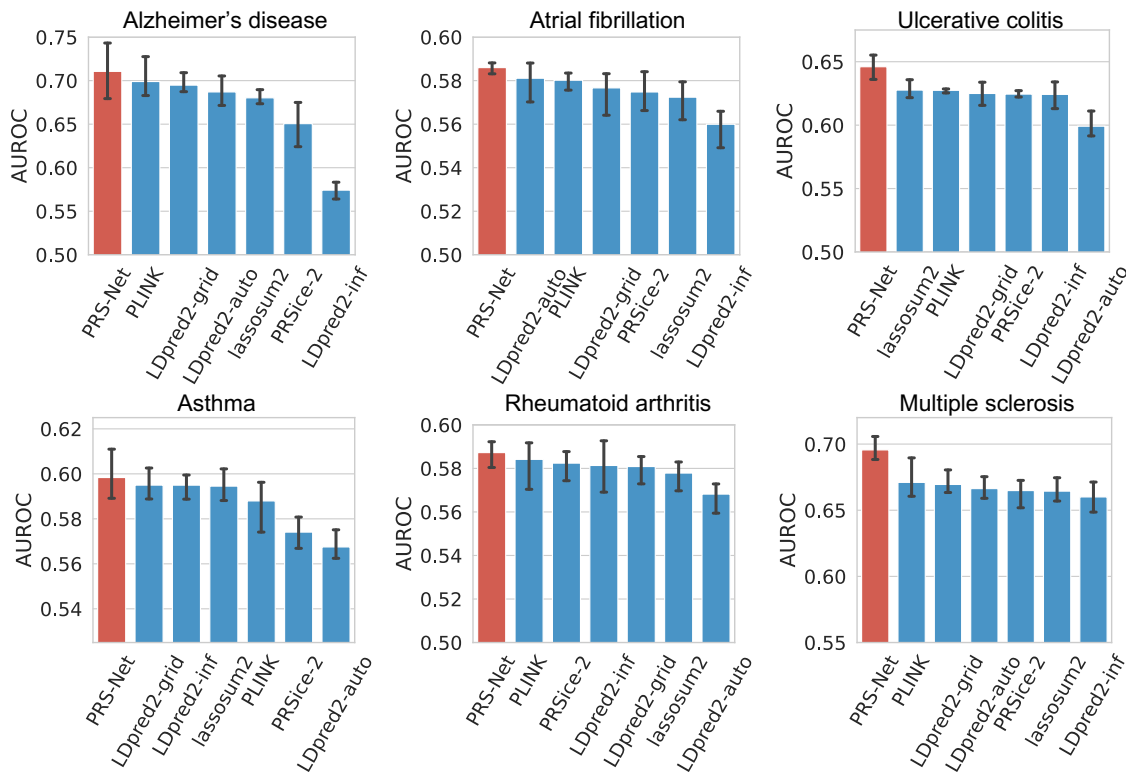


Fig. 2: The PRS prediction performance of PRS-Net compared to baseline methods across a range of complex diseases, including Alzheimer's disease, atrial fibrillation, ulcerative colitis, asthma, rheumatoid arthritis, and multiple sclerosis, measured in terms of the area under the receiver operating characteristic curve (AUROC). The bars are the estimated standard errors.

156 We extracted genotype-phenotype data from the UK Biobank database<sup>34</sup> for six different complex  
157 diseases, which encompassed Alzheimer's disease, atrial fibrillation, rheumatoid arthritis,  
158 multiple sclerosis, ulcerative colitis, and asthma. ICD-10 codes<sup>46</sup> were employed to define the dis-  
159 ease phenotypes (Supplementary Table 1). For our primary experiments, we focused exclusively  
160 on individuals of Western European ancestry due to the insufficient size of the non-European  
161 ancestry population, which did not provide an adequate amount of training data (Supplementary  
162 Table 2). Following a quality control process, each disease dataset consisted of roughly 411,000 in-  
163 dividuals (Supplementary Note 1.1). To prevent data leakage, we ensured that none of the GWAS  
164 were conducted on samples from the UK Biobank database (see Data availability). For each dis-  
165 ease dataset, we randomly partitioned it into training, validation, and test sets with a ratio of 8:1:1.  
166 To evaluate the performance of PRS-Net, we compared it against several previously proposed  
167 methods, such as C+T-based methods (PLINK<sup>21</sup> and PRSice<sup>214</sup>), lassosum<sup>217</sup>, and three vari-  
168 ations of LDpred2<sup>16</sup> (LDpred2-auto, LDpred2-grid, and LDpred2-inf), utilizing the area under the  
169 receiver operating characteristic curve (AUROC) as the metric. To ensure a rigorous and equi-  
170 table comparison, we utilized LD matrices estimated from European populations within the 1000

