

Bridging Heterogeneous Mutation Data to Enhance Disease-Gene Discovery

Kaiyin Zhou

Huazhong Agriculture University

Yuxing Wang

Huazhong Agriculture University

Kevin Bretonnel Cohen

University of Colorado Denver - Anschutz Medical Campus

Jin-Dong Kim

Research Organization of Information and Systems

Xiaohang Ma

Huazhong Agriculture University

Zhixue Shen

Huazhong Agriculture University

Xiangyu Meng

Wuhan University Zhongnan Hospital

Jingbo Xia (✉ xiajingbo.math@gmail.com)

Huazhong Agriculture University <https://orcid.org/0000-0002-7285-588X>

Research article

Keywords: Heterogeneous data, data fusion, generative model, GWAS, text mining

Posted Date: August 14th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-44127/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Briefings in Bioinformatics on April 13th, 2021. See the published version at <https://doi.org/10.1093/bib/bbab079>.

Bridging Heterogeneous Mutation Data to Enhance Disease-Gene Discovery

Kaiyin Zhou^{1#} zhoukaiyinhzau@gmail.com, Yuxing Wang^{1#} wang-yuxing@foxmail.com, Kevin Bretonnel Cohen² kevin.cohen@gmail.com, Jin-Dong Kim³ jdkim@dbcls.rois.ac.jp, Xiaohang Ma¹ mxh.hzau.edu.cn@webmail.hzau.edu.cn, Zhixue Shen¹ zhixue_shen@163.com, Xiangyu Meng^{4,5} mengxy_whu@163.com, Jingbo Xia^{1*} xiajingbo.math@gmail.com

¹Hubei Key Lab of Agricultural Bioinformatics, College of Informatics, Huazhong Agricultural University, Wuhan, Hubei Province, P.R. China

²School of Medicine, University of Colorado at Denver, Anschutz Medical Campus, Colorado, U.S

³Database Center for Life Science (DBCLS), Research Organization of Information and Systems (ROIS), Tokyo, Japan

⁴Department of Urology, Zhongnan Hospital of Wuhan University, Wuhan, Hubei Province, P.R. China

⁵Institut Curie, CNRS, Molecular Oncology Team, PSL Research University, Paris, France

#The authors have the same contributions

*Correspondence: xiajingbo.math@gmail.com, xjb@mail.hzau.edu.cn

Abstract

Background: Bridging heterogeneous mutation data fills in the gap between various data categories and propels discovery of disease-related genes. It is known that genome-wide association study (GWAS) infers significant *mutation associations* which link genotype and phenotype, and it is under-powered for pinpointing causal genes due to high false positive or negative rate. In the meantime, mutation events widely reported in literature unveil typical functional biological process, including *mutation types* like gain-of-function and loss-of-function.

Methods: To bring together the heterogeneous mutation data, we propose a pipeline, “Gene-Disease Association prediction by Mutation Data Bridging (GDAMDB)”, with a statistic generative model. The model learns the distribution parameters of *mutation associations* and *mutation types*, and recovers false negative GWAS mutations which fail to pass significant test but represent supportive evidences of functional biological process in literature.

Results: Eventually, GDAMDB is applied in Alzheimer’s disease which is a common inheritable neurodegenerative disorder with unknown pathological mechanism, and it predicted 79 AD-associated genes. Besides 12 of them come from the original GWAS study, 57 of them are supported to be AD-related by other GWAS or literature report.

34 **Conclusion:** Our model is capable of enhancing the GWAS-based gene association discovery by well
35 combining text mining results. The positive result indicates that bridging the heterogeneous mutation
36 data is contributory for the novel disease-related gene discovery.

37 **Key Words:** Heterogeneous data, data fusion, generative model, GWAS, text mining

38 1 Background

39 Genome-wide association study (GWAS) is helpful to identify disease-related genes through significant test on single
40 nucleotide polymorphisms (SNPs) across the entire genome, and the p-value of the SNPs is the *mutation association*
41 data that represent the relevance between the gene and disease. However, as generally recognized, the high false
42 negative rate of GWAS makes that not all *de facto* vital variations are able to pass the multiple testing, and the lack
43 of consideration about the biological mechanism between genes and phenotype makes GWAS insufficient to pinpoint
44 causal variations. Therefore, various researches considered to combine other Omics data with GWAS. For example,
45 Jia et al. [1] developed a dense module searching method, which combined protein-protein interaction (PPI) network
46 with GWAS data, to identify the candidate genes that was well studied in PPI but with low *p*-value in GWAS [2]. Wang
47 et al. [3] proposed that genes in the same functional pathway may work together to raise a disease but some of them
48 are difficult to reach the significant threshold in GWAS. Therefore, they combined the pathway knowledge with
49 existed GWAS data to retrieve the omitted genes. Though it has been attempted to combine various Omics data with
50 GWAS data, to combine another type of heterogeneous mutation data with *mutation association* GWAS is still a new
51 idea.

52 The *mutation type*, i.e., loss of function (LOF) or gain of function (GOF) [4], is the mutation categorical data which
53 bridges genes and diseases [5]. The LOF mutation in a gene results in the reduction or abolition of the gene function,
54 while GOF mutation in a gene results in the enhanced or new gene function. The changed genes function eventually
55 leads to the downstream molecular and cellular biological process, well supports the investigation of gene-disease
56 associations, and provides evidence to unveil the pathological mechanism of diseases.

57 To our best knowledge, there is no disease-related LOF/GOF database. Fortunately, the relevant biological
58 processes representing LOF or GOF are abundant in literature. Taking Alzheimer's disease (AD) as an example,
59 *mutation types* reported in literature represent supportive information in functional biological process. In 1997,
60 Citron et al. [6] proposed that the GOF mutation on PSEN1 and PSEN2 lead to an increase in $A\beta_{42}$ production. In 2013,

61 Guerreiro et al. [7] analyzed the genetic variability of TREM2 in 1,092 AD patient and 1,107 controls and found the
62 LOF mutation in TREM2 are associated with AD. In 2015, Steinberg et al. [8] identified that LOF variants in ABCA7
63 are related with AD in Icelanders. Though the *mutation type* information is usually well described in literature, most
64 of them are stated in an explicit semantics. For instance, the sentence “*In studies of cell lines transfected with beta-*
65 *amyloid precursor protein (beta APP) cDNAs, the beta APP mutation K670N/M671L found in a Swedish familial AD*
66 *(FAD) pedigree has previously been shown to cause a marked augmentation of A beta secretion.*”(pmid:7991571)
67 described a GOF *mutation type*, while “*In this issue of Neuron, describe two rare ADAM10 prodomain mutations that*
68 *cause late-onset Alzheimer’s disease by impairing prodomain chaperone function, attenuating alpha-secretase activity,*
69 *and reducing adult hippocampal neurogenesis.*”(pmid:24139026) described a LOF *mutation type*.

70 Owing to the rapid growth of algorithm strength, a combination of the state-of-the-art text mining strategy and
71 the customized corpus make it possible to extract *mutation type* from literature mining in a PubMed scale. An active
72 gene annotation corpus (AGAC) plays a role of a gold training set [9] in the area of LOF/GOF retrieval, which well
73 annotates the molecular and cellular events after mutation and captures the semantics of mutation events.

74 Though the availability of both *mutation type* and *mutation association* offered abundant evidences for gene
75 disease associations, neither single data is empirically rational. GWAS is known for high false positive or negative
76 rate, and the *mutation type* mined from literature only retrieves published knowledge. Therefore, it is illuminative
77 to use text mined *mutation type* to decrease the false positive or negative rate in GWAS, and connect these two
78 heterogeneous mutation data. Since generative model is capable of investigating the distribution of heterogeneous
79 mutation data, it helps to generate and discover the real *de facto* mutations based on both data. In this research, we
80 proposed a generative model, “Gene-Disease Association prediction by Mutation Data Bridging (GDAMDB)”, to
81 bridge these two heterogeneous mutation data, and enhance the gene-disease knowledge discovery. The model
82 learns the distribution parameters from *mutation association* data and *mutation type* data, and generating novel data
83 that contains the statistical characteristics of them. Hence, based on the novel data, the model works well to predict
84 novel disease-related genes, which in turn is used for loss/gain-of-function inference and novel gene-disease
85 association discovery.

86 As an application, GDAMDB was applied to Alzheimer’s disease (AD), which is a common neurodegenerative
87 disorder that threaten the elderly for a long time. Although the pathogenesis of AD is still unknown, it is widely
88 accepted that the accumulation of amyloid β forms the plaques on patients’ brain thus breaks the calcium
89 equilibrium of the neurons and finally leads the cell apoptosis [10]. Since AD is highly inheritable [11], identifying
90 the important genes sheds light on unveiling the mechanisms of the disease.

91 The *mutation association* data was downloaded from summary data of the International Genomics of Alzheimer’s
92 Project (IGAP) [12], which performed a two-stage GWAS on individuals of European ancestry on 7,055,881 SNPs, and
93 identified 11 new loci of AD. Furthermore, the *mutation types* data was extracted from a PubMed-scale event extraction.
94 After bridging the two heterogeneous mutation data with GDAMDB, we finally retrieved 79 AD associated genes.
95 Intuitively, 2 of them came from the original GWAS study, and 57 of them showed relevance with AD as supported by
96 other GWAS or literature evidences. Therefore, it is believed that GDAMDB successfully combined heterogeneous
97 GWAS and literature data and enhanced the disease-related gene discovery.

98

99 **2 Methods**

100 In this part, we will introduce the modules and models in the pipeline of GDAMDB model for gene-disease
101 association prediction.

102 **2.1 “Mutation Type Retrieval” module in GDAMDB**

103 We designed a “*Mutation Type Retrieval*” module to jointly complete the task of entity recognition and mutation type
104 classification. In this model, we transfer the parameters of BERT-base that released by Google to our model for fine-
105 tuning. Before the input, we use regex and SETH tool for filtering mutation unrelated sentences, the rest of the
106 sentences are thought to be related to *mutation type*. Then, those sentences are used as the input of our joint learning
107 model.

108 BERT-Base trains model on a large text corpus. For learning deeply bidirectional representation of words, it
109 masks 15% of the words in the input and run the entire sequence through a 12-layer deep bidirectional transformer,
110 each layer is made up of a multi-head self-attention, residual connection, batch normalize and fully connected layer.
111 It updates its parameters through optimizing two tasks, next sentence predicts and mask words predict.

112 BERT-Base model has been trained by Google, here we transfer its parameters to our joint model. In our model,
113 our input is abstract with sentences filtered, then the sentence is tokenized by WordPiece tool, besides, a
114 classification label ([CLS]) is added in the head of each abstract as classification encoding of each abstract.

115 After encoding by BERT we get the representation of each word, then a fully connected layer and Softmax are
116 used for normalize classification weights after that CRF loss function is employed to optimize entity recognition task.
117 At the same time, we use word-level attention for the output of the Softmax layer (The output of Softmax layer is
118 regarded as a priori knowledge of mutation type classification). Finally, the attention output vector is concatenate

119 with the classification label encoded vector, and then a multi-label classifier is used for mutation type classification.
120 Here, we train our model in AGAC train sets, and test out model in development sets. AGAC corpus is designed for
121 extracting gene-mutation type-disease triples from PubMed abstracts.

122 For filter out SNP related documents in the AD case study, we first use “Alzheimer disease” [MeSH Terms] OR
123 Alzheimer’s disease [Text Word]” as the search criteria for downloading 137,473 abstracts from the PubMed
124 database. Then regex and SNP recognize tool SETH was employed to filter out mutation related texts. SETH is an SNP
125 extraction tool for recognizing SNPs and other short sequence variations. Finally, we get 9,430 mutation related
126 abstracts.

127 The retrieved “mutation types” and part of the sentence evidences are presented in Figure 3. For instance, gene
128 APP is predicted to be related to a GOF *mutation type*, and the sentence evidence is from PubMed with ID equals to
129 16685645. The sentence, “...Promoter mutations that increase amyloid precursor-protein expression are associated
130 with Alzheimer disease...” clearly support the GOF prediction. Moreover, the full results are in the Supplementary file
131 S2 or in online data repository(<https://hzaubionlp.com/agac-on-alzheimers-disease/>)

132 **2.2 “Synchronization Filter” module in GDAMDB**

133 “Synchronization Filter” module uses a strategy to obtain a gene set with greatest size and most significant literature
134 support. It designs to optimize the probability that most genes in the gene set maintain not only literature
135 significance but also the GWAS significance. Actually, the whole idea of significance integration is an analogue of
136 signal synchronization”. Taking this concern, the module is named as “Synchronization Filter” module.