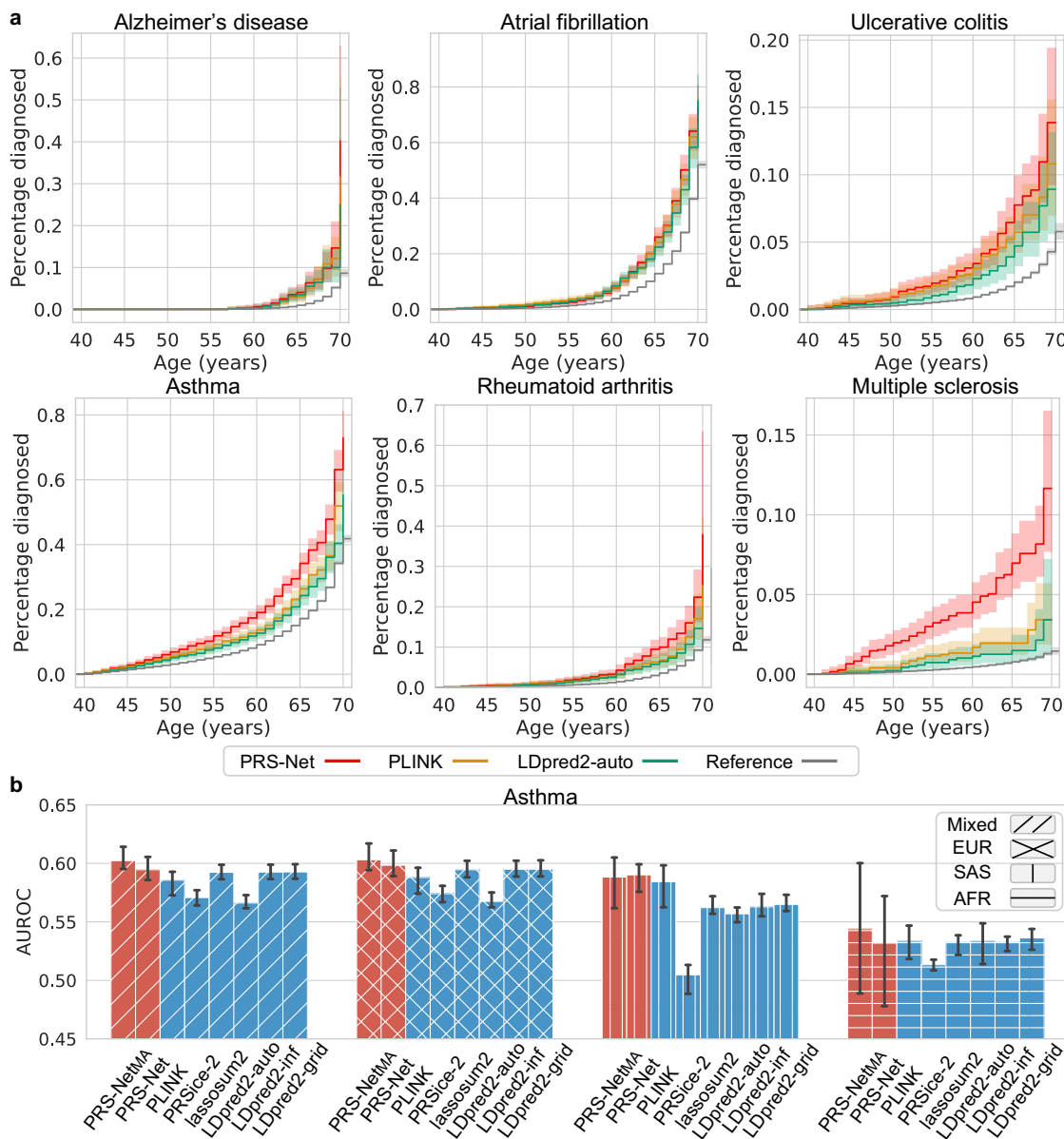
171 Genomes Project<sup>47</sup> across all methods in our study. Our results were based on three indepen-  
172 dent runs with different random seeds to ensure robustness and reliability. The results revealed  
173 that PRS-Net consistently outperformed all baseline methods on all disease datasets, resulting in  
174 relative improvements ranging from 0.5% to 3.7%. Interestingly, the largest improvements were  
175 obtained for two autoimmune diseases, i.e., ulcerative colitis (with a relative improvement of 3.0%)  
176 and multiple sclerosis (with a relative improvement of 3.7%), reinforcing the observed nonadditiv-  
177 ity of genomic factors underlying these diseases<sup>38,48–50</sup>. Altogether, our data demonstrates that  
178 PRS-Net possesses the capacity to capture more intricate associations between genotypes and  
179 phenotypes that are beyond the reach of previously proposed linear models.

180 We utilized the Aalen-Johansen estimator<sup>51</sup> to estimate the disease occurrence over a life-  
181 time for individuals categorized into high-risk and low-risk groups, as determined by the PRSs  
182 estimated by PRS-Net and baseline methods. High-risk individuals were defined as those with  
183 the highest 5% of PRSs, while low-risk individuals were identified as those with the lowest 5%  
184 of PRSs. The cumulative incidence plots revealed that individuals classified as high-risk by PRS-  
185 Net generally exhibited a heightened risk of disease throughout their lifetime compared to base-  
186 line methods, especially for ulcerative colitis, asthma, rheumatoid arthritis, and multiple sclerosis  
187 (Fig. 3a). Conversely, those categorized as low-risk by PRS-Net tended to maintain a lower risk of  
188 all diseases over their lifetime in comparison to baseline methods (Supplementary Fig. 1). These  
189 findings underscore the potential of PRS-Net as a powerful tool for individual risk stratification.

190 Next, we assessed the performance of PRS-Net and our multiple-ancestry model, PRS-Net<sub>MA</sub>,  
191 on a dataset comprising individuals from diverse ancestral backgrounds. Specifically, we curated a  
192 mixed-ancestry dataset encompassing Western European, South Asian, and African for asthma,  
193 which provides a reasonable number of asthma cases (over 1,000) for each ancestry (Supplemen-  
194 tary Table 2). The results revealed that PRS-Net outperformed baseline methods on the mixed  
195 ancestry and South Asian ancestry test sets, indicating that the PRS-Net trained solely on the  
196 Western European ancestry dataset captured the underlying disease biology independent of dif-  
197 ferent ancestries (Fig. 3b). Additionally, PRS-Net<sub>MA</sub> demonstrated superior performance when  
198 compared to PRS-Net on the mixed ancestry, Western European ancestry, and African ancestry  
199 test sets (Fig. 3b). These findings underscored the ability of PRS-Net<sub>MA</sub> to leverage the multi-  
200 ancestry dataset effectively, enhancing its portability in estimating PRS for individuals from diverse  
201 ancestral backgrounds.

### 202 **3.2 PRS-Net identifies disease-related genes and GGIs for Alzheimer’s disease and** 203 **multiple sclerosis**

204 Following the demonstration of the superior performance of PRS-Net in predicting PRS, we sought  
205 to explore its capability to identify risk genes and GGIs underlying complex diseases. Alzheimer’s  
206 disease, a progressively degenerative condition, has been the subject of extensive research for  
207 many years, leading to the identification of numerous genes associated with the disease<sup>52–56</sup>. We  
208 employed PRS-Net to identify disease-related genes and GGIs, with the expectation that our find-  
209 ings would align with prior research outcomes. Specifically, we first applied the Mann–Whitney U  
210 test<sup>57</sup> to each gene within our constructed GGI network, assessing whether the attention scores  
211 associated with the gene for individuals with Alzheimer’s disease were notably higher than those  
212 of the control group. This analysis yielded a gene set comprising 309 genes with compelling sta-  
213 tistical significance (P-value <0.001). Please refer to Supplementary Data 1 for the complete list  
214 of the genes. Subsequently, we conducted gene set enrichment analyses (GSEA)<sup>58</sup> utilizing the  
215 gene ontology (GO)<sup>59</sup> and Kyoto Encyclopedia of Genes and Genomes (KEGG)<sup>60</sup> datasets on  
216 the identified gene set. Notably, the GO terms related to lipoprotein particles emerged as sig-



**Fig. 3: a** The cumulative incidence plots of high-risk individuals (with the highest 5% PRSs) identified by PRS-Net and baseline methods. Each plot illustrates the estimated percentage of individuals diagnosed with a specific disease at different ages. We provide cumulative incidence plots for the original datasets as a reference. **b** The PRS prediction performance of PRS-Net compared to baseline methods on an asthma dataset encompassing multiple ancestries, including Western European (EUR), South Asian (SAS), and African (AFR) ancestry, measured in terms of the area under the receiver operating characteristic curve (AUROC). The results on the mixed ancestry test set are also reported. The bars are the estimated standard errors.

217 significantly enriched within the gene set (Supplementary Fig. 2a). This observation is in line with  
 218 prior studies that have implicated lipoprotein particles as significantly potential risk factors for  
 219 Alzheimer's disease<sup>61–63</sup> and have highlighted the role of metabolic dysregulation in the pro-  
 220 gression of Alzheimer's disease<sup>64,65</sup>. Notably, the exploration of high-density lipoprotein-inspired



8 H. Li et al.

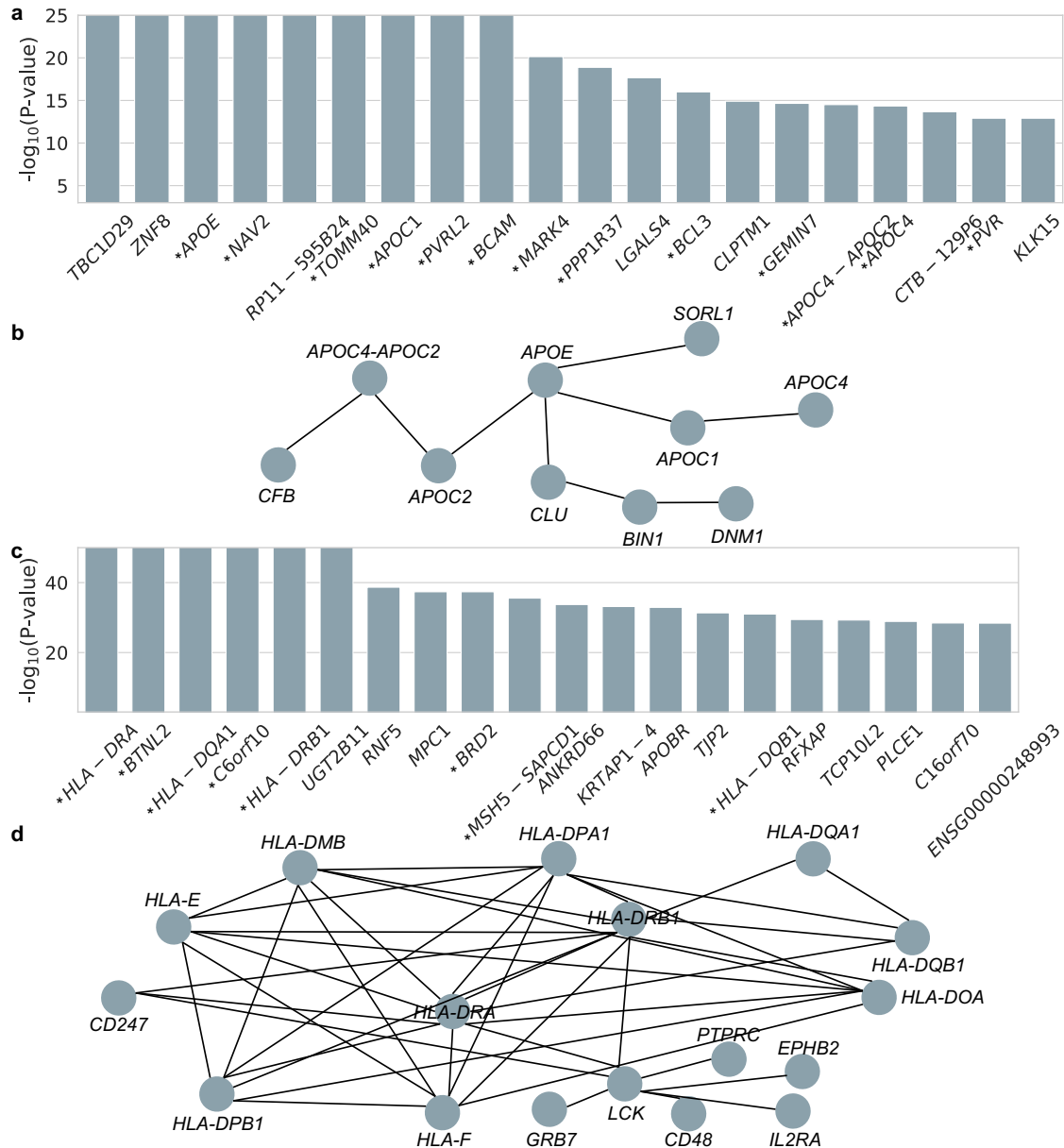


Fig. 4: PRS-Net identifies disease-related genes and GGIs for Alzheimer's disease and multiple sclerosis. **a** Top 20 genes with the highest statistical significance in the Mann-Whitney U test for Alzheimer's disease. The Mann-Whitney U test was utilized to assess whether the attention scores for a particular gene among the cases were significantly higher than those observed in the control group. An asterisk preceding the gene name signifies that the gene has been reported to be associated with Alzheimer's disease in previous studies. **b** Examples of interactions within the gene set with statistical significance ( $P$ -value  $< 0.001$ ) from the Mann-Whitney U test for Alzheimer's disease. **c** Top 20 genes with the highest statistical significance in the Mann-Whitney U test for multiple sclerosis. **d** Examples of interactions within the gene set with statistical significance ( $P$ -value  $< 0.001$ ) from the Mann-Whitney U test for multiple sclerosis.

222 illustrates the top 20 genes with the utmost statistical significance in the Mann-Whitney U test.  
223 Remarkably, 15 out of these 20 genes have been identified as potential risk factors for Alzheimer's  
224 disease in previous studies. One notable example is *APOE*, which is the most prevalent high-  
225 density lipoprotein in the central nervous system and has been consistently linked to Alzheimer's  
226 disease in numerous studies<sup>66–71</sup>. Fig. 4b illustrates examples of the interactions from the GGI  
227 network between genes within the identified gene set. Please refer to Supplementary Data 2 for  
228 the complete list of the GGIs. Interestingly, aside from *APOE*, other genes within the *APOE* gene  
229 cluster, including *APOC1*, *APOC2*, and *APOC4*, were also identified as disease-related genes.  
230 This finding aligns with previous studies that have shown interdependent or independent associ-  
231 ations of genes within the *APOE* gene cluster with Alzheimer's disease<sup>72–76</sup>. For instance, it has  
232 been shown that the variant *APOE* and *APOC2* exhibit interactive effects on metabolic pathways,  
233 potentially contributing to the risk of Alzheimer's disease<sup>72</sup>. *APOC1* also has been reported to  
234 serve as a risk factor for Alzheimer's disease in combination with *APOE*<sup>74</sup>. Furthermore, the com-  
235 bined effect of *APOE* and *CLU* on Alzheimer's disease has been observed<sup>77</sup>. *SORL1* is an *APOE*  
236 receptor gene, which has been recognized as a genetic risk factor in Alzheimer's disease. Recent  
237 research has elucidated the mechanistic connection between these two significant genetic factors  
238 in Alzheimer's disease<sup>78</sup>. A neuron-specific interaction between Alzheimer's disease risk factors  
239 *SORL1*, *APOE*, and *CLU* have also been shown in a recent study<sup>79</sup>. These observations highlight  
240 the proficiency of PRS-Net in not only identifying disease-related genes but also uncovering gene  
241 clusters that exhibit interactions contributing to the risk of Alzheimer's disease.