137 In a mathematical way, we assume there is a gene set for each g with f_{dg} , where f_{dg} is *mutation type* retrieved
138 by “Mutation Type Retrieval” module. Since p_{dg} of every g is traceable from GWAS summary data, we order all of the
139 g with f_{dg} with its p -value in a descending order. Generally, the topmost g has greater chance to has *mutation type*
140 info as reported in literature. However, speaking with probability, not all of the g with f_{dg} has greater significance
141 in GWAS. Therefore, from all genes with predicted f_{dg} value, we obtain the top n genes according to their p value,
142 and observe the GWAS significance of this gene set over random set with the same size. The hypothesis test method
143 introduced for this case is Wilcoxon test, where the zero hypothesis H_0 is:

144 H_0 : “The top n genes with literature significance f_{dg} ranked the same with the other genes in GWAS.”

145 Generally, if the p -value obtained by Wilcoxon test is least than a threshold significance value, the H_0 hypothesis
146 will be rejected, and it is accepted that the top n genes with f_{dg} are more significant in GWAS associations if
147 compared with other genes in GWAS. As we hope to get a seed gene set with greater size, the applied strategy is to

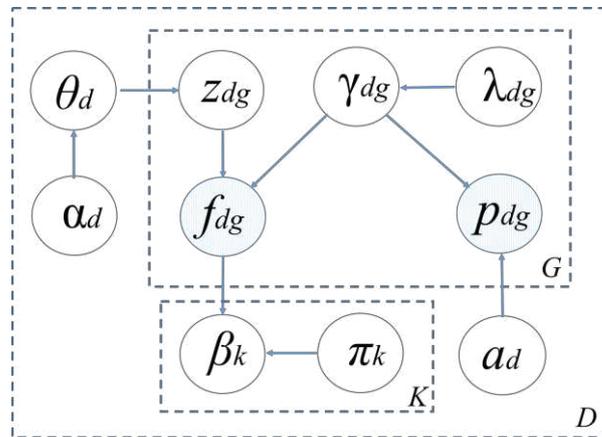
148 increase n gradually from 1 to the maximum. After increment of n , the size of the gene set increases, while p -value
 149 of Wilcoxon test specifies the advantage of the gene set over the whole. In most simulation tests, with the increase
 150 of n , the plot of $-\log p$ shows a bell shape. This plot suggests that there is a trade-off between the size of gene set and
 151 the overall significance over the whole. In the end, a peak value of $-\log p$ in Wilcoxon test corresponds a selection of
 152 proper value of n , which forms a synchronization filter.

153 2.3 “SNP-Gene Mapping” module in GDAMDB

154 For mapping SNP to specific genes. Firstly, we use Bedtools [15] to amplify SNP locations to left and right by 10kb
 155 base pairs. Bedtools is a fast, flexible tool set for genome arithmetic, and could be used for SNP flanking creating or
 156 calculating overlap of two sets of genomic features. Then, Bedtools is used to map those fragments onto the human
 157 genome by chromosomal location. Finally, we use the gene corresponding to the minimum p -value as the mapping
 158 result of the SNP.

159 2.4 Variational Inference on the solution of “Mutation Data Bridging” model in GDAMDB

160 The following mathematical setup defines the notations and symbols. For a disease d ($d = 1, \dots, D$) and a gene g
 161 ($g = 1, \dots, G$), $f_{dg} \in \{0,1\}^3$ encodes the associated *mutation type* of gene g for disease d , i.e., LOF/GOF/NA captured
 162 from literature, while $p_{dg} \in (0,1)$ refers to the p -value of the mapped *mutation association* of gene g for disease d in
 163 GWAS. Both of them are regarded as observations in the graphical model, and marked as dark circle in Figure 1. By
 164 introducing a latent variable $\gamma_{dg} \in \{0,1\}$, one switches the significance synchronization of the f_{dg} and p_{dg} . When γ_{dg}
 165 switch on, both *mutation type* and *mutation association* show significance; while it switches off, neither does. Thus,
 166 significance inconsistency is solved through this synchronization strategy.



167

168

Figure 1: The parameters setting of “Mutation Data Bridging” model.

169 Referred to Dai et al. [14], a Beta distribution with parameter a_d well cultivate a significant p -value. The greater
 170 the a_d , the least the p -value in GWAS. In the meantime, a uniform distribution leads to a non-significance of p -value.
 171 Thus, we have

$$172 \quad p_{dg} \sim \begin{cases} \text{Beta}(a_d, 1), & \text{if } \gamma_{dg} = 1; \\ U(0,1), & \text{if } \gamma_{dg} = 0. \end{cases} \quad (3)$$

173 We assume the switch variable $\gamma_{dg} \sim \text{Bernoulli}(\lambda_{dg})$, and $\lambda_{dg} \in (0,1)$. Intuitively, by considering the *mutation type*
 174 as event with full semantic augments, a sophisticated language model, Latent Dirichlet Allocation (LDA) [13], is
 175 introduced in this model which assumes the *mutation type* f_{dg} is generated by K latent topics, such as interaction,
 176 regulation, pathway, molecular or cell physiological activity, etc. Let $z_{dg} \in \{1, \dots, K\}$ index the latent topic, LDA
 177 assumes $z_{dg} \sim \text{Categorical}(\theta_d)$, $\theta_d \sim \text{Dir}(\alpha_d)$ where $\alpha_d \in \mathbb{R}^K$. Furthermore, the distribution of f_{dg} over the k -th
 178 topic is cultivated by a variable $\beta_k \sim \text{Dir}(\pi_k)$, $k = 1, 2, \dots, K$, where $\pi_k \in \mathbb{R}^3$. In the meantime, when γ_{dg} equals to
 179 zero, a zero vector is assigned to f_{dg} . So we have

$$180 \quad f_{dg} \sim \begin{cases} \text{Multi}(\beta_{z_{dg}}), & \text{if } \gamma_{dg} = 1; \\ \mathbf{0}, & \text{if } \gamma_{dg} = 0. \end{cases} \quad (4)$$

181 After applying variational inference in this model, the parameters of the distributions for p_{dg} and f_{dg} are
 182 computed by the iterative formulas in **Theorem 1** and **Theorem 2**, thus obtain the solutions of the model.

183 A brief introduction of the solutions is presented in Appendix B, while the complete proofs are shown in
 184 supplementary file: File S1. Proof of the generative model of Mutation Data Bridging. The solution of the ‘‘Mutation
 185 Data Bridging’’ model is based on **Theorem 1** and **Theorem 2**. For simplicity, all the parameters used in this model
 186 are shown as below:

187 Observation: $F = \{f_{dg}\}, P = \{p_{dg}\}$,

188 Latent variable: $\theta = \{\theta_d\}, \beta = \{\beta_k\}, \gamma = \{\gamma_{dg}\}, Z = \{z_{dg}\}$,

189 Model parameter: $\theta = \{\alpha_d, \pi_k, \lambda_{dg}, a_d\}$.

190 The first step of the variational inference is to derive an ELBO (Evidence lower bound) by using Jensen’s
 191 inequality to handle the logarithm of evidence $p(F, P)$.

$$\begin{aligned}
& \log p(\tilde{F}, P | \theta) \\
&= \log \int_{\theta} \int_{\beta} \sum_{\gamma} \sum_z p(\tilde{F}, P, \theta, \beta, \gamma, z | \theta) d\theta d\beta \\
192 \quad &= \log \int_{\theta} \int_{\beta} \sum_{\gamma} \sum_z \frac{p(\tilde{F}, P, \theta, \beta, \gamma, z | \theta)}{q(\theta, \beta, \gamma, z)} \cdot q(\theta, \beta, \gamma, z) d\theta d\beta \quad (5) \\
&\geq E_q[\log p(\tilde{F}, P, \theta, \beta, \gamma, z | \theta)] - E_q[q(\theta, \beta, \gamma, z)] \\
&:= \text{ELBO}.
\end{aligned}$$

193 The variational function 6, is represented by corresponding exponential family distributions.

$$194 \quad q(\theta, \beta, \gamma, z) = q(\theta)q(\beta)q(\gamma)q(z) \quad (6)$$

195 we assume $q(\theta_d) = \text{Dir}(\tilde{\alpha}_d)$ and $q(\beta_k) = \text{Dir}(\tilde{\pi}_k)$, both of which follow the Dirichlet distribution, while $q(\gamma_{dg}) =$
196 $\text{Bernoulli}(\tilde{\lambda}_{dg})$ follows Bernoulli distribution. In addition, $q(z_{dg}) = \text{Categorical}(\tilde{\theta}_{dg})$ follows a categorical
197 distribution.

198 Here, $\tilde{\alpha}$, $\tilde{\pi}$, $\tilde{\lambda}_{dg}$ and $\tilde{\theta}_{dg}$ are variational parameters under estimation. For simplicity, we denote $\tilde{\theta} =$
199 $\{\tilde{\alpha}, \tilde{\pi}, \tilde{\lambda}_{dg}, \tilde{\theta}_{dg} (d = 1, 2, \dots, D; g = 1, 2, \dots, G)\}$.

200 Variational parameters are estimated by differential computation of ELBO. And the solution of the model is given by
201 the following theorems.

202 **Theorem 1** Iteration formulas for variational parameters, $\tilde{\alpha}$, $\tilde{\pi}$, $\tilde{\lambda}$ and $\tilde{\theta}$, are:

$$\begin{aligned}
& \tilde{\alpha}_d^{(t+1)} = \alpha^{(t)} + \sum_{g=1}^G \tilde{\theta}_{dg}^{(t)} \\
& \tilde{\pi}_k^{(t+1)} = \pi_k^{(t)} + \sum_{g=1}^G \tilde{\theta}_{dg;k}^{(t)} \tilde{f}_{dg}^{(t)} \\
203 \quad & \tilde{\lambda}_{dg}^{(t+1)} = \text{sigmoid}(\log \frac{\lambda_{dg}^{(t)}}{1-\lambda_{dg}^{(t)}} a_d^{(t)} p_{dg}^{(a_d^{(t)}-1)}) \quad (7) \\
& \tilde{\theta}_{dg;k}^{(t+1)} \propto \exp(\sum_{f=1}^F \tilde{f}_{dg;f}^{(t)} (\psi(\tilde{\pi}_{k;f}^{(t)}) - \psi(\sum_{f=1}^F \tilde{\pi}_{k;f}^{(t)})) \\
& \quad + \psi(\tilde{\alpha}_{d;k}^{(t)}) - \psi(\sum_{f=1}^F \tilde{\alpha}_{d;k}^{(t)}))
\end{aligned}$$

204 **Theorem 2** Iteration for computing model parameters, α_d and π_k , is based on Newton's method, $\theta^{(n+1)} = \theta^{(n)} -$
205 $(\mathbf{H}f(\theta))^{-1} \cdot \nabla f(\theta)$, where $\mathbf{H}f(\theta)$ is the Hessian matrix and $\nabla f(\theta)$ is the gradient of $f(\theta)$. Then the Newton method
206 iteration of α_d and π_k is based on:

$$\begin{aligned}
\mathbf{HL}(\alpha_d)_{kj} &= D(\psi'(\sum_{k=1}^K \alpha_{d;k}) - \delta(k,j)\psi'(\alpha_{d;k})), \\
\nabla L(\alpha_d) &= D(\psi(\sum_{k=1}^K \alpha_{d;k}) - \psi(\alpha_{d;k})) \\
&\quad + \sum_{d=1}^D (\psi(\tilde{\alpha}_{d;k}) - \psi(\sum_{k=1}^K \tilde{\alpha}_{d;k})), \\
\mathbf{HL}(\pi_k)_{fj} &= \psi'(\sum_{f=1}^F \pi_{k;f}) - \delta(f,j)\psi'(\pi_{k;f}), \\
\nabla L(\pi_k) &= \psi(\sum_{f=1}^F \pi_{k;f}) - \psi(\pi_{k;f}) + \psi'(\tilde{\pi}_{k;f}) \\
&\quad - \psi'(\sum_{f=1}^F \tilde{\pi}_{k;f}).
\end{aligned} \tag{8}$$

$$\text{where } \delta(i,j) = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$$

In the meantime, the iteration of model parameters, λ_{dg} and a_d , is

$$\begin{aligned}
\lambda_{dg}^{(t+1)} &= \tilde{\lambda}_{dg}^{(t)}, \\
a_d^{(t+1)} &= -\sum_{g=1}^G \tilde{\lambda}_{dg}^{(t)} / \sum_{g=1}^G (\tilde{\lambda}_{dg}^{(t)} \log p_{dg})
\end{aligned} \tag{9}$$

2.5 Supplementary Files

There are three supplementary files with this research, the name of which are listed as below.