242 We also utilized PRS-Net to uncover genes and GGIs associated with multiple sclerosis. The  
243 Mann-Whitney U test identified a gene set with 456 potential risk genes (P-value <0.001). Please  
244 refer to Supplementary Data 3 for the complete list of the genes. The GSEA<sup>58</sup> using the KEGG<sup>60</sup>  
245 dataset on this gene set highlighted numerous immune-related pathways of statistical significance,  
246 such as antigen processing and presentation, allograft rejection, and graft-versus-host disease  
247 (Supplementary Fig. 4b). This finding aligns with the well-established understanding of multiple  
248 sclerosis as an autoimmune inflammatory disorder. The GSEA using the GO<sup>59</sup> dataset, unveiled  
249 significant enrichment of GO terms related to the major histocompatibility complex (MHC) protein  
250 complex within the identified gene set (Supplementary Fig. 4a), which can be supported by pre-  
251 vious studies that underscore the substantial genetic impact of MHCs on multiple sclerosis<sup>80–84</sup>.  
252 *HLA-DRA*, a subunit of *HLA-DR* which is a human MHC, was identified as the most significant  
253 gene in our analysis (Fig. 4c). Moreover, substantial *HLA* genes were identified as risk genes in  
254 our analysis (Fig. 4d). Please refer to Supplementary Data 4 for the complete list of the GGIs.  
255 This finding is in line with a previous study indicating that *HLA* interactions modulate genetic risk  
256 for multiple sclerosis<sup>85</sup>. Additionally, non-additive interactions between *HLAs* have been widely  
257 reported to significantly affect the risk of autoimmune diseases<sup>38,48–50</sup>. These discoveries col-  
258 lectively provide compelling evidence of the potential of PRS-Net to offer valuable insights that  
259 advance our understanding of diseases.

### 260 3.3 Ablation studies

261 To assess the effectiveness of specific design choices in PRS-Net, we conducted comprehensive  
262 ablation studies. We introduced various modified frameworks derived from PRS-Net, each with  
263 distinct constraints: PRS-Net-GGI (omitting the GGI network), PRS-Net-Att+Sum (replacing the  
264 attentive readout module with a sum readout module, which summarized the node feature to de-  
265 rive the global-level representations), PRS-Net-Att+Mean (replacing the attentive readout module  
266 with a mean readout module, which computes the average of node features to derive global-level  
267 representations), and PRS-Net-Att+Max (replacing the attentive readout module with a max read-

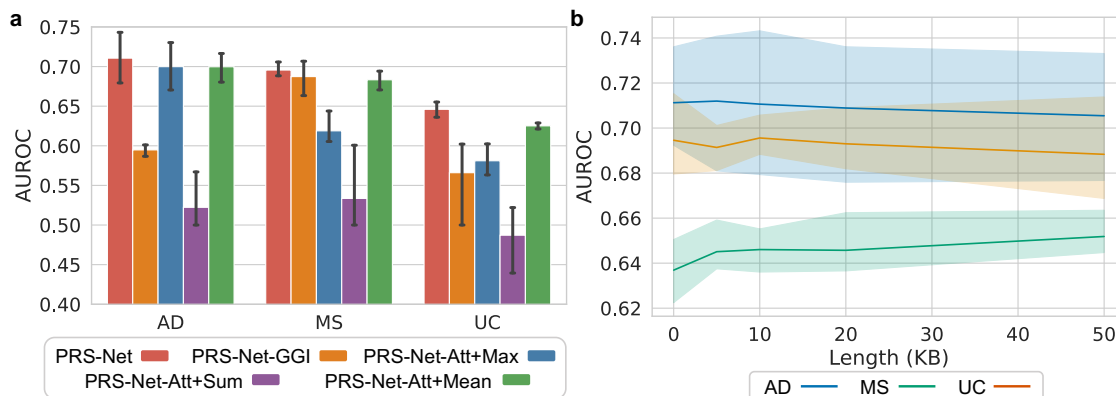


Fig. 5: The results of ablation studies on PRS-Net. **a** The comparison results of PRS-Net and its variations, including PRS-Net-GGI (omitting the GGI network), PRS-Net-Att+Sum (replacing the attentive readout module with a sum readout module, which summarized the node feature to derive the global-level representations), PRS-Net-Att+Mean (replacing the attentive readout module with a mean readout module, which computes the average of node features to derive global-level representations), and PRS-Net-Att+Max (replacing the attentive readout module with a max readout module, which extracts maximum values across node features to derive the global-level representations), conducted on the datasets of Alzheimer's disease (AD), multiple sclerosis (MS), and ulcerative colitis (UC). The bars are the estimated standard errors. **b** The PRS prediction performance of PRS-Net versus the extension lengths upstream and downstream of the transcription start and end sites.

268 out module, which extracts maximum values across node features to derive the global-level rep-  
269 resentations). We compared the performance of PRS-Net against these variants using datasets  
270 related to Alzheimer's disease, multiple sclerosis, and ulcerative colitis. The results showcased  
271 that PRS-Net surpassed PRS-Net-GGI by an average relative improvement of 11.6%, underscor-  
272 ing the significance of incorporating the GGI network to capture the intricate genetic interactions  
273 associated with diseases (Fig. 5a). Furthermore, PRS-Net outperformed PRS-Net-Att+Sum, PRS-  
274 Net-Att+Mean, and PRS-Net-Att+Max with average relative improvements of 33.0%, 2.2%, and  
275 8.4%, respectively, highlighting the effectiveness of the attentive readout module in summarizing  
276 node features (Fig. 5a).

277 Additionally, we explored the impact of varying extension lengths both upstream and down-  
278 stream of the transcription start and end sites when calculating gene-level PRSs. We assessed  
279 different length values, including 0, 5, 10, 20, and 50 KB, and subsequently evaluated their predic-  
280 tion performance. The results demonstrated that PRS-Net is generally robust to different extension  
281 lengths (Fig. 5b). However, it is noteworthy that the performance of PRS-Net on the multiple scler-  
282 osis dataset significantly declined when the extension length was set to 0 KB (Fig. 5b). This  
283 observation suggested that including SNPs from non-coding regions can indeed enhance the ac-  
284 curacy of PRS prediction.