S1. Supplementary file. The whole proof of the mutation data bridging model, Proof of Mutation Data Bridging Model.pdf.

S2. Supplementary data. The result of Mutation type retrieval model: genes with predicted *mutation type* information and the evidence sentence from PubMed texts, Mutation Type Data.xlsx.

S3. Supplementary data. GWAS or literature evidence of 79 predicted AD genes, 79 Predicted Genes.xlsx. The file contains 4 sheets: Sheet 1 is the 12 genes that repeated from IGPA GWAS result, and the three columns are gene id, gene symbol, p-value from IGPA GWAS summary statistics; Sheet 2 is the 24 genes supported by other GWAS, and the four columns are gene id, gene symbol, p-value from the evidence GWAS data, the accession of the evidence GWAS; Sheet 3 are the 33 genes supported by literature, and the 5 columns are gene id, gene symbol, p-value from IGPA GWAS summary statistics, the evidence sentence in literature, the evidence type; Sheet 4 are the newly predicted 10 genes, and the 3 columns are gene id, gene symbol, p-value from IGPA GWAS summary statistics.

226 **3 Results**

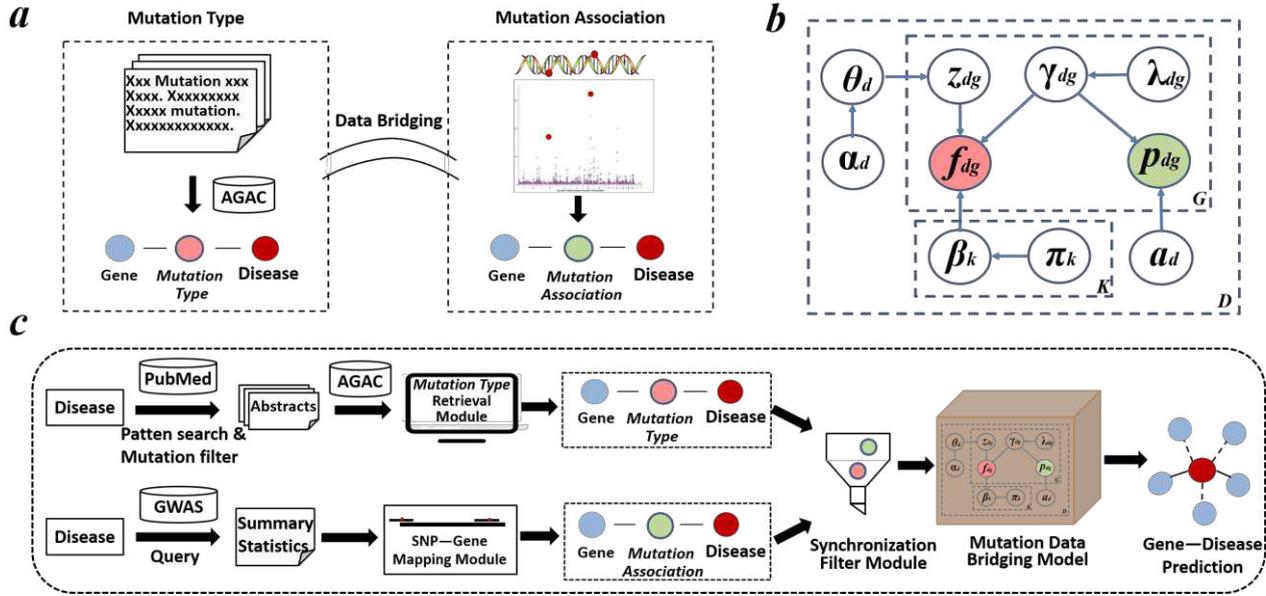
227 **3.1 Proposed Method “Gene-Disease Association Prediction by Mutation Data Bridging”**

228 **3.1.1 Main idea of bridging heterogeneous mutation data**

229 The two kinds of mutation data, *mutation association* and *mutation type*, share the commons but are with different
230 data characteristics, as shown in Figure 2 (a).

231 First, each mutation data comes from different resources. The *mutation association* is accessible from GWAS
232 summary statistics data, while the *mutation type* is available from literature mining. Second, they both support the
233 investigation of the gene-disease link, but the link is represented by p-value in *mutation association*, and by the
234 confidence value of LOF or GOF from the text mining module in *mutation type*. Third, the evidence of the link is under
235 different concerns. *Mutation association* data are from exacting and reliable experiments, and GWAS are widely
236 accepted as a powerful method to investigate the association between gene and disease. The *mutation type* is only
237 retrieved when a research report the gene are associated with the disease and also describe the mechanism between
238 the gene and disease in a published literature. Forth, both of them also has their weaknesses. For *mutation*
239 *association*, the high false negative rate and false negative rate in GWAS indicate that not all vital SNPs are able to
240 pass the multiple testing and not all passed genes are real important for the disease. In addition, the lack of
241 consideration about the biological mechanism between genes and phenotype makes GWAS insufficient to pinpoint
242 causal variations. Besides, it is difficult to conduct a GWAS on a large amount of case/control population. For
243 *mutation type*, since it comes from reported literature, it only represents part of the whole knowledge after text
244 mining.

245 Hence, considering the weaknesses and advantages of these two heterogeneous mutation data, we designed a
246 model to bridge *mutation association* and *mutation type* and achieves data fusion. Generally, a disease casual gene is
247 more likely to be identified by GWAS, and is also more likely to be discovered by other researches and described in
248 the literature. Therefore, bridging *mutation association* and *mutation type* is to integrate the mutation data in a
249 complementing way. In a simplified situation, if the *mutation association* of a significant SNP association failed to
250 pass or barely passed the threshold in GWAS, the *mutation type* of the gene helps to recover the association in the
251 manner of data fusion.



252

253 Figure 2: The heterogeneous mutation data and the pipeline of GDAMDB model. (a. The idea of bridging heterogeneous
 254 *mutation type* and *mutation association* by using data synchronization. b. Graphical model of the “Mutation Data
 255 Bridging” model. c. pipeline of GDAMDB model for gene-disease association prediction.)

256 3.1.2 Generative model bridges *mutation association* and *mutation type*

257 We designed a generative model by introducing a switch variable to bridge the *mutation association* and *mutation*
 258 *type* data. Here, the switch variable considers both the significance of *mutation association* mapped to the gene and
 259 the reported *mutation type* associated with the gene. Eventually, more reliable disease-related genes are predicted
 260 through the integration method.

261 As shown in Figure 2 (b), the parameter setting of f_{dg} follows the sophisticated language topic model, Latent
 262 Dirichlet Allocation (LDA) [13], which automatically organize words to form a readable text learns the probability
 263 distribution of each word in different latent topics and so as to generate a new text. By treating gene as the words,
 264 our model learns the probability distribution of each gene in different latent topics such as interaction, regulation,
 265 pathway, molecular or cell physiological activity, and so forth.

266 For a given disease d and a gene g , we denote *mutation type* and *mutation association* as f_{dg} and p_{dg} respectively,
 267 both of which are regarded as observations in a probability graph from a view of Bayesian statistics.

268 Referred to Dai et al. [14], a Beta distribution with parameter α_d well cultivate a significant p -value. The greater
 269 the α_d , the least the p -value in GWAS. In the meantime, a uniform distribution leads to a non-significance of p -value.

270 Thus we have

271
$$p_{dg} \sim \begin{cases} \text{Beta}(a_d, 1), & \text{if } \gamma_{dg} = 1; \\ U(0,1), & \text{if } \gamma_{dg} = 0. \end{cases} \quad (1)$$

272 We assume the switch variable $\gamma_{dg} \sim \text{Bernoulli}(\lambda_{dg})$, and $\lambda_{dg} \in (0,1)$.

273 As described above, f_{dg} can be generated by K latent topics in LDA. Let $z_{dg} \in \{1, \dots, K\}$ index the latent topic, LDA
 274 assumes $z_{dg} \sim \text{Categorical}(\theta_d)$, $\theta_d \sim \text{Dir}(\alpha_d)$ where $\alpha_d \in \mathbb{R}^K$. Furthermore, the distribution of f_{dg} over the k -th
 275 topic is cultivated by a variable $\beta_k \sim \text{Dir}(\pi_k)$, $k = 1, 2, \dots, K$, where $\pi_k \in \mathbb{R}^3$. In the meantime, when γ_{dg} equals to
 276 zero, a zero vector is assigned to f_{dg} . So we have

277
$$f_{dg} \sim \begin{cases} \text{Multi}(\beta_{z_{dg}}), & \text{if } \gamma_{dg} = 1; \\ \mathbf{0}, & \text{if } \gamma_{dg} = 0. \end{cases} \quad (2)$$

278 When a gene g mutated and the mutation related to disease d , γ_{dg} equals to 1. Decided by the value of γ_{dg} , the
 279 $p_{dg} \sim \text{Beta}(a_d, 1)$ which will mostly be a small value and the $f_{dg} \sim \text{Multi}(\beta_{z_{dg}})$ which will mostly be 1. The values
 280 are consisting with the facts that the gene should be significant in the GWAS result of this disease and the description
 281 about the gene mutation can be found in the literature of this disease.

282 After applying variational inference in this model, the parameters of the distributions for p_{dg} and f_{dg} are
 283 computed by the iterative formulas, thus obtain the solutions of the model. Definition of all the distribution
 284 parameters in the figure is provided in METHODS (Section 2.4), while a complete proof of the generative model and
 285 deduction can be found in the supplementary file, S1. Proof of Mutation Data Bridging Model.pdf.

286 3.1.3 Pipeline of “Gene-Disease Association prediction by Mutation Data Bridging”

287 The pipeline of “Gene-Disease Association prediction by Mutation Data Bridging (GDAMDB)” is shown in Figure 2(c),
 288 which consists of three data processing modules and one prediction model, i.e., “Mutation Type Retrieval” module,
 289 “SNP-Gene Mapping” module, “Synchronization Filter” module, and “Mutation Data Bridging” model.

290 (1) “Mutation Type Retrieval” is a text mining module, which is capable of jointly recognizing the gene and other
 291 entities from literature and classifying the *mutation type* of the genes based on the semantic in the sentence.
 292 The module is designed based on a Bidirectional Encoder Representations from Transformers (BERT) model
 293 which is released by Google and trained on a large text corpus. BERT contains 12 layers of deep bidirectional
 294 transformer, and each of them is fully connected and made up of a multi-head self-attention, residual connection,
 295 batch normalize. The complex construction makes BERT able to learn deep bidirectional representation of each

296 word. Therefore, we transfer and fine-tuned its parameters on our joint model. As shown in Figure 2 (c), after
297 paternal searching from PubMed and mutation filtering, the abstracts containing diseases and mutations are
298 input into BERT, then the presentations of each word in the abstracts are obtained. Subsequently, a fully
299 connected layer and softmax are used to normalize classification weights, and CRF loss function is employed to
300 optimize entity recognition task in the meantime. Finally, the model output the *mutation type* of genes in
301 abstracts.

302 (2) "SNP-Gene Mapping" module is to process the *mutation association* data. Since the pipeline is focus on gene, the
303 SNPs in GWAS data should be mapped on genes by bedtools [15]. The p-value of a gene is the p-value of its SNP
304 which is the lowest one.

305 (3) "Synchronization Filter" module designs to optimize the probability that most genes in the gene set maintain
306 not only literature significance but also the GWAS significance. Generally, the gene with great significant
307 *mutation association*, p_{dg} , in GWAS is likely to be described in literature with *mutation type*, f_{dg} , but not all of
308 genes satisfy the rule. Therefore, from all genes with predicted f_{dg} value, we obtain the top n genes according
309 to their p value, and observe the GWAS significance of this gene set over random set with the same size. The
310 hypothesis test method introduced for this case is Wilcoxon test.

311 (4) "Mutation Data Bridging" model is the generative model mentioned above which bridges *mutation association*
312 data and *mutation type* data by introducing a switch variable. The inputs of this model are the two processed
313 mutation data and the gene set selected by "Synchronization Filter" module. Then the distributions of the two
314 mutation data are computed by the model, and the predicted gene are obtained based on the switch variable of
315 each gene.