## 285 Discussion

286 In this study, we develop PRS-Net, a deep-learning framework that offers interpretable and im-  
287 proved PRS predictions. By constructing a GGI network and incorporating a graph neural net-  
288 work, PRS-Net fully takes advantage of the power of non-linear associations between genetic

289 factors and phenotypes. Additionally, the integration of an attentive readout module empowers  
290 PRS-Net to deliver interpretable predictions. Through comprehensive testing across six complex  
291 diseases, PRS-Net consistently achieved superior performance in comparison with baseline meth-  
292 ods in PRS prediction. Furthermore, we demonstrated the interpretability of PRS-Net by using it  
293 to identify specific genes and GGIs that significantly impact the risk of Alzheimer's disease and  
294 multiple sclerosis. In summary, PRS-Net provides a potent tool for accurate PRS prediction and  
295 biological discovery for complex diseases.

## 296 **Data availability**

297 The GWAS data for Alzheimer's disease can be accessed at [https://ctg.cncr.nl/software/summary\\_s](https://ctg.cncr.nl/software/summary_statistics/)  
298 [tatistics/](https://ctg.cncr.nl/software/summary_statistics/). The GWAS data for atrial fibrillation can be accessed at [https://cvd.hugeamp.org/download](https://cvd.hugeamp.org/downloads.html#summary/)  
299 [s.html#summary/](https://cvd.hugeamp.org/downloads.html#summary/). The GWAS data for ulcerative colitis can be accessed at [ftp://ftp.sanger.ac.uk/pub](ftp://ftp.sanger.ac.uk/pub/project/humgen/summary_statistics/human/2016-11-07/)  
300 [/project/humgen/summary\\_statistics/human/2016-11-07/](ftp://ftp.sanger.ac.uk/pub/project/humgen/summary_statistics/human/2016-11-07/). The GWAS data for asthma can be ac-  
301 cessed at <https://www.globalbiobankmeta.org/resources/>. The GWAS data for rheumatoid arthritis  
302 can be accessed at [https://data.cyverse.org/dav-anon/iplant/home/kazuyoshihigaki/ra\\_gwas/ra\\_g](https://data.cyverse.org/dav-anon/iplant/home/kazuyoshihigaki/ra_gwas/ra_gwas-10-28-2021.tar/)  
303 [was-10-28-2021.tar/](https://data.cyverse.org/dav-anon/iplant/home/kazuyoshihigaki/ra_gwas/ra_gwas-10-28-2021.tar/). The GWAS data for multiple sclerosis can be accessed at [https://imsgc.net/?](https://imsgc.net/?page_id=31/)  
304 [page\\_id=31/](https://imsgc.net/?page_id=31/). The UKBB dataset is available at <https://www.ukbiobank.ac.uk>.

## 305 **Code availability**

306 The source code of PRS-Net can be downloaded from the Github repository at [https://github.com/li](https://github.com/lihan97/PRS-Net)  
307 [han97/PRS-Net](https://github.com/lihan97/PRS-Net).

## 308 **Acknowledgments**

309 This work was supported in part by the National Natural Science Foundation of China (T2125007  
310 to J.Z.), the National Key Research and Development Program of China (2021YFF1201300 to  
311 J.Z.), the New Cornerstone Science Foundation through the XPLOER PRIZE (J.Z.), the Re-  
312 search Center for Industries of the Future (RCIF) at Westlake University (J.Z.), and the Westlake  
313 Education Foundation (J.Z.).

## 314 **Author contributions statement**

315 H.L. and S.Z. conceived the concept and designed the study. H.L. and S.Z. developed the method-  
316 ology and conducted data analysis. H.L., J.Z., M.S. and S.Z. are responsible for the data interpre-  
317 tation. S.Z., M.S. and J.Z. supervised the project. H.L. and S.Z. prepared the manuscript with the  
318 assistance from all other authors.

## 319 **Competing interests statement**

320 All authors declare no competing interests.

## 321 References

- 322 1. Robert Clarke, John F Peden, Jemma C Hopewell, Theodosios Kyriakou, Anuj Goel, Simon C Heath, Sarah Parish,  
323 Simona Barlera, Maria Grazia Franzosi, Stephan Rust, et al. Genetic variants associated with lipoprotein  
324 level and coronary disease. *New England Journal of Medicine*, 361(26):2518–2528, 2009.
- 325 2. CKDGen Consortium, KidneyGen Consortium, EchoGen consortium, CHARGE-HF consortium, Thor Aspelund,  
326 Melissa Garcia, Yen-Pei C Chang, Jeffrey R O’Connell, Nanette I Steinle, et al. Genetic variants in novel pathways  
327 influence blood pressure and cardiovascular disease risk. *Nature*, 478(7367):103–109, 2011.
- 328 3. Chris Gaiteri, Sara Mostafavi, Christopher J Honey, Philip L De Jager, and David A Bennett. Genetic variants in  
329 alzheimer disease—molecular and brain network approaches. *Nature Reviews Neurology*, 12(7):413–427, 2016.
- 330 4. Valeria Orrù, Maristella Steri, Gabriella Sole, Carlo Sidore, Francesca Virdis, Mariano Dei, Sandra Lai, Magdalena  
331 Zoledziewska, Fabio Busonero, Antonella Mulas, et al. Genetic variants regulating immune cell levels in health and  
332 disease. *Cell*, 155(1):242–256, 2013.
- 333 5. Ali Torkamani, Nathan E Wineinger, and Eric J Topol. The personal and clinical utility of polygenic risk scores.  
334 *Nature Reviews Genetics*, 19(9):581–590, 2018.
- 335 6. Cathryn M Lewis and Evangelos Vassos. Polygenic risk scores: from research tools to clinical instruments.  
336 *Genome medicine*, 12(1):1–11, 2020.
- 337 7. Shing Wan Choi, Timothy Shin-Heng Mak, and Paul F O’Reilly. Tutorial: a guide to performing polygenic risk score  
338 analyses. *Nature protocols*, 15(9):2759–2772, 2020.
- 339 8. Greg Gibson. On the utilization of polygenic risk scores for therapeutic targeting. *PLoS genetics*, 15(4):e1008060,  
340 2019.
- 341 9. Takahiro Konuma and Yukinori Okada. Statistical genetics and polygenic risk score for precision medicine.  
342 *Inflammation and regeneration*, 41(1):1–5, 2021.
- 343 10. Responsible use of polygenic risk scores in the clinic: potential benefits, risks and gaps. *Nature Medicine*,  
344 27(11):1876–1884, 2021.
- 345 11. Bjarni J Vilhjálmsson, Jian Yang, Hilary K Finucane, Alexander Gusev, Sara Lindström, Stephan Ripke, Giulio  
346 Genovese, Po-Ru Loh, Gaurav Bhatia, Ron Do, et al. Modeling linkage disequilibrium increases accuracy of  
347 polygenic risk scores. *The American journal of human genetics*, 97(4):576–592, 2015.
- 348 12. Jack Euesden, Cathryn M Lewis, and Paul F O’Reilly. Prsice: polygenic risk score software. *Bioinformatics*,  
349 31(9):1466–1468, 2015.
- 350 13. Timothy Shin Heng Mak, Robert Milan Porsch, Shing Wan Choi, Xueya Zhou, and Pak Chung Sham. Polygenic  
351 scores via penalized regression on summary statistics. *Genetic epidemiology*, 41(6):469–480, 2017.
- 352 14. Shing Wan Choi and Paul F O’Reilly. Prsice-2: Polygenic risk score software for biobank-scale data. *Gigascience*,  
353 8(7):giz082, 2019.
- 354 15. Luke R Lloyd-Jones, Jian Zeng, Julia Sidorenko, Loïc Yengo, Gerhard Moser, Kathryn E Kemper, Huanwei Wang,  
355 Zhili Zheng, Reedik Magi, Tõnu Esko, et al. Improved polygenic prediction by bayesian multiple regression on  
356 summary statistics. *Nature communications*, 10(1):5086, 2019.
- 357 16. Florian Privé, Julyan Arbel, and Bjarni J Vilhjálmsson. Ldpred2: better, faster, stronger. *Bioinformatics*, 36(22-  
358 23):5424–5431, 2020.
- 359 17. Florian Privé, Julyan Arbel, Hugues Aschard, and Bjarni J Vilhjálmsson. Identifying and correcting for misspecifi-  
360 cations in gwas summary statistics and polygenic scores. *Human Genetics and Genomics Advances*, 3(4), 2022.
- 361 18. William YS Wang, Bryan J Barratt, David G Clayton, and John A Todd. Genome-wide association studies: theoret-  
362 ical and practical concerns. *Nature Reviews Genetics*, 6(2):109–118, 2005.
- 363 19. Arthur Korte and Ashley Farlow. The advantages and limitations of trait analysis with gwas: a review. *Plant*  
364 *methods*, 9(1):1–9, 2013.
- 365 20. Emil Uffelmann, Qin Qin Huang, Nchangwi Syntia Munung, Jantina De Vries, Yukinori Okada, Alicia R Martin,  
366 Hilary C Martin, Tuuli Lappalainen, and Danielle Posthuma. Genome-wide association studies. *Nature Reviews*  
367 *Methods Primers*, 1(1):59, 2021.
- 368 21. Shaun Purcell, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel AR Ferreira, David Bender, Julian Maller,  
369 Pamela Sklar, Paul IW De Bakker, Mark J Daly, et al. Plink: a tool set for whole-genome association and population-  
370 based linkage analyses. *The American journal of human genetics*, 81(3):559–575, 2007.
- 371 22. Ying Wang, Kristin Tsuo, Masahiro Kanai, Benjamin M Neale, and Alicia R Martin. Challenges and opportunities  
372 for developing more generalizable polygenic risk scores. *Annual review of biomedical data science*, 5:293–320,  
373 2022.
- 374 23. Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document analysis and*  
375 *recognition*, volume 1, pages 278–282. IEEE, 1995.
- 376 24. Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
- 377 25. Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages  
378 1189–1232, 2001.

- 379 26. Jerome H Friedman. Stochastic gradient boosting. *Computational statistics & data analysis*, 38(4):367–378, 2002.
- 380 27. Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to  
381 boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- 382 28. Trevor Hastie, Saharon Rosset, Ji Zhu, and Hui Zou. Multi-class adaboost. *Statistics and its Interface*, 2(3):349–  
383 360, 2009.
- 384 29. Hadeel Alzoubi, Raid Alzubi, and Naeem Ramzan. Deep learning framework for complex disease risk prediction  
385 using genomic variations. *Sensors*, 23(9):4439, 2023.
- 386 30. Arnór I Sigurdsson, Ioannis Louloudis, Karina Banasik, David Westergaard, Ole Winther, Ole Lund, Sisse Rye  
387 Ostrowski, Christian Erikstrup, Ole Birger Vesterager Pedersen, Mette Nyegaard, et al. Deep integrative models  
388 for large-scale human genomics. *Nucleic Acids Research*, 51(12):e67–e67, 2023.
- 389 31. Pau Bellot, Gustavo de Los Campos, and Miguel Pérez-Enciso. Can deep learning improve genomic prediction of  
390 complex human traits? *Genetics*, 210(3):809–819, 2018.
- 391 32. Yu Xu, Dragana Vuckovic, Scott C Ritchie, Parsa Akbari, Tao Jiang, Jason Grealey, Adam S Butterworth, Willem H  
392 Ouwehand, David J Roberts, Emanuele Di Angelantonio, et al. Machine learning optimized polygenic scores for  
393 blood cell traits identify sex-specific trajectories and genetic correlations with disease. *Cell Genomics*, 2(1), 2022.
- 394 33. Saeed Masoudnia and Reza Ebrahimpour. Mixture of experts: a literature survey. *Artificial Intelligence Review*,  
395 42:275–293, 2014.
- 396 34. Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott,  
397 Jane Green, Martin Landray, et al. Uk biobank: an open access resource for identifying the causes of a wide range  
398 of complex diseases of middle and old age. *PLoS medicine*, 12(3):e1001779, 2015.
- 399 35. Heather J Cordell. Detecting gene–gene interactions that underlie human diseases. *Nature Reviews Genetics*,  
400 10(6):392–404, 2009.
- 401 36. Or Zuk, Eliana Hechter, Shamil R Sunyaev, and Eric S Lander. The mystery of missing heritability: Genetic inter-  
402 actions create phantom heritability. *Proceedings of the National Academy of Sciences*, 109(4):1193–1198, 2012.
- 403 37. Matthew B Taylor and Ian M Ehrenreich. Higher-order genetic interactions and their contribution to complex traits.  
404 *Trends in genetics*, 31(1):34–40, 2015.
- 405 38. Tobias L Lenz, Aaron J Deutsch, Buhm Han, Xinli Hu, Yukinori Okada, Stephen Eyre, Michael Knapp, Alexandra  
406 Zhernakova, Tom WJ Huizinga, Goncalo Abecasis, et al. Widespread non-additive and interaction effects within  
407 hla loci modulate the risk of autoimmune diseases. *Nature genetics*, 47(9):1085–1090, 2015.
- 408 39. Luis Varona, Andres Legarra, Miguel A Toro, and Zulma G Vitezica. Non-additive effects in genomic selection.  
409 *Frontiers in genetics*, 9:78, 2018.
- 410 40. Marta Guindo-Martínez, Ramon Amela, Silvia Bonàs-Guarch, Montserrat Puiggròs, Cecilia Salvo, Irene Miguel-  
411 Escalada, Caitlin E Carey, Joanne B Cole, Sina Rüeger, Elizabeth Atkinson, et al. The impact of non-additive  
412 genetic associations on age-related complex diseases. *Nature communications*, 12(1):2436, 2021.
- 413 41. James M Cheverud and Eric J Routman. Epistasis and its contribution to genetic variance components. *Genetics*,  
414 139(3):1455–1461, 1995.
- 415 42. Jason H Moore and Scott M Williams. Epistasis and its implications for personal genetics. *The American Journal*  
416 *of Human Genetics*, 85(3):309–320, 2009.
- 417 43. Ben Lehner. Molecular mechanisms of epistasis within and between genes. *Trends in Genetics*, 27(8):323–331,  
418 2011.
- 419 44. Damian Szklarczyk, Rebecca Kirsch, Mikaela Koutrouli, Katerina Nastou, Farrokh Mehryary, Radja Hachilif, An-  
420 nika L Gable, Tao Fang, Nadezhda T Doncheva, Sampo Pyysalo, et al. The string database in 2023: protein–  
421 protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic*  
422 *acids research*, 51(D1):D638–D646, 2023.
- 423 45. Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In  
424 *International Conference on Learning Representations*, 2018.
- 425 46. World Health Organization et al. Icd-10. international statistical classification of diseases and related health  
426 problems: Tenth revision 1992, volume 1= cim-10. classification statistique internationale des maladies et des  
427 problèmes de santé connexes: Dixième révision 1992, volume 1. 1992.
- 428 47. 1000 Genomes Project Consortium et al. A global reference for human genetic variation. *Nature*, 526(7571):68,  
429 2015.
- 430 48. Marc Lipsitch, Carl T Bergstrom, and Rustom Antia. Effect of human leukocyte antigen heterozygosity on infectious  
431 disease outcome: the need for allele-specific measures. *BMC Medical Genetics*, 4(1):1–9, 2003.
- 432 49. Sue Tsai and Pere Santamaria. Mhc class ii polymorphisms, autoreactive t-cells, and autoimmunity. *Frontiers in*  
433 *immunology*, 4:321, 2013.
- 434 50. Philippe Goyette, Gabrielle Boucher, Dermot Mallon, Eva Ellinghaus, Luke Jostins, Hailiang Huang, Stephan Ripke,  
435 Elena S Gusareva, Vito Anness, Stephen L Hauser, et al. High-density mapping of the mhc identifies a shared  
436 role for hla-drb1\* 01: 03 in inflammatory bowel diseases and heterozygous advantage in ulcerative colitis. *Nature*  
437 *genetics*, 47(2):172–179, 2015.