316 Briefly speaking, f_{dg} is retrieved from PubMed by using the "Mutation Type Retrieval" model and AGAC corpus,
317 and genes with significant f_{dg} are defined as "literature significant" genes. In the meantime, *mutation association*,
318 p_{dg} , is extracted from GWAS summary data by applying SNP inclusion criteria, and "GWAS significant" genes are
319 obtained by using a "SNP-Gene Mapping" module. In order to better synchronize the above heterogeneous mutation
320 data, a "Synchronization Filter" module creates a seed gene set consists of g with significant \hat{f}_{dg} and p_{dg} . After
321 feeding the observations, $\{\hat{f}_{dg}\}$ and $\{p_{dg}\}$, into the "Mutation Data Bridging" model, the model parameters are
322 obtained. Eventually, generative process is carried on to produce novel *mutation types* with significant \tilde{f}_{dg} . Thus,

323 new appeared gene g with \tilde{f}_{ag} is predicted with novel gene-disease association. All of the pipeline details are
324 elucidated in Online Method.

325 The purpose of GDAMDB is to accelerate the discovery rate of the gene associations of GWAS by integrating both
326 *mutation association* and *mutation type* information. A case study on Alzheimer's disease (AD) was carried on to
327 evaluate the performance of GDAMDB in the support of discovery of novel gene-disease associations.

328 **3.2 Application of GDAMDB on Alzheimer's Disease**

329 Alzheimer's Disease is a common neurodegenerative disorder, which impairs the memory, language and various
330 body behaviors. Till now, there are 116 AD studies in GWAS Catalog [16], three of which provide the summary
331 statistics data. We chose one of the studies that provided the most complete data. Although no database recording
332 the *mutation type* info of AD-related genes, there are lots of literatures that report the studies of AD pathogenesis.
333 The *mutation type* of the genes is widely implied in the description of the literatures. Since AD is a important diseases
334 and the mutation data of AD are available, we apply GDAMDB on AD to retrieve the genes that are undiscovered.

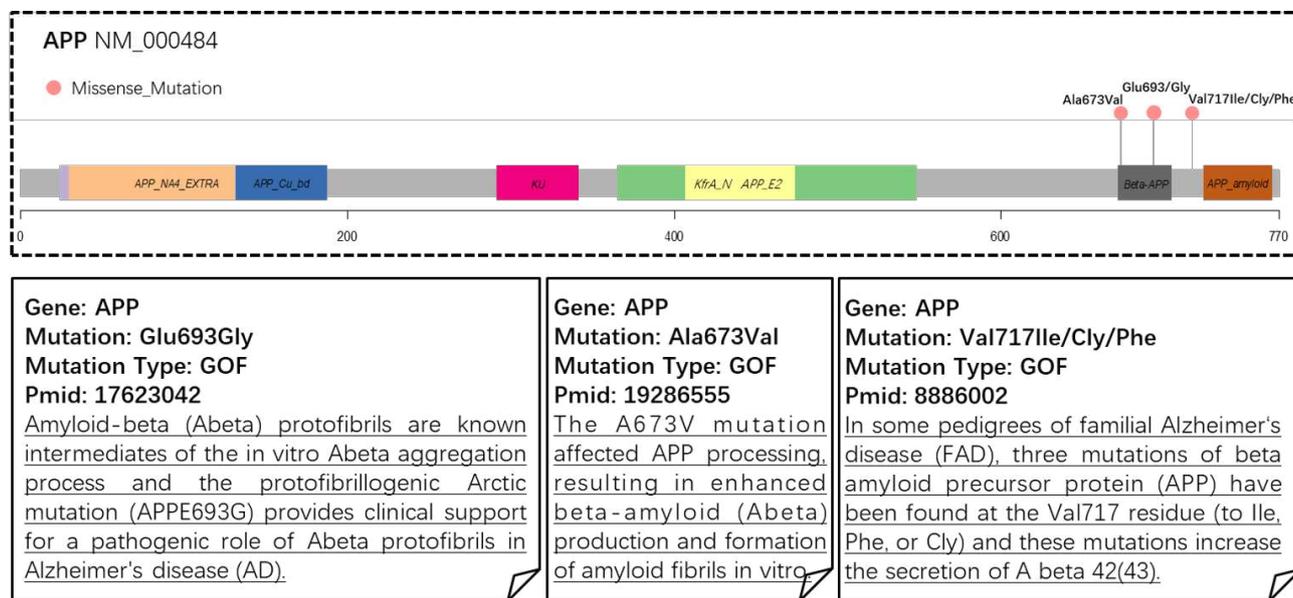
335 **3.2.1 Mutation association data of Alzheimer's disease**

336 In 2013, the International Genomics of Alzheimer's Project (IGAP) [12] performed a two-stage GWAS on individuals
337 of European ancestry on 7,055,881 SNPs. In stage 1, they meta-analyzed four previous AD GWAS datasets including
338 17,008 AD cases and 37,154 controls. In stage 2, they tested 211,632 SNPs on 8,572 AD cases and 11,312 controls.
339 The final result was obtained after meta-analysis of combining stage 1 and 2. We selected the summary statistics file
340 that combined stage 1 and 2, in which contains 1,513 genes. The p-value of each gene was same with the most
341 significant SNP in the gene.

342 **3.2.2 Mutation type data of Alzheimer's disease**

343 In the meantime, the MeSH term "Alzheimer's disease" was used as the key word to query PubMed database, and
344 137,473 abstracts were downloaded. To ensure that the literatures contain description about mutation, SETH [17]
345 was applied to filter the literature. SETH is able to recognize the SNP or other mutation semantic words in texts.
346 Thus, till this procedure, the abstracts that contains AD and mutations were left. After mutation filtering, 9,430
347 abstracts were input into "*Mutation Type Retrieval*" module. The module will compute the confidence value for each
348 abstract in each *mutation type*. The output of the module is the *mutation type* of genes in each abstract, of which the
349 confidence value passes the module threshold. Subsequently, we manually checked the result, and only preserved
350 the abstracts that clearly describe the mechanism of a mutated gene leading to AD. Finally, 65 genes with their

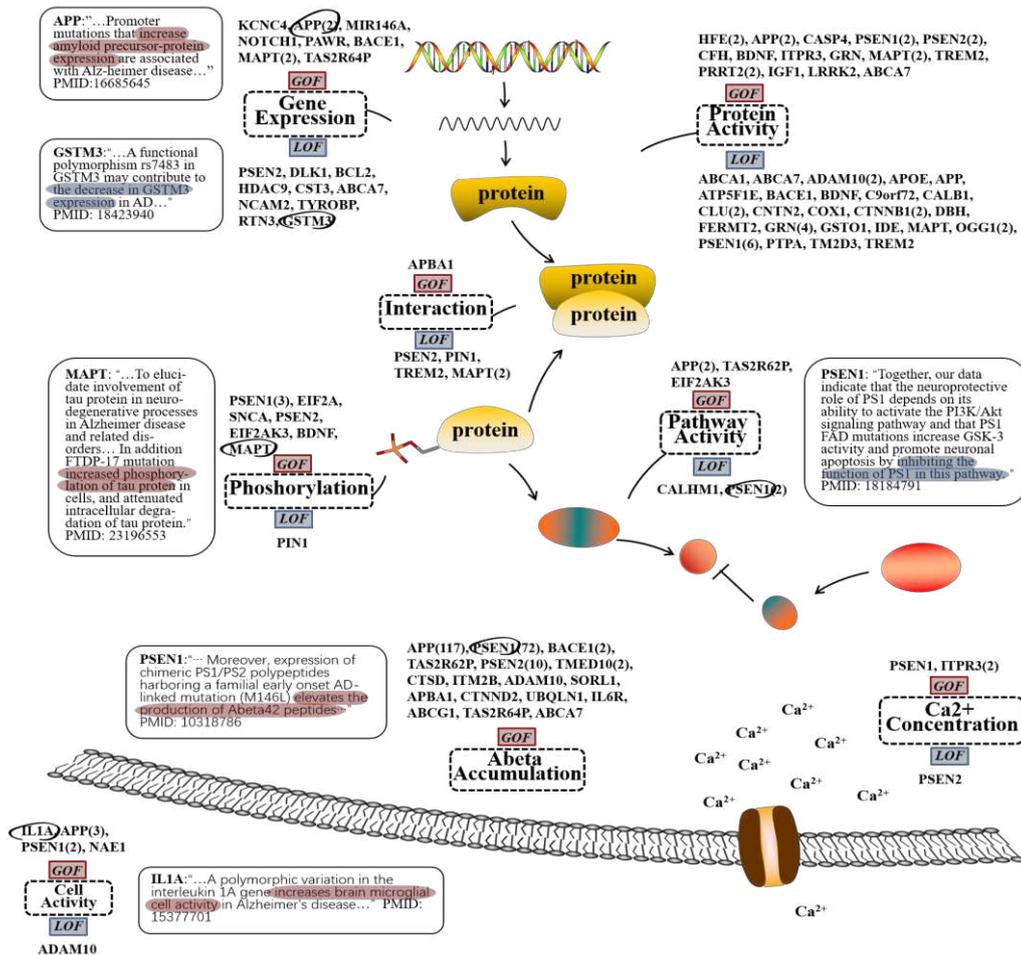
351 *mutation types* are firstly recognized from 325 abstracts where each abstract contains a conclusive sentence
 352 evidence leading to a *mutation type*. It is noted that the obtained plenty of AD-related LOF/GOF data is new to the
 353 AD community, while the full result is offered in supplementary data, S2: genes with predicted *mutation type*
 354 information and the evidence sentence from PubMed texts, Mutation Type Data.xlsx.



355
 356 Figure 3: The illustrated examples of *mutation types* retrieved from Glu693Gly, Ala673Val, and Val717Ile/Cly/Phe in
 357 APP.

358 56 abstracts clearly described the amino acid change of the mutations or the rs number of the SNPs, which totally
 359 report 64 mutations on 28 genes. Among the 64 mutations, 54 of them locate on the coding regions of the genes
 360 leading to the amino changes, and 10 of them locate on the non-coding region of the genes. Figure 3 is an example,
 361 which shows the specific mutations of APP recognized by the module. These five mutations are on three locations of
 362 APP, two of which locate on the Beta-APP domain of APP protein. The amino acid location 717 is found three
 363 mutations, and this location is between the sequence producing Beta-APP domain and sequence producing APP
 364 amyloid domain which form the beta-amyloid and is strongly implicated in the pathogenesis of AD. Moreover, the
 365 corresponding sentence evidences of these APP mutations are below. As introduced above, this module is able to
 366 recognize the entities and classify the *mutation type* of a gene. For example, the sentence in the middle, "The A673V
 367 mutation affected APP processing, resulting in enhanced beta amyloid (Abeta) production and formation of amyloid
 368 fibrils in vitro.", APP and enhanced will be recognized by the module. Based on enhanced, the confidence value of GOF
 369 will be higher than the value of LOF in this abstract, hence the *mutation type* of APP will be classified as GOF.
 370 Therefore, among the 325 abstracts, each one are recognized at least one gene and their *mutation types*. Besides that,
 371 all the 325 abstracts carry the clear semantic of the downstream biological processes after mutation, which can be

372 divided into 8 types after manual curation. As shown in figure 4, Gene Expression, Protein Activity, Interaction,
 373 Pathway Activity and Cell Activity are the fundamental biological processes which follows the central dogma and are
 374 from molecular level to cell level. In addition, the Phosphorylation, Abeta Accumulation and Ca²⁺ Concentration are
 375 frequently mentioned. Interestingly, these three biological processes are related to the known hypothesizes of AD
 376 pathogenesis. Abeta is the production of APP gene, the accumulation of which, especially Abeta42, forms the fibrillar
 377 amyloid plaques in brain and impair the ability of spatial learning and memory [18]. Phosphorylation related to
 378 another hypothesis of AD pathogenesis, especially the phosphorylation of Tau protein which encoded by MAPT gene.
 379 The hyperphosphorylation of Tau protein leads to neurofibrillary tangles in neurons and eventually results in the
 380 apoptosis of neurons [19]. Intracellular Ca²⁺ concentration is also thought as part of the cause of AD. The
 381 dysregulation of intracellular Ca²⁺ signaling disturbs many neural processes, which implicated in AD mechanism
 382 [20].



383

384

Figure 4: Biological process categories of 325 mutation types and the sentence evidences.