- 438 51. Odd O Aalen and Søren Johansen. An empirical transition matrix for non-homogeneous markov chains based on  
439 censored observations. *Scandinavian Journal of Statistics*, pages 141–150, 1978.
- 440 52. Dennis J Selkoe. Alzheimer's disease: genes, proteins, and therapy. *Physiological reviews*, 2001.
- 441 53. Rudolph E Tanzi and Lars Bertram. New frontiers in alzheimer's disease genetics. *Neuron*, 32(2):181–184, 2001.
- 442 54. Lars Bertram and Rudolph E Tanzi. Thirty years of alzheimer's disease genetics: the implications of systematic  
443 meta-analyses. *Nature Reviews Neuroscience*, 9(10):768–778, 2008.
- 444 55. Celeste M Karch, Carlos Cruchaga, and Alison M Goate. Alzheimer's disease genetics: from the bench to the  
445 clinic. *Neuron*, 83(1):11–26, 2014.
- 446 56. Ya-Ping Tang and Elliot S Gershon. Genetic studies in alzheimer's disease. *Dialogues in clinical neuroscience*,  
447 2022.
- 448 57. Henry B Mann and Donald R Whitney. On a test of whether one of two random variables is stochastically larger  
449 than the other. *The annals of mathematical statistics*, pages 50–60, 1947.
- 450 58. Aravind Subramanian, Pablo Tamayo, Vamsi K Mootha, Sayan Mukherjee, Benjamin L Ebert, Michael A Gillette,  
451 Amanda Paulovich, Scott L Pomeroy, Todd R Golub, Eric S Lander, et al. Gene set enrichment analysis:  
452 a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National  
453 Academy of Sciences*, 102(43):15545–15550, 2005.
- 454 59. Michael Ashburner, Catherine A Ball, Judith A Blake, David Botstein, Heather Butler, J Michael Cherry, Allan P  
455 Davis, Kara Dolinski, Selina S Dwight, Janan T Eppig, et al. Gene ontology: tool for the unification of biology.  
456 *Nature genetics*, 25(1):25–29, 2000.
- 457 60. Minoru Kanehisa and Susumu Goto. Kegg: kyoto encyclopedia of genes and genomes. *Nucleic acids research*,  
458 28(1):27–30, 2000.
- 459 61. Zhike Zhou, Yifan Liang, Xiaoqian Zhang, Junjie Xu, Jueying Lin, Rongwei Zhang, Kexin Kang, Chang Liu, Chuan-  
460 sheng Zhao, and Mei Zhao. Low-density lipoprotein cholesterol and alzheimer's disease: a systematic review and  
461 meta-analysis. *Frontiers in aging neuroscience*, 12:5, 2020.
- 462 62. Mirna Safieh, Amos D Korczyn, and Daniel M Michaelson. Apoe4: an emerging therapeutic target for alzheimer's  
463 disease. *BMC medicine*, 17:1–17, 2019.
- 464 63. Yi Jin, Kudzai Chifodya, Guochen Han, Wenxin Jiang, Yun Chen, Yang Shi, Qiao Xu, Yilong Xi, Jun Wang, Jianping  
465 Zhou, et al. High-density lipoprotein in alzheimer's disease: From potential biomarkers to therapeutics. *Journal of  
466 Controlled Release*, 338:56–70, 2021.
- 467 64. Suzanne Craft. The role of metabolic disorders in alzheimer disease and vascular dementia: two roads converged.  
468 *Archives of neurology*, 66(3):300–305, 2009.
- 469 65. Xu Yan, Yue Hu, Biyao Wang, Sijian Wang, and Xinwen Zhang. Metabolic dysregulation contributes to the progres-  
470 sion of alzheimer's disease. *Frontiers in neuroscience*, 14:530219, 2020.
- 471 66. MS Tsai, Eric G Tangalos, Ronald C Petersen, Glenn E Smith, Daniel J Schaid, Emre Kokmen, Robert J Ivnik, and  
472 Stephen N Thibodeau. Apolipoprotein e: risk factor for alzheimer disease. *American journal of human genetics*,  
473 54(4):643, 1994.
- 474 67. Warren J Strittmatter and Allen D Roses. Apolipoprotein e and alzheimer disease. *Proceedings of the National  
475 Academy of Sciences*, 92(11):4725–4727, 1995.
- 476 68. Marion R Meyer, JoAnn T Tschanz, Maria C Norton, Kathleen A Welsh-Bohmer, David C Steffens, Bonita W Wyse,  
477 and John Breitner. Apoe genotype predicts when—not whether—one is predisposed to develop alzheimer disease.  
478 *Nature genetics*, 19(4):321–322, 1998.
- 479 69. Robert C Green, J Scott Roberts, L Adrienne Cupples, Norman R Relkin, Peter J Whitehouse, Tamsen Brown,  
480 Susan LaRusse Eckert, Melissa Butson, A Dessa Sadovnick, Kimberly A Quaid, et al. Disclosure of apoe genotype  
481 for risk of alzheimer's disease. *New England Journal of Medicine*, 361(3):245–254, 2009.
- 482 70. Emmanuelle Genin, Didier Hannequin, David Wallon, Kristel Slegers, Mikko Hiltunen, Onofre Combarros,  
483 María Jesús Bullido, Sebastiaan Engelborghs, Peter De Deyn, Claudine Berr, et al. Apoe and alzheimer disease:  
484 a major gene with semi-dominant inheritance. *Molecular psychiatry*, 16(9):903–907, 2011.
- 485 71. Yu Yamazaki, Na Zhao, Thomas R Caulfield, Chia-Chen Liu, and Guojun Bu. Apolipoprotein e and alzheimer  
486 disease: pathobiology and targeting strategies. *Nature Reviews Neurology*, 15(9):501–518, 2019.
- 487 72. Robert A Hegele, W Carl Breckenridge, Diane W Cox, Graham F Maguire, J Alick Little, and Philip W Connelly. In-  
488 teraction between variant apolipoproteins c-ii and e that affects plasma lipoprotein concentrations. *Arteriosclerosis  
489 and Thrombosis: A Journal of Vascular Biology*, 11(5):1303–1309, 1991.
- 490 73. Sebastián Cervantes, Lluís Samarançh, José Manuel Vidal-Taboada, Isabel Lamet, María Jesús Bullido, Ana  
491 Frank-García, Francisco Coria, Albert Lleó, Jordi Clarimón, Elena Lorenzo, et al. Genetic variation in apoe cluster  
492 region and alzheimer's disease risk. *Neurobiology of aging*, 32(11):2107–e7, 2011.
- 493 74. Qin Zhou, Fan Zhao, Ze-ping Lv, Chen-guang Zheng, Wei-dong Zheng, Liang Sun, Na-na Wang, Shenghang Pang,  
494 Fabiana Michelsen de Andrade, Mian Fu, et al. Association between apoc1 polymorphism and alzheimer's disease:  
495 a case-control study and meta-analysis. *PloS one*, 9(1):e87017, 2014.