385

386

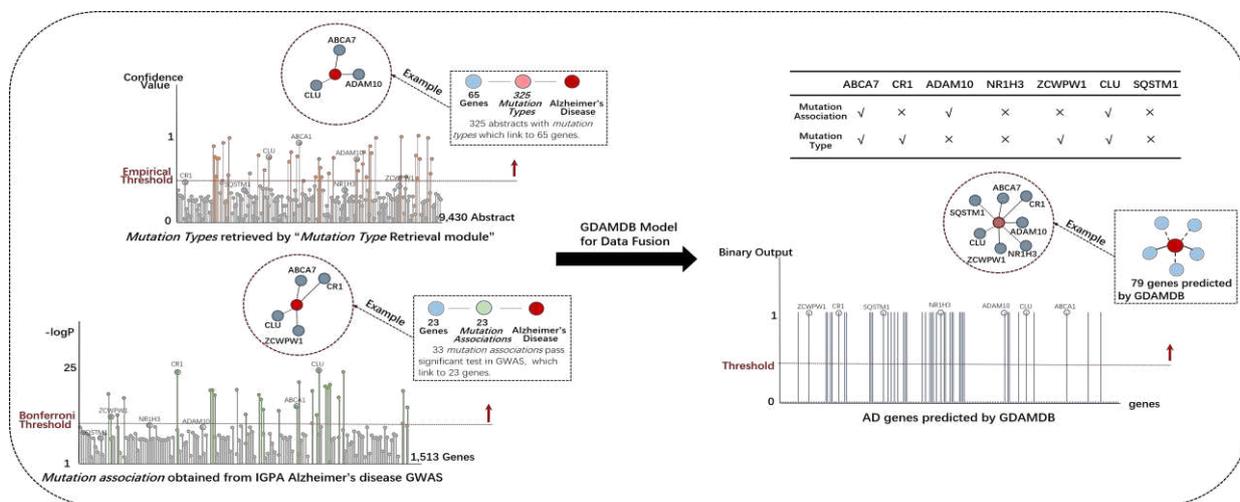
As shown in the upper left of the figure 4, it presents the genes that related to the gene expression biological process. the mutations in the 8 genes, KCNC4, APP, MIR146A, NOTCH1, PAWR, BACE1, MAPT, TAS2R64P, are GOF

387 mutation and leads to the increase of gene expression. In addition, the increase of expression of APP and MAPT are
 388 mentioned in two abstracts. 10 genes are mentioned LOF mutation that decrease the gene expression in abstracts.
 389 At the left side of the genes are two sentence examples. In the sentence, "...Promoter mutations that increase amyloid
 390 precursor-protein expression are associated with Alzheimer disease...", "increase" helps to confirm GOF and "amyloid
 391 precursor protein expression" helps to confirm that the biological process that effected by mutation is gene
 392 expression. Similarly, GSTM3 is grouped into the LOF of gene expression.

393 The biological process category of the genes and evidence sentence can be found in supplementary data S2, or
 394 in online data repository (<https://hzaubionlp.com/agac-on-alzheimers-disease/>)

395 3.2.3 Data fusion of heterogeneous mutation data

396 The data fusion by GDAMDB is shown on Figure 4. The left two graphs present confidence value of the gene *mutation*
 397 *type* in each abstract and the p-value of gene *mutation association*, both of the graphs are showing the rough
 398 distributions of data. In *mutation type* graph, the 325 gene-*mutation types*-AD are predicted and manually checked
 399 from 9,430 abstracts, and there are 65 unique genes since some of the genes are mentioned more than once in these
 400 abstracts. The empirical threshold represents the model parameters and human filtering. In *mutation association*
 401 graph, there are 23 *mutation associations* passed the final Bonferroni threshold, and they are mapped to 23 genes.
 402 The graph became bar graph after data fusion, since the output of the model is binary info representing the
 403 association between gene and disease or not. 79 genes are predicted to be the AD-related genes. The final prediction
 404 filtered some of the genes that passed the threshold in the single mutation data but recognized as the false negative
 405 genes by the model, and also retrieved the genes failed to pass the threshold in the single mutation data but
 406 recognized as the AD-related genes.



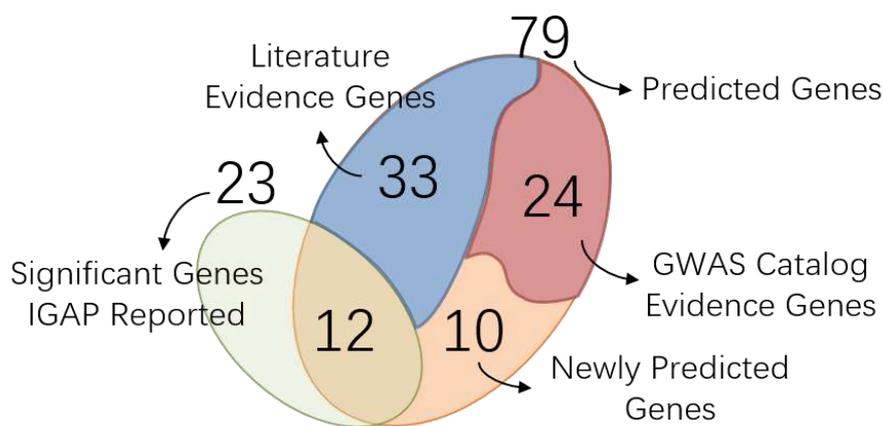
407
 408 Figure 5: Data fusion of heterogeneous AD mutation data improves the discovery of novel AD-related genes.

409 For example, as shown in the circle above *mutation types* graph, the *mutation types* of the three genes are
 410 retrieved by “mutation type retrieval module” and passed the empirical threshold, ABCA7, CLU and ADAM10. The
 411 circle above the *mutation association* graph contains four genes that passed the Bonferroni threshold, ABCA7, CLU,
 412 CR1 and ZCWPW1. There are different limitations make the information that mutation data contained is incomplete.
 413 Therefore, as marked on the graphs, ABCA7 and CLU both pass the threshold in two kinds of mutation data, but
 414 ADAM10, ZCWPW1 and CR1 only pass one. However, the -logp value of ADAM10 is close to the Bonferroni threshold,
 415 while the confidence value of ZCWPW1 and CR1 are close to the empirical threshold.

416 After data fusion, ABCA7, CLU, ADAM10, CR1 and ZCWPW1 are output by GDAMDB, which shows that GDAMDB
 417 is able to break the limitation of these two mutation data and save the important genes that are failed to pass the
 418 threshold. Besides, the genes, NR1H3 and SQSTM1, are retrieved in neither *mutation type* data nor *mutation*
 419 *association* data, but retrieved by GDAMDB after data fusion. It shows that GDAMDB is not simply merge the genes
 420 that are significant in one of the mutation data, but to learn the latent regularity of the mutation data distribution.

421 3.2.4 Novel discovery of AD-related genes after heterogeneous mutation data fusion

422 An encouraging result of AD-related gene discovery is shown in Figure 6. The left ellipse refers to the significant gene
 423 set which is reported in IGAP GWAS research [12]. Meanwhile, the right ellipse represents 79 genes that are
 424 predicted by GDAMDB.



425

426 Figure 6: 69 out of 79 predicted genes have supportive AD-related evidence.

427 As shown in Figure 6, 12 out of 79 predicted genes are reported in IGAP GWAS research, which are CR1, CD2AP,
 428 EPHA1, CLU, PICALM, ABCA7, HLA-DBR1, PTK2B, SORL1, INPP5D, ZCWPW1 and FERMT2. The red part of the big
 429 ellipse contains 24 genes with GWAS catalog evidences, which are reported to be AD-related in GWAS catalog

430 (https://www.ebi.ac.uk/gwas/efotraits/EFO_0000249). They are RNU6560P, GULOP, EPHA1-AS1, STAG3L5P-
431 PVRIG2P-PILRB, HLA-DQA1, GPR141, ADGRF2, SCARA3, SPI1, CSTF1, AP4M1, SCIMP, PILRA, EPDR1, RAPSIN,
432 MS4A6E, PSMC3, CASTOR3 and MS4A2 from GCST5922, HLA-DRB9, ADAM10 from GCST007320, MEF2C-AS1 from
433 GCST003427, NDUFAF6 from GCST009021, and MADD from GCST007825.

434 As for the other 43 genes which are not recorded in GWAS catalog as AD-related genes, 33 of them are reported
435 to be AD-related genes with confirmed literature evidence. In details, NYAP1, MYBPC3, BTNL2, HLA-DQB1, AGBL2,
436 COPS6, ZKSCAN1, MTCH2, MCM7, NDUFS3, PILRB, SLC39A13, USP50 and ACP2 [21] are identified as the additional
437 genes within the loci which contain AD-related SNP, while ADGRF4 [22] and RNU6-603P [23] are the gene adjacent
438 to AD-related SNP. MS4A4E [24], SLC25A1P1 [25], MBLAC1 [26], USP8 [27], TSPOAP1 [28], PGF [29], NR1H3 (also
439 named as LXR α) [30], TP53INP1 [31] and TRIP4 [32] are reported in AD GWAS study or the research that identified
440 new genes by performing further experiment on AD GWAS data. There are 5 microRNAs in the prediction gene set.
441 MicroRNA is known to regulate neuronal development, and the abnormal alteration of which contributes to
442 neurodegenerative disorders. The 5 microRNA are also found to be related to AD or neurodegenerative disease,
443 which are MIR-4487 [33], MIR142 [34], MIR25[35], MIR4736 [36] and MIR106B [37]. Besides, MAML1 [38], CR2
444 [39] and SQSTM1 [40] are reported to be related to neurodegenerative disease.

445 The rest of the 10 genes are all pseudogenes, which are IGHVIII-67-4, MTCO3P1, RNA5SP340, RN7SKP116,
446 SNORD3P2, SNORD3P3, SNORD3P1, RNU5E-10P, IGHV3-71 and IGHV2-70. Although there is no direct GWAS Catalog
447 evidence or literature evidence showing that these genes are associated with AD, the regulation function of the
448 pseudogenes in the pathogenesis of complex disease is worth attention. In 2011, Tay et al. [41] firstly proposed the
449 hypothesis that the competing endogenous RNAs (ceRNAs) may compete microRNA with the mRNA that contains
450 same microRNA response element with the ceRNAs. Costa et al. [42] suggested that ceRNAs may be involved in the
451 pathogenesis of neurodegenerative disease, such as AD.

452 Among the 11 genes that are reported in IGAP but aren't retrieved by GDAMDB, most of them are less reported
453 by GWAS Catalog. Especially, 4 of them are only reported once by IGAP in whole GWAS Catalog.

454 In summary, among the 79 predicted genes, the genes reported by IGAP, the genes with AD GWAS catalog
455 evidences, and the genes with AD literature evidences accounts for 86% (69/79), and the rest 14% (10/79)
456 predicted genes are suggested to be potentially involved in the pathogenesis mechanism of AD.

457 **4 Discussion**

458 In the case of data fusion in terms of knowledge discovery, the knowledge can be any form of data with different
459 format. In our research, the association relation between gene and disease can be the p-value, named as *mutation*
460 *association*, in GWAS, where the smaller p-value represents the more significant relevance between gene and disease.
461 Furthermore, this association relation can also be *mutation type* in literature, where the description about the
462 mechanism of mutations in disease pathogenesis directly indicate the details of the relation. When different data
463 reveal the relations in different aspects, taking both aspects into consideration leads to a more comprehensive
464 knowledge discovery. Besides, since the advantages and weakness vary from heterogeneous data, data fusion helps
465 to enhance the quality of both data. The relevance between a gene and a disease is adjusted after data fusion,
466 especially when a false negative *mutation association* of a gene fails to pass significant test in GWAS but is found to
467 be active with *mutation type* information in literature.

468 A generative model is capable of learning the data distribution of observations from two heterogeneous
469 categories, and generating novel data which represents the statistical characteristics of both observations, thus
470 achieves the data fusion of heterogeneous data. Therefore, by bridging *mutation type* data and *mutation association*
471 data, GDAMDB is capable of retrieving the important AD-related genes that are failed to pass the multiple testing in
472 GWAS or haven not been reported in literature. Eventually, our model retrieved 79 AD genes, and 57 of them are not
473 reported in the source GWAS study but 47 out of 57 are supported by convince evidences that are AD-related genes,
474 which positively shows the reliability of the model performance.

475 As a generative model, GDAMDB offers a way to enhance the disease-related gene discovery in a single mutation
476 data, and the implementation procedure of the model shows that the model is flexible to be adopted to each given
477 disease, in the case when the GWAS summary data and sufficiently abundant literature are available. All the results
478 in this research indicate that data fusion sheds light to the novel knowledge discovery.

479 **5 Conclusion**

480 This research drew a novel respective towards the data form of mutations, of which there are mutation associations
481 obtained from GWAS experiment and mutation type extracted from text mining. It is known that GWAS associations
482 are under-powered for pinpointing causal genes due to high false positive/negative rate, and integration of other
483 mutation information is possibly an effective addition. Thus, we used a PubMed-wide text mining strategy to
484 pinpoint vital genes which carry core semantics of the mutation effect, and came up with the mining of *mutation*
485 *type*, which associated gene and disease in an interpretable LOF/GOF way.

486 GDAMDB is a model to bridge the two heterogeneous types of mutation data from a same disease. This model
487 designs a switch variable γ_{dg} to synchronize the importance of two types of mutation data, learns the distribution
488 of both data, and discover novel significant genes that potentially related to disease.

489 The case study in AD made use of real GWAS data from IGAP and went through a thorough data integration via
490 applying GDAMDB. Finally, 57 out of 79 predicted genes are closely related to AD. The results obtained in this
491 research fully showed that bridging the heterogeneous mutation data integrate information from GWAS and
492 literature, thus shed lights on novel disease-related gene discovery.

493 **6 Abbreviations**

494 Abbreviations used in this paper are listed as below.

- 495
- 496 AD, Alzheimer's disease
- 497 AGAC, active gene annotation corpus
- 498 APP, amyloid precursor protein
- 499 BERT, bidirectional encoder representations from transformers
- 500 GDAMDB, gene-disease association prediction by mutation data bridging
- 501 GOF, gain of function
- 502 GWAS, genome-wide association study
- 503 IGAP, international genomics of Alzheimer's project
- 504 LDA, latent Dirichlet allocation
- 505 LOF, loss of function
- 506 SNP, single nucleotide polymorphism

507

508 **7 Declaration**

509 **7.1 Ethics approval and consent to participate**

510 This research does not involve human or animals.

511 **7.2 Consent for publication**

512 All the authors have consented for the publication.

513 **7.3 Availability of data and material**

514 Data and materials are available in the supplementary files.

515 **7.4 Competing interests**

516 None of the authors have any competing interests.

517 **7.5 Funding**

518 The research is supported by Hubei Province Funds for Natural Science (No. 2019CFB552).

519 **7.6 Authors' contributions**

520 KZ developed the algorithm and coding. YW checked the biological knowledge and wrote the manuscript. KBC and
521 JDK took part in the algorithm discussion and the manuscript writing. XM took part in the algorithm development.
522 ZS curated the biological data. XM checked the biological data. JX designed the whole pipeline and took
523 responsibility of the whole research. All authors read and approved the final manuscript.

524 **7.7 Acknowledgements**

525 The authors would like to express their gratitude to all HZAU BioNLP teams members who joined many discussions
526 related to this research.

527

528 **References**

- 529 [1] P. Jia, S. Zheng, J. Long, W. Zheng, and Z. Zhao. dm-gwas: dense module searching for genome-wide association
530 studies in protein-protein interaction networks. *Bioinformatics*, 27(1):95–102.
- 531 [2] Andrea Califano, Atul J Butte, Stephen Friend, Trey Ideker, and Eric Schadt. Leveraging models of cell regulation
532 and gwas data in integrative network-based association studies. *Nature genetics*, 44(8):841, 2012.
- 533 [3] Kai Wang, Mingyao Li, and Hakon Hakonarson. Analysing biological pathways in genome-wide association
534 studies. *Nature Reviews Genetics*, 11(12):843–854.
- 535 [4] Yongsheng Li, Yunpeng Zhang, Xia Li, Song Yi, and Juan Xu. Gain-of-function mutations: an emerging advantage
536 for cancer biology. *Trends in biochemical sciences*, 2019.
- 537 [5] Zhong Yi Wang and Hong Yu Zhang. Rational drug repositioning by medical genetics. *Nature biotechnology*,
538 31(12):1080, 2013.
- 539 [6] Martin Citron, David Westaway, Weiming Xia, George Carlson, Thekla Diehl, Georges Levesque, Kelly Johnson-
540 Wood, Michael Lee, Peter Seubert, Angela Davis, et al. Mutant presenilins of alzheimer's disease increase
541 production of 42-residue amyloid β -protein in both transfected cells and transgenic mice. *Nature medicine*,
542 3(1):67–72, 1997.
- 543 [7] Rita Guerreiro, Aleksandra Wojtas, Jose Bras, Minerva Carrasquillo, Ekaterina Rogava, Elisa Majounie, Carlos
544 Cruchaga, Celeste Sassi, John SK Kauwe, Steven Younkin, et al. Trem2 variants in alzheimer's disease. *New
545 England Journal of Medicine*, 368(2):117–127, 2013.

- 546 [8] Stacy Steinberg, Hreinn Stefansson, Thorlakur Jonsson, Hrefna Johannsdottir, Andres Ingason, Hannes Helgason,
547 Patrick Sulem, Olafur Th Magnusson, Sigurjon A Gudjonsson, Unnur Unnsteinsdottir, et al. Loss-of-function
548 variants in *abca7* confer risk of alzheimer's disease. *Nature genetics*, 47(5):445–447, 2015.
- 549 [9] Yuxing Wang, Kaiyin Zhou, Jin Dong Kim, Kevin Bretonnel Cohen, Mina Gachloo, Yuxin Ren, Shanghui Nie, Xuan
550 Qin, Panzhong Lu, and Jingbo Xia. An active gene annotation corpus and its application on anti-epilepsy drug
551 discovery. *BIBM 2019: International Conference on Bioinformatics & Biomedicine, San Diego, U.S.*, 2019.
- 552 [10] Dennis J Selkoe. Translating cell biology into therapeutic advances in alzheimer's disease. *Nature*,
553 399(6738):A23–A31, 1999.
- 554 [11] Robert S Wilson, Sandra Barral, Joseph H Lee, Sue E Leurgans, Tatiana M Foroud, Robert A Sweet, Neill Graff-
555 Radford, Thomas D Bird, Richard Mayeux, David A Bennett, et al. Heritability of different forms of memory in
556 the late onset alzheimer's disease family study. *Journal of Alzheimer's Disease*, 23(2):249–255, 2011.
- 557 [12] Jean Charles Lambert, Carla A Ibrahim Verbaas, Denise Harold, Adam C Naj, Rebecca Sims, C'eline Bellenguez,
558 Gyungah Jun, Anita L DeStefano, Joshua C Bis, Gary W Beecham, et al. Meta-analysis of 74,046 individuals
559 identifies 11 new susceptibility loci for alzheimer's disease. *Nature genetics*, 45(12):1452, 2013.
- 560 [13] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning*
561 *research*, 3(Jan):993–1022, 2003.
- 562 [14] Mingwei Dai, Jingsi Ming, Mingxuan Cai, Jin Liu, Can Yang, Xiang Wan, and Zongben Xu. Iguess: a statistical
563 approach to integrating individual-level genotype data and summary statistics in genome-wide association
564 studies. *Bioinformatics*, 33(18):2882–2889, 2017.
- 565 [15] Aaron R Quinlan. Bedtools: the swiss-army tool for genome feature analysis. *Current protocols in bioinformatics*,
566 47(1):11–12, 2014.
- 567 [16] Annalisa Buniello, Jacqueline A L MacArthur, Maria Cerezo, Laura W Harris, James Hayhurst, Cinzia Malangone,
568 Aoife McMahon, Joannella Morales, Edward Mountjoy, Elliot Sollis, et al. The nhgri-ebi gwas catalog of published
569 genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic acids research*,
570 47(D1):D1005–D1012, 2019.
- 571 [17] Philippe Thomas, Tim Rocktäschel, Jörg Hakenberg, Yvonne Lichtblau, and Ulf Leser. Seth detects and
572 normalizes genetic variants in text. *Bioinformatics*, 32(18):2883–2885, 2016.
- 573 [18] John Hardy and David Allsop. Amyloid deposition as the central event in the aetiology of alzheimer's disease.
574 *Trends in pharmacological sciences*, 12:383–388, 1991.

- 575 [19] Wanjoo Chun and GV Johnson. The role of tau phosphorylation and cleavage in neuronal cell death. *Front Biosci*,
576 12(733):e56, 2007.
- 577 [20] Frank M LaFerla. Calcium dyshomeostasis and intracellular signalling in alzheimer's disease. *Nature Reviews*
578 *Neuroscience*, 3(11):862–872, 2002.
- 579 [21] Anastasia G Efthymiou and Alison M Goate. Late onset alzheimer's disease genetics implicates microglial
580 pathways in disease risk. *Molecular neurodegeneration*, 12(1):43, 2017.
- 581 [22] Anna A Pimenova, Towfique Raj, and Alison M Goate. Untangling genetic risk for alzheimer's disease. *Biological*
582 *psychiatry*, 83(4):300–310, 2018.
- 583 [23] Brian W Kunkle, Benjamin Grenier-Boley, Rebecca Sims, Joshua C Bis, Adam C Naj, Anne Boland, Maria
584 Vronskaya, Sven J van der Lee, Alex Amlie-Wolf, Celine Bellenguez, et al. Meta-analysis of genetic association
585 with diagnosed alzheimer's disease identifies novel risk loci and implicates abeta, tau, immunity and lipid
586 processing. *bioRxiv*, page 294629, 2018.
- 587 [24] Paul Hollingworth, Denise Harold, Rebecca Sims, Amy Gerrish, Jean Charles Lambert, Minerva M Carrasquillo,
588 Richard Abraham, Marian L Hamshere, Jaspreet Singh Pahwa, Valentina Moskva, et al. Common variants at
589 *abca7*, *ms4a6a/ms4a4e*, *epha1*, *cd33* and *cd2ap* are associated with alzheimer's disease. *Nature genetics*,
590 43(5):429, 2011.
- 591 [25] Shubhabrata Mukherjee, Joshua C Russell, Daniel T Carr, Jeremy D Burgess, Mariet Allen, Daniel J Serie, Kevin L
592 Boehme, John SK Kauwe, Adam C Naj, David W Fardo, et al. Systems biology approach to lateonset alzheimer's
593 disease genome-wide association study identifies novel candidate genes validated using brain expression data
594 and *caenorhabditis elegans* experiments. *Alzheimer's & Dementia*, 13(10):1133–1142, 2017.
- 595 [26] Iris J Broce, Chin Hong Tan, Chun Chieh Fan, Iris Jansen, Jeanne E Savage, Aree Witoelar, Natalie Wen,
596 Christopher P Hess, William P Dillon, Christine M Glastonbury, et al. Dissecting the genetic relationship between
597 cardiovascular risk factors and alzheimer's disease. *Acta neuropathologica*, 137(2):209–226, 2019.
- 598 [27] Eniola Funmilayo Aduke Yeates and Giuseppina Tesco. The endosome-associated deubiquitinating enzyme
599 *usp8* regulates *bace1* enzyme ubiquitination and degradation. *Journal of Biological Chemistry*, 291(30):15753–
600 15766, 2016.
- 601 [28] Gyungah R Jun, Jaeyoon Chung, Jesse Mez, Robert Barber, Gary W Beecham, David A Bennett, Joseph D Buxbaum,
602 Goldie S Byrd, Minerva M Carrasquillo, Paul K Crane, et al. Transethnic genome-wide scan identifies novel
603 alzheimer's disease loci. *Alzheimer's & Dementia*, 13(7):727–738, 2017.