- 496 75. Yvonne Shao, McKenzie Shaw, Kaitlin Todd, Maria Khrestian, Giana D'Aleo, P John Barnard, Jeff Zahratka, Jagan  
497 Pillai, Chang-En Yu, C Dirk Keene, et al. Dna methylation of tomm40-apoe-apoc2 in alzheimer's disease. Journal  
498 of human genetics, 63(4):459–471, 2018.
- 499 76. Alexander M Kulminski, Ethan Jain-Washburn, Ian Philipp, Liang He, Yury Loika, Elena Loiko, Olivia Bagley, Svet-  
500 lana Ukraintseva, Anatoliy Yashin, Konstantin Arbeev, et al. Apoe epsilon4 allele and tomm40-apoc1 variants jointly  
501 contribute to survival to older ages. Aging Cell, 21(12):e13730, 2022.
- 502 77. Florence F Roussotte, Boris A Gutman, Sarah K Madsen, John B Colby, Paul M Thompson, Alzheimer's Dis-  
503 ease Neuroimaging Initiative, et al. Combined effects of alzheimer risk variants in the clu and apoe genes on  
504 ventricular expansion patterns in the elderly. Journal of Neuroscience, 34(19):6537–6545, 2014.
- 505 78. Hyo Lee, Aimee J Aylward, Richard V Pearse, Yi-Chen Hsieh, Zachary M Augur, Courtney R Benoit, Vicky Chou,  
506 Allison Knupp, Cheryl Pan, Srilakshmi Goberdhan, et al. Cell-type-specific regulation of apoe levels in human  
507 neurons by the alzheimer's disease risk gene sorl1. bioRxiv, pages 2023–02, 2023.
- 508 79. Pranav Preman and Amaia M Arranz. A neuron-specific interaction between alzheimer's disease risk factors sorl1,  
509 apoe, and clu. Cell Reports, 42(9), 2023.
- 510 80. Ute Traugott. Multiple sclerosis: relevance of class i and class ii mhc-expressing cells to lesion development.  
511 Journal of neuroimmunology, 16(2):283–302, 1987.
- 512 81. Sunhee C Lee, GR Wayne Moore, George Golenwsky, and Cedric S Raine. Multiple sclerosis: a role for astroglia  
513 in active demyelination suggested by class ii mhc expression and ultrastructural study. Journal of neuropathology  
514 and experimental neurology, 49(2):122–136, 1990.
- 515 82. Multiple Sclerosis Genetics Group, Jonathan L Haines, Henry A Terwedow, Katie Burgess, Margaret A Pericak-  
516 Vance, Jackie B Rimmler, Eden R Martin, Jorge R Oksenberg, Robin Lincoln, David Y Zhang, et al. Linkage of  
517 the mhc to familial multiple sclerosis suggests genetic heterogeneity. Human molecular genetics, 7(8):1229–1234,  
518 1998.
- 519 83. Matthew R Lincoln, Alexandre Montpetit, M Zameel Cader, Janna Saarela, David A Dymment, Milvi Tiislar, Vincent  
520 Ferretti, Pentti J Tienari, A Dessa Sadovnick, Leena Peltonen, et al. A predominant role for the hla class ii region  
521 in the association of the mhc region with multiple sclerosis. Nature genetics, 37(10):1108–1112, 2005.
- 522 84. David A Dymment, Blanca M Herrera, M Zameel Cader, Cristen J Willer, Matthew R Lincoln, A Dessa Sadovnick,  
523 Neil Risch, and George C Ebers. Complex interactions among mhc haplotypes in multiple sclerosis: susceptibility  
524 and resistance. Human molecular genetics, 14(14):2019–2026, 2005.
- 525 85. Class ii hla interactions modulate genetic risk for multiple sclerosis. Nature genetics, 47(10):1107–1113, 2015.