- 604 [29] Pu-Ting Xu, Yiju Li, Xue-Jun Qin, Clemens R Scherzer, Hong Xu, Donald E Schmechel, Christine M Hulette, John
605 Ervin, Steven R Gullans, Jonathan Haines, et al. Differences in apolipoprotein e3/3 and e4/4 allele-specific gene
606 expression in hippocampus in alzheimer disease. *Neurobiology of disease*, 21(2):256–275, 2006.
- 607 [30] Noam Zelcer, Negar Khanlou, Ryan Clare, Qingguang Jiang, Erin G Reed-Geaghan, Gary E Landreth, Harry V
608 Vinters, and Peter Tontonoz. Attenuation of neuroinflammation and alzheimer’s disease pathology by liver x
609 receptors. *Proceedings of the National Academy of Sciences*, 104(25):10601–10606, 2007.
- 610 [31] Valentina Escott-Price, C’eline Bellenguez, Li-San Wang, Seung-Hoan Choi, Denise Harold, Lesley Jones, Peter
611 Holmans, Amy Gerrish, Alexey Vedernikov, Alexander Richards, et al. Gene-wide analysis detects two new
612 susceptibility genes for alzheimer’s disease. *PloS one*, 9(6):e94661, 2014.
- 613 [32] A Ruiz, S Heilmann, Tim Becker, Isabel Hern’andez, Hermann Wagner, M Thelen, A Mauleon, M RosendeRoca,
614 C’eline Bellenguez, JC Bis, et al. Follow-up of loci from the international genomics of alzheimer’s disease project
615 identifies trip4 as a novel susceptibility gene. *Translational psychiatry*, 4(2):e358, 2014.
- 616 [33] Ling Hu, Rong Zhang, Qiong Yuan, Yinping Gao, Mary Q Yang, Chunxiang Zhang, Jiankun Huang, Yufei Sun,
617 William Yang, Jack Y Yang, et al. The emerging role of microRNA-4487/6845-3p in alzheimer’s disease
618 pathologies is induced by a β 25–35 triggered in sh-sy5y cell. *BMC systems biology*, 12(7):119, 2018.
- 619 [34] Juhyun Song and Young-Kook Kim. Identification of the role of mir-142-5p in alzheimer’s disease by
620 comparative bioinformatics and cellular analysis. *Frontiers in molecular neuroscience*, 10:227, 2017.
- 621 [35] Linlin Wang, Li Min, Qingdong Guo, Junxia Zhang, Hailun Jiang, Shuai Shao, Jianguo Xing, Linlin Yin, Jianghong
622 Liu, Rui Liu, et al. Profiling microRNA from brain by microarray in a transgenic mouse model of alzheimer’s
623 disease. *BioMed research international*, 2017, 2017.
- 624 [36] Masataka Katsu, Yuka Hama, Jun Utsumi, Ken Takashina, Hiroshi Yasumatsu, Fumiaki Mori, Koichi Wakabayashi,
625 Mikio Shoji, and Hidenao Sasaki. MicroRNA expression profiles of neuron-derived extracellular vesicles in
626 plasma from patients with amyotrophic lateral sclerosis. *Neuroscience letters*, 2019.
- 627 [37] Jaekwang Kim, Hyejin Yoon, Cristina M Ram’irez, Sang-Mi Lee, Hyang-Sook Hoe, Carlos Fern’andezHernando,
628 and Jungsu Kim. Mir-106b impairs cholesterol efflux and increases a β levels by repressing abca1 expression.
629 *Experimental neurology*, 235(2):476–483, 2012.
- 630 [38] Mariana Saint Just Ribeiro, Magnus L Hansson, Mikael J Lindberg, Anita E Popko-Scibor, and Annika E’ Wallberg.
631 Gsk3 β is a negative regulator of the transcriptional coactivator mam11. *Nucleic acids research*, 37(20):6691–
632 6700, 2009.

- 633 [39] Maiko Moriyama, Takeshi Fukuhara, Markus Britschgi, Yingbo He, Ramya Narasimhan, Saul Villeda, Hector
634 Molina, Brigitte T Huber, Mike Holers, and Tony Wyss-Coray. Complement receptor 2 is expressed in neural
635 progenitor cells and regulates adult hippocampal neurogenesis. *Journal of Neuroscience*, 31(11):3981–3989,
636 2011.
- 637 [40] Yifeng Du, Michael C Wooten, and Marie W Wooten. Oxidative damage to the promoter region of sqstm1/p62 is
638 common to neurodegenerative disease. *Neurobiology of disease*, 35(2):302–310, 2009.
- 639 [41] Yvonne Tay, Lev Kats, Leonardo Salmena, Dror Weiss, Shen Mynn Tan, Ugo Ala, Florian Karreth, Laura Poliseno,
640 Paolo Provero, Ferdinando Di Cunto, et al. Coding-independent regulation of the tumor suppressor pten by
641 competing endogenous mrnas. *Cell*, 147(2):344–357, 2011.
- 642 [42] Valerio Costa, Roberta Esposito, Marianna Aprile, and Alfredo Ciccodicola. Non-coding rna and pseudogenes in
643 neurodegenerative diseases: “the (un) usual suspects”. *Frontiers in genetics*, 3:231, 2012.

644

645 Figure legend:

646 Figure 1: The heterogeneous mutation data and the pipeline of GDAMDB model. (a. The idea of bridging
647 heterogeneous *mutation type* and *mutation association* by using data synchronization. b. Graphical model of the
648 “Mutation Data Bridging” model. c. pipeline of GDAMDB model for gene-disease association prediction.)

649

650 Figure 2: The illustrated examples of *mutation types* retrieved from Glu693Gly, Ala673Val, and Val717Ile/Cly/Phe in
651 APP.

652

653 Figure 3: Biological process categories of 325 mutation types and the sentence evidences.

654

655 Figure 4: Data fusion of heterogeneous AD mutation data improves the discovery of novel AD-related genes.

656 Figure 5: 69 out of 79 predicted genes have supportive AD-related evidence.

657

658 Figure 6: The parameters setting of “Mutation Data Bridging” model.

Figures

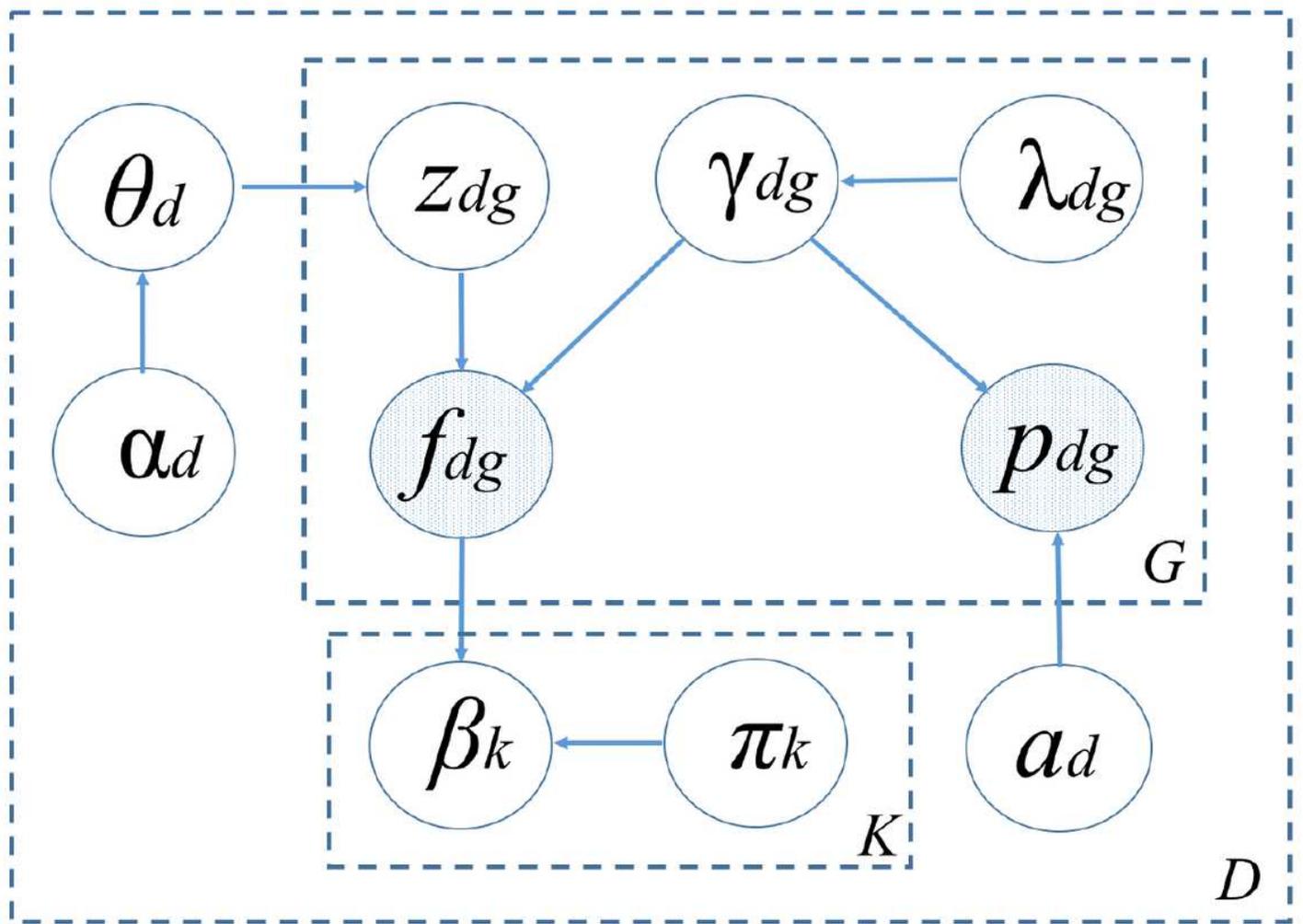


Figure 1

The parameters setting of "Mutation Data Bridging" model.

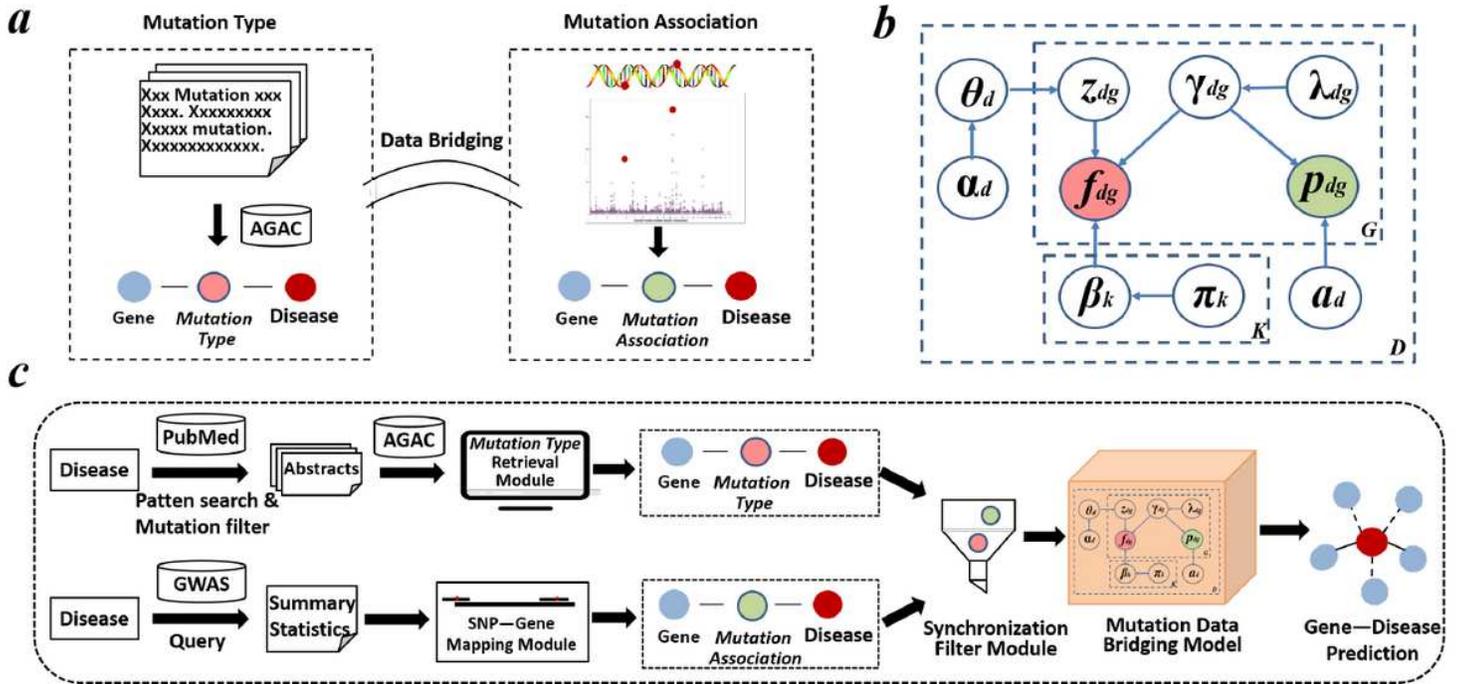


Figure 2

The heterogeneous mutation data and the pipeline of GDAMDB model. (a. The idea of bridging heterogeneous mutation type and mutation association by using data synchronization. b. Graphical model of the “Mutation Data Bridging” model. c. pipeline of GDAMDB model for gene-disease association prediction.)

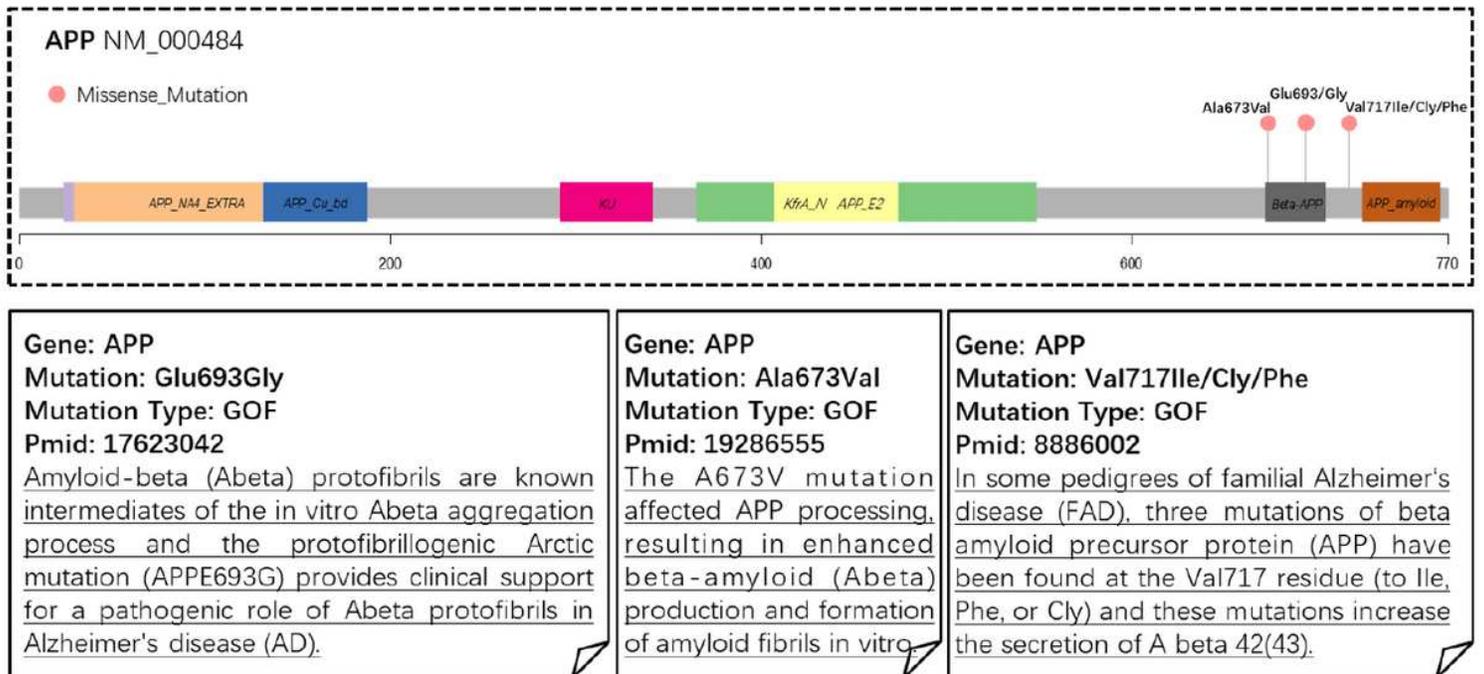


Figure 3

The illustrated examples of mutation types retrieved from Glu693Gly, Ala673Val, and Val717Ile/Cly/Phe in APP.

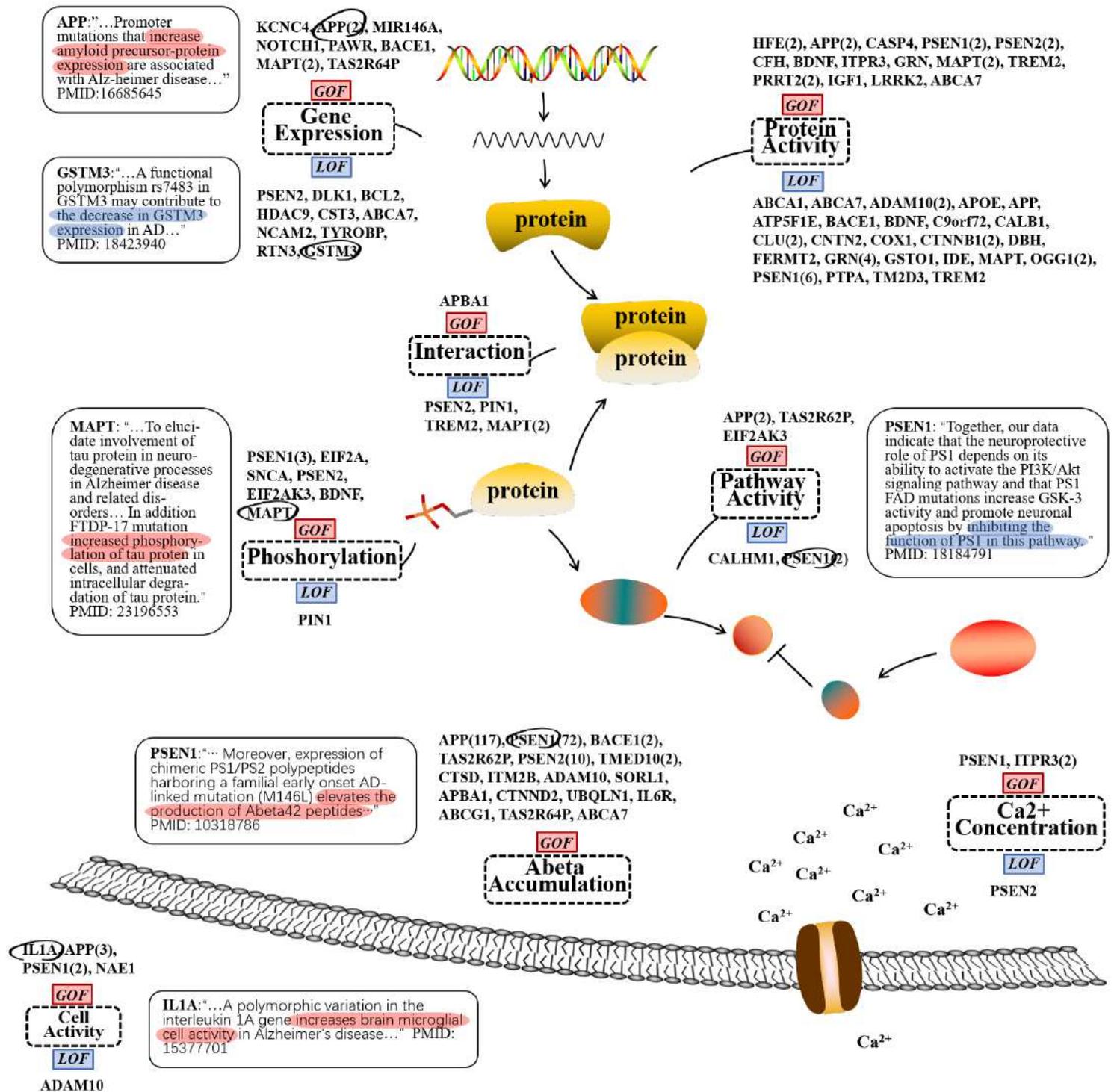


Figure 4

Biological process categories of 325 mutation types and the sentence evidences.

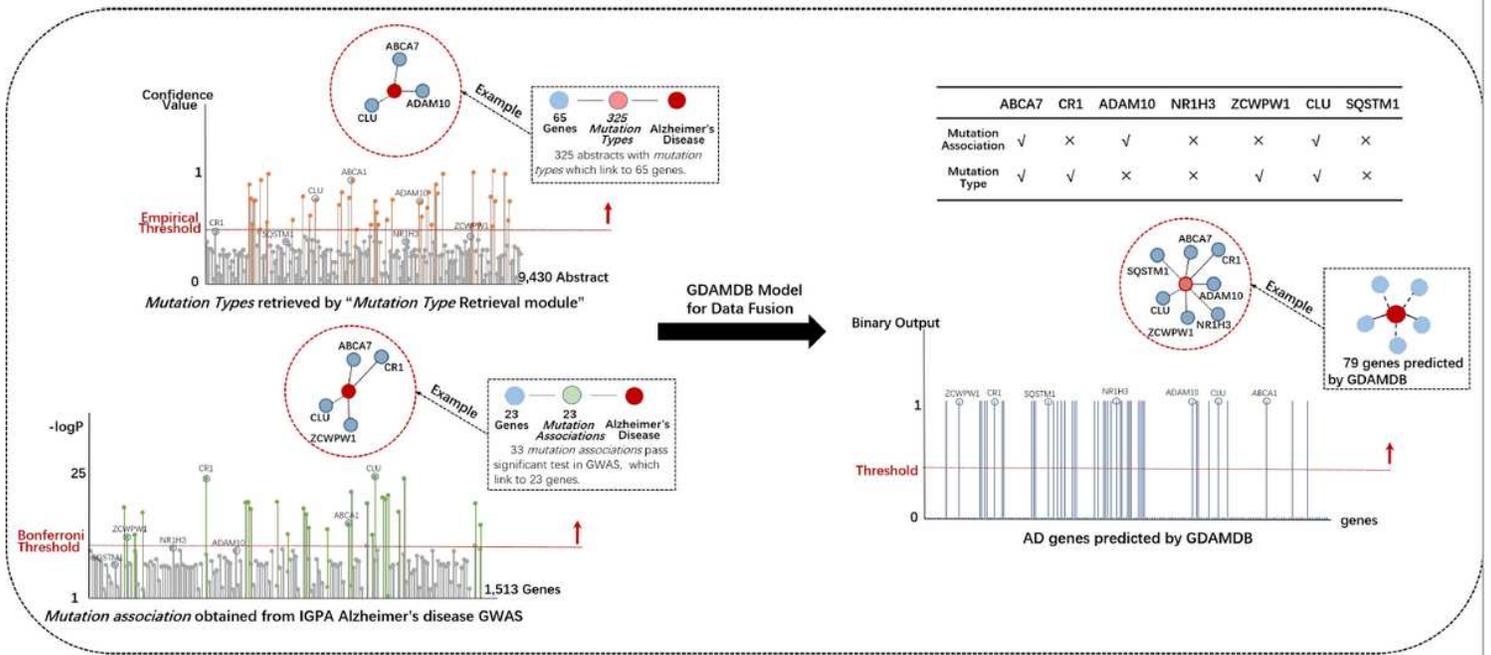


Figure 5

Data fusion of heterogeneous AD mutation data improves the discovery of novel AD-related genes.

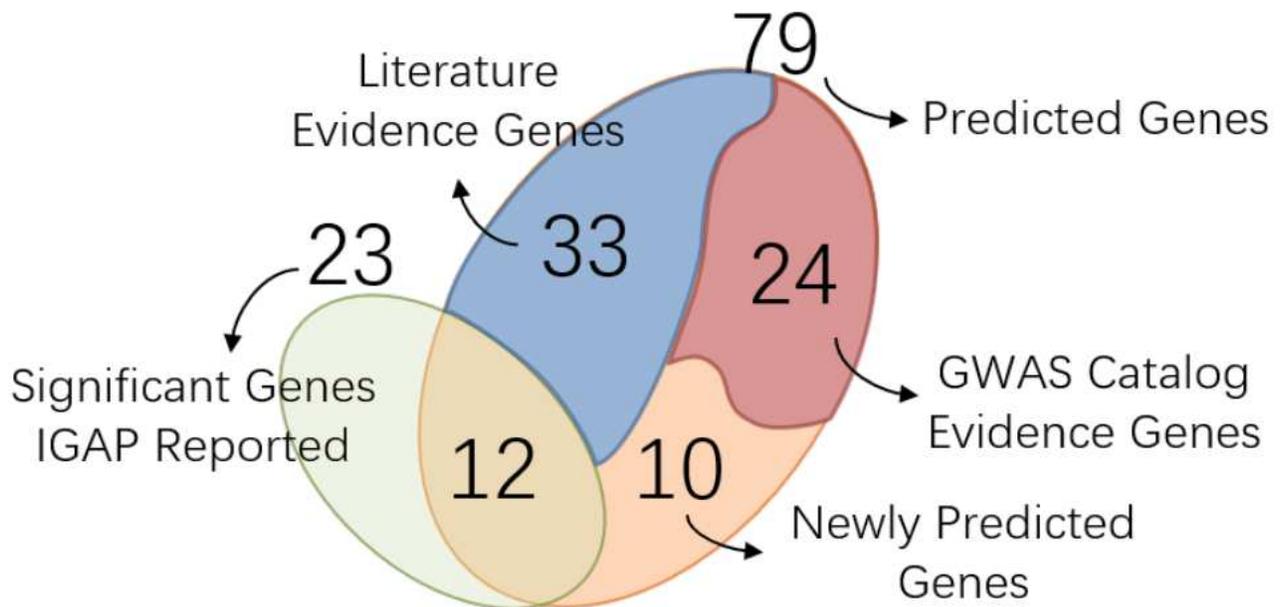


Figure 6

69 out of 79 predicted genes have supportive AD-related evidence.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [S379PredictedGenes.xlsx](#)
- [S2MutationTypeData.xlsx](#)
- [S1ProofofMutationDataBridgingModel.pdf](#)